# Phylogeny and evolution of Adephaga and Neuropterida (Insecta: Holometabola) as inferred from analyses of next-generation sequencing data

Cumulative dissertation submitted in partial fulfilment of

the requirements for the doctoral degree

(Dr. rer. nat.)

by

**Alexandros Vasilikopoulos**

from Athens, Greece

Bonn, February 2021

Faculty of Mathematics and Natural Sciences

Rheinische Friedrich-Wilhelms-Universität Bonn

Angefertigt mit Genehmigung der Mathematisch-Naturwissenschaftlichen Fakultät der
Rheinischen Friedrich-Wilhelms-Universität Bonn.

Angefertigt am Zoologischen Forschungsmuseum Alexander Koenig, Bonn.

Erstgutachter:                               Prof. Dr. Bernhard Misof

Zweitgutachter:                              Prof. Dr. Oliver Niehuis

Fachnahes Kommissionsmitglied:               Prof. Dr. Albert Haas

Fachfremdes Kommissionsmitglied:             Prof. Dr. Diana Imhof

Tag der Promotion: 31.08.2021

Erscheinungsjahr: 2021

# Summary

Knowing the evolutionary relationships of species is fundamental for comparative studies in biology as well as for biodiversity conservation. The main topic of this thesis is the investigation of phylogenetic relationships of two groups of holometabolous insects: Adephaga and Neuropterida. The phylogenetic analyses performed here were based on extensive genomic data that were obtained using two genome-reduction approaches: 1) transcriptomics and 2) hybrid enrichment of protein-coding exons. Several methods were applied to alleviate potential errors in the phylogenetic inferences such as: a) different data-subsampling strategies and b) application of methods and models that take into account different types of heterogeneity in the data. Evaluation of inferred evolutionary hypotheses was performed using a combination of methods such as: 1) different quartet-based measures of phylogenetic incongruence applied together with commonly used branch support measures, 2) congruency tests with morphology-based phylogenies, 3) comparisons of results from analyses of different data types (i.e., amino acids and nucleotides) and 4) comparison of results between best-fitting and less-fitting models of molecular evolution. This combination of methods and data was applied for the first time here to infer the phylogeny of Adephaga and Neuropterida. In general, these integrative phylogenomic approaches for species-tree inference and for evaluating inferred evolutionary hypotheses result in a better understanding of the phylogeny Adephaga and Neuropterida but they also help to identify unresolved or difficult phylogenetic

questions in the backbone phylogeny of these groups and potentially also in other groups of species. Additionally, the presented approaches for critically evaluating results of phylogenomic analyses constitute a valuable resource for future studies focusing on the reconciliation of molecular and morphological phylogenies.

In chapter 1, I provide a general introduction to the field of molecular systematics and phylogenomics and a brief introduction to the phylogeny of Adephaga and Neuropterida. In chapter 2, I assemble a large transcriptomic data matrix to infer the phylogeny and divergence times of Neuropterida and I evaluate the inferred results using different measures of phylogenomic incongruence. In chapter 3, I focus on the relationships in the adephagan superfamily Dytiscoidea using transcriptomes and apply data subsampling strategies in order to reduce deviation from model assumptions. I also evaluate results based on concatenation-based and gene tree-based measures of phylogenomic incongruence. In chapter 4, a combination of transcriptomes with new hybrid-enrichment data is performed to generate the most species-rich phylogenomic taxon sampling for Adephaga presented to date. I use this dataset to infer the phylogeny of Adephaga and perform evaluation-focused exploratory analyses with different evolutionary models and with quartet-based measures of phylogenetic support. Lastly in chapter 5, I discuss the most important results of this thesis in the historical context of knowledge on the phylogeny of these groups and I provide useful directions and advice for overcoming the limits of phylogenomic data and methods in future molecular systematic studies.

**Chapter 2.** The latest advancements in DNA sequencing technologies have facilitated the resolution of the phylogeny of insects, yet parts of the tree of Holometabola remain unresolved. The phylogeny of Neuropterida has been extensively studied, but no strong

consensus exists concerning the phylogenetic relationships within the order Neuroptera. Here, we assembled a novel transcriptomic dataset to address previously unresolved issues in the phylogeny of Neuropterida and to infer divergence times within the group. We tested the robustness of our phylogenetic estimates by comparing summary coalescent and concatenation-based phylogenetic approaches and by employing different quartet-based measures of phylogenomic incongruence, combined with data permutations. Our results suggest that the order Raphidioptera is sister to Neuroptera + Megaloptera. Coniopterygidae is inferred as sister to all remaining neuropteran families suggesting that larval cryptonephry could be a ground plan feature of Neuroptera. A clade that includes Nevrorthidae, Osmylidae, and Sisyridae (i.e. Osmyloidea) is inferred as sister to all other Neuroptera except Coniopterygidae, and Dilaridae is placed as sister to all remaining neuropteran families. Ithonidae is inferred as the sister group of monophyletic Myrmeleontiformia. The phylogenetic affinities of Chrysopidae and Hemerobiidae were dependent on the data type analyzed, and quartet-based analyses showed only weak support for the placement of Hemerobiidae as sister to Ithonidae + Myrmeleontiformia. Our molecular dating analyses suggest that most families of Neuropterida started to diversify in the Jurassic and our ancestral character state reconstructions suggest a primarily terrestrial environment of the larvae of Neuropterida and Neuroptera. Our extensive phylogenomic analyses consolidate several key aspects in the backbone phylogeny of Neuropterida, such as the basal placement of Coniopterygidae within Neuroptera and the monophyly of Osmyloidea. Furthermore, they provide new insights into the timing of diversification of Neuropterida. Despite the vast amount of analyzed molecular data, we found that certain nodes in the tree of Neuroptera are not robustly resolved. Therefore, we emphasize the importance of integrating the results of morphological analyses with

those of sequence-based phylogenomics. We also suggest that comparative analyses of genomic metacharacters should be incorporated into future phylogenomic studies of Neuropterida.


**Chapter 3.** The beetle superfamily Dytiscoidea, placed within the suborder Adephaga, comprises six families. The phylogenetic relationships of these families, whose species are aquatic, remain highly contentious. In particular the monophyly of the geographically disjunct Aspidytidae (China and South Africa) remains unclear. Here we use a phylogenomic approach to demonstrate that Aspidytidae are indeed monophyletic, as we inferred this phylogenetic relationship from analyzing nucleotide sequence data filtered for compositional heterogeneity and from analyzing amino-acid sequence data. Our analyses suggest that Aspidytidae are the sister group of Amphizoidae, although the support for this relationship is not unequivocal. A sister group relationship of Hygrobiidae to a clade comprising Amphizoidae, Aspidytidae, and Dytiscidae is supported by analyses in which model assumptions are violated the least. In general, we find that both concatenation and the applied coalescent method are sensitive to the effect of among-species compositional heterogeneity. Four-cluster likelihood-mapping suggests that despite the substantial size of the dataset and the use of advanced analytical methods, statistical support is weak for the inferred phylogenetic placement of Hygrobiidae. These results indicate that other kinds of data (e.g. genomic meta-characters) are possibly required to resolve the above-specified persisting phylogenetic uncertainties. Our study illustrates various data-driven confounding effects in phylogenetic reconstructions and highlights the need for careful monitoring of model violations prior to phylogenomic analysis.

**Chapter 4.** Adephaga is the second largest suborder of Coleoptera and contains aquatic and terrestrial groups of insects that are sometimes classified as Hydradephaga and Geadephaga respectively. Phylogenetic relationships of Adephaga have been extensively studied, but the relationshisps of the families of Geadephaga and some relationships within Dytiscoidea, such as the placement of Hygrobiidae, remain obscure. Here, we generate new DNA-hybridization baits for exon-capture phylogenomics and we combine the new hybrid-capture sequence data with transcriptomes to generate the largest phylogenomic taxon sampling for Adephaga presented to date. Our analyses show that the new baits can be successfully applied to recover the target loci in across divergent lineages of Adephaga. Concatenated analyses of moderately trimmed supermatrices strongly support the paraphyly of "Hydradephaga" with Gyrinidae placed as sister to all other families as in morphology-based phylogenies. All analyses under the site-heterogeneous models suggest Trachypachidae as sister to a clade Carabidae + Cicindelidae in congruence with previous morphological studies. Haliplidae is inferred as sister to Dytiscoidea, while a clade of Noteridae (+ most likely Meruidae) is inferred as sister to all remaining Dytiscoidea. A strongly supported clade Hygrobiidae + (Amphizoidae + Aspidytidae) is inferred in most analyses of the site-heterogeneous C60, PMSF and CAT+GTR models of moderately trimmed supermatrices under full taxon sampling. In general, we find that very stringent trimming of supermatrices results in reduced deviation from model assumptions but at the same time in reduction of phylogenetic information and in reduced phylogenetic resolution of Adephaga. We also find that site-heterogeneous C60 models provide greater stability of phylogenetic relationships of Adephaga across analyses of different amino-acid supermatrices than site-homogeneous models. Therefore site-heterogeneous C60 models can potentially reduce incongruence in

phylogenomics. Lastly, we show that gene-tree errors are prominent in the data, even after subsampling genes to potentially reduce these errors but we also show that subsampling genes based on the likelihood mapping criterion results in higher topological congruence to the concatenation-based tree. Overall, our analyses demonstrate that moderate alignment trimming strategies, application of site-heterogeneous models and mitigation of gene-tree errors should be routinely included in the phylogenomic pipeline in order to more accurately infer the phylogeny of species.

# List of Figures

# List of Tables

# Contents

# 1. General introduction

Biologists have always been interested in discovering, describing and measuring biological diversity. A key aspect of describing biological diversity includes the deciphering of the phylogeny (i.e., evolutionary relationships) of species and populations of species that constitutes the basis for conservation science (Lean and Maclaurin, 2016). Additionally, deciphering the timing of the origins of species through evolutionary analyses provides a window into the biological history of Earth. How are different species related to each other? When did their common ancestor inhabit the Earth? Is it possible to explain present geographical distributions of species by examining phylogenetic patterns? What can we say about the evolution of morphological and ecological traits when looking into the phylogeny of species? How can we assess the reliability of evolutionary hypotheses? These are all questions that are subjects of the fields of biological systematics and phylogenetics. In my dissertation, I focus on the higher-level phylogeny of the holometabolous insects in the superorder Neuropterida and those in the beetle suborder Adephaga by using analyses of next-generation sequencing data. Here, I will first introduce the common concepts, methods and problems in the field of molecular systematics and phylogenetics, in order to facilitate their discussion in the context of the phylogenetic inference of Adephaga and Neuropterida. In addition, I will provide a brief introduction to the current knowledge on the biology and phylogeny of these insect groups in order to provide a framework for discussing my phylogenetic results in the context of existing knowledge.

## 1.1. What is molecular systematics?

The biological field of "systematics" is broadly defined as the study of detecting, describing and explaining biological diversity (Moritz and Hillis, 1996). One

of the tasks of systematists is to infer the evolutionary relationships of species and use these inferences to classify them (Yang and Rannala, 2012). In the early years of systematics scientists relied on subjective criteria for classifying biological diversity. However, since Willi Hennig's notion of "phylogenetic systematics", scientists have increasingly been using more objective criteria for biological classifications that are based on strictly phylogenetic (i.e., evolutionary) concepts (Hennig, 1966). One example of such concepts is that monophyletic groups (i.e., clades, meaning groups of species that include all descendants of a common ancestor) should only be defined by shared derived characters that are unique for these groups (i.e., synapomorphies). These concepts have revolutionized systematic biology ever since and they were also later used to allow evolutionary inferences of species relationships (Felsenstein, 1981; Fitch, 1971). Nowadays, most modern systematists use evidence from evolutionary analyses of species relationships in order to classify them.

Molecular systematics is a modern branch of systematics that tries to address the above-mentioned questions by the analysis of genetic markers (Moritz and Hillis, 1996). In general, molecular systematics is broadly defined to include among others: inferences of species relationships and their classifications, studies of population structure and hybridization as well as studies of species boundaries (Moritz and Hillis, 1996). Molecular systematists use molecular data to infer phylogenetic trees. Phylogenetic trees are diagrams depicting ancestor-descendant relationships between organisms or gene sequences (Holder and Lewis, 2003). The branching pattern of phylogenetic trees provides information on the relationships of species or genes and is called the "topology" of the tree. The inference of phylogenetic trees is the subject of

the field of "phylogenetics". Accordingly, "molecular phylogenetics" is the study of evolutionary relationships among species or genes with the use of molecular data.


## 1.2. The concept of homology in molecular systematics

In order to infer the evolutionary relationships among species by any method a set of homologous characters has to be used. Homologous characters are characters that derive from a common ancestral character (Fitch, 2000). In molecular systematics and phylogenetics the most common types of data used for evolutionary inferences are homologous genetic fragments (e.g., genes derived from a common ancestral gene). In particular, phylogenetic analyses of species relationships are typically based on the analyses of orthologous genes (Smith and Hahn, 2020). Orthology is a specific type of homology which defines genes that derive from an ancestral speciation event, in contrast to paralogy which defines genes that derive from an ancestral duplication event (Fitch, 2000, 1970). It is logical that inference of the phylogeny of species should be based on the analyses of genes that reflect the speciation history of the species under investigation (Moritz and Hillis, 1996). For the purpose of inferring species trees, orthologous genes are first aligned before phylogenetic reconstructions, to produce a hypothesis on the positional homology of characters within the orthologous sequences (i.e., a multiple sequence alignment, MSA). During this process, gaps are inserted in the sequences to account for insertions or deletions when the sequences among different species are of unequal length. This results in an aligned set of sequences whereby the residues of the same column derive from the same ancestral residue in the common ancestor of the investigated species (Kapli et al., 2020). The main advantages of using molecular data for inferring species trees are 1) vast number of available

characters in comparison to morphology and 2) the existence of sophisticated statistical models that are used to describe various evolutionary processes at the molecular level (e.g., Crotty et al., 2020; Lartillot and Philippe, 2004; Tavaré, 1986; Yang, 1996).

## 1.3. A brief overview of molecular phylogenetic methods

Phylogenetic inference from molecular sequence data is based on statistical methods that use a set of aligned homologous sequences as a basis to infer a phylogenetic tree (or set of trees). In general, phylogenetic inference methods are divided into: 1) distance-based and 2) character-based methods (Kapli et al., 2020; Van de Peer, 2009). Distance-based methods derive evolutionary distances among pairs of taxa in the MSAs and use these distances to infer a phylogenetic tree. They may or may not use a model of sequence evolution to infer the true evolutionary distances (Van de Peer, 2009). After distances have been inferred, the sequence alignments are not used anymore in distance-based methods (Yang and Rannala, 2012). Some distance-based methods (e.g., neighbor joining) try to fit a phylogenetic tree to a pairwise distance matrix by using a clustering approach (Felsenstein, 1988; Van de Peer, 2009; Yang and Rannala, 2012). On the other hand, character-based methods use a specific optimality criterion to evaluate all possible phylogenetic trees (Bleidorn, 2017). Maximum parsimony is the oldest of the character-based methods and it is based on the criterion that the phylogenetic tree that explains all of the observed data by invoking the fewest character state changes is the most likely (Baum and Smith, 2013; Futuyma, 2013). Maximum parsimony has mostly been used in analyses of morphological data whereas its application in molecular phylogenetics has been replaced by the application of more sophisticated statistical inference methods.

Molecular sequence data are commonly analyzed with two other character-based methods that always use an explicit model of sequence evolution to account for multiple substitutions between character states: maximum likelihood (ML) or Bayesian inference (BI) (Felsenstein, 1981; Rannala and Yang, 1996). Because of the nature of molecular sequence data (i.e., a few character states), homoplasy can be prominent in the analyses of DNA or amino-acid sequences (Boore and Fuerstenberg, 2008; Jeffroy et al., 2006). Multiple substitutions at homologous alignment sites are therefore modelled through models of sequence evolution that are used to quantify the amount of change between aligned homologous sequence data (e.g., Tavaré, 1986). These models have been expanded to include different types of heterogeneity such as the heterogeneity of evolutionary rates among different sites of the alignment (Yang, 1996, 1994). Such models are the basis for phylogenetic reconstruction under ML and BI methods. ML tries to find the tree that maximizes the likelihood of observing the data given a model of sequence evolution (Baum and Smith, 2013), whereas BI uses probability distributions of parameters to evaluate trees based on their posterior probability, the probability that the tree is true given the data and the model (Baum and Smith, 2013). The commonly used models of molecular evolution assume that substitution of nucleotides or amino acids over time is a stochastic process (i.e., Markov process) in which the change from one state to another is only dependent on the last character state and not on any of the previous states (Galtier et al., 2005; Strimmer and von Haeseler, 2009). Other assumptions of most commonly used models of sequence evolution are that the DNA or amino-acid sequences have evolved under stationary (i.e., the relative frequencies of amino acids or nucleotides are at equilibrium), homogeneous (i.e., substitution rates do not change over time), and

reversible conditions (SRH conditions). Reversibility means that the probability of sampling nucleotide i and going to nucleotide j is the same as that of sampling nucleotide j and going to nucleotide i (Jermiin et al., 2008). If the assumptions of the models are violated by the data, the results of phylogenetic analyses of model-based phylogenetic reconstructions can be misleading (Ababneh et al., 2006; Jermiin et al., 2008, 2004).

When the data include two or more genetic loci, two main phylogenetic approaches for inferring species trees from these data exist: 1) supermatrix- or concatenation-based approaches (de Queiroz and Gatesy, 2007) and 2) supertree approaches (Bininda-Emonds, 2004; Kapli et al., 2020). In supermatrix-based approaches all independent genetic loci are concatenated to generate a supermatrix. This supermatrix is then analyzed by assuming a common topology across loci and using a model of sequence evolution that incorporates substitutional heterogeneity across the alignment sites of the supermatrix (e.g., Lanfear et al., 2012; Lartillot and Philippe, 2004). In the case of amino-acid supermatrices, the most common ways to model heterogeneity of substitution processes across sites are 1) site-homogeneous partition models and 2) site-heterogeneous mixture models (Kapli et al., 2020). Both types of models use different substitution matrices for different groups of alignment sites but they do so in different ways. The most complex site-heterogeneous mixture models use probabillity distributions and Markov chain Monte Carlo algorithms (MCMC) to assign sites to a specific class (or category) depending on stationary equilibrium frequencies (Lartillot and Philippe, 2004), whereas most partition models require an *a priori* defined partitioning scheme that is used as a basis to identify groups of partitions (or genes) presumably evolving under the same substitution processes

(Frandsen et al., 2015). Site-heterogeneous mixture models take into account the heterogeneity in amino-acid propensities across different alignment sites whereas partition models do not explicitly model this type of heterogeneity and therefore they are called site-homogeneous models. Despite this, partition models or less complex site-heterogeneous mixture models (e.g., Wang et al., 2019) are the only models scalable to the analyses of very large supermatrices (Frandsen et al., 2015).

Supertree methods combine information from multiple trees inferred under different sets of data to infer a species tree (Bininda-Emonds, 2014). According to some authors, summary multispecies coalescent methods (MSC methods) constitute a special case of supertree methods that are explicitly robust to genealogical heterogenity due to incomplete lineage sorting (ILS) (Bininda-Emonds, 2014). In general, coalescent-based methods for species tree reconstruction are divided into "full" coalescent methods (e.g., Drummond and Rambaut, 2007) and summary coalescent methods (e.g., Mirarab et al., 2014). Full coalescent Bayesian methods co-estimate gene trees and species trees by using MCMC algorithms to average over gene trees and other parameters but are computationally demanding and not applicable to large phylogenomic datasets (Xu and Yang, 2016). Summary MSC methods take a two-step approach to species tree inference (Liu et al., 2019). First a set of gene trees has to be inferred by any method and this set is then used by the summary method to infer a species tree (e.g., Liu et al., 2010; Mirarab et al., 2014). Summary MSC methods can be statistically consistent in some cases in which concatenation fails to provide statistically consistent results due to ILS (Roch and Steel, 2015), but they are sensitive to gene-tree estimation errors (Kapli et al., 2020; Roch and Warnow, 2015).

## 1.4. Inferring the divergence times of species

Biologists often want to assign dates to specific nodes of a tree which provide information on the timing of origin of the analyzed species. The idea of using molecular sequence data to date species divergences relied on the assumption that if proteins have similar rates of evolution among different lineages (e.g., Zuckerkandl and Pauling, 1965), then these rates can be used to estimate species divergence times, because the amount of difference between the sequences will be proportional to the time since the species diverged. This is commonly referred to as the molecular clock hypothesis (Zuckerkandl and Pauling, 1965). However, nowadays it is clear that the the assumption of uniform substitution rates among different lineages may hold for closely related species but is generally unrealistic for distantly related species (Yang, 2014; Yoder and Yang, 2000). Because of this, estimation of species divergence times is usually performed with relaxed-clock models that use distances among aligned sequence data and fossil calibrations to estimate species divergence times (Hasegawa et al., 1985; Rannala and Yang, 2007). In relaxed-clock models, different lineages (or branches) are allowed to have different rates of substitution that are either independent or autocorrelated (Ho and Duchêne, 2014). Furthermore, calibration information in the form of dates of fossils (or other geological events) is necessary in order to have absolute and not relative times of divergence (Ho and Duchêne, 2014). The most commonly used methods for estimating species divergence times are Bayesian methods that either co-estimate species divergence times and the phylogeny of species or use a fixed tree topology to infer species divergence times (Drummond and Rambaut, 2007; Rannala and Yang, 2007; Thorne et al., 1998).

## 1.5. Evaluating inferred phylogenetic relationships

Because there are various reasons that the topology of the inferred tree might deviate from the true species tree (some of them described in a following subsection) the field of phylogenetics should not be seemed as an attempt to "build" or "reconstruct" the true tree, but rather to examine alternative hypotheses and to quantify the extent to which the results support or exclude certain hypotheses (Baum and Smith, 2013). Accordingly, phylogenetic relationships of species as inferred from a specific method should be regarded as hypotheses or estimates of the evolutionary relationships of species. Ideally, after having inferred a phylogenetic tree scientists might wish to assign some confidence to the inferred phylogenetic hypothesis. There are two major ways to test the reliability of the inferred phylogenetic tree by any method: 1) statistical tests, and 2) congruency or plausibility tests (Futuyma, 2013; Wägele, 2005).

A first type of simple statistical tests can be applied to assess the reliability of the inferred tree by testing whether or not a phylogenetic hypothesis fits the data better than an alternative hypothesis (Futuyma, 2013). In a maximum likelihood framework, this test is simply the comparison the log-likelihood scores of the data under the two hypotheses and the model (Futuyma, 2013). For example, it is common practice that multiple independent maximum likelihood tree searches are performed for the same molecular dataset and some of them result in trees with different topologies due to local optima in the likelihood surface (Money and Whelan, 2012). In these cases the tree with the best score is selected as the "maximum-likelihood" tree (Money and Whelan, 2012).

Another class of statistical tests assess whether or not the likelihoods of two models or trees, are significantly different or could be explained by random effects

(Church et al., 2015; Goldman et al., 2000; Schmidt, 2009; Shimodaira, 2002). One example is the approximately unbiased test (AU test) for tree topologies (Shimodaira, 2002). These types of statistical tests compare the inferred phylogenetic tree and one or more additional phylogenetic trees which differ from the inferred tree in having one or more mutually exclusive clades (e.g., Shimodaira and Hasegawa, 1999). Subsequently, the tests are used to reject or accept alternative trees with a certain degree of confidence but some of them require the investigator to make subjective decisions on specific steps of the analyses and are also sensitive to model misspecification (Church et al., 2015).

Finally, a third class of statistical methods exists which makes it is possible to quantify support or phylogenetic signal in favour of specific branches or clades on a phylogenetic tree. This can be done by using either 1) conventional measures of branch support (e.g., bootstrap, jacknifing and Bayesian posterior probabilities) (Felsenstein, 1985; Källersjö et al., 1998; Rannala and Yang, 1996), 2) by using approaches that examine number of sites, genes or partitions supporting a specific branch in the inferred phylogenetic tree (e.g., Ané et al., 2007; Baum, 2007; Shen et al., 2017) and 3) by using measures of support that are based on the analyses of quartets of taxa (e.g., Minh et al., 2020; Pease et al., 2018; Strimmer and von Haeseler, 1997; Zhou et al., 2020). Many of these methods (e.g., bootstrap, posterior probabilities, four-cluster likelihood mapping, site- and gene-wise likelihoods) enable us to measure support (or signal) in favour of alternative hypotheses (or branches) that may or may not be present in the inferred tree.

The second category of tests for testing the reliability of phylogenetic trees are the congruency tests that leverage information from multiple independent sources of data to assess confidence in a particular phylogenetic result. This type of test differs from the above-described tests in that it is based on comparing results from completely

independent sources of data (Wägele, 2005). Examples of such independent sources of data can be non-overlapping sets of genetic sequence data (such as coding versus non-coding regions of the genome or different genes), morphological versus molecular sequence data, or genomic rearrangements compared with morphological data.

## 1.6. The transition from molecular phylogenetics to phylogenomics

In the early years of molecular systematics scientists have used the information from a few or single genes to reconstruct the evolutionary relationships of species and to infer species divergence times and biogeographic patterns (Hillis et al., 1996). These early analyses relied on the assumption that phylogenetic trees inferred from analysing the sequences of single or a few genes are representative of the true phylogeny of the species. Later in the early 2000s, it became clear that increasing the number of genes in a phylogenetic analysis results in higher branch support values of the inferred phylogenetic relationships (Rokas et al., 2003). It also became evident that including more genes in a phylogenetic analysis results in the accumulation of phylogenetic signal (Delsuc et al., 2005; Simion et al., 2020). Although it is now known that high branch support values are not necessarily due to accumulated signal (e.g., Hoang et al., 2018), these observations resulted in the addition of more data becoming the standard procedure in phylogenetic research, in an attempt to overcome biased phylogenetic estimates due to insufficient phylogenetic signal (Philippe et al., 2017), and in order to end phylogenetic incongruence (Gee, 2003; Rokas et al., 2003). Because of this, the science of molecular systematics shifted from using a few genes to using information from entire genomes or large portions of genomes to infer phylogenetic relationships of species (Eisen and Fraser, 2003; Young and Gillung, 2020). Such evolutionary

reconstructions based on analyses of genomes are the subject of the scientific field of phylogenomics (Delsuc et al., 2005). Although the term "phylogenomics" was originally used to describe prediction of gene function at the genomic level (Eisen, 1998), it is is now a composite field of research that includes: 1) the utilization of genome-scale data for deciphering the evolutionary relationships of species (Delsuc et al., 2005; Eisen and Fraser, 2003; Young and Gillung, 2020), 2) the prediction of gene functions based on phylogenetic analyses of gene families (Brown and Sjölander, 2006; Eisen, 1998; Eisen and Fraser, 2003) and 3) the study of gene-repertoire evolution by looking into the phylogenetic histories of genes in the genomes of different species (e.g., Fernández and Gabaldón, 2020; Julca et al., 2020).

## 1.7. Next-generation sequencing techniques and strategies in phylogenomics

The utilization of genomic data in molecular systematics is now possible due to the advancements in next-generation sequencing technologies (NGS) (Lemmon and Lemmon, 2013). Specifically, the invention of massively parallel DNA sequencing techniques has enabled the simultaneous sequencing of hundreds of thousands of DNA reads at relatively high accuracy and within a short amount of time (Bentley et al., 2008; Margulies et al., 2005). Sequence tagging of the desired DNA fragments (i.e. multiplexing) has allowed sequencing genomic data of several individuals on the same sequencing run and therefore has dramatically reduced the sequencing costs for molecular systematic studies (Glenn, 2011; McCormack et al., 2013). Several sequencing techniques (and platforms) exist for sequencing genomic data from several

specimens in parallel (Bentley et al., 2008; Eid et al., 2009) and each of them has specific advantages and limitations (see Bleidorn, 2017; Glenn, 2011).

There are multiple sequencing strategies (i.e., data collection approaches) that are commonly used to sequence multiple genetic loci across the genomes of interest (Bleidorn, 2017). These sequencing strategies are divided into two main categories: 1) sequencing strategies that target the whole genome of the organism of interest non-specifically (see Bleidorn, 2017), and 2) genome-reduction (or genome-partitioning) sequencing strategies that target only a fraction of the genome of interest in a specific manner (Jones and Good, 2016; McCormack et al., 2013). Within the first category whole-genome shotgun sequencing (WGS) describes the strategy in which the total genomic DNA is extracted, sheared rendomly into fragments and sequenced non-specifically without first establishing a physical map (Bleidorn, 2017). The sequenced reads are then assembled into contigs and used for downstream detection of suitable loci for phylogenetic reconstruction. A specific type of WGS sequencing is "genome skimming" that refers to non-specific and shallow sequencing (i.e. low-depth or low-coverage sequencing) of the genome of interest. This strategy is mostly useful for obtaining high-copy parts of the genome such as mitogenomes, plastomes and repetitive elements (Dodsworth, 2015; Straub et al., 2012). Genome skimming has been successfully applied to infer phylogenetic relationships in animals (e.g., Richter et al., 2015) and plants (e.g., Malé et al., 2014). WGS at high coverage depths is increasingly becoming the standard procedure in many phylogenetic studies due to reduction of the sequencing costs (e.g., Árnason et al., 2018; Edelman et al., 2019; Jarvis et al., 2014).

The most common genome-reduction strategies that are used in molecular systematics are 1) RNA sequencing or transcriptome sequencing (e.g. Misof et al.,

2014), 2) target-DNA enrichment, sequence capture or hybrid enrichment (e.g. Bragg et al., 2016; Lemmon et al., 2012) and 3) restriction-site-associated DNA sequencing (RAD-seq) (e.g., Díaz-Arce et al., 2016). Transcriptome sequencing refers to sequencing the genes that are expressed in the organism or tissue of interest at the time of isolating the tissue samples. With this approach the extracted mRNA is used as template to construct the reverse complement DNA and this complementary DNA then undergoes NGS library preparation and high-throughput sequencing (Young and Gillung, 2020). Hybrid enrichment or sequence capture is a specific type of DNA-target enrichment (Mamanova et al., 2010). Hybrid enrichment strategies use available genomic or transcriptomic data to design oligonucleotide probes (or baits) that bind specifically to the regions of interest, and remove non-target regions before amplification and high-throughput sequencing of the targets (Lemmon and Lemmon, 2013). Two common hybrid-enrichment sequencing strategies are: 1) hybrid enrichment of ultraconserved elements (UCEs, Faircloth et al., 2012) and 2) exon capture (Bi et al., 2012; Bragg et al., 2016). Transcriptomic and hybrid-enrichment approaches have been used both in addressing both shallow- and deep-level phylogenetic questions (e.g., Hamilton et al., 2016; Rodríguez et al., 2017; Wickett et al., 2014; Zhang et al., 2020). The third type of genome partitioning strategies (i.e., RAD-seq) starts with shearing the extracted genomic DNA into fragments. Subsequently, a selection of those restriction-based fragments is made based on their size and the selected fragments undergo library preparation and high-throughput sequencing (Baird et al., 2008; Lemmon and Lemmon, 2013). The RAD-seq sequencing strategy is mostly effective for investigating shallow-level phylogenetic and phylogeographic questions (Lemmon and Lemmon, 2013).

## 1.8. Types and sources of phylogenetic error

Phylogenetic inference of species relationships can be biased due to different reasons therefore resulting in phylogenetic estimation error. One of the reasons is because of the poor quality of the data used to infer the species tree. There are three main sources of error related to poor data quality: 1) the presence of cross-contaminated sequences, 2) homology errors and 3) multiple sequence alignment (MSA) errors (Philippe et al., 2017, 2011; Simion et al., 2018). For example, it has been previously shown that cross-species contamination may result in distantly related species being erroneously inferred as closely related (Laurin-Lemay et al., 2012; Simion et al., 2018). Therefore it is important that the data used for phylogenetic inference are screened for potential contaminations. Another issue related to data quality is errors in orthology inference (Philippe et al., 2011). Processes such as horizontal gene transfer, high rates of sequence evolution, and differential gene duplication and loss might result in the erroneous identification of orthologs (Altenhoff et al., 2019; Dalquen et al., 2013; Natsidis et al., 2020). In general, there are two types of *de novo* orthology inference approaches: graph-based and tree-based and they both suffer from specific limitations that make the identification of orthologs a non-trivial task (Gabaldón, 2008; Kapli et al., 2020; Smith and Pease, 2017). MSAs represent homology hypotheses for the characters of the aligned sequences, and therefore MSA quality is of paramount importance to accurate phylogenetic inference. There are multiple algorithms for inferring multiple sequence alignments, and each of them has a different degree of accuracy and computational speed (e.g., Bradley et al., 2009; Yamada et al., 2016). In addition, much effort has been made on developing methods for removing unreliable alignment blocks (e.g. Castresana, 2000; Kück et al., 2010;

Sela et al., 2015) or unreliable sequence fragments (e.g., Di Franco et al., 2019) in an attempt to improve the quality of MSAs.

Even when systematists have data of good quality, phylogenetic analyses can still result in biased phylogenetic estimates. There are two main types of error that are not related to the data quality problems described above. In some cases, phylogenetic analyses might result in poorly resolved or conflicting phylogenetic results due to insufficient phylogenetic signal that is bound to the length of the analyzed sequences (Delsuc et al., 2005; Jeffroy et al., 2006). This is especially common when inferring phylogenetic relationships from single genes or a few genes that do not contain enough information to reliably infer the phylogeny of species. This type of phylogenetic error is commonly referred to as "stochastic error" or "sampling error" (Kapli et al., 2020; Yang and Rannala, 2012). Another reason that the inferred phylogenetic trees might be incorrect is due to violation of model assumptions in model-based phylogenetic inference methods (Duchêne et al., 2017; Ho and Jermiin, 2004). This type of error is due to insufficient modelling and is referred to as "systematic error" (Kapli et al., 2020). Overall, this second type of error is more important when investigating ancient divergences because of the higher degree of erosion of phylogenetic signal and because of various heterogeneous processes potentially present in the analyzed data (Kapli et al., 2020). It has been pointed out that increasing the amount of data (i.e., number of genes) in the phylogenetic analyses results in reduction of stochastic errors whereas systematic errors might actually increase (Rodríguez-Ezpeleta et al., 2007).

Most commonly, systematic errors occur because the model is too simplistic and fails to describe heterogeneous processes in the data (Kapli et al., 2020). Examples of such heterogeneous processes are: 1) among-site rate variation (Yang, 1996), 2) among-

lineage (or across-taxa) rate variation (Felsenstein, 1978), 3) heterogeneity of substitution processes across the sites of the alignment (Lanfear et al., 2017; Lartillot and Philippe, 2004), 4) heterotachy (i.e. site-specific heterogeneity of evolutionary rate over time) (Lopez et al., 2002), 5) compositional heterogeneity across taxa (Foster, 2004), and 6) heterogeneity of genealogical relationships among different loci (Degnan and Rosenberg, 2009).

Genealogical heterogeneity of different loci (or genes) could be due to many different biological processes that cause gene trees to differ from the species tree (Maddison, 1997). Examples of such processes are: horizontal gene transfer (HGT, Marcet-Houben and Gabaldón, 2016), genomic introgression (Fontaine et al., 2014) and ILS (also referred to as deep coalescence, Degnan and Rosenberg, 2009). Such processes are not accounted for by the most commonly applied phylogenetic models in a supermatrix framework. On the other hand, methods that account for genealogical heterogeneity across genes usually only take into account one of these biological processes and do not incorporate other types of gene-tree heterogeneity (de Queiroz and Gatesy, 2007; Liu et al., 2019). This in turn makes the different methods sensitive to species-tree estimation errors when their assumptions are violated (e.g., Jiao et al., 2020). Overall, existing phylogenetic approaches and the applied models of sequence evolution are designed to tackle some of the heterogeneity in the data but not all possible types of heterogeneity. For this reason, the relative importance of different misleading factors and the selection of most appropriate phylogenetic methods and approaches is under the judgement of the research investigator (Simion et al., 2020). In addition to selecting a method that is suited for the analyses of the data at hand, there are three objective strategies to overcome biased estimates of phylogeny due to

systematic errors: 1) selecting datasets that are less likely to deviate from the model assumptions, 2) selecting best-fit substitution models based on objective statistical criteria for downstream phylogenetic reconstruction (Sullivan and Joyce, 2005), 3) evaluating model adequacy after the phylogenetic trees are inferred (Feuda et al., 2017; Jermiin et al., 2020; Shepherd and Klaere, 2019). These strategies are not mutually exclusive and a combination of them can be applied.

## 1.9. An introduction to the biology and phylogeny of Adephaga and Neuropterida

Insects are the most species-rich group of animals and their phylogeny has been extensively studied (Hennig, 1969; Kristensen, 1999; Misof et al., 2014; Wheeler et al., 2001). Within the larger phylogenetic clade of Insecta (i.e., insects *sensu stricto*) the species that undergo complete metamorphosis form a well-supported monophylum: the Holometabola (or Endopterygota) (e.g., Meusemann et al., 2010). Adephaga and Neuropterida, which are the focus of this dissertation, are phylogenetically placed within Holometabola and more specifically in the larger clade Neuropteroidea (or Neuropteriformia) (Peters et al., 2014). The clade Neuropteroidea includes the megadiverse Coleoptera (suborders: Adephaga, Archostemata, Myxophaga, Polyphaga), the Strepsiptera and the Neuropterida.

Neuropterida is a superorder of Holometabola with relatively few species (a little more than 6500 extant described species) and includes the orders Megaloptera, Neuroptera and Raphidioptera (Oswald, 2019). In general, Neuropterida are considered a relict group of insects (Engel et al., 2018), mainly due to the small number of species, the heterogeneity of the taxa, the vicariant geographical distributions and the rich fossil

records (Aspöck, 2002). Most members of Neuropterida are predators both as adults and larvae while there are some remarkable ecological adaptations within Neuroptera (e.g., larvae of the family Sisyridae use bryozoans and sponges as hosts, Winterton et al., 2010). The majority of adult insects in Neuropterida have terrestrial lifestyles, although there are some groups whose adults are always found in close proximity to water (e.g., Neuroptera: Nevrorthidae, Aspöck et al., 2017). Likewise, most larvae of Neuropterida are terrestrial but the larvae of some species of Neuropterida are strictly aquatic (i.e., Megaloptera, and members of neuropteran families Nevrorthidae and Sisyridae, Wang et al., 2017; Winterton et al., 2010). Hypotheses about the interordinal phylogeny of Neuropterida have reached a stable consensus among scientists in the last years (Wang et al., 2017; Y. Wang et al., 2019; Winterton et al., 2018, 2010). For example, a sister group relationship of Megaloptera and Neuroptera is considered a relatively robust hypothesis based on latest phylogenomic analyses (Misof et al., 2014; Wang et al., 2017; Winterton et al., 2018). Despite this, open questions in the phylogeny of the group remain, especially concerning the family relationships within the species-rich and ecologically diverse Neuroptera (Wang et al., 2017; Winterton et al., 2018). In addition, the timing of diversification of the major lineages of Neuropterida differs among previous molecular studies (Misof et al., 2014; Montagna et al., 2019; Winterton et al., 2018, 2010) and the pattern of evolution of larval ecologies has yet to be conclusively determined (Wang et al., 2017). For example, some analyses have suggested that the common ancestor of Neuroptera might have had either aquatic or terrestrial larvae depending on the analytical method used (Wang et al., 2017).

Adephaga is the second most species-rich suborder of Coleoptera after the suborder Polyphaga, and mostly includes predatory species whereas most species of Polyphaga are primarily phytophagous. The phylogenetic position of Adephaga within Coleoptera has been a matter of controversy because Polyphaga as sister to a clade Adephaga + (Archostemata + Myxophaga), which is suggested by morphology-based analyses, is not corroborated by phylogenomic analyses (Beutel et al., 2019; McKenna et al., 2019). The suborder Adephaga includes insects with either aquatic or terrestrial lifestyles, although there exist some adephagan families with species living in semi-aquatic or hygropetric habits (e.g., Aspidytidae, Balke et al., 2003). The phylogeny of the families of Adephaga have been extensively studied and these relationships have been more stable across analyses of different types of data (i.e., morphology and molecules) in comparison to the familial relationships in Neuroptera. Despite this, some open questions in the phylogeny of Adephaga remain.

An open question in the evolutionary history of Adephaga concerns the monophyly or not of the aquatic groups of Adephaga. Species in the aquatic family Gyrinidae have very different adaptations to life in water than species in other aquatic families of Adephaga and are generally morphologically distinct from species in other families of Adephaga (Beutel et al., 2020; Beutel and Roughley, 1988). This in turn suggests that the aquatic lineages do not constitute a monophyletic group and that more than one transition from terrestrial to aquatic habitats took place in the evolution of Adephaga (Beutel and Roughley, 1988). This hypothesis has been corroborated based on concatenation-based phylogenomic analyses of transcriptomes and UCEs (Gustafson et al., 2020; McKenna et al., 2019). Despite this, some recent reanalyses of phylogenomic data based on summary coalescent methods tentatively suggest the

monophyly of the aquatic groups (Freitas et al., 2020). In addition, previous phylogenomic studies did not employ complex site-heterogeneous models of sequence evolution in order to assess the placement of Gyrinidae in the tree of Adephaga. Other open phylogenetic questions concern the placement of the families Trachypachidae and Hygrobiidae the position of which has not been reconciled with morphology based on analyses of different types of data (e.g., Beutel et al., 2020; Gustafson et al., 2020). In general, one limitation of previous phylogenomic analyses of Adephaga is that their taxon sampling was limited in order to robustly test hypotheses concerning all familial relationships of Adephaga and to provide intra-familial hypotheses on the evolution of these groups. Lastly, there are no phylogenomic studies to date that examine both the effects of model misspecification and data-subsampling strategies on the phylogenetic inference of Adephaga as a whole.


## 1.10. Research focus and aims of the thesis

One of the aims of my thesis is to generate new genomic resources and tools for studying the evolutionary relationships of Adephaga and Neuropterida. To that end, I focus on generating two new ortholog sets (sets of clusters of orthologous and single-copy genes, sets of COGs, Tatusov et al., 1997) that are appropriate for transcriptome-based orthology assignment in these groups (see Petersen et al., 2017). Subsequently, my aim is to use these ortholog sets for identifying the genes of interest in the transcriptomes of Neuropterida and Adephaga and use them for downstream phylogeny reconstruction. It should be noted that newly generated and assembled transcriptomes of Neuropterida and Adephaga had already been generated in the frame of the 1KITE consortium (https://1kite.org/, accessed on 18.12.2020) and were made available for

analyses. Some additional transcriptomes of Adephaga, which were generated specifically for this project, had been sequenced and were also readily available but had not been previously processed or assembled with bioinformatic methods (see chapter 3). A few transcriptomes of Adephaga were obtained from corresponding databases and research collaborators and incorporated into my analyses (see chapter 4). In the frame of generating new genomic resources for Adephaga, another goal of my dissertation is to infer a new set of DNA-hybridization baits that are applicable for capturing phylogenetically informative genes in different lineages of the suborder.

In chapter 2, my main goal is to generate a large and carefully curated dataset for inferring the relationships within Neuropterida based on transcriptomes in order to potentially reconciliate molecular and morphological phylogenies of Neuropterida. In order to avoid biased inference of phylogeny my goal is to generate a large informative supermatrix and analyze it by taking into account possible violations of model assumptions by 1) sub-sampling the data to minimize deviation from model assumptions, 2) by using complex site-heterogeneous models in addition to partitioned site-homogeneous models and 3) by using summary coalescent phylogenetic approaches to assess potential incongruence due to ILS. Another goal is to evaluate the reliability of inferred phylogenetic hypotheses using conventional measures of branch support (e.g., bootstrap) combined with alternative measures of phylogenomic incongruence (e.g. four-cluster likelihood mapping, FcLM). This unique combination of methods and data for studying the phylogeny of Neuropterida is applied here for the first time. Lastly, another goal of my dissertation is to use the most reliable estimate of the phylogeny to infer the temporal pattern of diversification of the major lineages of

Neuropterida by using molecular dating methods and to also infer the pattern of evolution of larval ecologies using ancestral character state reconstruction methods.

In chapter 3, my main focus is to use a targeted sampling of new and old transcriptomes to infer the phylogeny of the superfamily Dytiscoidea, which is placed within the suborder Adephaga, and particularly of the small families Amphizoidae, Aspidytidae and Hygrobiidae. In this context, I will focus on minimizing the effects of potential model violations by 1) subsampling the data to reduce potential deviation from model assumptions, and 2) by using site-heterogeneous mixture models. Furthermore, my aim is again to assess the reliability of inferences using alternative and conventional measures of phylogenetic support. This new transcriptomic dataset is the largest dataset ever compiled to address phylogeny of the superfamily Dytiscoidea within the suborder Adephaga.

In chapter 4, my focus is to to generate a new widely applicable toolkit of DNA-hybridization baits for studying the phylogeny of the beetle suborder Adephaga as a whole. Firstly, my goal is to utilize all available genomic resources and the previously generated ortholog set used for Dytiscoidea in order to infer a novel set of DNA-hybridization probes and test its applicability for locus recovery in various lineages of Adephaga. Secondly, the aim is to efficiently combine newly generated sequence capture data with transcriptomes to infer the phylogeny of the suborder. Thirdly, another aim is to compare complex and less complex models in the phylogenetic inference of Adephaga, and to assess the effects of data-subsampling, performed here to reduce deviations from model assumptions, on the phylogenetic inferences. Specifically, by using a broad selection of models for model selection and by applying both site-heterogeneous and less fitting site-homogeneous models, I want to investigate

what the effect of misspecifying the model and overly trimming the data is on the results of phylogenetic reconstructions. This is a particularly interesting investigation concerning a few relationships of Adephaga that are generally well established based on analyses of other types of data, such as analyses morphological data. Another purpose is to investigate whether or not the disagreement between summary coalescent and concatenation methods could be reduced by careful selection of genes and whether or not gene-tree discordance in the data can be mainly explained by technical factors that do not have a biological basis.

One last overarching goal that is applicable to both groups of taxa I studied here is to reconcile disagreements between different studies specifically concerning controversial phylogenetic hypotheses that are conflicting between analyses of morphological and molecular sequence data. The potential for reconciliation of phylogenetic results between different types of data and methods constitutes the basis for further comparative evolutionary studies in these fascinating groups of insects.

## 1.11. References

Ababneh, F., Jermiin, L.S., Ma, C., Robinson, J., 2006. Matched-pairs tests of homogeneity with applications to homologous nucleotide sequences. Bioinformatics 22, 1225–1231.

Altenhoff, A.M., Glover, N.M., Dessimoz, C., 2019. Inferring orthology and paralogy, in: Anisimova, M. (Ed.), Evolutionary genomics: statistical and computational methods. Humana Press, New York, pp. 149–175.

Ané, C., Larget, B., Baum, D.A., Smith, S.D., Rokas, A., 2007. Bayesian estimation of concordance among gene trees. Mol. Biol. Evol. 24, 412–426.

Árnason, Ú., Lammers, F., Kumar, V., Nilsson, M.A., Janke, A., 2018. Whole-genome sequencing of the blue whale and other rorquals finds signatures for introgressive gene flow. Sci. Adv. 4, eaap9873.

Aspöck, U., 2002. Phylogeny of the Neuropterida (Insecta: Holometabola). Zool. Scr. 31, 51–55.

Aspöck, U., Aspöck, H., Liu, X., 2017. The Nevrorthidae, mistaken at all times: phylogeny and review of present knowledge (Holometabola, Neuropterida, Neuroptera). Dtsch. Entomol. Zeitschrift 64, 77–110.

Baird, N.A., Etter, P.D., Atwood, T.S., Currey, M.C., Shiver, A.L., Lewis, Z.A., Selker, E.U., Cresko, W.A., Johnson, E.A., 2008. Rapid SNP discovery and genetic mapping using sequenced RAD markers. PLoS One 3, e3376.

Balke, M., Ribera, I., Beutel, R.G., 2003. ASPIDYTIDAE: on the discovery of a new beetle family: detailed morphological analysis, description of a second species, and key to fossil and extant adephagan families (Coleoptera), in: Jach, M.A., Ji, L.

(Eds.), Water beetles of China. Zoologisch-Botanische Gesellschaft & Wiener Coleopterologenverein, pp. 53–66.

Baum, D.A., 2007. Concordance trees, concordance factors, and the exploration of reticulate genealogy. Taxon 56, 417–426.

Baum, D.A., Smith, S.D., 2013. Tree thinking: an introduction to phylogenetic biology, 1st ed. Roberts and Company Publishers, Inc., Greenwood village, Colorado.

Bentley, D.R., Balasubramanian, S., Swerdlow, H.P., Smith, G.P., Milton, J., Brown, C.G., Hall, K.P., Evers, D.J., Barnes, C.L., Bignell, H.R., Boutell, J.M., Bryant, J., Carter, R.J., Keira Cheetham, R., Cox, A.J., Ellis, D.J., Flatbush, M.R., Gormley, N.A., Humphray, S.J., Irving, L.J., Karbelashvili, M.S., Kirk, S.M., Li, H., Liu, X., Maisinger, K.S., Murray, L.J., Obradovic, B., Ost, T., Parkinson, M.L., Pratt, M.R., Rasolonjatovo, I.M.J., Reed, M.T., Rigatti, R., Rodighiero, C., Ross, M.T., Sabot, A., Sankar, S. V., Scally, A., Schroth, G.P., Smith, M.E., Smith, V.P., Spiridou, A., Torrance, P.E., Tzonev, S.S., Vermaas, E.H., Walter, K., Wu, X., Zhang, L., Alam, M.D., Anastasi, C., Aniebo, I.C., Bailey, D.M.D., Bancarz, I.R., Banerjee, S., Barbour, S.G., Baybayan, P.A., Benoit, V.A., Benson, K.F., Bevis, C., Black, P.J., Boodhun, A., Brennan, J.S., Bridgham, J.A., Brown, R.C., Brown, A.A., Buermann, D.H., Bundu, A.A., Burrows, J.C., Carter, N.P., Castillo, N., Catenazzi, M.C.E., Chang, S., Neil Cooley, R., Crake, N.R., Dada, O.O., Diakoumakos, K.D., Dominguez-Fernandez, B., Earnshaw, D.J., Egbujor, U.C., Elmore, D.W., Etchin, S.S., Ewan, M.R., Fedurco, M., Fraser, L.J., Fuentes Fajardo, K. V., Scott Furey, W., George, D., Gietzen, K.J., Goddard, C.P., Golda, G.S., Granieri, P.A., Green, D.E., Gustafson, D.L., Hansen, N.F., Harnish, K., Haudenschild, C.D., Heyer, N.I., Hims, M.M., Ho, J.T., Horgan, A.M., Hoschler, K., Hurwitz, S., Ivanov, D. V.,

Johnson, M.Q., James, T., Huw Jones, T.A., Kang, G.D., Kerelska, T.H., Kersey, A.D., Khrebtukova, I., Kindwall, A.P., Kingsbury, Z., Kokko-Gonzales, P.I., Kumar, A., Laurent, M.A., Lawley, C.T., Lee, S.E., Lee, X., Liao, A.K., Loch, J.A., Lok, M., Luo, S., Mammen, R.M., Martin, J.W., McCauley, P.G., McNitt, P., Mehta, P., Moon, K.W., Mullens, J.W., Newington, T., Ning, Z., Ling Ng, B., Novo, S.M., O'Neill, M.J., Osborne, M.A., Osnowski, A., Ostadan, O., Paraschos, L.L., Pickering, L., Pike, Andrew C., Pike, Alger C., Chris Pinkard, D., Pliskin, D.P., Podhasky, J., Quijano, V.J., Raczy, C., Rae, V.H., Rawlings, S.R., Chiva Rodriguez, A., Roe, P.M., Rogers, John, Rogert Bacigalupo, M.C., Romanov, N., Romieu, A., Roth, R.K., Rourke, N.J., Ruediger, S.T., Rusman, E., Sanches-Kuiper, R.M., Schenker, M.R., Seoane, J.M., Shaw, R.J., Shiver, M.K., Short, S.W., Sizto, N.L., Sluis, J.P., Smith, M.A., Ernest Sohna Sohna, J., Spence, E.J., Stevens, K., Sutton, N., Szajkowski, L., Tregidgo, C.L., Turcatti, G., Vandevondele, S., Verhovsky, Y., Virk, S.M., Wakelin, S., Walcott, G.C., Wang, J., Worsley, G.J., Yan, J., Yau, L., Zuerlein, M., Rogers, Jane, Mullikin, J.C., Hurles, M.E., McCooke, N.J., West, J.S., Oaks, F.L., Lundberg, P.L., Klenerman, D., Durbin, R., Smith, A.J., 2008. Accurate whole human genome sequencing using reversible terminator chemistry. Nature 456, 53–59.

Beutel, R.G., Pohl, H., Yan, E. V., Anton, E., Liu, S.-P., Ślipiński, A., McKenna, D., Friedrich, F., 2019. The phylogeny of Coleopterida (Hexapoda) – morphological characters and molecular phylogenies. Syst. Entomol. 44, 75–102.

Beutel, R.G., Ribera, I., Fikáček, M., Vasilikopoulos, A., Misof, B., Balke, M., 2020. The morphological evolution of the Adephaga (Coleoptera). Syst. Entomol. 45, 378–395.

Beutel, R.G., Roughley, R.E., 1988. On the systematic position of the family Gyrinidae (Coleoptera: Adephaga). J. Zool. Syst. Evol. Res. 26, 380–400.

Bi, K., Vanderpool, D., Singhal, S., Linderoth, T., Moritz, C., Good, J.M., 2012. Transcriptome-based exon capture enables highly cost-effective comparative genomic data collection at moderate evolutionary scales. BMC Genomics 13, 403.

Bininda-Emonds, O.R.P., 2014. An introduction to supertree construction (and partitioned phylogenetic analyses) with a view toward the distinction between gene trees and species trees, in: Garamszegi, L.Z. (Ed.), Modern phylogenetic comparative methods and their application in evolutionary biology. Springer-Verlag Berlin Heidelberg, Berlin, Heidelberg, pp. 49–76.

Bininda-Emonds, O.R.P., 2004. The evolution of supertrees. Trends Ecol. Evol. 19, 315–322.

Bleidorn, C., 2017. Phylogenomics - an introduction, 1st ed. Springer International Publishing.

Boore, J.L., Fuerstenberg, S.I., 2008. Beyond linear sequence comparisons: the use of genome-level characters for phylogenetic reconstruction. Philos Trans R Soc L. B Biol Sci. 363, 1445–1451.

Bradley, R.K., Roberts, A., Smoot, M., Juvekar, S., Do, J., Dewey, C., Holmes, I., Pachter, L., 2009. Fast statistical alignment. PLoS Comput. Biol. 5, e1000392.

Bragg, J.G., Potter, S., Bi, K., Moritz, C., 2016. Exon capture phylogenomics: efficacy across scales of divergence. Mol. Ecol. Resour. 16, 1059–1068.

Brown, D., Sjölander, K., 2006. Functional classification using phylogenomic inference. PLoS Comput. Biol. 2, e77.

Castresana, J., 2000. Selection of conserved blocks from multiple alignments for their use in phylogenetic analysis. Mol. Biol. Evol. 17, 540–552.

Church, S.H., Ryan, J.F., Dunn, C.W., 2015. Automation and evaluation of the SOWH test with SOWHAT. Syst. Biol. 64, 1048–1058.

Crotty, S.M., Minh, B.Q., Bean, N.G., Holland, B.R., Tuke, J., Jermiin, L.S., von Haeseler, A., 2020. GHOST: recovering historical signal from heterotachously evolved sequence alignments. Syst. Biol. 69, 249–264.

Dalquen, D.A., Altenhoff, A.M., Gonnet, G.H., Dessimoz, C., 2013. The impact of gene duplication, insertion, deletion, lateral gene transfer and sequencing error on orthology inference: a simulation study. PLoS One 8, e56925.

de Queiroz, A., Gatesy, J., 2007. The supermatrix approach to systematics. Trends Ecol. Evol. 22, 34–41.

Degnan, J.H., Rosenberg, N.A., 2009. Gene tree discordance, phylogenetic inference and the multispecies coalescent. Trends Ecol. Evol. 24, 332–340.

Delsuc, F., Brinkmann, H., Philippe, H., 2005. Phylogenomics and the reconstruction of the tree of life. Nat. Rev. Genet. 6, 361–375.

Di Franco, A., Poujol, R., Baurain, D., Philippe, H., 2019. Evaluating the usefulness of alignment filtering methods to reduce the impact of errors on evolutionary inferences. BMC Evol. Biol. 19, 21.

Díaz-Arce, N., Arrizabalaga, H., Murua, H., Irigoien, X., Rodríguez-Ezpeleta, N., 2016. RAD-seq derived genome-wide nuclear markers resolve the phylogeny of tunas. Mol. Phylogenet. Evol. 102, 202–207.

Dodsworth, S., 2015. Genome skimming for next-generation biodiversity analysis. Trends Plant Sci. 20, 525–527.

Drummond, A.J., Rambaut, A., 2007. BEAST: Bayesian evolutionary analysis by sampling trees. BMC Evol. Biol. 7, 214.

Duchêne, D.A., Duchêne, S., Ho, S.Y.W., 2017. New statistical criteria detect phylogenetic bias caused by compositional heterogeneity. Mol. Biol. Evol. 34, 1529–1534.

Edelman, N.B., Frandsen, P.B., Miyagi, M., Clavijo, B., Davey, J., Dikow, R.B., García-Accinelli, G., Van Belleghem, S.M., Patterson, N., Neafsey, D.E., Challis, R., Kumar, S., Moreira, G.R.P., Salazar, C., Chouteau, M., Counterman, B.A., Papa, R., Blaxter, M., Reed, R.D., Dasmahapatra, K.K., Kronforst, M., Joron, M., Jiggins, C.D., McMillan, W.O., Di Palma, F., Blumberg, A.J., Wakeley, J., Jaffe, D., Mallet, J., 2019. Genomic architecture and introgression shape a butterfly radiation. Science. 366, 594–599.

Eid, J., Fehr, A., Gray, J., Luong, K., Lyle, J., Otto, G., Peluso, P., Rank, D., Baybayan, P., Bettman, B., Bibillo, A., Bjornson, K., Chaudhuri, B., Christians, F., Cicero, R., Clark, S., Dalal, R., DeWinter, A., Dixon, J., Foquet, M., Gaertner, A., Hardenbol, P., Heiner, C., Hester, K., Holden, D., Kearns, G., Kong, X., Kuse, R., Lacroix, Y., Lin, S., Lundquist, P., Ma, C., Marks, P., Maxham, M., Murphy, D., Park, I., Pham, T., Phillips, M., Roy, J., Sebra, R., Shen, G., Sorenson, J., Tomaney, A., Travers,

K., Trulson, M., Vieceli, J., Wegener, J., Wu, D., Yang, A., Zaccarin, D., Zhao, P., Zhong, F., Korlach, J., Turner, S., 2009. Real-time DNA sequencing from single polymerase molecules. Science. 323, 133–138.

Eisen, J.A., 1998. Phylogenomics: improving functional predictions for uncharacterized genes by evolutionary analysis. Genome Res. 8, 163–167.

Eisen, J.A., Fraser, C.M., 2003. Phylogenomics: intersection of evolution and genomics. Science. 300, 1706–1707.

Engel, M.S., Winterton, S.L., Breitkreuz, L.C.V., 2018. Phylogeny and evolution of Neuropterida: where have wings of lace taken us? Annu. Rev. Entomol. 63, 531–551.

Faircloth, B.C., McCormack, J.E., Crawford, N.G., Harvey, M.G., Brumfield, R.T., Glenn, T.C., 2012. Ultraconserved elements anchor thousands of genetic markers spanning multiple evolutionary timescales. Syst. Biol. 61, 717–726.

Felsenstein, J., 1988. Phylogenies from molecular sequences: inference and reliability. Annu. Rev. Genet. 22, 521–565.

Felsenstein, J., 1985. Confidence limits on phylogenies: an approach using the bootstrap. Evolution. 39, 783–791.

Felsenstein, J., 1981. Evolutionary trees from DNA sequences: a maximum likelihood approach. J. Mol. Evol. 17, 368–376.

Felsenstein, J., 1978. Cases in which parsimony or compatibility methods will be positively misleading. Syst. Biol. 27, 401–410.

Fernández, R., Gabaldón, T., 2020. Gene gain and loss across the metazoan tree of life. Nat. Ecol. Evol. 4, 524–533.

Feuda, R., Dohrmann, M., Pett, W., Philippe, H., Rota-Stabelli, O., Lartillot, N., Wörheide, G., Pisani, D., 2017. Improved modeling of compositional heterogeneity supports sponges as sister to all other animals. Curr. Biol. 27, 3864–3870.

Fitch, W.M., 2000. Homology: a personal view on some of the problems. Trends Genet. 16, 227–231.

Fitch, W.M., 1971. Toward defining the course of evolution: minimum change for a specific tree topology. Syst. Zool. 20, 406–416.

Fitch, W.M., 1970. Distinguishing homologous from analogous proteins. Syst. Zool. 19, 99–113.

Fontaine, M.C., Pease, J.B., Steele, A., Waterhouse, R.M., Neafsey, D.E., Sharakhov, I. V, Jiang, X., Hall, A.B., Catteruccia, F., Kakani, E., Mitchell, S.N., Wu, Y.-C., Smith, H.A., Love, R.R., Lawniczak, M.K., Slotman, M.A., Emrich, S.J., Hahn, M.W., Besansky, N.J., 2014. Extensive introgression in a malaria vector species complex revealed by phylogenomics. Science. 347, 1258524.

Foster, P., 2004. Modeling compositional heterogeneity. Syst. Biol. 53, 485–495.

Frandsen, P.B., Calcott, B., Mayer, C., Lanfear, R., 2015. Automatic selection of partitioning schemes for phylogenetic analyses using iterative k-means clustering of site rates. BMC Evol. Biol. 15, 13.

Freitas, F.V., Branstetter, M.G., Griswold, T., Almeida, E.A.B., 2020. Partitioned gene-tree analyses and gene-based topology testing help resolve incongruence in a phylogenomic study of host-specialist bees (Apidae: Eucerinae). Mol. Biol. Evol. doi: 10.1093/molbev/msaa277.

Futuyma, D.J., 2013. Evolution, 3rd ed. Sinauer Associates, Inc., Sunderland, Massachusetts.

Gabaldón, T., 2008. Large-scale assignment of orthology: back to phylogenetics? Genome Biol. 9, 235.

Galtier, N., Gascuel, O., Jean-Marie, A., 2005. Markov models in molecular evolution, in: Nielsen, R. (Ed.), Statistical methods in molecular evolution. Springer New York, New York, pp. 3–24.

Gee, H., 2003. Ending incongruence. Nature 425, 782.

Glenn, T.C., 2011. Field guide to next-generation DNA sequencers. Mol. Ecol. Resour. 11, 759–769.

Goldman, N., Anderson, J.P., Rodrigo, A.G., 2000. Likelihood-based tests of topologies in phylogenetics. Syst. Biol. 49, 652–670.

Gustafson, G.T., Baca, S.M., Alexander, A.M., Short, A.E.Z., 2020. Phylogenomic analysis of the beetle suborder Adephaga with comparison of tailored and generalized ultraconserved element probe performance. Syst. Entomol. 45, 552–570.

Hamilton, C.A., Lemmon, A.R., Lemmon, E.M., Bond, J.E., 2016. Expanding anchored hybrid enrichment to resolve both deep and shallow relationships within the spider tree of life. BMC Evol. Biol. 16, 212.

Hasegawa, M., Kishino, H., Yano, T., 1985. Dating of the human-ape splitting by a molecular clock of mitochondrial DNA. J. Mol. Evol. 22, 160–174.

Hennig, W., 1969. Die Stammesgeschichte der Insekten, 1st ed. Waldemar Kramer, Frankfurt am Main.

Hennig, W., 1966. Phylogenetic Systematics, 1st ed. University of Illinois Press, Urbana.

Hillis, D.M., Mable, B.K., Moritz, C., 1996. Applications of molecular systematics, in: Hillis, D.M., Mable, B.K., Moritz, C. (Eds.), Molecular systematics. Sinauer Associates, Inc., Sunderland, Massachusetts, pp. 515–543.

Ho, S.Y.W., Duchêne, S., 2014. Molecular-clock methods for estimating evolutionary rates and timescales. Mol. Ecol. 23, 5947–5965.

Ho, S.Y.W., Jermiin, L.S., 2004. Tracing the decay of the historical signal in biological sequence data. Syst. Biol. 53, 623–637.

Hoang, D.T., Chernomor, O., von Haeseler, A., Minh, B.Q., Le, S.V., 2018. UFBoot2: improving the ultrafast bootstrap approximation. Mol. Biol. Evol. 35, 518–522.

Holder, M., Lewis, P.O., 2003. Phylogeny estimation: traditional and Bayesian approaches. Nat. Rev. Genet. 4, 275–284.

Jarvis, E.D., Mirarab, S., Aberer, A.J., Li, B., Houde, P., Li, C., Ho, S.Y.W., Faircloth, B.C., Nabholz, B., Howard, J.T., Suh, A., Weber, C.C., da Fonseca, R.R., Li, J.,

Zhang, F., Li, H., Zhou, L., Narula, N., Liu, L., Ganapathy, G., Boussau, B., Bayzid, M.S., Zavidovych, V., Subramanian, S., Gabaldón, T., Capella-Gutiérrez, S., Huerta-Cepas, J., Rekepalli, B., Munch, K., Schierup, M., Lindow, B., Warren, W.C., Ray, D., Green, R.E., Bruford, M.W., Zhan, X., Dixon, A., Li, S., Li, N., Huang, Y., Derryberry, E.P., Bertelsen, M.F., Sheldon, F.H., Brumfield, R.T., Mello, C. V, Lovell, P. V, Wirthlin, M., Schneider, M.P.C., Prosdocimi, F., Samaniego, J.A., Velazquez, A.M.V., Alfaro-Núñez, A., Campos, P.F., Petersen, B., Sicheritz-Ponten, T., Pas, A., Bailey, T., Scofield, P., Bunce, M., Lambert, D.M., Zhou, Q., Perelman, P., Driskell, A.C., Shapiro, B., Xiong, Z., Zeng, Y., Liu, S., Li, Z., Liu, B., Wu, K., Xiao, J., Yinqi, X., Zheng, Q., Zhang, Y., Yang, H., Wang, J., Smeds, L., Rheindt, F.E., Braun, M., Fjeldsa, J., Orlando, L., Barker, F.K., Jønsson, K.A., Johnson, W., Koepfli, K.-P., O'Brien, S., Haussler, D., Ryder, O.A., Rahbek, C., Willerslev, E., Graves, G.R., Glenn, T.C., McCormack, J., Burt, D., Ellegren, H., Alström, P., Edwards, S. V., Stamatakis, A., Mindell, D.P., Cracraft, J., Braun, E.L., Warnow, T., Jun, W., Gilbert, M.T.P., Zhang, G., 2014. Whole-genome analyses resolve early branches in the tree of life of modern birds. Science. 346, 1320–1331.

Jeffroy, O., Brinkmann, H., Delsuc, F., Philippe, H., 2006. Phylogenomics: the beginning of incongruence? Trends Genet. 22, 225–231.

Jermiin, L.S., Catullo, R.A., Holland, B.R., 2020. A new phylogenetic protocol: dealing with model misspecification and confirmation bias in molecular phylogenetics. NAR Genomics Bioinforma. 2. doi: 10.1093/nargab/lqaa041.

Jermiin, L.S., Ho, S.Y.W., Ababneh, F., Robinson, J., Larkum, A.W.D., 2004. The biasing effect of compositional heterogeneity on phylogenetic estimates may be underestimated. Syst. Biol. 53, 638–643.

Jermiin, L.S., Jayaswal, V., Ababneh, F., Robinson, J., 2008. Phylogenetic model evaluation, in: Keith, J.M. (Ed.), Bioinformatics. Methods in molecular Biology™, Vol 452. Humana Press, Totowa, pp. 331–363.

Jiao, X., Flouri, T., Rannala, B., Yang, Z., 2020. The impact of cross-species gene flow on species tree estimation. Syst. Biol. 69, 830–847.

Jones, M.R., Good, J.M., 2016. Targeted capture in evolutionary and ecological genomics. Mol. Ecol. 25, 185–202.

Julca, I., Marcet-Houben, M., Cruz, F., Vargas-Chavez, C., Johnston, J.S., Gómez-Garrido, J., Frias, L., Corvelo, A., Loska, D., Cámara, F., Gut, M., Alioto, T., Latorre, A., Gabaldón, T., 2020. Phylogenomics identifies an ancestral burst of gene duplications predating the diversification of Aphidomorpha. Mol. Biol. Evol. 37, 730–756.

Källersjö, M., Farris, J.S., Chase, M.W., Bremer, B., Fay, M.F., Humphries, C.J., Petersen, G., Seberg, O., Bremer, K., 1998. Simultaneous parsimony jackknife analysis of 2538rbcL DNA sequences reveals support for major clades of green plants, land plants, seed plants and flowering plants. Plant Syst. Evol. 213, 259–287.

Kapli, P., Yang, Z., Telford, M.J., 2020. Phylogenetic tree building in the genomic age. Nat. Rev. Genet. 21, 428–444.

Kristensen, N.P., 1999. Phylogeny of endopterygote insects, the most successful lineage of living organisms. Eur. J. Entomol. 96, 237–253.

Kück, P., Meusemann, K., Dambach, J., Thormann, B., von Reumont, B.M., Wägele, J.W., Misof, B., 2010. Parametric and non-parametric masking of randomness in sequence alignments can be improved and leads to better resolved trees. Front. Zool. 7, 10.

Lanfear, R., Calcott, B., Ho, S.Y.W., Guindon, S., 2012. PartitionFinder: combined selection of partitioning schemes and substitution models for phylogenetic analyses. Mol. Biol. Evol. 29, 1695–1701.

Lanfear, R., Frandsen, P.B., Wright, A.M., Senfeld, T., Calcott, B., 2017. Partitionfinder 2: new methods for selecting partitioned models of evolution for molecular and morphological phylogenetic analyses. Mol. Biol. Evol. 34, 772–773.

Lartillot, N., Philippe, H., 2004. A Bayesian mixture model for across-site heterogeneities in the amino-acid replacement process. Mol. Biol. Evol. 21, 1095–1109.

Laurin-Lemay, S., Brinkmann, H., Philippe, H., 2012. Origin of land plants revisited in the light of sequence contamination and missing data. Curr. Biol. 22, R593–R594.

Lean, C., Maclaurin, J., 2016. The value of phylogenetic diversity, in: Pellens, R., Grandcolas, P. (Eds.), Biodiversity conservation and phylogenetic systematics: preserving our evolutionary heritage in an extinction crisis. Springer International Publishing, Open access book, pp. 19–37.

Lemmon, A.R., Emme, S.A., Lemmon, E.M., 2012. Anchored hybrid enrichment for massively high-throughput phylogenomics. Syst. Biol. 61, 727–744.

Lemmon, E.M., Lemmon, A.R., 2013. High-throughput genomic data in systematics and phylogenetics. Annu. Rev. Ecol. Evol. Syst 44, 99–121.

Liu, L., Anderson, C., Pearl, D., Edwards, S. V., 2019. Modern phylogenomics: building phylogenetic trees using the multispecies coalescent model, in: Anisimova, M. (Ed.), Evolutionary genomics: statistical and computational methods. Humana Press, New York, pp. 211–239.

Liu, L., Yu, L., Edwards, S. V., 2010. A maximum pseudo-likelihood approach for estimating species trees under the coalescent model. BMC Evol. Biol. 10, 302.

Lopez, P., Casane, D., Philippe, H., 2002. Heterotachy, an important process of protein evolution. Mol. Biol. Evol. 19, 1–7.

Maddison, W.P., 1997. Gene trees in species trees. Syst. Biol. 46, 523–536.

Malé, P.J.G., Bardon, L., Besnard, G., Coissac, E., Delsuc, F., Engel, J., Lhuillier, E., Scotti-Saintagne, C., Tinaut, A., Chave, J., 2014. Genome skimming by shotgun sequencing helps resolve the phylogeny of a pantropical tree family. Mol. Ecol. Resour. 14, 966–975.

Mamanova, L., Coffey, A.J., Scott, C.E., Kozarewa, I., Turner, E.H., Kumar, A., Howard, E., Shendure, J., Turner, D.J., 2010. Target-enrichment strategies for next-generation sequencing. Nat. Methods 7, 111–118.

Marcet-Houben, M., Gabaldón, T., 2016. Horizontal acquisition of toxic alkaloid synthesis in a clade of plant associated fungi. Fungal Genet. Biol. 86, 71–80.

Margulies, M., Egholm, M., Altman, W.E., Attiya, S., Bader, J.S., Bemben, L.A., Berka, J., Braverman, M.S., Chen, Y.J., Chen, Z., Dewell, S.B., Du, L., Fierro, J.M.,

Gomes, X. V., Godwin, B.C., He, W., Helgesen, S., Ho, C.H., Irzyk, G.P., Jando, S.C., Alenquer, M.L.I., Jarvie, T.P., Jirage, K.B., Kim, J.B., Knight, J.R., Lanza, J.R., Leamon, J.H., Lefkowitz, S.M., Lei, M., Li, J., Lohman, K.L., Lu, H., Makhijani, V.B., McDade, K.E., McKenna, M.P., Myers, E.W., Nickerson, E., Nobile, J.R., Plant, R., Puc, B.P., Ronan, M.T., Roth, G.T., Sarkis, G.J., Simons, J.F., Simpson, J.W., Srinivasan, M., Tartaro, K.R., Tomasz, A., Vogt, K.A., Volkmer, G.A., Wang, S.H., Wang, Y., Weiner, M.P., Yu, P., Begley, R.F., Rothberg, J.M., 2005. Genome sequencing in microfabricated high-density picolitre reactors. Nature 437, 376–380.

McCormack, J.E., Hird, S.M., Zellmer, A.J., Carstens, B.C., Brumfield, R.T., 2013. Applications of next-generation sequencing to phylogeography and phylogenetics. Mol. Phylogenet. Evol. 66, 526–538.

McKenna, D.D., Shin, S., Ahrens, D., Balke, M., Beza-Beza, C., Clarke, D.J., Donath, A., Escalona, H.E., Friedrich, F., Letsch, H., Liu, S., Maddison, D., Mayer, C., Misof, B., Murin, P.J., Niehuis, O., Peters, R.S., Podsiadlowski, L., Pohl, H., Scully, E.D., Yan, E. V, Zhou, X., Ślipiński, A., Beutel, R.G., 2019. The evolution and genomic basis of beetle diversity. Proc. Natl. Acad. Sci. U. S. A. 116, 24729–24737.

Meusemann, K., von Reumont, B.M., Simon, S., Roeding, F., Strauss, S., Kück, P., Ebersberger, I., Walzl, M., Pass, G., Breuers, S., Achter, V., von Haeseler, A., Burmester, T., Hadrys, H., Wägele, J.W., Misof, B., 2010. A phylogenomic approach to resolve the arthropod tree of life. Mol. Biol. Evol. 27, 2451–2464.

Minh, B.Q., Hahn, M.W., Lanfear, R., 2020. New methods to calculate concordance factors for phylogenomic datasets. Mol. Biol. Evol. 37, 2727–2733.

Mirarab, S., Reaz, R., Bayzid, M.S., Zimmermann, T., S. Swenson, M., Warnow, T., 2014. ASTRAL: genome-scale coalescent-based species tree estimation. Bioinformatics 30, 541–548.

Misof, B., Liu, S., Meusemann, K., Peters, R.S., Donath, A., Mayer, C., Frandsen, P.B., Ware, J., Flouri, T., Beutel, R.G., Niehuis, O., Petersen, M., Izquierdo-Carrasco, F., Wappler, T., Rust, J., Aberer, A.J., Aspöck, U., Aspöck, H., Bartel, D., Blanke, A., Berger, S., Bohm, A., Buckley, T.R., Calcott, B., Chen, J., Friedrich, F., Fukui, M., Fujita, M., Greve, C., Grobe, P., Gu, S., Huang, Y., Jermiin, L.S., Kawahara, A.Y., Krogmann, L., Kubiak, M., Lanfear, R., Letsch, H., Li, Y., Li, Z., Li, J., Lu, H., Machida, R., Mashimo, Y., Kapli, P., McKenna, D.D., Meng, G., Nakagaki, Y., Navarrete-Heredia, J.L., Ott, M., Ou, Y., Pass, G., Podsiadlowski, L., Pohl, H., von Reumont, B.M., Schutte, K., Sekiya, K., Shimizu, S., Slipinski, A., Stamatakis, A., Song, W., Su, X., Szucsich, N.U., Tan, M., Tan, X., Tang, M., Tang, J., Timelthaler, G., Tomizuka, S., Trautwein, M., Tong, X., Uchifune, T., Walzl, M.G., Wiegmann, B.M., Wilbrandt, J., Wipfler, B., Wong, T.K.F., Wu, Q., Wu, G., Xie, Y., Yang, S., Yang, Q., Yeates, D.K., Yoshizawa, K., Zhang, Q., Zhang, R., Zhang, W., Zhang, Y., Zhao, J., Zhou, C., Zhou, L., Ziesmann, T., Zou, S., Xu, X., Yang, H., Wang, J., Kjer, K.M., Zhou, X., 2014. Phylogenomics resolves the timing and pattern of insect evolution. Science. 346, 763–767.

Money, D., Whelan, S., 2012. Characterizing the phylogenetic tree-search problem. Syst. Biol. 61, 228.

Montagna, M., Tong, K.J., Magoga, G., Strada, L., Tintori, A., Ho, S.Y.W., Lo, N., 2019. Recalibration of the insect evolutionary time scale using Monte San Giorgio fossils suggests survival of key lineages through the End-Permian extinction. Proc. R. Soc. B Biol. Sci. 286, 20191854.

Moritz, C., Hillis, D.M., 1996. Molecular systematics: context and controversies, in: Hillis, D.M., Mable, B.K., Moritz, C. (Eds.), Molecular systematics. Sinauer Associates, Inc., Sunderland, Massachusetts, pp. 1–13.

Natsidis, P., Kapli, P., Schiffer, P.H., Telford, M.J., 2020. Systematic errors in orthology inference: a bug or a feature for evolutionary analyses? BioRxiv. doi: 10.1101/2020.11.03.366625.

Oswald, J.D., 2019. LDL Neuropterida Species of the World (version July 2018) [WWW Document]. Species 2000 ITIS Cat. Life, 26th Febr. 2019. URL www.catalogueoflife.org/col (accessed 3.12.19).

Pease, J.B., Brown, J.W., Walker, J.F., Hinchliff, C.E., Smith, S.A., 2018. Quartet Sampling distinguishes lack of support from conflicting support in the green plant tree of life. Am. J. Bot. 105, 385–403.

Peters, R.S., Meusemann, K., Petersen, M., Mayer, C., Wilbrandt, J., Ziesmann, T., Donath, A., Kjer, K.M., Aspöck, U., Aspöck, H., Aberer, A., Stamatakis, A., Friedrich, F., Hünefeld, F., Niehuis, O., Beutel, R.G., Misof, B., 2014. The evolutionary history of holometabolous insects inferred from transcriptome-based phylogeny and comprehensive morphological data. BMC Evol. Biol. 14, 52.

Petersen, M., Meusemann, K., Donath, A., Dowling, D., Liu, S., Peters, R.S., Podsiadlowski, L., Vasilikopoulos, A., Zhou, X., Misof, B., Niehuis, O., 2017.

Orthograph: a versatile tool for mapping coding nucleotide sequences to clusters of orthologous genes. BMC Bioinformatics 18, 111.

Philippe, H., Brinkmann, H., Lavrov, D. V., Littlewood, D.T.J., Manuel, M., Wörheide, G., Baurain, D., 2011. Resolving difficult phylogenetic questions: why more sequences are not enough. PLoS Biol. 9, e1000602.

Philippe, H., de Vienne, D.M., Ranwez, V., Roure, B., Baurain, D., Delsuc, F., 2017. Pitfalls in supermatrix phylogenomics. Eur. J. Taxon. 283, 1–25.

Rannala, B., Yang, Z., 2007. Inferring speciation times under an episodic molecular clock. Syst. Biol. 56, 453–466.

Rannala, B., Yang, Z., 1996. Probability distribution of molecular evolutionary trees: a new method of phylogenetic inference. J. Mol. Evol. 43, 304–311.

Richter, S., Schwarz, F., Hering, L., Böggemann, M., Bleidorn, C., 2015. The utility of genome skimming for phylogenomic analyses as demonstrated for glycerid relationships (Annelida, Glyceridae). Genome Biol. Evol. 7, 3443–3462.

Roch, S., Steel, M., 2015. Likelihood-based tree reconstruction on a concatenation of aligned sequence data sets can be statistically inconsistent. Theor. Popul. Biol. 100, 56–62.

Roch, S., Warnow, T., 2015. On the robustness to gene tree estimation error (or lack thereof) of coalescent-based species tree methods. Syst. Biol. 64, 663–676.

Rodríguez-Ezpeleta, N., Brinkmann, H., Roure, B., Lartillot, N., Lang, B.F., Philippe, H., 2007. Detecting and overcoming systematic errors in genome-scale phylogenies. Syst. Biol. 56, 389–399.

Rodríguez, A., Burgon, J.D., Lyra, M., Irisarri, I., Baurain, D., Blaustein, L., Göçmen, B., Künzel, S., Mable, B.K., Nolte, A.W., Veith, M., Steinfartz, S., Elmer, K.R., Philippe, H., Vences, M., 2017. Inferring the shallow phylogeny of true salamanders (Salamandra) by multiple phylogenomic approaches. Mol. Phylogenet. Evol. 115, 16–26.

Rokas, A., Williams, B.L., King, N., Carroll, S.B., 2003. Genome-scale approaches to resolving incongruence in molecular phylogenies. Nature 425, 798–804.

Schmidt, H.A., 2009. Testing tree topologies, in: Lemey, P., Salemi, M., Vandamme, A.-M. (Eds.), The phylogenetic handbook: a practical approach to phylogenetic analysis and hypothesis testing. Cambridge University Press, New York, pp. 381–403.

Sela, I., Ashkenazy, H., Katoh, K., Pupko, T., 2015. GUIDANCE2: accurate detection of unreliable alignment regions accounting for the uncertainty of multiple parameters. Nucleic Acids Res. 43, W7–W14.

Shen, X.-X., Hittinger, C.T., Rokas, A., 2017. Contentious relationships in phylogenomic studies can be driven by a handful of genes. Nat. Ecol. Evol. 1, 0126.

Shepherd, D.A., Klaere, S., 2019. How well does your phylogenetic model fit your data? Syst. Biol. 68, 157–167.

Shimodaira, H., 2002. An approximately unbiased test of phylogenetic tree selection. Syst. Biol. 51, 492–508.

Shimodaira, H., Hasegawa, M., 1999. Multiple comparisons of log-likelihoods with applications to phylogenetic inference. Mol. Biol. Evol. 16, 1114–1116.

Simion, P., Belkhir, K., François, C., Veyssier, J., Rink, J.C., Manuel, M., Philippe, H., Telford, M.J., 2018. A software tool "CroCo" detects pervasive cross-species contamination in next generation sequencing data. BMC Biol. 16, 28.

Simion, P., Delsuc, F., Philippe, H., 2020. To what extent current limits of phylogenomics can be overcome?, in: Scornavacca, C., Delsuc, F., Galtier, N. (Eds.), Phylogenetics in the genomic era. Authors' open access book, pp. 2.1:1–2.1:34. hal-02535651.

Smith, M.L., Hahn, M.W., 2020. New approaches for inferring phylogenies in the presence of paralogs. Trends Genet. doi:10.1016/j.tig.2020.08.012.

Smith, S.A., Pease, J.B., 2017. Heterogeneous molecular processes among the causes of how sequence similarity scores can fail to recapitulate phylogeny. Brief. Bioinform. 18, 451–457.

Straub, S.C.K., Parks, M., Weitemier, K., Fishbein, M., Cronn, R.C., Liston, A., 2012. Navigating the tip of the genomic iceberg: next-generation sequencing for plant systematics. Am. J. Bot. 99, 349–364.

Strimmer, K., von Haeseler, A., 2009. Genetic distances and nucleotide substitution models, in: Lemey, P., Salemi, M., Vandamme, A.-M. (Eds.), The phylogenetic handbook: a practical approach to phylogenetic analysis and hypothesis testing. Cambridge University Press, New York, pp. 111–140.

Strimmer, K., von Haeseler, A., 1997. Likelihood-mapping: a simple method to visualize phylogenetic content of a sequence alignment. Proc. Natl. Acad. Sci. U. S. A. 94, 6815–6819.

Sullivan, J., Joyce, P., 2005. Model selection in phylogenetics. Annu. Rev. Ecol. Evol. Syst. 36, 445–466.

Tatusov, R.L., Koonin, E. V., Lipman, D.J., 1997. A genomic perspective on protein families. Science. 278, 631–637.

Tavaré, S., 1986. Some probabilistic and statistical problems in the analysis of DNA sequences. Lect. Math. Life Sci. 17, 57–86.

Thorne, J.L., Kishino, H., Painter, I.S., 1998. Estimating the rate of evolution of the rate of molecular evolution. Mol. Biol. Evol. 15, 1647–1657

Van de Peer, Y., 2009. Phylogenetic inference based on distance methods, in: Lemey, P., Salemi, M., Vandamme, A.-M. (Eds.), The phylogenetic handbook: a practical approach to phylogenetic analysis and hypothesis testing. Cambridge University Press, New York, pp. 142–180.

Wägele, J.W., 2005. Foundations of phylogenetic systematics, 2nd ed. Verlag Dr. Friedrich Pfeil, Munich.

Wang, H.-C., Susko, E., Roger, A.J., 2019. The relative importance of modeling site pattern heterogeneity versus partition-wise heterotachy in phylogenomic inference. Syst. Biol. 68, 1003–1019.

Wang, Y., Liu, X., Garzón-Orduña, I.J., Winterton, S.L., Yan, Y., Aspöck, U., Aspöck, H., Yang, D., 2017. Mitochondrial phylogenomics illuminates the evolutionary history of Neuropterida. Cladistics 33, 617–636.

Wang, Y., Zhou, X., Wang, L., Liu, X., Yang, D., Rokas, A., 2019. Gene selection and evolutionary modeling affect phylogenomic inference of Neuropterida based on transcriptome data. Int. J. Mol. Sci. 20, 1072.

Wheeler, W.C., Whiting, M., Wheeler, Q.D., Carpenter, J.M., 2001. The phylogeny of the extant hexapod orders. Cladistics 17, 113–169.

Wickett, N.J., Mirarab, S., Nguyen, N., Warnow, T., Carpenter, E., Matasci, N., Ayyampalayam, S., Barker, M.S., Burleigh, J.G., Gitzendanner, M.A., Ruhfel, B.R., Wafula, E., Der, J.P., Graham, S.W., Mathews, S., Melkonian, M., Soltis, D.E., Soltis, P.S., Miles, N.W., Rothfels, C.J., Pokorny, L., Shaw, A.J., DeGironimo, L., Stevenson, D.W., Surek, B., Villarreal, J.C., Roure, B., Philippe, H., DePamphilis, C.W., Chen, T., Deyholos, M.K., Baucom, R.S., Kutchan, T.M., Augustin, M.M., Wang, J., Zhang, Y., Tian, Z., Yan, Z., Wu, X., Sun, X., Wong, G.K.-S., Leebens-Mack, J., 2014. Phylotranscriptomic analysis of the origin and early diversification of land plants. Proc. Natl. Acad. Sci. U. S. A. 111, E4859–E4868.

Winterton, S.L., Hardy, N.B., Wiegmann, B.M., 2010. On wings of lace: phylogeny and Bayesian divergence time estimates of Neuropterida (Insecta) based on morphological and molecular data. Syst. Entomol. 35, 349–378.

Winterton, S.L., Lemmon, A.R., Gillung, J.P., Garzon, I.J., Badano, D., Bakkes, D.K., Breitkreuz, L.C.V., Engel, M.S., Lemmon, E.M., Liu, X., Machado, R.J.P., Skevington, J.H., Oswald, J.D., 2018. Evolution of lacewings and allied orders using anchored phylogenomics (Neuroptera, Megaloptera, Raphidioptera). Syst. Entomol. 43, 330–354.

Xu, B., Yang, Z., 2016. Challenges in species tree estimation under the multispecies coalescent model. Genetics 204, 1353–1368.

Yamada, K.D., Tomii, K., Katoh, K., 2016. Application of the MAFFT sequence alignment program to large data - reexamination of the usefulness of chained guide trees. Bioinformatics 32, 3246–3251.

Yang, Z., 2014. Molecular evolution: a statistical approach, 1st ed. Oxford University Press, Oxford.

Yang, Z., 1996. Among-site rate variation and its impact on phylogenetic analyses. Trends Ecol. Evol. 11, 367–372.

Yang, Z., 1994. Maximum likelihood phylogenetic estimation from DNA sequences with variable rates over sites: approximate methods. J. Mol. Evol. 39, 306–314.

Yang, Z., Rannala, B., 2012. Molecular phylogenetics: principles and practice. Nat. Rev. Genet. 13, 303–314.

Yoder, A.D., Yang, Z., 2000. Estimation of primate speciation dates using local molecular clocks. Mol. Biol. Evol. 17, 1081–1090.

Young, A.D., Gillung, J.P., 2020. Phylogenomics — principles, opportunities and pitfalls of big-data phylogenetics. Syst. Entomol. 45, 225–247.

Zhang, Y.M., Buffington, M.L., Looney, C., László, Z., Shorthouse, J.D., Ide, T., Lucky, A., 2020. UCE data reveal multiple origins of rose gallers in North America: global phylogeny of *Diplolepis* Geoffroy (Hymenoptera: Cynipidae). Mol. Phylogenet. Evol. 153, 106949.

Zhou, X., Lutteropp, S., Czech, L., Stamatakis, A., von Looz, M., Rokas, A., 2020. Quartet-based computations of internode certainty provide robust measures of phylogenetic incongruence. Syst. Biol. 69, 308–324.

Zuckerkandl, E., Pauling, L., 1965. Evolutionary divergence and convergence in proteins, in: Bryson, V., Vogel, H.J. (Eds.), Evolving genes and proteins. Academic Press, New York, pp. 97–166.

# 2. An integrative phylogenomic approach to elucidate the evolutionary history and divergence times of Neuropterida (Insecta: Holometabola)

This chapter is published in the following article (included here with minor formatting edits such as numbering of headings):

Vasilikopoulos, A., Misof, B., Meusemann, K., Lieberz, D., Flouri, T., Beutel, R.G., Niehuis, O., Wappler, T., Rust, J., Peters, R.S., Donath, A., Podsiadlowski, L., Mayer, C., Bartel, D., Böhm, A., Liu, S., Kapli, P., Greve, C., Jepson, J.E., Liu, X., Zhou, X., Aspöck, H., Aspöck, U., 2020. An integrative phylogenomic approach to elucidate the evolutionary history and divergence times of Neuropterida (Insecta: Holometabola). BMC Evol. Biol. 20, 64. https://doi.org/10.1186/s12862-020-01631-6


A correction to this article is published in:

Vasilikopoulos, A., Misof, B., Meusemann, K., Lieberz, D., Flouri, T., Beutel, R.G., Niehuis, O., Wappler, T., Rust, J., Peters, R.S., Donath, A., Podsiadlowski, L., Mayer, C., Bartel, D., Böhm, A., Liu, S., Kapli, P., Greve, C., Jepson, J.E., Liu, X., Zhou, X., Aspöck, H., Aspöck, U., 2020. Correction to: An integrative phylogenomic approach to elucidate the evolutionary history and divergence times of Neuropterida (Insecta: Holometabola). BMC Evol. Biol. 20, 133. https://doi.org/10.1186/s12862-020-01695-4


Authors' contributions to the original article:

AV, BM, KM, RGB, ON, RSP, XZ, HA and UA conceived and designed the study. BM and XZ contributed to funding acquisition. AV, BM, KM, DL, TF, RGB, TW, JR, AD, LP, DB, AB, SL, PL and JEJ performed research. KM, ON, RSP, CG, XL, HA, UA collected and processed samples. BM, KM, RSP, SL, and CG contributed to biosample processing, SL and XZ organized transcriptome sequencing and transcriptome assembly. AD, LP and KM organized sequence submissions to NCBI. AV, BM, KM,

DL, RGB, LP, DB, AB, and SL analyzed data. TF and CM developed new bioinformatic workflows. AV, BM, KM, DL, RGB, ON, TW, JR, XL, HA and UA wrote the manuscript draft. All authors contributed with comments and suggestions to the final manuscript version.

## 2.1. Background

The insect superorder Neuropterida contains more than 6500 described and extant species that are classified into three holometabolous insect orders: Megaloptera (alderflies, dobsonflies and fishflies), Neuroptera (lacewings, antlions and relatives) and Raphidioptera (snakeflies). Among these three, Neuroptera is by far the most species-rich order with 5917 species, in comparison to the much less diverse Megaloptera and Raphidioptera (386 and 253 species respectively) (Oswald, 2019). Within Holometabola, Neuropterida is considered the sister group of Coleopterida, and both together form the clade Neuropteroidea (or Neuropteriformia) (Misof et al., 2014; Peters et al., 2014; Wiegmann et al., 2009). Overall, the monophyly of Neuropterida is well established but morphological evidence in support of this monophyly is only based on a small number of inconspicuous characters (summarized by Aspöck, 2002 and by Aspöck et al., 1980). The phylogenetic relationships of neuropterid insects have received considerable attention based on the analyses of different types of data such as the anatomy of adults (Aspöck et al., 2001; Aspöck and Aspöck, 2008; Beutel et al., 2010b; Randolf et al., 2017, 2014, 2013), or the anatomy of larvae (Aspöck et al., 2001; Badano et al., 2017; Beutel et al., 2010a; MacLeod, 1964). Other studies have combined morphological and molecular evidence in a phylogenetic framework (Winterton et al., 2010; Yang et al., 2012), and recently several studies have analyzed

genome-scale molecular datasets (Cameron et al., 2009; Song et al., 2019; Wang et al., 2017; Y. Wang et al., 2019; Winterton et al., 2018; Zhao et al., 2013). These phylogenomic studies have included analyses of different types of data such as hybrid enrichment data (Machado et al., 2019; Winterton et al., 2018), mitochondrial genome sequences (Cameron et al., 2009; Song et al., 2019; Wang et al., 2017; Zhao et al., 2013), and transcriptomic data (Y. Wang et al., 2019). Analyses of these types of data did not reach a full consensus on the phylogenetic relationships of Neuropterida, specifically concerning the backbone tree of Neuroptera. Here, we present the largest dataset of phylogenetically informative molecular characters compiled to date, across a large number of neuropterid and outgroup species, in an attempt to resolve the existing phylogenetic uncertainties in the phylogeny of Neuropterida and infer the temporal pattern of diversification within the group. A further important goal of this study is to identify sources of phylogenetic signal in the data and assess the effects of confounding factors on the phylogenetic reconstructions, in order to identify methodological problems behind open questions or conflicting phylogenetic results.

Recent phylogenetic investigations of Neuropterida have converged on the hypothesis that the order Raphidioptera is sister to a clade comprising Megaloptera and Neuroptera (Aspöck and Aspöck, 2008; Cameron et al., 2009; Haring et al., 2011; Haring and Aspöck, 2004; Wang et al., 2017; Winterton et al., 2018, 2010; Zhao et al., 2014). Raphidioptera is a relict group of holometabolous insects with most of its species geographically distributed over small areas in the northern hemisphere (except eastern North America) (H. Aspöck, 2002; Haring et al., 2011). Owing to their distinctly higher species diversity in the Mesozoic, and their very limited morphological divergence since then, some authors refer to them as "living fossils"

(Aspöck, 2000, 1998; H. Aspöck, 2002; Aspöck and Aspöck, 2007; Winterton et al., 2018). The order is divided into two extant families: Raphidiidae (209 described extant species) and Inocelliidae (44 described extant species) (Oswald, 2019). The monophyly of Raphidioptera and of each raphidiopteran family is well established. However, previous phylogenomic analyses of Neuropterida have suffered from taxon-sampling limitations within the order (Wang et al., 2017; Y. Wang et al., 2019; Winterton et al., 2018). Therefore, a comprehensive phylogenetic analysis of snakeflies based on the analysis of genomic sequence data has yet to be performed. The order Megaloptera comprises two extant families: Corydalidae (Corydalinae: dobsonflies and Chauliodinae: fishflies with 303 described extant species in total) and Sialidae (alderflies: 83 described extant species) (Oswald, 2019). This order includes the oldest known holometabolous insects with an aquatic lifestyle of the larvae (Rivera-Gasperín et al., 2019). The monophyly of Megaloptera has been questioned before (Achtelig, 1967; Beutel et al., 2011; Winterton et al., 2010; Yang et al., 2012), as has been the monophyly of the family Corydalidae (Contreras-Ramos, 2004). Nevertheless, recent morphological and molecular evidence suggests that Corydalidae and Sialidae are monophyletic sister taxa within the monophyletic Megaloptera (Aspöck and Aspöck, 2008; Liu et al., 2016; Wang et al., 2017; Winterton et al., 2018).

The order Neuroptera comprises 16 extant families. In comparison to the adults, the larvae of Neuroptera have evolved a very broad spectrum of morphological adaptations to very different habitats and lifestyles (Winterton et al., 2018, 2010). Only two neuropteran families contain species with strictly aquatic larvae (i.e., Nevrorthidae, Sisyridae) (Aspöck et al., 2017; Winterton et al., 2018). The larvae of Sisyridae (spongillaflies) use freshwater bryozoans and sponges as hosts, whereas the larvae of

Nevrorthidae (mermaids) are generalist benthic predators (Aspöck et al., 2017; Winterton et al., 2010). Other remarkable adaptations of the larvae within Neuroptera include predators of termites (some Berothidae) (Brushwein, 1987; Komatsu, 2014; Tauber and Tauber, 1968), parasitoids of bees and wasps (Mantispidae: some Symphrasinae) (Dejean and Canard, 1990), predators of spider eggs (Mantispidae: Mantispinae) (Redborg, 1998; Schremmer, 1983), fossorial pit-trap builders (some Myrmeleontidae) (Badano et al., 2017; Engel et al., 2018; X. Liu et al., 2015; Winterton et al., 2018, 2010), and possibly also phytophagous root suckers (Ithonidae, *Oliarces*) (Faulkner, 1990). The monophyly of Neuroptera has never been questioned and is strongly supported by the unique and complex sucking tubes of the larvae (Aspöck and Aspöck, 2007; Winterton et al., 2018). However, there is currently a lack of consensus on the phylogeny of neuropteran families mainly because analyses of different types of phylogenomic data have suggested conflicting topologies. In addition, the morphological characters of the adults are affected by homoplasy (Beutel et al., 2010a; Randolf et al., 2017) and although larval morphology yields important information, the phylogenetic signal from analyzing larval characters appears to be partly eroded (Wang et al., 2017; Winterton et al., 2018, 2010), probably due to far-reaching specialization, especially in the case of the miniaturized Coniopterygidae (dustywings).

Concerning the phylogeny of neuropteran families, conflicting phylogenetic results have emerged both among different molecular studies (Wang et al., 2017; Winterton et al., 2018) as well as among different datasets or methods applied within the same study (Winterton et al., 2018). One example of conflicting hypotheses concerns the monophyly, or non-monophyly, of the suborder Myrmeleontiformia (Wang et al., 2017; Winterton et al., 2018). Myrmeleontiformia contains the five

families Ascalaphidae (owlflies), Myrmeleontidae (antlions), Nemopteridae (thread-winged lacewings), Nymphidae (split-footed lacewings) and Psychopsidae (silky lacewings). The family Psychopsidae is most likely the sister group to all remaining Myrmeleontiformia, as suggested by analyses of morphological characters (Badano et al., 2018; Beutel et al., 2010b, 2010a; Engel et al., 2018; Jandausch et al., 2018). It should, however, be noted that similar complex male genital sclerites of Psychopsidae and Nemopteridae have been interpreted as synapomorphies indicating a possible sister group relationship of these two families (Aspöck and Aspöck, 2008). Recently, target DNA enrichment-based phylogenomic analyses suggested a clade of Ithonidae (moth lacewings) + Nymphidae, implying paraphyletic Myrmeleontiformia (Machado et al., 2019; Winterton et al., 2018). In contrast, phylogenetic analyses of mitochondrial genomes did not corroborate this result but suggested monophyletic Myrmeleontiformia (Song et al., 2019; Wang et al., 2017). Other conflicting hypotheses among previous phylogenomic studies include the disruption, or not, of a clade comprising Chrysopidae (green lacewings) and Hemerobiidae (brown lacewings) and the exact affinities of these two families to a clade of Ithonidae + Myrmeleontiformia (Song et al., 2019; Wang et al., 2017; Winterton et al., 2018). A clade comprising Mantispidae (mantid lacewings), Berothidae (beaded lacewings), and Rhachiberothidae (thorny lacewings), collectively referred to as Mantispoidea (Aspöck et al., 2001; Winterton et al., 2018), was recovered in all previous phylogenomic studies, but the exact placement of this clade within Neuroptera remains elusive. Lastly, the inter-relationships of Osmylidae (lance lacewings), Nevrorthidae, and Sisyridae also remain unresolved. All previous phylogenomic studies suggested that these three families branch off close to the base of the neuropteran tree, but reconstructed different

topologies among these groups (Song et al., 2019; Wang et al., 2017; Winterton et al., 2018).

Despite the above-outlined discrepancies among phylogenomic studies, some results seem to be robust across phylogenomic studies, but they are in conflict with the results of morphological studies. Such conflicts include the phylogenetic placement of Coniopterygidae as sister to the remaining families of Neuroptera, as suggested by previous analyses of genomic sequence data, but also by analyses of a small number of molecular markers (Winterton et al., 2010), or by total evidence analyses (Yang et al., 2012). Most cladistic analyses of morphological characters instead suggest that Nevrorthidae is the sister group to all other neuropteran families (Aspöck et al., 2001; Aspöck and Aspöck, 2008; Beutel et al., 2010b, 2010a; Jandausch et al., 2018). The family Sisyridae has also been proposed as sister to all other Neuroptera based on the analysis of morphological characters (Randolf et al., 2013). A consensus on the basal splitting patterns within Neuroptera is essential for inferring the ancestral lifestyle of the neuropteran larvae, and also for tracing morphological character evolution within the order (Winterton et al., 2018). Most importantly, the paraphyly of Myrmeleontiformia as suggested by target DNA enrichment-based phylogenomic studies, was a surprising result especially given the long-lasting (Aspöck et al., 2012a) and strong support of morphological studies in favor of monophyletic Myrmeleontiformia. Hence, a reevaluation of the previously proposed paraphyly of Myrmeleontiformia based on other kinds of data or methods is needed (Badano et al., 2018).

Previous molecular studies of the phylogeny of Neuropterida have mostly relied on conventional measures of branch support, such as the non-parametric bootstrap

(Felsenstein, 1985) and the Bayesian posterior probabilities (Rannala and Yang, 1996). However, the usage of these measures alone has often proven insufficient for the purpose of estimating the robustness of the inferred molecular phylogenies (Evangelista et al., 2018; Johnson et al., 2018; Salichos and Rokas, 2013; Simmons et al., 2004; Simmons and Norton, 2014; Wägele et al., 2009), especially when the size of the dataset increases (Cloutier et al., 2019; Dell'Ampio et al., 2014; Gadagkar et al., 2005; Seo, 2008; Simmons, 2012), or when overly simplified evolutionary models are used (Huelsenbeck and Rannala, 2004; Lemmon and Moriarty, 2004). A plethora of quartet-based approaches for estimating phylogenomic incongruence and node certainty in molecular phylogenies has been proposed lately (Johnson et al., 2018; Kück et al., 2017; Misof et al., 2014; Pease et al., 2018; Sayyari and Mirarab, 2016; Zhou et al., 2020). These approaches rely on the calculation of phylogenetic signal from quartets of taxa and they can be used to identify conflicting signals and potentially inflated support for certain phylogenetic clades, but have not yet been applied to the phylogeny of Neuropterida. Given the putatively misleading nature of the existing branch support measures in a maximum likelihood or Bayesian phylogenetic framework, combined with the incongruent results of previous phylogenomic studies, a thorough evaluation of the conflicts in the phylogenetic tree of Neuropterida is currently needed.

The purpose of this study is to provide: 1) a phylogenomic framework and updated divergence time estimates of Neuropterida, 2) an evaluation of conflicting phylogenetic signals in the backbone phylogeny of the group, and 3) a discussion of the implications for morphological character evolution within Neuropterida based on the results of the present contribution and those of other studies. In an effort to resolve the existing incongruencies we assembled a novel transcriptomic dataset of Neuropterida

and of suitable outgroup species, and assessed the robustness of our phylogenetic estimates with concatenation-based quartet approaches combined with data permutations and with gene tree-based quartet approaches. We additionally estimated divergence times of the major lineages of Neuropterida by using an approach that enables monitoring the effect of data selection on the Bayesian posterior divergence times of Neuropterida.

## 2.2. Results

### 2.2.1. Orthology assignment, alignment refinement, protein domain identification and supermatrix evaluation

On average, 3292 sequences per transcriptome or official gene set (OGS) passed the reciprocal best-hit criterion during the orthology assignment step (max. = 3909, min. = 1935). We excluded a total number of 21 transcriptomes and OGSs from our dataset because we found too few target genes (orthologs) within them (Additional file 1:Table S1). The majority of the excluded transcriptomes and OGSs refer to outgroup taxa (17 outgroup and four ingroup species). Alignment masking resulted in removal of a total number of 1,307,572 alignment sites at the amino-acid sequence level (~ 45% of alignment sites). Concatenation of the masked amino-acid sequence alignments resulted in a supermatrix composed of 6869 domain-based partitions spanning more than 1.5 million amino-acid alignment sites (supermatrix A, Table 2.1). Supermatrices E and F did not significantly differ in their overall completeness, data coverage in terms of presence/absence of partitions (i.e., saturation, Table 2.1), information content and deviation from stationary, (time-) reversible and homogeneous (SRH) conditions (Table 2.1). We selected supermatrix E for downstream analyses due to its larger size in terms

of total alignment length and number of partitions (see Additional file 2). The optimization of the partitioning scheme of supermatrix E with the software PartitionFinder resulted in a total number of 1825 meta-partitions.

**Table 2.1**: Descriptive statistics for each of the analyzed amino-acid supermatrices that were partitioned according to protein-domain clans, protein families and to single protein domains. Information content calculated with the software MARE is a relative measure of phylogenetic informativeness and data coverage. Completeness scores calculated with AliStat indicate the proportion of non-ambiguous characters.

| Amino-acid supermatrix | No. of alignment sites | No. of domain-based partitions or meta-partitions | No. of species | Inform-ation content (MARE) | Saturation (MARE) | Completeness score ($C_a$) (AliStat) | Median pairwise *p*-value for the Bowker's test (SymTest) |
|---|---|---|---|---|---|---|---|
| A | 1,550,004 | 6869 partitions | 121 | 0.432 | 0.804 | 0.628 | 2.22e-141 |
| B | 1,087,525 | 4261 partitions | 119 | 0.636 | 0.909 | 0.659 | 8.22e-092 |
| C | 1,506,256 | 5353 partitions | 121 | 0.554 | 0.820 | 0.628 | 4.46e-137 |
| D | 1,506,256 | 5353 partitions | 119 | 0.557 | 0.826 | 0.635 | 8.68e-137 |
| E | 931,450 | 3635 partitions | 119 | 0.667 | 0.923 | 0.657 | 8.13e-068 |
| F | 920,182 | 3603 partitions | 119 | 0.669 | 0.923 | 0.657 | 1.40e-066 |
| E (RCFV-corrected) | 383,656 | 314 (meta-partitions) | 119 | 0.662 | 0.997 | 0.713 | 9.33e-018 |
| E (Decisive) | 228,933 | 209 (meta-partitions) | 119 | 0.619 | 1.000 | 0.796 | 3.29e-013 |

## 2.2.2. Phylogeny of Neuropterida: concatenation-based and summary coalescent phylogenetic analyses

Phylogenetic analyses of the domain-based partitioned amino-acid sequence data yielded congruent topologies (with respect to the phylogenetic relationships of major lineages) with those obtained when analyzing the second codon positions of the nucleotide sequence data (Fig. 2.1, Additional file 3: Figures S1–S5). In addition, the phylogenetic trees yielded by the analyses of the reduced amino-acid supermatrices (decisive and RCFV-corrected versions of supermatrix E, Table 2.1) are topologically congruent with trees that resulted from the analyses of the above-mentioned datasets, concerning the phylogenetic relationships within Neuropterida (Additional file 3: Figures S6–S9). Analyses with the site-heterogeneous mixture models also delivered topologies congruent to the analyses of the above-mentioned datasets (Additional file 3: Figures S10–S14). All these analyses support Coleopterida (Coleoptera + Strepsiptera) as sister to Neuropterida, the monophyly of all neuropterid orders and families, and the sister group relationship between Raphidioptera and Megaloptera + Neuroptera (Fig. 2.1, Additional file 3: Figures S1–S14).

The inferred relationships within Raphidioptera suggest the monophyly of the family Raphidiidae, placement of the Nearctic genus *Agulla* as sister to a clade comprising all the Palearctic Raphidiidae. These relationships received maximum bootstrap and maximum bootstrap by transfer (TBE) support (Fig. 2.1, Additional file 3: Figure S2). Within the Palearctic Raphidiidae the genus Mongoloraphidia was inferred as the sister taxon to all remaining Raphidiidae. Within Neuroptera, a sister group relationship between Coniopterygidae and all remaining neuropteran families received maximum bootstrap and maximum TBE support (Fig. 2.1, Additional file 3:

Figure S2). A clade comprising Osmylidae, Sisyridae, and Nevrorthidae (i.e., Osmyloidea, Winterton et al., 2018) was inferred as sister to all neuropteran families except Coniopterygidae. Dilaridae was placed as the sister group to all other Neuroptera except Coniopterygidae and Osmyloidea. A clade comprising Mantispidae and Berothidae (i.e., Mantispoidea excluding Rhachiberothidae for which transcriptomic data were not available) received high statistical branch support in all analyses of the above-mentioned analyzed datasets (Fig. 2.1, Additional file 3: Figures S1–S5). A sister group relationship between Ithonidae and Myrmeleontiformia (excluding Psychopsidae for which transcriptomic data were not available) was inferred with maximum bootstrap and maximum TBE support. Furthermore, analyses of concatenated domain-partitioned amino-acid data and those of second codon positions suggest Chrysopidae as sister to Mantispidae + Berothidae, and Hemerobiidae as the sister group of Ithonidae + Myrmeleontiformia (Fig. 2.1, Additional file 3: Figures S1–S5). Within Myrmeleontiformia, Nemopteridae is placed as sister to a clade of Ascalaphidae + Myrmeleontidae. Even though non-parametric bootstrap and TBE support for the monophyly of Myrmeleontidae + Ascalaphidae is high, non-parametric bootstrap support for the monophyly of Myrmeleontidae is very low (Fig. 2.1). These results were congruent with the results of the summary coalescent analyses of gene partitions at the amino-acid sequence level, except for the sister group relationship of *Mongoloraphidia* to the remaining Palearctic Raphidiidae (Fig. 2.2a, Additional file 3: Fig. S15–S17, see also Additional file 2). Within Neuroptera, the results of the phylogenetic analyses of domain-based partitioned amino-acid sequence data are also congruent with the concatenation-based analyses of genes at the amino-acid sequence level, except for the disruption of the clade Mantispoidea + Chrysopidae in the

concatenated analyses of genes with increased species coverage (Additional file 3: Fig.

S18–S21).



(**Fig. 2.1**: *see caption on next page*)

**Fig. 2.1** (see figure on previous page): Phylogenetic relationships of Neuropterida based on the analyses of the concatenated amino-acid sequence data of supermatrix E. Colored circles depict phylogenetic branch support values based on 100 non-parametric bootstrap replicates. Bars on the individual nodes show the 95% confidence intervals (equal-tail CI) of the posterior divergence time estimates. Blue squares indicate the time-calibrated nodes. Divergence time estimates were calculated from a single summarized MCMC chain (first independent analysis, run 1) that included all parameter values from each individual meta-partition analysis when including all fossil calibrations. Insect photos from top to bottom: *Dichrostigma flavipes*, *Sialis lutaria*, *Chrysopa perla* (all photos by O. Niehuis).

The summary coalescent analyses and the concatenation-based analyses of gene partitions when analyzing codon-based nucleotide sequence data (with all codon positions included) suggest different topologies concerning the inter-familiar phylogenetic relationships of Neuroptera (Additional file 3: Figures S22–S29, see also Additional file 2). Specifically, analyses of the codon-based nucleotide sequence data with both methods yielded paraphyletic Myrmeleontiformia and further suggest a sister group relationship of Chrysopidae with a clade of Ithonidae + paraphyletic Myrmeleontiformia (Additional file 3: Figures S22–S29). Additional topological differences concern the inferred relationships within Osmyloidea depending on the method and the data type analyzed (e.g. Figures 2.1 and 2.2 and Additional file 3: Figures S1–S29, see also Additional file 2). Overall, the topological differences inferred from the different analyses mainly concern the inter-relationships of the four monophyletic groups: Chrysopidae, Hemerobiidae, Mantispoidea, Ithonidae + Myrmeleontiformia. The different hypotheses concerning the relationships of these four groups (e.g., Hemerobiidae vs. Chrysopidae as sister to Ithonidae + Myrmeleontiformia), are characteristic of the different types of data that were analyzed

(i.e., amino-acid vs. codon-based nucleotide sequence data with all codon positions included, see Additional file 2). The family Hemerobiidae was inferred as sister to Ithonidae + monophyletic Myrmeleontiformia when analyzing amino-acid sequences or second-codon positions of nucleotide sequences, irrespective of the applied phylogenetic method (i.e., concatenation vs. summary coalescent phylogenetic analysis, Fig. 2.1–2.2, Additional file 3: Fig. S1–S14, S15, S18), or partitioning strategy (i.e., domain-based partitioning vs. gene-based partitioning, Additional file 3: Fig. S1–S2, S10–14, S18–S21).

### 2.2.3. Tests for the presence of confounding signal via four-cluster likelihood mapping and data permutations

The four-cluster likelihood mapping (FcLM) approach delivered strong statistical support for most inferred phylogenetic relationships (Additional file 1: Table S2). For example, a clade Megaloptera + Neuroptera is strongly supported by the FcLM analyses with no detectable confounding signal (Fig. 2.2b). Support for Coniopterygidae instead of Nevrorthidae as the sister group to the remaining Neuroptera also received strong FcLM support without detectable confounding signal (Fig. 2.2b, Hypothesis 5: 99.40% of quartets). The monophyly of Osmyloidea is also strongly supported without detectable confounding signal (99.70% of quartets, Hypothesis 8, Additional file 1: Table S2, see also Hypothesis 4a). A potential sister group relationship of Osmylidae and Chrysopidae, as suggested by some previous morphological studies, is not supported by the FcLM branch support tests (Hypotheses 4a and 4b, Fig. 2.2b and Additional file 1: Table S2). The monophyly of Myrmeleontiformia (Nymphidae, Nemopteridae, Ascalaphidae, Myrmeleontidae) is

strongly supported by our FcLM tests without detectable confounding signal (Fig. 2.2b, Hypothesis 7: 94.70% of quartets).

Nevertheless, the results of FcLM analyses showed conflicting signal for some splits in the backbone tree of Neuroptera (Fig. 2.2b, Additional file 1: Table S2). For example, the FcLM analyses do not unequivocally support the sister group relationship of Sisyridae and Nevrorthidae (i.e., 51.80% of quartets support Nevrorthidae + Sisyridae, Fig. 2.2b, Additional file 1: Table S2, Hypotheses 2 and 3). Moreover, FcLM analyses do not unequivocally support a clade Mantispoidea + Chrysopidae (46.10% of quartets, Hypothesis 9, Additional file 1: Table S2). The sister group relationship of Hemerobiidae to Ithonidae + Myrmeleontiformia received only moderate support in FcLM analyses (72.40% of quartets in Hypothesis 6a). FcLM analyses on the permuted matrices showed that there was no substantial contribution of confounding factors for this sister group relationship, although there exists some weak signal (43.30% of quartets) possibly originating from non-random distribution of missing data in support of the results of tree reconstructions (Hypothesis 6a, permutations I and II, Additional file 1: Table S2). When using a different definition of groups of taxa, the placement of Hemerobiidae as sister to Ithonidae + Myrmeleontiformia was supported by only 36.60% of the analyzed quartets (Hypothesis 6b, Additional file 1: Table S2).

**Fig. 2.2:** Gene tree-based and concatenation-based quartet analyses of the phylogenetic relationships of Neuropterida. **a)** Phylogenetic relationships of Neuropterida, as they resulted from the summary coalescent phylogenetic analysis with ASTRAL, when analyzing the full set of gene trees (3983 gene trees inferred at the amino-acid sequence level). Pie charts on branches show ASTRAL quartet support (quartet-based frequencies of alternative quadripartition topologies around a given internode). Arrows indicate the numbers of the

corresponding tree nodes in Fig. 2.1, and the corresponding hypotheses in the FcLM analyses.
**b)** Results of FcLM analyses for a selection of phylogenetic hypotheses applied at the amino-acid sequence level (supermatrix E). The first column shows the results of FcLM when the original data of supermatrix E were analyzed. The second column shows the results of FcLM after phylogenetic signal had been eliminated from supermatrix E (i.e., permutation no. I, see Additional file 2)

### 2.2.4. Divergence times of Neuropterida

Our molecular-dating analyses illustrate that most meta-partitions contained enough signal to overrule the prior assumptions (i.e., marginal prior distributions) on the divergence times of Neuropterida (Fig. 2.3), except for the ancient splits within the outgroup taxa. Given a fixed topology and node-age calibrations, the distribution of median posterior divergence times among meta-partitions when compared with the distribution of the median values of the marginal prior distributions, constitutes evidence for the dominant influence of signal in the datasets (Fig. 2.3). It does however also show extensive variation in signal among meta-partitions. This variation in signal is more prominent for certain nodes (e.g., crown Raphidioptera, Fig. 2.3), whereas the individual median posterior age estimates are less dispersed compared to the overall median for others (e.g., crown Ithonidae + Myrmeleontiformia).

The combined dating analysis of the meta-partitions from the first run in MCMCTree (Fig. 2.1, Additional file 1: Table S3) suggests that the phylogenetic split between Coleopterida and Neuropterida (i.e., Neuropteroidea) occurred in the end of the Devonian period (median = 364.3 Mya, CI = 392.9–325.9, Additional file 1: Tables S3, S4). Crown Neuropterida started to diversify in the middle of Carboniferous (median = 321.7 Mya, CI = 362.0–282.4 Mya). Although Raphidioptera was inferred as

the earliest branching lineage within Neuropterida, the most recent common ancestor of crown Raphidioptera was estimated to have lived at the beginning of the Cretaceous period (median = 132.1, CI = 238.2–61.7 Mya). There is extensive variation in signal among metapartitions for this particular split (Fig. 2.3) that is reflected in the very wide confidence intervals (95% equal-tail and 95% higher posterior density CI, Fig. 2.1, Additional file 1: Tables S3, S4). The split between the Nearctic *Agulla* and all remaining Raphidiidae in the dataset was estimated to have occurred in the middle of the Eocene (median = 44.1, CI = 103.6–21.1 Mya). The split of crown Megaloptera was estimated to have occurred at the beginning of the Triassic period (median = 238.9, CI = 303.4–180.8 Mya), while crown Neuroptera started to diversify much earlier at the beginning of the Permian (median = 280.8, CI = 327.4–241.7 Mya). The crown group of Osmyloidea started to diversify at the beginning of the Jurassic (median = 197.4, CI = 266.7–121.7 Mya). Many consecutive deep splits in the phylogeny of Neuroptera (e.g. crown Osmyloidea, crown Coniopteryginae, and the split between Hemerobiidae, Mantispoidea, Chrysopidae, and Myrmeleontiformia) were estimated to have occurred at the end of the Triassic or the beginning of the Jurassic (Figs. 1 and 3). Lastly, most crown groups of the different neuropterid families (e.g. the crown groups of Chrysopidae, Hemerobiidae, Nemopteridae, Ithonidae, and the common ancestor of Ascalaphidae + Myrmeleontidae) started to diversify during the Cretaceous (Fig. 2.1). Posterior node-age estimates and confidence intervals that resulted from the combined analysis of the second independent run (run 2) with MCMCTree are very similar (Additional file 1: Table S4), which suggests that the two independent chains (each composed of the combined parameter values of the individual meta-partitions) have

converged to very similar posterior node-age estimates (Additional file 3: Figures S30, S31).



(**Fig. 2.3**: *see caption on next page*)

**Fig. 2.3** (see figure on previous page): Distribution of the median posterior node ages among the different meta-partitions. Arrows indicate the corresponding crown groups of Neuropterida and outgroups. Numbers on x-axis correspond to the node number IDs of the tree in Fig. 2.1. The distribution of the median posterior age estimates of the individual meta-partitions from the first independent dating analysis ($\alpha = 2$, run 1) is shown in blue. The distribution of the median age estimates when running the analyses without data (i.e., marginal prior) is shown in red

### 2.2.5. Evolution of larval characters and lifestyles within Neuropterida

We traced the evolution of larval characters within Neuroptera based on the best topology (overall best maximum likelihood tree, ML tree, Fig. 2.1) that resulted from the analysis of domain-based partitioned amino-acid sequence data. The implications for the evolution of larval characters in Neuroptera under parsimony are outlined in Additional file 1: Table S5. Autapomorphies of Neuroptera, Myrmeleontiformia and Coniopterygidae (two terminals included in the studies by Beutel et al., 2010a) and Jandausch et al., 2018) are not affected by the phylogenetic pattern obtained in the present study. With the parsimony approach the reconstruction of ancestral states remained ambiguous with respect to the larval habitat of Neuroptera (terrestrial versus aquatic, Additional file 1: Table S5). In contrast, our Bayesian stochastic character mapping (SCM) analyses suggest a primarily terrestrial larval habitat in the last common ancestor of Neuroptera but also in the last common ancestor of the entire Neuropterida (Fig. 2.4). This result is recovered irrespective of the inferred relationships within Osmyloidea (Additional file 3: Figures S32–S34). Additionally, the parsimony-based analysis remained ambiguous with respect to the ancestral character state of the larval gula in Neuroptera. A large posterior sclerotized plate as it is present in Nevrorthidae (and also in Raphidioptera and Megaloptera) may be ancestral, with a

small posterior rectangular sclerite preserved as vestige in Polystoechotinae, and a small anteromedian triangular sclerite as a *de novo* formation in Myrmeleontiformia. Following the principle of parsimony, the "maxillary head" as defined by Aspöck et al., 2001) (i.e., the complete absence of a gula) could be a ground plan apomorphy of Neuroptera, and the secondary gain of a gula consequently an apomorphy of Nevrorthidae, Polystoechotinae and Myrmeleontiformia. The specialized terminal seta of the flagellum is interpreted as secondarily absent in Nevrorthidae on the one hand, and in Ithonidae and Myrmeleontiformia on the other, in the latter case as a potentially synapomorphic feature of these two groups. The poison channel and the intrinsic musculature of the maxillary stylets are secondarily absent in Sisyridae (Jandausch et al., 2018). The trumpet-shaped empodium is likely an apomorphy of Neuroptera excluding Coniopterygidae and Osmyloidea, and the secondary loss of this feature is a synapomorphy of Ithonidae and Myrmeleontiformia (Jandausch et al., 2018). The ground plan of Neuroptera with respect to the larval cryptonephry is ambivalent. This feature could represent an apomorphy of Neuroptera (Additional file 1: Table S5).

## 2.3. Discussion

### 2.3.1. Statistical robustness of phylogenomic results and potential pitfalls in phylogenetic reconstructions

Previously published phylogenomic analyses have suggested robustly resolved backbone trees of Neuropterida (Song et al., 2019; Wang et al., 2017; Winterton et al., 2018, 2010) 20–22]that were in part incongruent to inferred phylogenetic relationships based on analyses of morphological characters. The most recent molecular analyses at odds with morphological analyses were based on extensive genomic data (Machado et

al., 2019; Wang et al., 2017; Winterton et al., 2018) and therefore the incongruences between these molecular and morphological phylogenies cannot be easily dismissed. Since the accumulation and characterization of extensive genomic data is now the standard procedure in phylogenetics, as it is also true for the analyses of the phylogeny of Neuropterida, the evaluation of statistical robustness of the inferred phylogenies is becoming a complex yet essential task (Kumar et al., 2012). It is obvious that conventional analyses of statistical robustness, in most cases performed with the classical non-parametric bootstrap, might not scale well with the quantity of the data (Cloutier et al., 2019; Phillips et al., 2004; Rodríguez-Ezpeleta et al., 2007; Smith et al., 2015). This is because bootstrap support values provide an assessment of the sampling effects and repeatability of the analyses but cannot assess the accuracy of the inferred phylogenetic trees (Phillips et al., 2004). Alternative or complementary measures of phylogenomic incongruence are warranted to identify phylogenetic relationships with potentially inflated support (Johnson et al., 2018; Pease et al., 2018; Salichos et al., 2014; Salichos and Rokas, 2013). In order to identify potentially inflated branch support of the inferred relationships within Neuropterida, we have used a combination of gene tree-based and concatenation-based quartet methods and compared results with those of the classical non-parametric bootstrapping approach and with those of the newly described bootstrap by transfer support measure (TBE). We observed that a few seemingly well supported phylogenetic relationships assessed by bootstrapping are in fact inflated due to potentially confounding factors in the data. In most instances, concatenation-based and gene tree-based quartet methods deliver congruent pictures, that are in several cases in stark contrast to the classical resampling approaches. We conclude from these observations that at least parts of the backbone tree of

Neuropterida should still not be considered robustly resolved. Below we discuss two examples from the backbone tree of Neuroptera that do not receive unequivocal support from our quartet analyses:

*Phylogenetic relationships within Osmyloidea* – We observed incongruent topologies between concatenation and the summary coalescent phylogenetic analyses concerning the splits within Osmyloidea. Summary coalescent phylogenetic analyses at the amino-acid sequence level suggest a clade of Sisyridae + (Osmylidae + Nevrorthidae), whereas all concatenated analyses of amino-acid sequence data suggested a clade of Osmylidae + (Nevrorthidae + Sisyridae). This incongruence between methods was only present when analyzing aminoacid sequence alignments. The analyses of the codon-based nucleotide sequence alignments (with all codon positions included) resulted in phylogenetic relationships congruent to the summary coalescent approach. Despite the high bootstrap and high TBE support from the concatenated analyses of amino-acid sequence data for a sister group relationship of Sisyridae and Nevrorthidae, our FcLM analyses do not unequivocally support the inferred phylogenetic relationships within Osmyloidea. Specifically, quartet support calculated with ASTRAL and FcLM analyses show almost equal proportions of quartets supporting each of the two above-mentioned prevalent phylogenetic hypotheses. Moreover, the FcLM analyses suggest substantial influence from taxon sampling and possibly from non-random distribution of missing data for this particular phylogenetic relationship. Putting the results of the concatenation-based, summary coalescent and FcLM analyses together, we conclude that the phylogenetic relationships of the three families in Osmyloidea should be considered for now unresolved.

*Phylogenetic position of Hemerobiidae* – Our analyses of amino-acid sequence data and those of second codon positions of the nucleotide sequence data, suggest Hemerobiidae as sister to Ithonidae + monophyletic Myrmeleontiformia, whereas analyses of the complete codon-based nucleotide sequence alignments suggest Chrysopidae as sister to Ithonidae + paraphyletic Myrmeleontiformia. These incongruencies again warrant a detailed examination of potentially confounding signals. The FcLM analyses do not unequivocally support Hemerobiidae as sister to Ithonidae + Myrmeleontiformia (72.40 and 36.60% of quartets), despite the maximum bootstrap and maximum TBE support for this relationship (100%). The FcLM analyses also show some weak putatively misleading signal in support of this relationship that possibly originates from non-random distribution of missing data. Since the FcLM and ASTRAL quartet analyses do not unequivocally support Hemerobiidae as sister to Ithonidae + Myrmeleontiformia, we consider this part of the neuropteran tree as statistically not robustly resolved.

## 2.3.2. Different data types and not different tree-inference methods are responsible for some of the phylogenomic incongruences

Although many previous phylogenomic studies have focused on the biological causes of incongruence that results from analyzing the data with coalescent-based or concatenation-based phylogenetic methods (Cloutier et al., 2019; Edwards et al., 2016; Kubatko and Degnan, 2007; Song et al., 2012), little attention has been given to the effects of the different analyzed data types on phylogenetic inference (Jeffroy et al., 2006). Such data-type effects have been discussed before either in the context of analyzing different genomic regions, such as analyzing introns vs. analyzing coding

sequences (Reddy et al., 2017), or in the context of analyzing the same coding regions at different levels (i.e., nucleotides vs. amino acids) (Gillung et al., 2018; Jeffroy et al., 2006; Vasilikopoulos et al., 2019). Here, we find that some of the inferred relationships within Neuroptera (i.e., the monophyly of Myrmeleontiformia and the position of Chrysopidae, Hemerobiidae and Mantispoidea) are characteristic of the data type that was analyzed (i.e., amino acids vs. codon-based nucleotide sequences with all codon positions included) irrespective of the tree-inference method. Given sufficient phylogenetic signal, the expectation is that the analyses of the same genomic regions at the nucleotide sequence level and the translational level should reflect the same evolutionary history. If the analyses of different data types result in discrepancies, this is most likely due to the failure of the applied substitution models to accommodate the evolutionary history in the analyzed data. Thus, the above-mentioned data-type effects probably stem from violations of the model assumptions by the analyzed data. Additionally, the observation that these data-type effects are quite robust across different tree-inference methods further suggests that both concatenation and summary coalescent methods are sensitive to these violations of model assumptions. An important open question is why some branches in the tree of Neuroptera may be more prone to data-type effects than others. Ancient rapid radiations have been proposed as candidates for such data-driven effects in phylogenetic reconstructions (Reddy et al., 2017).

**Fig. 2.4**: Summarized results of stochastic character mapping analyses (SCM) for the evolution of larval ecologies based on 10,000 sampled character histories. Stochastic character maps were generated under the ER model and by using the topology and branch lengths of the chronogram of Fig. 2.1. Colored circles at the tips show the coded state for each species. Pie charts on internal tree nodes show posterior probabilities of states at each node under the model used. Internal nodes with a posterior probability lower than 1.00 are depicted in larger size (note: for the SCM analyses we assumed that larval ecologies remain constant within the same family).

## 2.3.3. Implications of our phylogenetic reconstructions concerning the evolution of Neuropterida

*Inter-ordinal phylogenetic affinities of Neuroptera, Megaloptera, and Raphidioptera*

Within holometabolous insects, Neuropterida is inferred as the sister group to Coleopterida, a phylogenetic hypothesis that is in accordance with the latest views on the phylogeny of Holometabola (Misof et al., 2014; Peters et al., 2014). The monophyly of Neuropteroidea (Coleopterida + Neuropterida) is supported by the presence of a prognathous or slightly inclined head in the adults of this group (Peters et al., 2014). We estimated the most recent common ancestor of Neuropteroidea to have lived in the late Devonian (~ 363 Mya), an estimate that is earlier than what has been suggested (Misof et al., 2014; Tong et al., 2015), and during a time interval that coincides with the appearance of the first tetrapod vertebrates and the formation of the first land forests.

In our study, the order Raphidioptera is placed as sister to Megaloptera + Neuroptera, in agreement with the results of most previous molecular studies (Cameron et al., 2009; Kjer et al., 2006; Misof et al., 2014; Peters et al., 2014; Wang et al., 2017; Winterton et al., 2018; Zhao et al., 2013). The notion that Megaloptera is the sister group to Neuroptera was first introduced by Boudreaux (1979), on the premise of common wing venation characters. This idea was revived later with the argument that aquatic larvae represent a synapomorphic feature for Neuroptera and Megaloptera, with secondary terrestrialization in Neuroptera (Aspöck, 1995). Our phylogenetic results and FcLM analyses are in agreement with the results of those morphological studies and with recent phylogenomic analyses of mitochondrial genomes or target DNA enrichment data concerning the inter-ordinal relationships of Neuropterida (Wang et al., 2017; Winterton et al., 2018; Zhao et al., 2013). Hence, the traditional hypothesis that

Neuroptera is the sister group to Megaloptera + Raphidioptera (Achtelig, 1978, 1975; Beutel et al., 2011; Beutel and Gorb, 2001; Hennig, 1969; Kristensen, 1991), that was suggested by a few studies based on the analyses of a few genes (McKenna and Farrell, 2010; Wheeler and Hayashi, 2001; Whiting et al., 1997; Wiegmann et al., 2009), is highly unlikely. We inferred the first split among the crown Neuropterida to have occurred in the middle of the Carboniferous (~ 321 Mya). This node-age estimate is slightly older than the age inferred in previously published phylogenomic studies, that proposed a common origin of the extant Neuropterida in the late Carboniferous or the early Permian (Wang et al., 2017; Winterton et al., 2018).

*Evolutionary history of Raphidioptera*

Within Raphidioptera, both Raphidiidae and Inocelliidae are recovered as monophyletic in all of our analyses and with high statistical support. We estimated the common ancestor of extant Raphidioptera to have lived during the early Cretaceous (~136 Mya), although it is evident from the fossil record that stem lineages of Raphidioptera were distinctly diverse much earlier in the Mesozoic (Aspöck and Aspöck, 2007). Our results suggest the placement of the Nearctic genus *Agulla* as sister to the Palearctic Raphidiidae. Although the Nearctic genus *Alena* is not included in our analyses, the above-mentioned relationship suggests the monophyly of the Palearctic Raphidiidae and corroborates previous molecular phylogenetic analyses of Raphidiidae (Haring et al., 2011). Furthermore, the results of the analyses of domain-based partitioned data are in agreement with previous molecular phylogenetic analyses of the Raphidiidae, that suggested the division of the Palearctic Raphidiidae into an Eastern Palearctic (*Mongoloraphidia* clade) and a Western Palearctic (*Ohmella*, *Puncha* and

*Phaeostigma* clades) radiation (Haring et al., 2011). Biogeographical aspects of the phylogeny of extant Raphidioptera are discussed in more detail by Aspöck et al. (2012b).

### *Evolutionary history of Megaloptera*

The order Megaloptera is inferred as monophyletic in all analyses and the family Corydalidae is also inferred as monophyletic. These results are congruent with the results of target DNA enrichment-based phylogenomic analyses of Neuropterida (Winterton et al., 2018). In addition, these results are in agreement with morphological analyses of genital and non-genital characters and with most morphology-based phylogenies of Neuropterida (Aspöck et al., 2001; Aspöck and Aspöck, 2008; Zhao et al., 2014). There are only few morphological autapomorphies of Megaloptera such as the shift of the bases of the male gonocoxites 9 to the base of tergum 9 (Liu et al., 2016). Morphological characters supporting the monophyly of Corydalidae are scarce and they concern mostly genital characters and wing-base structures (Liu et al., 2016; Zhao et al., 2014). Our taxon sampling does not allow further assessment of the monophyly of the corydalid subfamilies Corydalinae and Chauliodinae, but recent phylogenetic investigations have shown that the current taxonomic classification is supported by the analyses of molecular or morphological characters (Liu et al., 2016; Winterton et al., 2018; Zhao et al., 2014). We estimated the common ancestor of extant Megaloptera to have lived in the early Triassic (~ 239 Mya), an estimate that is younger than estimates derived from analyses of target DNA enrichment data (Winterton et al., 2018), but in agreement to the results of analyses of mitochondrial genomes (Wang et al., 2017).

*Evolutionary history of Neuroptera*

The order Neuroptera is inferred as monophyletic and our divergence time estimates suggest that its members started to diverge in the end of the Carboniferous (~ 301 Mya), while the common ancestor of the extant Neuroptera is estimated to have lived in the early Permian (~ 281 Mya). Our inferred phylogenetic trees corroborate the results of previous phylogenomic studies that suggested the family Coniopterygidae as sister to all remaining neuropteran families (Song et al., 2019; Wang et al., 2017; Winterton et al., 2018). The idea that the dustywings are the sister group of the remaining families of Neuroptera is very old (Withycombe, 1925) and was originally based on a number of characters that this family shares with Megaloptera, such as the reduced number of Malpighian tubules (six in Coniopterygidae instead of eight in other Neuroptera) and the reduced number of abdominal ganglia of their larvae (Withycombe, 1925). However, it should be noted that these features could be the result of miniaturization in the dustywings. Moreover, the alternative character states would be plesiomorphic, and therefore they constitute no arguments for monophyletic Neuroptera excluding Coniopterygidae. In our study Coniopterygidae is inferred as an ancient lacewing group that started to diversify in the middle of the Permian (~ 281 Mya). This result is in agreement with the findings of recent molecular dating analyses of Neuropterida (Wang et al., 2017; Winterton et al., 2018).

The phylogenetic placement of Coniopterygidae as sister to all remaining Neuroptera is in contrast with the majority of morphological analyses that have instead suggested Nevrorthidae as the most ancient lineage within the order (Aspöck et al., 2001; Aspöck and Aspöck, 2008; Beutel et al., 2010a). The monophyly of Neuroptera with the exclusion of Nevrorthidae is morphologically supported by the formation of an

undivided postmentum, the far-reaching modification or loss of the larval gula and the presence of cryptonephric Malpighian tubules of the larvae (Beutel et al., 2010a). Specifically, in all terrestrial neuropteran larvae (including Coniopterygidae) the distal parts of the Malpighian tubules are connected with the colon, a phenomenon referred to as larval cryptonephry. In the aquatic larva of *Nevrorthus* all Malpighian tubules are free, while the aquatic larvae of Sisyridae have one cryptonephric tubule. The phenomenon of cryptonephry results in an improved water re-absorption mechanism and is apparently an adaptation to terrestrial environment, especially to a more exposed lifestyle and life in drier habitats. The original idea concerning the evolution of cryptonephry within Neuroptera is in contrast with the herewith presented phylogenetic relationships and with other molecular phylogenies (Wang et al., 2017; Winterton et al., 2018, 2010), that suggest cryptonephry might be an apomorphic feature of Neuroptera with a putative secondary loss in Nevrothidae and secondary modification in Sisyridae. Despite the lack of morphological autapomorphies for a clade comprising Neuroptera excluding Coniopterygidae, this robust result across molecular analyses and methods suggests that a sister group relationship of Nevrorthidae to all other neuropteran families is unlikely.

A clade of Nevrorthidae, Sisyridae and Osmylidae (i.e., Osmyloidea) is inferred as sister to all remaining neuropteran families except Coniopterygidae and this clade is stable across analyses of different datasets and methods. This clade was also strongly supported in all quartet analyses, which in turn suggests that the placement of these three families in a monophyletic group is robust. This result is also in agreement with the results of analyses of target DNA enrichment data (Winterton et al., 2018). Potential synapomorphies of Osmyloidea are the semi-aquatic or aquatic larval ecologies and the

secondarily multi-segmented antennae of the larvae (Jandausch et al., 2019). Within

Osmyloidea, a sister group relationship of Nevrorthidae and Sisyridae is congruent with

the analyses of mitochondrial genomes (Wang et al., 2017) and with older studies based

on the analysis of a few genes (Winterton et al., 2010). Moreover, a single shift to an

aquatic lifestyle conforms to a branching pattern of Nevrorthidae and Sisyridae as sister

clades. It should, however, be noted that the larvae of Nevrorthidae and Sisyridae have

very different breathing and feeding adaptations, an observation that contrasts their

sister group relationship (Jandausch et al., 2019). The recent discovery of a complex

submental gland with a multiporous opening in adults of *Nevrorthus* and *Osmylus*

(Randolf et al., 2014) could corroborate the monophyly of Osmylidae + Nevrorthidae

as revealed by our summary coalescent analyses and by previous analyses of target

DNA enrichment sequence data (Winterton et al., 2018). In the context of our best ML

tree (Fig. 2.1), either the stem species of Neuroptera must have evolved this gland, with

subsequent multiple losses, or it must have evolved in the stem species of Osmylidae +

(Nevrorthidae + Sisyridae) and was then secondarily lost in Sisyridae. A clade of

Osmylidae + Nevrorthidae has been presented elsewhere: e.g., by Zwick (1967) (based

on macrochaete of the neck, and the size of the palps), by Yang et al. (2012) (mainly

based on fossils), and in the recent target DNA enrichment-based phylogenomic study

of Neuropterida (Winterton et al., 2018). Another interesting observation in this context

is that the adults of Osmylidae are the only neuropterans with ocelli. Given that the

possession of ocelli is most likely a plesiomorphic feature, as they are present in the

adults of Raphidiidae and Corydalidae, we can hypothesize that the median eyes must

have been reduced several times independently within Neuroptera, with possible

vestiges still preserved in several groups.

A robust inference of the most archaic phylogenetic events within Neuroptera is essential for deciphering the evolution of lifestyle transitions of their larvae. Aquatic versus terrestrial habits of ancestral neuropteran larvae as well as a possible ancestral aquatic larvae of Neuropterida have been discussed in detail by authors of previous studies (Wang et al., 2017; Winterton et al., 2018). Specifically, previous ancestral character state reconstructions (ACSR) of the larval ecologies of Neuropterida have suggested that the common ancestor of Neuroptera might have had aquatic larvae (Wang et al., 2017; Winterton et al., 2018). Under the scenario of primarily aquatic neuropteran larvae, the results of our transcriptomic analysis would imply that the larvae of Coniopterygidae acquired terrestrial habits secondarily. In a second step Osmylidae must also have acquired terrestrial larvae independently, and finally in a third step the stem species of the remaining Neuroptera must also have acquired terrestrial larvae. Although three independent transitions to terrestrial lifestyle within Neuroptera is a possible scenario, it is not the most parsimonious. In an alternative scenario, with the stem species of Neuroptera being primarily terrestrial in the larval stages, the larvae of Sisyridae and Nevrorthidae would be secondarily aquatic as assumed by Gaumont (1976). Our parsimony-based ACSR of larval ecologies do not provide unequivocal support for either aquatic or terrestrial larvae in the last common ancestor of Neuroptera. In contrast, our SCM analyses unequivocally support primarily terrestrial larvae of Neuroptera and Neuropterida. However, it should be noted that parsimony-based ACSRs suffer from a number of limitations (Bollback, 2006; Huelsenbeck et al., 2003) and that our parsimony-based analysis is based on a less extensive taxon sampling (Jandausch et al., 2019). For these reasons we consider the estimates of SCM analyses as more reliable. The hypothesis of primarily terrestrial

larvae of Neuropterida and Neuroptera suggests either two or three independent shifts to aquatic larval lifestyles within Neuropterida depending on the inferred topology within Osmyloidea. Interestingly, this hypothesis implies that the stem species of Megaloptera + Neuroptera had terrestrial larvae and that the larvae of Megaloptera are secondarily aquatic. We conclude from these observations that at least two shifts to aquatic habitats must have occurred in the early evolution of Neuropterida.

The family Dilaridae (pleasing lacewings) has been traditionally considered to form a clade with the families Mantispidae, Berothidae and Rhachiberothidae. The unofficial term "dilarid clade" has been used to describe this phylogenetic assemblage (Aspöck et al., 2001; Beutel et al., 2010a, 2010b; Jandausch et al., 2018). We could not corroborate a clade that includes these four families as suggested by other authors (Jandausch et al., 2018; Randolf et al., 2014). All analyses place Dilaridae as sister to all remaining Neuroptera except Coniopterygidae and Osmyloidea. This result is in accordance with previous sequenced-based phylogenomic analyses (Wang et al., 2017; Winterton et al., 2018). Most importantly, the monophyly of the neuropteran families except Coniopterygidae and Osmyloidea is strongly supported by previous analyses of mitochondrial genomic rearrangements (Wang et al., 2017; Zhao et al., 2013).

Mantispidae and Berothidae were recovered as sister taxa with strong statistical branch support in all phylogenetic analyses, but the placement of this clade within Neuroptera is not robustly resolved. Concatenation-based and summary coalescent phylogenetic analyses of amino-acid sequences suggest a sister group relationship of Mantispoidea with Chysopidae. However, the different quartet analyses did not unequivocally support this sister group relationship. Our results corroborate previous views suggesting a close phylogenetic affinity of Berothidae and Mantispidae (Aspöck

et al., 2001; Jandausch et al., 2018). Despite the fact that the family Rhachiberothidae is not included in our analyses, the monophyly of Mantispoidea is strongly supported by the presence of overlapping scales on antennae and maxillae, the presence of thoracic "trichobothria",and by their hypermetamorphic development (Aspöck et al., 2001; Jandausch et al., 2018). The phylogenetic relationships within Mantispoidea, as well as the monophyly of Mantispidae, have remained unresolved (Winterton et al., 2018), yet our taxon sampling does not allow testing any hypothesis concerning the phylogeny of Mantispoidea.

A clade Chrysopidae + Hemerobiidae, suggested by analyses of mitochondrial genomes (Song et al., 2019; Wang et al., 2017) and morphological characters (Aspöck and Aspöck, 2008), is not corroborated in our study. The conflicting phylogenetic hypotheses between the analyses of different data types presented here corroborate the results of Winterton et al. (2018) concerning the affinities of Chrysopidae and Hemerobiidae. In their analyses of amino-acid sequence alignments Mantispoidea was inferred as sister to Chrysopidae, while Hemerobiidae was inferred as sister to Ithonidae + Myrmeleontiformia. These results are identical to our own results based on analyses of amino-acid sequence data. However, it should be noted that there is presently no morphological support in favor of these phylogenetic relationships. Morphological apomorphies shared by Hemerobiidae and Chrysopidae (Aspöck and Aspöck, 2008; Wang et al., 2017) and the results of our quartet-based analyses show that the above-mentioned relationships require further scrutiny. The previously suggested clade Chrysopidae + Osmylidae that was based on analyses of larval head characters (Aspöck et al., 2001) is also not supported by our FcLM analyses. The main argument for this sister group relationship was based on length of the cardines, and the

possession of special prothoracic glands (Gusten and Dettner, 1992). However, varying lengths of the cardines are gradual modifications rather than discrete character states. Additionally, data on the prothoracic glands are missing for most neuropteran families. Therefore, the arguments for a clade Chrysopidae + Osmylidae are not convincing.

The family Ithonidae is inferred as monophyletic and sister to monophyletic Myrmeleontiformia. The monophyly of Myrmeleontiformia is also strongly supported by our FcLM analyses and by previous analyses of morphological characters (Badano et al., 2018, 2017). The synapomorphies supporting the monophyly of Myrmeleontiformia, including the Psychopsidae, have already been documented by MacLeod (1964), by Beutel et al. (2010a), and more recently by Badano et al. (2017). Overall, the larval cephalic morphology of Myrmeleontiformia differs profoundly from that of other groups of Neuroptera (Beutel et al., 2010a; Jandausch et al., 2018), including among others the anterior shift of the tentorium and the greatly enlarged muscles of the paired mouthparts to handle the huge sucking tubes. Although Psychopsidae is not included in our study, we expected that if there is phylogenetic signal supporting a clade Ithonidae + Nymphidae, as suggested by other authors (Winterton et al., 2018), the FcLM analyses would support this clade. Our phylogenetic analyses of amino-acid sequence alignments are in contrast with the results of the analyses of target DNA enrichment data that suggested paraphyletic Myrmeleontiformia in relation to Ithonidae (Machado et al., 2019; Winterton et al., 2018). Interestingly, when we analyzed codon-based nucleotide sequences with all three codon positions included, Myrmeleontiformia was rendered paraphyletic in relation to Ithonidae similarly to the results of Winterton et al. (2018). The study of Winterton et al. (2018) was the first molecular study to challenge the clade

Myrmeleontiformia. In contrast, we received high statistical support in most phylogenetic analyses and in FcLM analyses in favor of the monophyly of this group.

Within Myrmeleontiformia (excluding Psychopsidae), Nymphidae is inferred as the earliest diverging lineage. Larval synapomorphies of Myrmeleontiformia excluding Psychopsidae are the conspicuously raised ocular region, a sensory pit on the apical labial palpomere, a strongly developed mid-dorsal cervical apodeme, a distinctly widened body posterior to the prothorax, and a compact and laterally rounded abdomen (Beutel et al., 2010a; Jandausch et al., 2018). The monophyly of the family Nemopteridae has been questioned before (Monserrat, 1996), but has been corroborated later (Badano et al., 2017). We inferred Nemopteridae as monophyletic with strong statistical support and sister to a clade of Ascalaphidae + monophyletic Myrmeleontidae. These results are congruent with those of most recent cladistic analyses of Myrmeleontiformia based on analyses of larval characters (Badano et al., 2018). However, non-parametric bootstrap support for the monophyly of Myrmeleontidae in the analyses of amino-acid sequence alignments was very low, and the same applies for the gene tree-based quartet support for this particular phylogenetic relationship. Previous phylogenomic analyses of the owlflies and antlions have suggested that Myrmeleontidae are polyphyletic with respect to Ascalaphidae (Machado et al., 2019; Winterton et al., 2018). Based on that premise, it has been suggested that Ascalaphidae should be placed in a subfamily of Myrmeleontidae together with the antlion tribes Palparini, Dimarini and Stilbopterygini (Machado et al., 2019). Since we did not recover Ascalaphidae nested within Myrmeleontidae, we retain the taxonomic status of Ascalaphidae as a separate family. The monophyly of the

Myrmeleontidae has been corroborated based on several fossorial habits of their larvae and specific features linked with them (Badano et al., 2018, 2017).

It is essential to mention that the different phylogenetic relationships of neuropteran families presented here corroborate previous results on the evolution of the larval gula-like sclerite within Neuroptera (Winterton et al., 2018). Winterton et al. (2018) interpreted a pattern of evolution of the larval gula in Neuropterida according the results of their analyses. The result showed that the presence of gula is the ancestral state of the entire Neuropterida clade. As such, the presence of gula in the larvae of Nevrorthidae, Ithonidae, and Myrmeleontiformia could be formed either by numerous multiple losses in other lacewings, or could have at least two independent gains in these groups. When considering the larval gula in Myrmeleontiformia, this sclerite is usually reduced to a narrow sclerite medially dividing the two greatly enlarged genal sclerites, a structure that appears different from the gula in Megaloptera and Raphidioptera. Accordingly, the gula of Neuroptera is called "gula-like sclerite" by Winterton et al. (2018) due to its likely non-homologous origin but contrary to the hypothesis of its homologous origin within Neuropterida implied by U. Aspöck (2002). Our parsimony-based character mapping analysis suggested an independent gain of the gula-like sclerite in the members of Ithonidae and Myrmeleontiformia similarly to the suggestion by Winterton et al. (2018). Because the herewith presented phylogenetic incongruencies mainly concern the phylogenetic position of Hemerobiidae, Chrysopidae and Mantispoidea and because the larvae of these groups lack a gula-like sclerite, the previously suggested pattern for the evolution of this morphological feature is unaffected by our results. Hence, an independent gain or reinvention of this gula-like sclerite in Ithonidae and in Myrmeleontiformia appears very likely.

## 2.4. Conclusions

We draw four major conclusions from our analyses: (1) Part of the backbone tree of Neuropterida receives strong statistical support in several independent phylogenetic analyses and should be considered for now the most likely scenario of neuropterid evolution. One such scenario is the early split between Raphidioptera and Megaloptera + Neuroptera. Within Neuroptera, all analyses support an early split between Coniopterygidae and the remaining Neuroptera which cannot be corroborated with morphological analyses. The families Nevrorthidae, Sisyridae and Osmylidae form a monophyletic group sister to all other Neuroptera except Coniopterygidae. The family Dilaridae is the sister group to all remaining Neuroptera except Coniopterygidae and Osmyloidea. Despite these seemingly robust phylogenetic results, the phylogenetic relationships between the most species rich groups of Neuroptera (i.e., Chrysopidae, Ithonidae + Myrmeleontiformia, Hemerobiidae, Mantispoidea) are still not robustly resolved. For several branches in the neuropteran tree, the seemingly high branch support appears to be inflated and should be taken with caution. (2) Comparing concatenation versus summary coalescent approaches, and additional quartet-based measures of phylogenomic incongruence such as the FcLM approach, illustrates the potential of inflated branch support particularly derived from non-parametric resampling methods. Scientists are therefore advised to critically evaluate branch support in phylogenomic analyses and assume a conservative position. (3) The analyses of neuropterid relationships have received a lot of attention in the past and an extensive amount of phylogenomic data has been generated. However, parts of the backbone tree of Neuropterida can still not be robustly resolved which is disappointing, but reflecting a picture seen in other analyses of ancient phylogenetic splits as well. It will be

necessary to invest molecular data beyond primary gene sequence information, for example structural genomic data (Cloutier et al., 2019; Niehuis et al., 2012). (4) Without an interplay of molecular and detailed morphological analyses, we will not be able to spot the major problems in biased results of any kind. Morphological analyses are critically needed to deliver a complete picture of the evolution of Neuropterida.

## 2.5. Methods

### 2.5.1. Taxon sampling

We sequenced and de novo assembled 88 whole-body transcriptomes of 85 species of Neuropterida (Raphidioptera: 18 species, Megaloptera: seven species, Neuroptera: 60 species, Additional file 1: Table S6), comprising representatives of all extant families of Neuropterida except Rhachiberothidae and Psychopsidae. For the species *Parvoraphidia microstigma*, *Palpares libelluloides*, *Peyerimhoffina gracilis*, two transcript libraries of separate specimens were generated respectively, sequenced and assembled (Table S1). RNA isolation, RNA library preparation, transcriptome sequencing, transcriptome assembly, and transcriptome quality assessment were performed according to the procedures described by Misof et al. (2014) and by Peters et al. (2017) (see Additional file 2). We complemented our dataset with publicly available transcriptomic and genomic (official gene sets, OGS) sequence data of eight neuropterid and 41 outgroup species, representing all currently recognized holometabolous insect orders (Additional file 1: Table S7). In total, our sampling comprised 96 transcriptomes of Neuropterida (from 92 species) and 45 transcriptomes and official gene-sets of non-neuropterid insects (from 41 species, see Additional file 2).

## 2.5.2. Orthology assignment, multiple sequence alignment, alignment refinement and alignment masking

We identified a set of 3983 clusters of orthologous single-copy genes (COGs) at the hierarchical level "Endopterygota" (i.e., Holometabola), based on a custom profile query in OrthoDB7 (Waterhouse et al., 2013) (see Additional file 2 for details). The custom query allowed COGs only to be included in the ortholog set if single-copy genes of all selected reference taxa were present in a given COG. As reference genomes, we selected *Acromyrmex echinatior* v. 3.8 (Nygaard et al., 2011), *Tribolium castaneum* v. 3.0 (Richards et al., 2008), *Bombyx mori* v. 2.0 (Xia et al., 2004), and *Drosophila melanogaster* v. 5.51 (Adams et al., 2000) (see Additional file 1: Table S8).

Mapping of putative orthologous transcripts to each COG, at the translational (amino-acid, aaCOGs) and at the transcriptional level (nucleotide, nCOGs), was performed with the software package Orthograph v. 0.5 (Petersen et al., 2017) (see Additional file 2). Subsequently, we selected a subset of outgroup and ingroup species with a high number of assigned orthologs for downstream analyses (Additional file 1: Table S1). Specifically, if more than one transcriptome/OGS were processed from the same outgroup or ingroup species, the dataset with the highest number of identified orthologs was included in downstream analyses. We did not exclude ingroup taxa based on their completeness (measured by the number of assigned orthologs), except in those cases in which more than one transcriptome from the same species were used in the orthology assignment step. Overall, we considered transcriptomes of the outgroup species to be of high completeness when putative orthologous transcripts from these datasets were assigned to at least 3000 COGs (Additional file 1:Table S1, with the exception of *Mengenilla moldrzyki*). The filtered dataset consisted of 124 species (92

neuropterid species and 32 outgroup species) including the four reference species of the ortholog set.

Orthologous amino-acid sequences were aligned with MAFFT v. 7.123 (Katoh and Standley, 2013) and by applying the L-INS-i algorithm. We followed the procedures outlined by Misof et al. (2014) for identifying potentially non-orthologous and misaligned sequences. Details on the applied alignment-refinement procedure, the removal of putative outliers, and the generation of codon-based alignments (corresponding to the amino-acid alignments) are given in Additional file 2. Based on the rationale of previous phylogenomic studies employing various alignment masking (i.e., alignment-column filtering) methods (Dell'Ampio et al., 2014; Fernandez et al., 2017, 2016; Laumer et al., 2015; Li et al., 2017; Meusemann et al., 2010; Misof et al., 2014; Schwentner et al., 2017; Sharma et al., 2014; von Reumont et al., 2012) we used ALISCORE v. 1.2 (Kück et al., 2010; Misof and Misof, 2009), to identify and mask putatively randomly similar aligned sections at the amino-acid sequence level and also masked the corresponding nucleotide sequence codons.

### 2.5.3. Concatenation of supermatrices

We combined the results of alignment masking and protein-domain identification (see Additional file 2) to generate amino-acid and nucleotide sequence supermatrices partitioned according to protein-domain clans, families and single domains following the procedure described by Misof et al. (2014). Subsequently, we generated subsets of the original concatenated supermatrix to improve data coverage and information content, and to assess any putative effects of violations of the SRH conditions assumed by the substitution models in our phylogenetic analyses (Table 2.1,

Additional file 2). For each amino-acid supermatrix, we calculated the overall alignment completeness scores and generated heatmaps of pairwise completeness scores with AliStat v. 1.6 (current version available from: https://github.com/thomaskf/ AliStat) (Wong et al., 2020). Overall deviation from SRH conditions within each supermatrix (Jermiin et al., 2004) was measured with the Bowker's test of symmetry (Bowker, 1948) and by generating heatmaps as implemented in SymTest v. 2.0.47 (current version available from: https://github.com/ottmi/symtest, see Misof et al., 2014).

## 2.5.4. Phylogenetic analyses of amino-acid sequence data partitioned according to protein-domain clans and families, and to single protein domains

We selected the amino-acid supermatrix E (Table 2.1, Additional file 1: Table S9, details in Additional file 2) for downstream analyses, because it showed increased phylogenetic information content and data coverage compared to the supermatrices A, B, C, and D, while being only slightly less informative and larger than supermatrix F (Table 2.1) (Misof et al., 2013; Wong et al., 2020). We used PartitionFinder v. 2.0.0-pre11 (Lanfear et al., 2017) to identify the optimal combination of partitions into meta-partitions, and to infer the respective amino-acid substitution models for each meta-partition prior to tree reconstructions (Additional file 2). The resulting partitioning scheme with the best AICc and the accompanying selected models for each meta-partition were used as input for IQ-TREE v. 1.3.13 (Nguyen et al., 2015) to conduct 100 independent maximum likelihood tree searches (see Additional file 2). We selected the tree with the highest log-likelihood score among all tree searches as the maximum likelihood tree (best ML tree).

Based on the best ML tree, we calculated branch support from 100 non-parametric bootstrap replicates as well as from 10,000 replicates of the SH-like approximate likelihood ratio test (SH-aLRT) (Guindon et al., 2010) with IQ-TREE v. 1.3.13. We assessed whether or not the number of bootstrap replicates was sufficient to accurately infer branch support by running the *a posteriori* bootstop test in RAxML v. 8.2.8 (Pattengale et al., 2010; Stamatakis, 2014) and by doing ten independent tests with different random seeds (see Additional file 2). We calculated an additional branch support metric by applying the bootstrap by transfer support measure based on our calculated bootstrap trees (Lemoine et al., 2018). We also tested for the presence of rogue taxa in our dataset with RogueNaRok v. 1.0 (Aberer et al., 2013). Finally, we rooted the presented tree (Fig. 2.1) by selecting the split between Hymenoptera and all remaining holometabolous taxa using the software Seaview v. 4.5.4 (Gouy et al., 2010).

Modeling site-heterogeneous processes of amino-acid substitutions by incorporating site specific amino-acid profiles into phylogenetic reconstruction can potentially alleviate phylogenetic artifacts due to model misspecification (Lartillot et al., 2007; Le et al., 2008; H.-C. Wang et al., 2019). We therefore performed an additional tree search on supermatrix E with the PMSF mixture model implemented in IQ-TREE v. 1.5.5 (Wang et al., 2018) (Additional file 2) and compared results of this phylogenetic reconstruction with those described above. In order to control for the effects of missing data, we generated two reduced versions of supermatrix E by keeping only those alignment sites with at least 90% or 95% of the total number of species present (207,582 and 110, 708 amino-acid alignment sites respectively). For each of these two reduced matrices, we conducted two additional tree searches with the rapid

approximation to the PMSF model in IQ-TREE v. 1.5.5 (see Additional file 2 for details).

Heterogeneous amino-acid composition among species in the dataset can severely bias phylogenetic reconstructions due to violation of substitution model assumptions (Ababneh et al., 2006; Feuda et al., 2017; Foster, 2004; Jermiin et al., 2008, 2004; Vasilikopoulos et al., 2019). We therefore controlled for among-species compositional heterogeneity in the analyzed amino-acid supermatrix E by masking subsets with a relative composition frequency variation (RCFV) value greater than or equal to 0.1 (Fernandez et al., 2016; Zhong et al., 2011), calculated with BaCoCa v. 1.01 (Kück and Struck, 2014). We monitored the effect of this masking by applying the Bowker's symmetry tests across taxa with SymTest v. 2.0.47. With this RCFV-corrected dataset we conducted five ML tree searches with IQ-TREE v. 1.6.6 by specifying the previously estimated most-fitted substitution models for each meta-partition. We calculated 1000 ultrafast bootstraps (UFB) (Hoang et al., 2018) and 10,000 SH-aLRT replicates for the RCFV-corrected dataset with IQ-TREE v. 1.6.6 (see Additional file 2).

We studied the effect of potentially confounding signal, like non-random distribution of data coverage and violations of SRH conditions, on our phylogenetic reconstructions with the FcLM approach (Strimmer and von Haeseler, 1997) as described by Misof et al. (2014). We formulated nine phylogenetic hypotheses, that are in part based on the results of our tree reconstructions and partly on published alternative phylogenetic hypotheses. For each of the nine tested hypotheses (Additional file 1: Table S2), we used a permutation approach to assess signal originating from non-random distribution of data coverage and violations of SRH conditions in supermatrix E. Accompanying the FcLM approach, we generated a decisive subset of supermatrix E

(Table 2.1) (Dell'Ampio et al., 2014), and which included only meta-partitions with 1) data for all species, 2) less than 30% ambiguous sites (< 30% of X/−), and 3) an alignment length of at least 500 amino-acid sites. The selected meta-partitions were concatenated into a decisive supermatrix (209 meta-partitions, 228,933 aligned amino-acid sites) with FASconCAT-G v. 1.02 (Kück and Longo, 2014). The phylogenetic analyses of this decisive supermatrix followed the scheme of the previous analyses (Additional file 2).

### 2.5.5. Concatenation-based phylogenetic analyses of the second codon positions

We compared the results of tree reconstructions based on data at the amino-acid and nucleotide sequence levels. Substitutions at the nucleotide sequence level follow different processes than substitutions at the amino-acid sequence level, and thus the analyses at the nucleotide level can be considered an independent test of the results based on the amino-acid sequence data. Published investigations have consistently demonstrated that the base composition of second codon positions of protein-coding nucleotide sequences are the most homogeneous across taxa and thus least violate assumptions of the applied nucleotide substitution models (Misof et al., 2014; Timmermans et al., 2016; Vasilikopoulos et al., 2019). We therefore selected the nucleotide supermatrix corresponding to the amino-acid supermatrix D (Table 2.1) and evaluated the degree of deviation from SRH conditions on different subsets of this matrix (Ababneh et al., 2006; Jermiin et al., 2008). We performed the pairwise symmetry tests of homogeneity, by selecting the Bowker's test in SymTest v. 2.0.47, on the following datasets: 1) the entire nucleotide supermatrix, 2) only first codon positions of the nucleotide supermatrix 3) only third codon positions of the nucleotide

supermatrix, and 4) only second codon positions of the nucleotide sequence supermatrix. Since the second codon positions showed the least deviation from the SRH conditions, we masked all first and third codon positions and further proceeded by analyzing a dataset composed exclusively of second codon positions. We calculated the most appropriate partitioning scheme to analyze the second codon positions of supermatrix D, with the *k*-means algorithm (Frandsen et al., 2015) in PartitionFinder v. 2.0.0-pre11, and conducted 100 independent maximum likelihood searches with IQ-TREE v. 1.3.13 (details in Additional file 2). We calculated branch support values from 100 non-parametric bootstraps and 100 TBE replicates and mapped them onto the tree with the highest log-likelihood among all tree searches.

## 2.5.6. Concatenation-based vs. summary coalescent phylogenetic analyses of gene partitions

The concatenation approach has been criticized for being ignorant against gene tree discordance due to ILS and thus for being susceptible to tree reconstruction biases caused by these effects (Edwards, 2009; Edwards et al., 2016; Kubatko and Degnan, 2007; Xu and Yang, 2016). Currently it is unclear which approach delivers the most reliable topological estimates when analyzing empirical data (de Queiroz and Gatesy, 2007; Edwards et al., 2016; Gatesy and Springer, 2014; L. Liu et al., 2015; Sayyari et al., 2017; Simmons and Gatesy, 2015; Springer and Gatesy, 2016; Tonini et al., 2015; Xu and Yang, 2016). To explore the sensitivity of our supermatrix-based analyses to the putative effects of gene tree discordance we used the 3983 alignments of COGs to conduct summary coalescent analyses with ASTRAL III v. 5.6.1 (Zhang et al., 2018). We first removed ambiguous-only sites (X, N, −) from each amino-acid and nucleotide

sequence alignment. Subsequently, we used ModelFinder in IQ-TREE v. 1.6.3 (Kalyaanamoorthy et al., 2017) to infer the best fitting substitution model for each gene separately at the translational level and the transcriptional level (see Additional file 2 for details) based on the BIC criterion. We considered all combinations of modelling ASRV. At the nucleotide sequence level all three codon positions for each gene were included in the phylogenetic analyses. We performed ten independent ML tree searches for each gene with the respective best fitting model and selected the best ML gene tree among these searches to be used for the summary coalescent analyses. Coalescent-based species trees were inferred separately at the amino-acid and the nucleotide sequence levels. The resulting species trees were then scored and annotated by comparing the gene trees with the inferred species tree (Sayyari and Mirarab, 2016). We considered the quartet support values of the summary coalescent analyses (q1, q2, q3) complementary to our FcLM analyses for assessing the conflict in our dataset (Fig. 2a; note that the coalescent method does not test for putative confounding signal *per se*). It has been suggested that low data coverage may have a negative impact on summary coalescent methods (Sayyari et al., 2017). In order to account for this negative effect, we selected only these gene partitions with at least 95% species coverage (min. = 115 leaf terminals, 2083 genes) and repeated coalescent species tree analyses both at the amino-acid and nucleotide sequence levels. Finally, results of the different coalescent analyses were compared to those based on domain-based partitioned and gene-based partitioned concatenated supermatrices (see Additional file 2 for details). We used ETE v. 3.0 (Huerta-Cepas et al., 2016) to visualize quartet support, as an indication of gene tree conflict, on the species trees that were inferred with ASTRAL (e.g., Fig. 2.2a).

## 2.5.7. Estimation of divergence times of Neuropterida

We used 129 meta-partitions of the decisive amino-acid supermatrix (supermatrix E-Decisive, see Additional file 2 and Table 1) to estimate the divergence times of the major lineages of Neuropterida based on 12 fossil calibrations (Additional file 1: Table S10). The fossil calibrations were selected according to the criteria described by Parham et al. (2012) (see Additional file 2). We extracted the 129 meta-partitions from the decisive supermatrix and reestimated the most suitable substitution models for each individual meta-partition using IQ-TREE v. 1.6.6 (with the AICc criterion), by restricting model selection to a set of amino-acid substitution matrices available in the PAML package (Yang, 2007, JTT+G, LG+G,WAG +G, DAYHOFF + G, JTTDCMUT + G, DCMUT + G) and by using the fixed topology of the best ML tree. Subsequently, substitution rates per time unit for each meta-partition were estimated with codeml v. 4.9e (part of the PAML software suite) under the assumption of a strict clock (clock = 1), and by using the fixed topology of the best ML tree and the above-selected substitution models. The age of the root was fixed at 362.35 million years ago (Mya) in each ML analysis. This root age was derived as the average between the oldest known hexapod fossils at 411 Mya and the minimum age 313.7 Mya for Aparaglossata (Wolfe et al., 2016) (i.e., Holometabola without Hymenoptera, see Peters et al., 2014). The purpose of these analyses was to calculate a rough estimate of the mean rate prior for each meta-partition to be used for estimating the divergence times in MCMCTree v. 4.9e (part of the PAML software suite (Yang, 2007).

Calculation of the Hessian matrices followed the standard procedure, applying the fitted substitution models (+ G with four rate categories) for each meta-partition (Additional file 2). Similarly to the approach proposed by Misof et al. (2014),

divergence time estimation was performed for each of the 129 meta-partitions separately with the approximate likelihood method (dos Reis and Yang, 2011). We used the same set of calibration points (Additional file 1: Table S10), the independent-rates model (Rannala and Yang, 2007) and the topology of the best ML tree for each separate analysis. The estimated substitution rate of each meta-partition was used as the mean ($\mu$)ofthe Dirichlet-gamma prior (rgene_gamma) in MCMCTree v. 4.9e. We specified a hard maximum bound for the age of the root at 411 Mya in all analyses and ran each MCMCTree chain for 550,000 generations, sampling every 10th generation and discarding the first 50,000 samples as a burn-in (Additional file 2). For each meta-partition, three different analyses were performed:

1. Two independent analyses (run 1 and run 2) with the same calibrations and diffuse rate priors ($\alpha$ =2) to check for repeatability of the analyses (Additional file 2).

2. One calibration without data (usedata = 0) to assess whether or not the results without data were significantly different, implying that the data harbor sufficient information for reliably estimating divergence times.

For each of the three separate analyses (two analyses with data and one without) parameter outputs of the separate analyses of the meta-partitions were combined in a single MCMC summarized file. We mapped the posterior mean node ages and 95% confident intervals (equal-tail CI) on the overall best ML tree (Fig. 2.1). The branch lengths of the resulting chronogram were calculated as the posterior mean node-age difference between two nodes. The posterior node-age estimates from the 129 meta-

partitions were used to calculate median posterior node-age estimates in R v. 3.4.3 (R Core Team, 2015) (Fig. 2.3, Additional file 2). The datasets used for estimation of divergence times and all the analyzed supermatrices are deposited in the Dryad repository (see availability of Supplementary Materials).

### 2.5.8. Tracing the evolution of larval characters within Neuropterida

In addition to the informal discussion of implications of the proposed phylogeny for our understanding of the evolution of neuropterid insects in general, we also formally analyzed character transformations (Fig. 2.1) with Mesquite v. 3.2 (Maddison and Maddison, 2001). For this analysis we selected a data matrix comprising 86 larval characters from a previously published morphological study with focus on Neuroptera (Jandausch et al., 2019) (see also Beutel et al., 2010a and Jandausch et al., 2018). We analyzed this character matrix under the constrained topology of our best ML tree (Fig. 2.1) using maximum parsimony (see Additional file 4). A summary of the interpretation of results for the most important characters is provided in Additional file 1: Table S5.

Previous ACSRs of the larval ecologies of Neuropterida have suggested that ancestral Neuropterida most likely had an aquatic larva (Wang et al., 2017; Winterton et al., 2018). However, a clade of Nevrorthidae + Sisyridae as sister to Osmylidae has not been inferred in previous phylogenomic studies, and the taxon sampling of outgroup species was not as extensive as in our study (Wang et al., 2017; Winterton et al., 2018). Therefore, we additionally used a Bayesian approach to reconstruct the ancestral states of larval ecologies of Neuropterida. Specifically, we used the stochastic character mapping method (SCM) (Bollback, 2006; Huelsenbeck et al., 2003), as implemented in the R package phytools v. 0.6.99 (Revell, 2012) (see Additional file 2 for details and

additional sensitivity analyses). We simulated 10,000 character histories conditioned on

the topology and branch lengths of the best ML tree (Fig. 2.1), and by using the best

fitted model of character evolution. The results of the SCM analyses were visualized

using ape v. 5.3 (Paradis and Schliep, 2018).

## 2.6. References

Ababneh, F., Jermiin, L.S., Ma, C., Robinson, J., 2006. Matched-pairs tests of homogeneity with applications to homologous nucleotide sequences. Bioinformatics 22, 1225–1231.

Aberer, A.J., Krompass, D., Stamatakis, A., 2013. Pruning rogue taxa improves phylogenetic accuracy: an efficient algorithm and webservice. Syst. Biol. 62, 162–166.

Achtelig, M., 1978. Entwicklung und Morphologie der inneren und ausseren weiblichen Genitalorgane der Kamelhalsfliegen (Neuropteroidea: Raphidioptera). Entomol. Ger. 4, 140–163.

Achtelig, M., 1975. Die Abdomenbasis der Neuropteroidea (Insecta, Holometabola). Eine vergleichend anatomische Untersuchung des Skeletts und der Muskulatur. Zoomorphologie 82, 201–242.

Achtelig, M., 1967. Über die Anatomie des Kopfes von *Raphidia flavipes* Stein und die Verwandtschaftsbeziehungen der Raphidiidae zu den Megaloptera. Zool. J. Abt. Anat. Ontog. Tiere 84, 249–312.

Adams, M.D., Celniker, S.E., Holt, R.A., Evans, C.A., Gocayne, J.D., Amanatides, P.G., Scherer, S.E., Li, P.W., Hoskins, R.A., Galle, R.F., George, R.A., Lewis, S.E., Richards, S., Ashburner, M., Henderson, S.N., Sutton, G.G., Wortman, J.R., Yandell, M.D., Zhang, Q., Chen, L.X., Brandon, R.C., Rogers, Y.C., Blazej, R.G., Champe, M., Pfeiffer, B.D., Wan, K.H., Doyle, C., Baxter, E.G., Helt, G., Nelson, C.R., Miklos, G.L.G., Abril, J.F., Agbayani, A., An, H., Andrews-pfannkoch, C., Baldwin, D., Ballew, R.M., Basu, A., Baxendale, J., Bayraktaroglu, L., Beasley,

E.M., Beeson, K.Y., Benos, P. V, Berman, B.P., Bhandari, D., Bolshakov, S., Borkova, D., Botchan, M.R., Bouck, J., Brokstein, P., Brottier, P., Burtis, K.C., Busam, D. a, Butler, H., Cadieu, E., Center, A., Chandra, I., Cherry, J.M., Cawley, S., Dahlke, C., Davenport, L.B., Davies, P., Pablos, B. De, Delcher, A., Deng, Z., Mays, A.D., Dew, I., Dietz, S.M., Dodson, K., Doup, L.E., Downes, M., Dugan-rocha, S., Dunkov, B.C., Dunn, P., Durbin, K.J., Evangelista, C.C., Ferraz, C., Ferriera, S., Fleischmann, W., Fosler, C., Gabrielian, A.E., Garg, N.S., Gelbart, W.M., Glasser, K., Glodek, A., Gong, F., Gorrell, J.H., Gu, Z., Guan, P., Harris, M., Harris, N.L., Harvey, D., Heiman, T.J., Hernandez, J.R., Houck, J., Hostin, D., Houston, K. a, Howland, T.J., Wei, M., Ibegwam, C., Jalali, M., Kalush, F., Karpen, G.H., Ke, Z., Kennison, J. a, Ketchum, K. a, Kimmel, B.E., Kodira, C.D., Kraft, C., Kravitz, S., Kulp, D., Lai, Z., Lasko, P., Lei, Y., Levitsky, A. a, Li, J., Li, Z., Liang, Y., Lin, X., Liu, X., Mattei, B., McIntosh, T.C., Mcleod, M.P., Mcpherson, D., Merkulov, G., Milshina, N. V, Mobarry, C., Morris, J., Moshre, A., Mount, S.M., Moy, M., Murphy, B., Murphy, L., Muzny, D.M., Nelson, D.L., Nelson, D.R., Nelson, K. a, Nixon, K., Nusskern, D.R., Pacleb, J.M., Palazzolo, M., Pittman, G.S., Pan, S., Pollard, J., Puri, V., Reese, M.G., Reinert, K., Remington, K., Saunders, R.D.C., Scheeler, F., Shen, H., Shue, B.C., Side, I., Simpson, M., Skupski, M.P., Smith, T., Spier, E., Spradling, A.C., Stapleton, M., Strong, R., Sun, E., Svirskas, R., Tector, C., Turner, R., Venter, E., Wang, A.H., Wang, X., Wang, Z., Wassarman, D. a, Weinstock, G.M., Weissenbach, J., Williams, S.M., Woodage, T., Worley, K.C., Wu, D., Yang, S., Yao, Q.A., Ye, J., Yeh, R., Zaveri, J.S., Zhan, M., Zhang, G., Zhao, Q., Zheng, L., Zheng, X.H., Zhong, F.N., Zhong, W., Zhou, X., Zhu, S., Zhu, X., Smith, H.O., Gibbs, R. a,

Myers, E.W., Rubin, G.M., Venter, J.C., 2000. The genome sequence of *Drosophila melanogaster*. Science. 287, 2185–2195.

Aspöck, H., 2002. The biology of Raphidioptera: a review of present knowledge. Acta Zool. Acad. Sci. Hungaricae 48, 35–50.

Aspöck, H., 2000. Der endkreidezeitliche Impakt und das Überleben der Raphidiopteren., in: Int. Entomol. Tag. 1999, Entomologica Basiliensia. pp. 223–233.

Aspöck, H., 1998. Distribution and biogeography of the order Raphidioptera: updated facts and a new hypothesis. Acta Zool. Fenn. 209, 33–44.

Aspöck, H., Aspöck, U., Hölzel, H., 1980. Die Neuropteren Europas. Eine zusammenfassende Darstellung der Systematik, Ökologie und Chorologie der Neuropteroidea (Megaloptera, Raphidioptera, Planipennia) Europas. Goecke & Evers, Krefeld.

Aspöck, U., 2002. Phylogeny of the Neuropterida (Insecta: Holometabola). Zool. Scr. 31, 51–55.

Aspöck, U., 1995. Neue Hypothesen zum System der Neuropterida. Mitteilungen der Dtsch. Gesellschaft fur Allg. und Angew. Entomol. 10, 633–636.

Aspöck, U., Aspöck, H., 2008. Phylogenetic relevance of the genital sclerites of Neuropterida (Insecta: Holometabola). Syst. Entomol. 33, 97–127.

Aspöck, U., Aspöck, H., 2007. Verbliebene Vielfalt vergangener Blüte. Zur Evolution, Phylogenie und Biodiversität der Neuropterida (Insecta: Endopterygota). Denisia 20, 451–516.

Aspöck, U., Aspöck, H., Liu, X., 2017. The Nevrorthidae, mistaken at all times:

  phylogeny and review of present knowledge (Holometabola, Neuropterida,

  Neuroptera). Dtsch. Entomol. Zeitschrift 64, 77–110.

Aspöck, U., Haring, E., Aspöck, H., 2012a. The phylogeny of the Neuropterida: long

  lasting and current controversies and challenges (Insecta: Endopterygota).

  Arthropod Syst. Phylogeny 70, 119–129.

Aspöck, U., Haring, E., Aspöck, H., 2012b. Biogeographical implications of a

  molecular phylogeny of the Raphidiidae (Raphidioptera). Mitteilungen der Dtsch.

  Gesellschaft für Allg. und Angew. Entomol. 18, 575–582.

Aspöck, U., Plant, J.D., Nemeschkal, H.L., 2001. Cladistic analysis of Neuroptera and

  their systematic position within Neuropterida (Insecta: Holometabola:

  Neuropterida: Neuroptera). Syst. Entomol. 26, 73–86.

Badano, D., Aspöck, U., Aspöck, H., Cerretti, P., 2017. Phylogeny of

  Myrmeleontiformia based on larval morphology (Neuropterida: Neuroptera). Syst.

  Entomol. 42, 94–117.

Badano, D., Engel, M.S., Basso, A., Wang, B., Cerretti, P., 2018. Diverse Cretaceous

  larvae reveal the evolutionary and behavioural history of antlions and lacewings.

  Nat. Commun. 9, 3257.

Beutel, R.G., Friedrich, F., Aspöck, U., 2010a. The larval head of Nevrorthidae and the

  phylogeny of Neuroptera (Insecta). Zool. J. Linn. Soc. 158, 533–562.

Beutel, R.G., Friedrich, F., Hörnschemeyer, T., Pohl, H., Hünefeld, F., Beckmann, F.,

  Meier, R., Misof, B., Whiting, M.F., Vilhelmsen, L., 2011. Morphological and

molecular evidence converge upon a robust phylogeny of the megadiverse Holometabola. Cladistics 27, 341–355.

Beutel, R.G., Gorb, S.N., 2001. Ultrastructure of attachment specializations of hexapods (Arthropoda): evolutionary patterns inferred from a revised ordinal phylogeny. J. Zool. Syst. Evol. Res. 39, 177–207.

Beutel, R.G., Zimmermann, D., Krauß, M., Randolf, S., Wipfler, B., 2010b. Head morphology of *Osmylus fulvicephalus* (Osmylidae, Neuroptera) and its phylogenetic implications. Org. Divers. Evol. 10, 311–329.

Bollback, J.P., 2006. SIMMAP: stochastic character mapping of discrete traits on phylogenies. BMC Bioinformatics 7, 88.

Boudreaux, H.B., 1979. Arthropod phylogeny with special reference to insects. John Wiley & Sons, New York.

Bowker, A.H., 1948. A test for symmetry in contingency tables. J. Am. Stat. Assoc. 43, 572–574.

Brushwein, J.R., 1987. Bionomics of *Lomamyia hamata* (Neuroptera: Berothidae). Ann. Entomol. Soc. Am. 80, 671–679.

Cameron, S.L., Sullivan, J., Song, H., Miller, K.B., Whiting, M.F., 2009. A mitochondrial genome phylogeny of the Neuropterida (lace-wings, alderflies and snakeflies) and their relationship to the other holometabolous insect orders. Zool. Scr. 38, 575–590.

Cloutier, A., Sackton, T.B., Grayson, P., Clamp, M., Baker, A.J., Edwards, S. V., 2019. Whole-genome analyses resolve the phylogeny of flightless birds (Palaeognathae) in the presence of an empirical anomaly zone. Syst. Biol. 68, 937–955.

Contreras-Ramos, A., 2004. Is the family Corydalidae (Neuropterida, Megaloptera) a monophylum? Denisia 13, 135–140.

de Queiroz, A., Gatesy, J., 2007. The supermatrix approach to systematics. Trends Ecol. Evol. 22, 34–41.

Dejean, A., Canard, M., 1990. Reproductive behaviour of *Trichoscelia santreni* (Navas) (Neuroptera: Mantispidae) and parasitization of the colonies of *Polybia diguetana* R. Du Buysson (Hymenoptera: Vespidae). Neuroptera Int. 6, 19–26.

Dell'Ampio, E., Meusemann, K., Szucsich, N.U., Peters, R.S., Meyer, B., Borner, J., Petersen, M., Aberer, A.J., Stamatakis, A., Walzl, M.G., Minh, B.Q., von Haeseler, A., Ebersberger, I., Pass, G., Misof, B., 2014. Decisive data sets in phylogenomics: lessons from studies on the phylogenetic relationships of primarily wingless insects. Mol. Biol. Evol. 31, 239–249.

dos Reis, M., Yang, Z., 2011. Approximate likelihood calculation on a phylogeny for Bayesian estimation of divergence times. Mol. Biol. Evol. 28, 2161–2172.

Edwards, S. V., 2009. Is a new and general theory of molecular systematics emerging? Evolution 63, 1–19.

Edwards, S. V., Xi, Z., Janke, A., Faircloth, B.C., McCormack, J.E., Glenn, T.C., Zhong, B., Wu, S., Lemmon, E.M., Lemmon, A.R., Leache, A.D., Liu, L., Davis, C.C., 2016. Implementing and testing the multispecies coalescent model: a valuable paradigm for phylogenomics. Mol. Phylogenet. Evol. 94, 447–462.

Engel, M.S., Winterton, S.L., Breitkreuz, L.C.V., 2018. Phylogeny and evolution of Neuropterida: where have wings of lace taken us? Annu. Rev. Entomol. 63, 531–551.

Evangelista, D., Thouzé, F., Kohli, M.K., Lopez, P., Legendre, F., 2018. Topological support and data quality can only be assessed through multiple tests in reviewing Blattodea phylogeny. Mol. Phylogenet. Evol. 128, 112–122.

Faulkner, D.K., 1990. Current knowledge of the biology of the moth-lacewing *Oliarces clara* Banks (Insecta: Neuroptera: Ithonidae), in: Mansell, M.W., Aspöck, H. (Eds.), Advances in Neuropterology, Third International Symposium, Kruger National Park, South Africa. 3–4 February 1988. Department of Agricultural Development, Pretoria, pp. 197–203.

Felsenstein, J., 1985. Confidence limits on phylogenies: an approach using the bootstrap. Evolution 39, 783–791.

Fernandez, R., Edgecombe, G.D., Giribet, G., 2016. Exploring phylogenetic relationships within Myriapoda and the effects of matrix composition and occupancy on phylogenomic reconstruction. Syst. Biol. 65, 871–889.

Fernandez, R., Sharma, P., Tourinho, A.L., Giribet, G., 2017. The Opiliones Tree of Life: shedding light on harvestmen relationships through transcriptomics. Proc. R. Soc. B 284, 20162340.

Feuda, R., Dohrmann, M., Pett, W., Philippe, H., Rota-Stabelli, O., Lartillot, N., Wörheide, G., Pisani, D., 2017. Improved modeling of compositional heterogeneity supports sponges as sister to all other animals. Curr. Biol. 27, 3864–3870.

Foster, P., 2004. Modeling compositional heterogeneity. Syst. Biol. 53, 485–495.

Frandsen, P.B., Calcott, B., Mayer, C., Lanfear, R., 2015. Automatic selection of
partitioning schemes for phylogenetic analyses using iterative k-means clustering
of site rates. BMC Evol. Biol. 15, 13.

Gadagkar, S.R., Rosenberg, M.S., Kumar, S., 2005. Inferring species phylogenies from
multiple genes: concatenated sequence tree versus consensus gene tree. J. Exp.
Zool. Part B Mol. Dev. Evol. 304B, 64–74.

Gatesy, J., Springer, M.S., 2014. Phylogenetic analysis at deep timescales: unreliable
gene trees, bypassed hidden support, and the coalescence/concatalescence
conundrum. Mol. Phylogenet. Evol. 80, 231–266.

Gaumont, J., 1976. L'appareil digestif des larves de Planipennes. Ann. des Sci. Nat. –
Zool. Biol. Anim. 18, 145–250.

Gillung, J.P., Winterton, S.L., Bayless, K.M., Khouri, Z., Borowiec, M.L., Yeates, D.,
Kimsey, L.S., Misof, B., Shin, S., Zhou, X., Mayer, C., Petersen, M., Wiegmann,
B.M., 2018. Anchored phylogenomics unravels the evolution of spider flies
(Diptera, Acroceridae) and reveals discordance between nucleotides and amino
acids. Mol. Phylogenet. Evol. 128, 233–245.

Gouy, M., Guindon, S., Gascuel, O., 2010. Sea view version 4: a multiplatform
graphical user interface for sequence alignment and phylogenetic tree building.
Mol. Biol. Evol. 27, 221–224.

Guindon, S., Dufayard, J.F., Lefort, V., Anisimova, M., Hordijk, W., Gascuel, O., 2010.
New algorithms and methods to estimate maximum-likelihood phylogenies:
assessing the performance of PhyML 3.0. Syst. Biol. 59, 307–321.

Gusten, R., Dettner, K., 1992. The prothoracic gland of the Chrysopidae (Neuropteroidea: Planipennia), in: Zombori, L., Peregovits, L. (Eds.), Proceedings of the Fourth European Congress of Entomology and the XIII Internationale Symposium Fur Die Entomofaunistik Mitteleuropas. Gödölö. Hungary, 1991. Hungarian Natural History Museum, Budapest, Hungary, pp. 56–60.

Haring, E., Aspöck, H., Bartel, D., Aspöck, U., 2011. Molecular phylogeny of the Raphidiidae (Raphidioptera). Syst. Entomol. 36, 16–30.

Haring, E., Aspöck, U., 2004. Phylogeny of the Neuropterida: a first molecular approach. Syst. Entomol. 29, 415–430.

Hennig, W., 1969. Die Stammesgeschichte der Insekten, 1st ed. Waldemar Kramer, Frankfurt am Main.

Hoang, D.T., Chernomor, O., von Haeseler, A., Minh, B.Q., Le, S.V., 2018. UFBoot2: improving the ultrafast bootstrap approximation. Mol. Biol. Evol. 35, 518–522.

Huelsenbeck, J.P., Nielsen, R., Bollback, J.P., 2003. Stochastic mapping of morphological characters. Syst. Biol. 52, 131–158.

Huelsenbeck, J.P., Rannala, B., 2004. Frequentist properties of Bayesian posterior probabilities of phylogenetic trees under simple and complex substitution models. Syst. Biol. 53, 904–913.

Huerta-Cepas, J., Serra, F., Bork, P., 2016. ETE 3: Reconstruction, analysis, and visualization of phylogenomic data. Mol. Biol. Evol. 33, 1635–1638.

Jandausch, K., Beutel, R.G., Bellstedt, R., 2019. The larval morphology of the

   spongefly *Sisyra nigra* (Retzius, 1783) (Neuroptera: Sisyridae). J. Morphol. 280,

   1742–1758.

Jandausch, K., Pohl, H., Aspöck, U., Winterton, S.L., Beutel, R.G., 2018. Morphology

   of the primary larva of *Mantispa aphavexelte* Aspöck & Aspöck, 1994

   (Neuroptera: Mantispidae) and phylogenetic implications to the order of

   Neuroptera. Arthropod Syst. Phylogeny 76, 529–560.

Jeffroy, O., Brinkmann, H., Delsuc, F., Philippe, H., 2006. Phylogenomics: the

   beginning of incongruence? Trends Genet. 22, 225–231.

Jermiin, L.S., Ho, S.Y.W., Ababneh, F., Robinson, J., Larkum, A.W.D., 2004. The

   biasing effect of compositional heterogeneity on phylogenetic estimates may be

   underestimated. Syst. Biol. 53, 638–643.

Jermiin, L.S., Jayaswal, V., Ababneh, F., Robinson, J., 2008. Phylogenetic model

   evaluation, in: Keith, J.M. (Ed.), Bioinformatics. Methods in molecular biology™,

   vol 452. Humana Press, Totowa, pp. 331–363.

Johnson, K.P., Dietrich, C.H., Friedrich, F., Beutel, R.G., Wipfler, B., Peters, R.S.,

   Allen, J.M., Petersen, M., Donath, A., Walden, K.K.O., Kozlov, A.M.,

   Podsiadlowski, L., Mayer, C., Meusemann, K., Vasilikopoulos, A., Waterhouse,

   R.M., Cameron, S.L., Weirauch, C., Swanson, D.R., Percy, D.M., Hardy, N.B.,

   Terry, I., Liu, S., Zhou, X., Misof, B., Robertson, H.M., Yoshizawa, K., 2018.

   Phylogenomics and the evolution of hemipteroid insects. Proc. Natl. Acad. Sci.

   115, 12775–12780.

Kalyaanamoorthy, S., Minh, B.Q., Wong, T.K.F., von Haeseler, A., Jermiin, L.S., 2017. ModelFinder: fast model selection for accurate phylogenetic estimates. Nat. Methods 14, 587–589.

Katoh, K., Standley, D.M., 2013. MAFFT multiple sequence alignment software version 7: improvements in performance and usability. Mol. Biol. Evol. 30, 772–780.

Kjer, K.M., Carle, F.L., Litman, J., Ware, J., 2006. A molecular phylogeny of Hexapoda. Arthropod Syst. Phylogeny 64, 35–44.

Komatsu, T., 2014. Larvae of the Japanese termitophilous predator *Isoscelipteron okamotonis* (Neuroptera, Berothidae) use their mandibles and silk web to prey on termites. Insectes Soc. 61, 203–205.

Kristensen, N.P., 1991. Phylogeny of extant hexapods, in: Naumann, I.D., Carne, P.B., Lawrence, J.F., Nielsen, E.S., Spradberry, J.P., Taylor, R.W., Whitten, M.J., Littlejohn, M.J. (Eds.), The insects of Australia: a textbook for students and research workers. Melbourne University Press, Melbourne, pp. 125–140.

Kubatko, L.S., Degnan, J.H., 2007. Inconsistency of phylogenetic estimates from concatenated data under coalescence. Syst. Biol. 56, 17–24.

Kück, P., Longo, G.C., 2014. FASconCAT-G: extensive functions for multiple sequence alignment preparations concerning phylogenetic studies. Front. Zool. 11, 81.

Kück, P., Meusemann, K., Dambach, J., Thormann, B., von Reumont, B.M., Wägele, J.W., Misof, B., 2010. Parametric and non-parametric masking of randomness in sequence alignments can be improved and leads to better resolved trees. Front. Zool. 7, 10.

Kück, P., Struck, T.H., 2014. BaCoCa - a heuristic software tool for the parallel

assessment of sequence biases in hundreds of gene and taxon partitions. Mol.

Phylogenet. Evol. 70, 94–98.

Kück, P., Wilkinson, M., Groß, C., Foster, P.G., Wägele, J.W., 2017. Can quartet

analyses combining maximum likelihood estimation and Hennigian logic

overcome long branch attraction in phylogenomic sequence data? PLoS One 12,

e0183393.

Kumar, S., Filipski, A.J., Battistuzzi, F.U., Kosakovsky Pond, S.L., Tamura, K., 2012.

Statistics and truth in phylogenomics. Mol. Biol. Evol. 29, 457–472.

Lanfear, R., Frandsen, P.B., Wright, A.M., Senfeld, T., Calcott, B., 2017. Partitionfinder

2: New methods for selecting partitioned models of evolution for molecular and

morphological phylogenetic analyses. Mol. Biol. Evol. 34, 772–773.

Lartillot, N., Brinkmann, H., Philippe, H., 2007. Suppression of long-branch attraction

artefacts in the animal phylogeny using a site-heterogeneous model. BMC Evol.

Biol. 7, S4.

Laumer, C.E., Bekkouche, N., Kerbl, A., Goetz, F., Neves, R.C., Sørensen, M. V,

Kristensen, R.M., Hejnol, A., Dunn, C.W., Giribet, G., Worsaae, K., 2015.

Spiralian phylogeny informs the evolution of microscopic lineages. Curr. Biol. 25,

2000–2006.

Le, S.Q., Lartillot, N., Gascuel, O., 2008. Phylogenetic mixture models for proteins.

Philos. Trans. R. Soc. Lond. B. Biol. Sci. 363, 3965–3976.

Lemmon, A.R., Moriarty, E.C., 2004. The importance of proper model assumption in

Bayesian phylogenetics. Syst. Biol. 53, 265–277.

Lemoine, F., Domelevo Entfellner, J.-B., Wilkinson, E., Correia, D., Dávila Felipe, M.,
De Oliveira, T., Gascuel, O., 2018. Renewing Felsenstein's phylogenetic bootstrap
in the era of big data. Nature 556, 452–456.

Li, Z., De La Torre, A.R., Sterck, L., Cánovas, F.M., Avila, C., Merino, I., Cabezas,
J.A., Cervera, M.T., Ingvarsson, P.K., Van de Peer, Y., 2017. Single-copy genes as
molecular markers for phylogenomic studies in seed plants. Genome Biol. Evol. 9,
1130–1147.

Liu, L., Xi, Z., Wu, S., Davis, C.C., Edwards, S. V., 2015. Estimating phylogenetic
trees from genome-scale data. Ann. N. Y. Acad. Sci. 1360, 36–53.

Liu, X., Lü, Y., Aspöck, H., Yang, D., Aspöck, U., 2016. Homology of the genital
sclerites of Megaloptera (Insecta: Neuropterida) and their phylogenetic relevance.
Syst. Entomol. 41, 256–286.

Liu, X., Winterton, S.L., Wu, C., Piper, R., Ohl, M., 2015. A new genus of mantidflies
discovered in the Oriental region, with a higher-level phylogeny of Mantispidae
(Neuroptera) using DNA sequences and morphology. Syst. Entomol. 40, 183–206.

Machado, R.J.P., Gillung, J.P., Winterton, S.L., Garzón-Orduña, I.J., Lemmon, A.R.,
Lemmon, E.M., Oswald, J.D., 2019. Owlflies are derived antlions: anchored
phylogenomics supports a new phylogeny and classification of Myrmeleontidae
(Neuroptera). Syst. Entomol. 44, 418–450.

MacLeod, E.G., 1964. Comparative morphological studies on the head capsule and
cervix of the larval Neuroptera (Insecta). Harvard University.

Maddison, W.P., Maddison, D.R., 2001. Mesquite: a modular system for evolutionary
analysis.

McKenna, D.D., Farrell, B.D., 2010. 9-genes reinforce the phylogeny of Holometabola
    and yield alternate views on the phylogenetic placement of Strepsiptera. PLoS One
    5, e11887.

Meusemann, K., von Reumont, B.M., Simon, S., Roeding, F., Strauss, S., Kück, P.,
    Ebersberger, I., Walzl, M., Pass, G., Breuers, S., Achter, V., von Haeseler, A.,
    Burmester, T., Hadrys, H., Wägele, J.W., Misof, B., 2010. A phylogenomic
    approach to resolve the arthropod tree of life. Mol. Biol. Evol. 27, 2451–2464.

Misof, B., Liu, S., Meusemann, K., Peters, R.S., Donath, A., Mayer, C., Frandsen, P.B.,
    Ware, J., Flouri, T., Beutel, R.G., Niehuis, O., Petersen, M., Izquierdo-Carrasco,
    F., Wappler, T., Rust, J., Aberer, A.J., Aspöck, U., Aspöck, H., Bartel, D., Blanke,
    A., Berger, S., Bohm, A., Buckley, T.R., Calcott, B., Chen, J., Friedrich, F., Fukui,
    M., Fujita, M., Greve, C., Grobe, P., Gu, S., Huang, Y., Jermiin, L.S., Kawahara,
    A.Y., Krogmann, L., Kubiak, M., Lanfear, R., Letsch, H., Li, Y., Li, Z., Li, J., Lu,
    H., Machida, R., Mashimo, Y., Kapli, P., McKenna, D.D., Meng, G., Nakagaki, Y.,
    Navarrete-Heredia, J.L., Ott, M., Ou, Y., Pass, G., Podsiadlowski, L., Pohl, H., von
    Reumont, B.M., Schutte, K., Sekiya, K., Shimizu, S., Slipinski, A., Stamatakis, A.,
    Song, W., Su, X., Szucsich, N.U., Tan, M., Tan, X., Tang, M., Tang, J.,
    Timelthaler, G., Tomizuka, S., Trautwein, M., Tong, X., Uchifune, T., Walzl, M.G.,
    Wiegmann, B.M., Wilbrandt, J., Wipfler, B., Wong, T.K.F., Wu, Q., Wu, G., Xie,
    Y., Yang, S., Yang, Q., Yeates, D.K., Yoshizawa, K., Zhang, Q., Zhang, R., Zhang,
    W., Zhang, Y., Zhao, J., Zhou, C., Zhou, L., Ziesmann, T., Zou, S., Xu, X., Yang,
    H., Wang, J., Kjer, K.M., Zhou, X., 2014. Phylogenomics resolves the timing and
    pattern of insect evolution. Science. 346, 763–767.

Misof, B., Meyer, B., von Reumont, B.M., Kück, P., Misof, K., Meusemann, K., 2013. Selecting informative subsets of sparse supermatrices increases the chance to find correct trees. BMC Bioinformatics 14, 348.

Misof, B., Misof, K., 2009. A Monte Carlo approach successfully identifies randomness in multiple sequence alignments: a more objective means of data exclusion. Syst. Biol. 58, 21–34.

Monserrat, V.J., 1996. Larval stages of European Nemopterinae, with systematic considerations on the family Nemopteridae (Insecta, Neuroptera). Dtsch. Entomol. Zeitschrift 43, 99–121.

Nguyen, L.-T., Schmidt, H.A., Von Haeseler, A., Minh, B.Q., 2015. IQ-TREE: A fast and effective stochastic algorithm for estimating maximum-likelihood phylogenies. Mol. Biol. Evol. 32, 268–274.

Niehuis, O., Hartig, G., Grath, S., Pohl, H., Lehmann, J., Tafer, H., Donath, A., Krauss, V., Eisenhardt, C., Hertel, J., Petersen, M., Mayer, C., Meusemann, K., Peters, R.S., Stadler, P.F., Beutel, R.G., Bornberg-Bauer, E., McKenna, D.D., Misof, B., 2012. Genomic and morphological evidence converge to resolve the enigma of Strepsiptera. Curr. Biol. 22, 1309–1313.

Nygaard, S., Zhang, G., Schiøtt, M., Li, C., Wurm, Y., Hu, H., Zhou, J., Ji, L., Qiu, F., Rasmussen, M., Pan, H., Hauser, F., Krogh, A., Grimmelikhuijzen, C.J.P., Wang, J., Boomsma, J.J., 2011. The genome of the leaf-cutting ant *Acromyrmex echinatior* suggests key adaptations to advanced social life and fungus farming. Genome Res. 21, 1339–1348.

Oswald, J.D., 2019. LDL Neuropterida species of the world (version July 2018)

    [WWW Document]. Species 2000 ITIS Cat. Life, 26th Febr. 2019. URL

    www.catalogueoflife.org/col (accessed 3.12.19).

Paradis, E., Schliep, K., 2018. Ape 5.0: an environment for modern phylogenetics and

    evolutionary analyses in R. Bioinformatics 35, 526–528.

Parham, J.F., Donoghue, P.C.J., Bell, C.J., Calway, T.D., Head, J.J., Holroyd, P.A.,

    Inoue, J.G., Irmis, R.B., Joyce, W.G., Ksepka, D.T., Patané, J.S.L., Smith, N.D.,

    Tarver, J.E., Van Tuinen, M., Yang, Z., Angielczyk, K.D., Greenwood, J.M.,

    Hipsley, C.A., Jacobs, L., Makovicky, P.J., Müller, J., Smith, K.T., Theodor, J.M.,

    Warnock, R.C.M., Benton, M.J., 2012. Best practices for justifying fossil

    calibrations. Syst. Biol. 61, 346–359.

Pattengale, N.D., Alipour, M., Bininda-Emonds, O.R.P., Moret, B.M.E., Stamatakis, A.,

    2010. How many bootstrap replicates are necessary? J. Comput. Biol. 17, 337–

    354.

Pease, J.B., Brown, J.W., Walker, J.F., Hinchliff, C.E., Smith, S.A., 2018. Quartet

    sampling distinguishes lack of support from conflicting support in the green plant

    tree of life. Am. J. Bot. 105, 385–403.

Peters, R.S., Krogmann, L., Mayer, C., Donath, A., Gunkel, S., Meusemann, K.,

    Kozlov, A., Podsiadlowski, L., Petersen, M., Lanfear, R., Diez, P.A., Heraty, J.,

    Kjer, K.M., Klopfstein, S., Meier, R., Polidori, C., Schmitt, T., Liu, S., Zhou, X.,

    Wappler, T., Rust, J., Misof, B., Niehuis, O., 2017. Evolutionary history of the

    Hymenoptera. Curr. Biol. 27, 1013–1018.

Peters, R.S., Meusemann, K., Petersen, M., Mayer, C., Wilbrandt, J., Ziesmann, T., Donath, A., Kjer, K.M., Aspöck, U., Aspöck, H., Aberer, A., Stamatakis, A., Friedrich, F., Hünefeld, F., Niehuis, O., Beutel, R.G., Misof, B., 2014. The evolutionary history of holometabolous insects inferred from transcriptome-based phylogeny and comprehensive morphological data. BMC Evol. Biol. 14, 52.

Petersen, M., Meusemann, K., Donath, A., Dowling, D., Liu, S., Peters, R.S., Podsiadlowski, L., Vasilikopoulos, A., Zhou, X., Misof, B., Niehuis, O., 2017. Orthograph: a versatile tool for mapping coding nucleotide sequences to clusters of orthologous genes. BMC Bioinformatics 18, 111.

Phillips, M.J., Delsuc, F., Penny, D., 2004. Genome-scale phylogeny and the detection of systematic biases. Mol. Biol. Evol. 21, 1455–1458.

R Core Team, 2015. R: A language and environment for statistical computing.

Randolf, S., Zimmermann, D., Aspöck, U., 2017. Head anatomy of adult *Coniopteryx pygmaea*: effects of miniaturization and the systematic position of Coniopterygidae (Insecta: Neuroptera). Arthropod Struct. Dev. 46, 304–322.

Randolf, S., Zimmermann, D., Aspöck, U., 2014. Head anatomy of adult *Nevrorthus apatelios* and basal splitting events in Neuroptera (Neuroptera: Nevrorthidae). Arthropod Syst. Phylogeny 72, 111–136.

Randolf, S., Zimmermann, D., Aspöck, U., 2013. Head anatomy of adult *Sisyra terminalis* (Insecta: Neuroptera: Sisyridae) - functional adaptations and phylogenetic implications. Arthropod Struct. Dev. 42, 565–582.

Rannala, B., Yang, Z., 2007. Inferring speciation times under an episodic molecular clock. Syst. Biol. 56, 453–466.

Rannala, B., Yang, Z., 1996. Probability distribution of molecular evolutionary trees: a new method of phylogenetic inference. J. Mol. Evol. 43, 304–311.

Redborg, K.E., 1998. Biology of the Mantispidae. Annu. Rev. Entomol. 43, 175–194.

Reddy, S., Kimball, R.T., Pandey, A., Hosner, P.A., Braun, M.J., Hackett, S.J., Han, K.-L., Harshman, J., Huddleston, C.J., Kingston, S., Marks, B.D., Miglia, K.J., Moore, W.S., Sheldon, F.H., Witt, C.C., Yuri, T., Braun, E.L., 2017. Why do phylogenomic data sets yield conflicting trees? Data type influences the avian Tree of Life more than taxon sampling. Syst. Biol. 66, 857–879.

Revell, L.J., 2012. Phytools: an R package for phylogenetic comparative biology (and other things). Methods Ecol. Evol. 3, 217–223.

Richards, S., Gibbs, R.A., Weinstock, G.M., Brown, S.J., Denell, R., Beeman, R.W., Gibbs, R., Bucher, G., Friedrich, M., Grimmelikhuijzen, C.J.P., Klingler, M., Lorenzen, M., Roth, S., Schröder, R., Tautz, D., Zdobnov, E.M., Muzny, D., Attaway, T., Bell, S., Buhay, C.J., Chandrabose, M.N., Chavez, D., Clerk-Blankenburg, K.P., Cree, A., Dao, M., Davis, C., Chacko, J., Dinh, H., Dugan-Rocha, S., Fowler, G., Garner, T.T., Garnes, J., Gnirke, A., Hawes, A., Hernandez, J., Hines, S., Holder, M., Hume, J., Jhangiani, S.N., Joshi, V., Khan, Z.M., Jackson, L., Kovar, C., Kowis, A., Lee, S., Lewis, L.R., Margolis, J., Morgan, M., Nazareth, L. V, Nguyen, N., Okwuonu, G., Parker, D., Ruiz, S.-J., Santibanez, J., Savard, J., Scherer, S.E., Schneider, B., Sodergren, E., Vattahil, S., Villasana, D., White, C.S., Wright, R., Park, Y., Lord, J., Oppert, B., Brown, S., Wang, L., Weinstock, G., Liu, Y., Worley, K., Elsik, C.G., Reese, J.T., Elhaik, E., Landan, G., Graur, D., Arensburger, P., Atkinson, P., Beidler, J., Demuth, J.P., Drury, D.W., Du,

Y.-Z., Fujiwara, H., Maselli, V., Osanai, M., Robertson, H.M., Tu, Z., Wang, J., Wang, S., Song, H., Zhang, L., Werner, D., Stanke, M., Morgenstern, B., Solovyev, V., Kosarev, P., Brown, G., Chen, H.-C., Ermolaeva, O., Hlavina, W., Kapustin, Y., Kiryutin, B., Kitts, P., Maglott, D., Pruitt, K., Sapojnikov, V., Souvorov, A., Mackey, A.J., Waterhouse, R.M., Wyder, S., Kriventseva, E. V, Kadowaki, T., Bork, P., Aranda, M., Bao, R., Beermann, A., Berns, N., Bolognesi, R., Bonneton, F., Bopp, D., Butts, T., Chaumot, A., Denell, R.E., Ferrier, D.E.K., Gordon, C.M., Jindra, M., Lan, Q., Lattorff, H.M.G., Laudet, V., von Levetsow, C., Liu, Z., Lutz, R., Lynch, J.A., da Fonseca, R.N., Posnien, N., Reuter, R., Schinko, J.B., Schmitt, C., Schoppmeier, M., Shippy, T.D., Simonnet, F., Marques-Souza, H., Tomoyasu, Y., Trauner, J., Van der Zee, M., Vervoort, M., Wittkopp, N., Wimmer, E.A., Yang, X., Jones, A.K., Sattelle, D.B., Ebert, P.R., Nelson, D., Scott, J.G., Muthukrishnan, S., Kramer, K.J., Arakane, Y., Zhu, Q., Hogenkamp, D., Dixit, R., Jiang, H., Zou, Z., Marshall, J., Elpidina, E., Vinokurov, K., Oppert, C., Evans, J., Lu, Z., Zhao, P., Sumathipala, N., Altincicek, B., Vilcinskas, A., Williams, M., Hultmark, D., Hetru, C., Hauser, F., Cazzamali, G., Williamson, M., Li, B., Tanaka, Y., Predel, R., Neupert, S., Schachtner, J., Verleyen, P., Raible, F., Walden, K.K.O., Angeli, S., Forêt, S., Schuetz, S., Maleszka, R., Miller, S.C., Grossmann, D., 2008. The genome of the model beetle and pest *Tribolium castaneum*. Nature 452, 949–955.

Rivera-Gasperín, S., Ardila-Camacho, A., Contreras-Ramos, A., 2019. Bionomics and ecological services of Megaloptera larvae (dobsonflies, fishflies, alderflies). Insects 10, 86.

Rodríguez-Ezpeleta, N., Brinkmann, H., Roure, B., Lartillot, N., Lang, B.F., Philippe, H., 2007. Detecting and overcoming systematic errors in genome-scale phylogenies. Syst. Biol. 56, 389–399.

Salichos, L., Rokas, A., 2013. Inferring ancient divergences requires genes with strong phylogenetic signals. Nature 497, 327–331.

Salichos, L., Stamatakis, A., Rokas, A., 2014. Novel information theory-based measures for quantifying incongruence among phylogenetic trees. Mol. Biol. Evol. 31, 1261–1271.

Sayyari, E., Mirarab, S., 2016. Fast coalescent-based computation of local branch support from quartet frequencies. Mol. Biol. Evol. 33, 1654–1668.

Sayyari, E., Whitfield, J.B., Mirarab, S., 2017. Fragmentary gene sequences negatively impact gene tree and species tree reconstruction. Mol. Biol. Evol. 34, 3279–3291.

Schremmer, F., 1983. Beitrag zur Entwicklungsgeschichte und zum Kokonbau von *Mantispa styriaca*. Zeitschrift der Arbeitsgemeinschaft Österreichischer Entomol. 35, 21–26.

Schwentner, M., Combosch, D.J., Pakes Nelson, J., Giribet, G., 2017. A phylogenomic solution to the origin of insects by resolving crustacean-hexapod relationships. Curr. Biol. 26, 1569–1571.

Seo, T.K., 2008. Calculating bootstrap probabilities of phylogeny using multilocus sequence data. Mol. Biol. Evol. 25, 960–971.

Sharma, P.P., Kaluziak, S.T., Pérez-Porro, A.R., González, V.L., Hormiga, G., Wheeler, W.C., Giribet, G., 2014. Phylogenomic interrogation of Arachnida reveals systemic conflicts in phylogenetic signal. Mol. Biol. Evol. 31, 2963–2984.

Simmons, M.P., 2012. Misleading results of likelihood-based phylogenetic analyses in the presence of missing data. Cladistics 28, 208–222.

Simmons, M.P., Gatesy, J., 2015. Coalescence vs. concatenation: sophisticated analyses vs. first principles applied to rooting the angiosperms. Mol. Phylogenet. Evol. 91, 98–122.

Simmons, M.P., Norton, A.P., 2014. Divergent maximum-likelihood-branch-support values for polytomies. Mol. Phylogenet. Evol. 73, 87–96.

Simmons, M.P., Pickett, K.M., Miya, M., 2004. How meaningful are Bayesian support values? Mol. Biol. Evol. 21, 188–199.

Smith, S.A., Moore, M.J., Brown, J.W., Yang, Y., 2015. Analysis of phylogenomic datasets reveals conflict, concordance, and gene duplications with examples from animals and plants. BMC Evol. Biol. 15, 150.

Song, N., Li, X., Zhai, Q., Bozdoğan, H., Yin, X.-M., 2019. The mitochondrial genomes of neuropteridan insects and implications for the phylogeny of Neuroptera. Genes 10, 108.

Song, S., Liu, L., Edwards, S. V., Wu, S., 2012. Resolving conflict in eutherian mammal phylogeny using phylogenomics and the multispecies coalescent model. Proc. Natl. Acad. Sci. 109, 14942–14947.

Springer, M.S., Gatesy, J., 2016. The gene tree delusion. Mol. Phylogenet. Evol. 94, 1–
33.

Stamatakis, A., 2014. RAxML version 8: a tool for phylogenetic analysis and post-
analysis of large phylogenies. Bioinformatics 30, 1312–1313.

Strimmer, K., von Haeseler, A., 1997. Likelihood-mapping: a simple method to
visualize phylogenetic content of a sequence alignment. Proc. Natl. Acad. Sci. U.
S. A. 94, 6815–6819.

Tauber, C.A., Tauber, M.J., 1968. *Lomamyia latipennis* (Neuroptera, Berothidae) life
history and larval descriptions. Can. Entomol. 100, 623–629.

Timmermans, M.J.T.N., Barton, C., Haran, J., Ahrens, D., Culverwell, C.L., Ollikainen,
A., Dodsworth, S., Foster, P.G., Bocak, L., Vogler, A.P., 2016. Family-level
sampling of mitochondrial genomes in Coleoptera: compositional heterogeneity
and phylogenetics. Genome Biol. Evol. 8, 161–175.

Tong, K.J., Duchêne, S., Ho, S.Y.W., Lo, N., 2015. Comment on "Phylogenomics
resolves the timing and pattern of insect evolution". Science. 349, 487.

Tonini, J., Moore, A., Stern, D., Shcheglovitova, M., Ortí, G., 2015. Concatenation and
species tree methods exhibit statistically indistinguishable accuracy under a range
of simulated conditions. PLOS Curr. Tree Life 7:
ecurrents.tol.34260cc27551a527b124ec5f6334b6be.

Vasilikopoulos, A., Balke, M., Beutel, R.G., Donath, A., Podsiadlowski, L., Pflug, J.M.,
Waterhouse, R.M., Meusemann, K., Peters, R.S., Escalona, H.E., Mayer, C., Liu,
S., Hendrich, L., Alarie, Y., Bilton, D.T., Jia, F., Zhou, X., Maddison, D.R.,
Niehuis, O., Misof, B., 2019. Phylogenomics of the superfamily Dytiscoidea

(Coleoptera: Adephaga) with an evaluation of phylogenetic conflict and systematic error. Mol. Phylogenet. Evol. 135, 270–285.

von Reumont, B.M., Jenner, R.A., Wills, M.A., Dell'Ampio, E., Pass, G., Ebersberger, I., Meyer, B., Koenemann, S., Iliffe, T.M., Stamatakis, A., Niehuis, O., Meusemann, K., Misof, B., 2012. Pancrustacean phylogeny in the light of new phylogenomic data: support for Remipedia as the possible sister group of Hexapoda. Mol. Biol. Evol. 29, 1031–1045.

Wägele, J.W., Letsch, H., Klussmann-Kolb, A., Mayer, C., Misof, B., Wägele, H., 2009. Phylogenetic support values are not necessarily informative: the case of the Serialia hypothesis (a mollusk phylogeny). Front. Zool. 6, 12.

Wang, H.-C., Minh, B.Q., Susko, E., Roger, A.J., 2018. Modeling site heterogeneity with posterior mean site frequency profiles accelerates accurate phylogenomic estimation. Syst. Biol. 67, 216–235.

Wang, H.-C., Susko, E., Roger, A.J., 2019. The relative importance of modeling site pattern heterogeneity versus partition-wise heterotachy in phylogenomic inference. Syst. Biol. 68, 1003–1019.

Wang, Y., Liu, X., Garzón-Orduña, I.J., Winterton, S.L., Yan, Y., Aspöck, U., Aspöck, H., Yang, D., 2017. Mitochondrial phylogenomics illuminates the evolutionary history of Neuropterida. Cladistics 33, 617–636.

Wang, Y., Zhou, X., Wang, L., Liu, X., Yang, D., Rokas, A., 2019. Gene selection and evolutionary modeling affect phylogenomic inference of Neuropterida based on transcriptome data. Int. J. Mol. Sci. 20, 1072.

Waterhouse, R.M., Tegenfeldt, F., Li, J., Zdobnov, E.M., Kriventseva, E. V., 2013. OrthoDB: a hierarchical catalog of animal, fungal and bacterial orthologs. Nucleic Acids Res. 41, D358–D365.

Wheeler, W.C., Hayashi, C.Y., 2001. The phylogeny of the extant hexapod orders. Cladistics 17, 173–192.

Whiting, M.F., Carpenter, J.C., Wheeler, Q.D., Wheeler, W.C., 1997. The Strepsiptera problem: phylogeny of the holometabolous insect orders inferred from 18S and 28S ribosomal DNA sequences and morphology. Syst. Biol. 46, 1–68.

Wiegmann, B.M., Trautwein, M.D., Kim, J.-W., Cassel, B.K., Bertone, M. A., Winterton, S.L., Yeates, D.K., 2009. Single-copy nuclear genes resolve the phylogeny of the holometabolous insects. BMC Biol. 7, 34.

Winterton, S.L., Hardy, N.B., Wiegmann, B.M., 2010. On wings of lace: phylogeny and Bayesian divergence time estimates of Neuropterida (Insecta) based on morphological and molecular data. Syst. Entomol. 35, 349–378.

Winterton, S.L., Lemmon, A.R., Gillung, J.P., Garzon, I.J., Badano, D., Bakkes, D.K., Breitkreuz, L.C.V., Engel, M.S., Lemmon, E.M., Liu, X., Machado, R.J.P., Skevington, J.H., Oswald, J.D., 2018. Evolution of lacewings and allied orders using anchored phylogenomics (Neuroptera, Megaloptera, Raphidioptera). Syst. Entomol. 43, 330–354.

Withycombe, C.L., 1925. XV. Some Aspects of the Biology and Morphology of the Neuroptera. With special reference to the immature stages and their possible phylogenetic significance. Trans. R. Entomol. Soc. London 72, 303–411.

Wolfe, J.M., Daley, A.C., Legg, D.A., Edgecombe, G.D., 2016. Fossil calibrations for the arthropod Tree of Life. Earth-Science Rev. 160, 43–110.

Wong, T.K.F., Kalyaanamoorthy, S. Meusemann, K., Yeates, D.K., Misof, B., Jermiin, L.S., 2020. A minimum reporting standard for multiple sequence alignments. NAR Genomics Bioinforma. 2. doi: 10.1093/nargab/lqaa024

Xia, Q., Zhou, Z., Lu, C., Cheng, D., Dai, F., Li, B., Zhao, P., Zha, X., Cheng, T., Chai, C., Pan, G., Xu, Jinshan, Liu, C., Lin, Y., Qian, J., Hou, Y., Wu, Z., Li, Guanrong, Pan, M., Li, Chunfeng, Shen, Y., Lan, X., Yuan, L., Li, T., Xu, H., Yang, G., Wan, Y., Zhu, Y., Yu, M., Shen, W., Wu, D., Xiang, Z., Yu, J., Wang, Jun, Li, R., Shi, J., Li, H., Li, Guangyuan, Su, J., Wang, X., Li, Guoqing, Zhang, Zengjin, Wu, Q., Li, Jun, Zhang, Q., Wei, N., Xu, Jianzhe, Sun, H., Dong, L., Liu, D., Zhao, S., Zhao, X., Meng, Q., Lan, F., Huang, X., Li, Y., Fang, L., Li, Changfeng, Li, D., Sun, Y., Zhang, Zhenpeng, Yang, Z., Huang, Y., Xi, Y., Qi, Q., He, D., Huang, H., Zhang, X., Wang, Z., Li, W., Cao, Y., Yu, Y., Yu, H., Li, Jinhong, Ye, Jiehua, Chen, H., Zhou, Y., Liu, B., Wang, Jing, Ye, Jia, Ji, H., Li, Shengting, Ni, P., Zhang, J., Zhang, Y., Zheng, H., Mao, B., Wang, W., Ye, C., Li, Songgang, Wang, Jian, Wong, G.K.-S., Yang, H., 2004. A draft sequence for the genome of the domesticated silkworm (*Bombyx mori*). Science. 306, 1937–1940.

Xu, B., Yang, Z., 2016. Challenges in species tree estimation under the multispecies coalescent model. Genetics 204, 1353–1368.

Yang, Q., Makarkin, V.N., Winterton, S.L., Khramov, A. V., Ren, D., 2012. A remarkable new family of Jurassic insects (Neuroptera) with primitive wing venation and its phylogenetic position in Neuropterida. PLoS One 7, e44762.

Yang, Z., 2007. PAML 4: Phylogenetic analysis by maximum likelihood. Mol. Biol. Evol. 24, 1586–1591.

Zhang, C., Rabiee, M., Sayyari, E., Mirarab, S., 2018. ASTRAL-III: polynomial time species tree reconstruction from partially resolved gene trees. BMC Bioinformatics 19, 153.

Zhao, C., Liu, X., Yang, D., 2014. Wing base structural data support the sister relationship of Megaloptera and Neuroptera (Insecta: Neuropterida). PLoS One 9, e114695.

Zhao, J., Li, H., Winterton, S.L., Liu, Z., 2013. Ancestral gene organization in the mitochondrial genome of *Thyridosmylus langii* (McLachlan, 1870) (Neuroptera: Osmylidae) and implications for lacewing evolution. PLoS One 8, e62943.

Zhong, M., Hansen, B., Nesnidal, M., Golombek, A., Halanych, K.M., Struck, T.H., 2011. Detecting the symplesiomorphy trap: a multigene phylogenetic analysis of terebelliform annelids. BMC Evol. Biol. 11, 369.

Zhou, X., Lutteropp, S., Czech, L., Stamatakis, A., von Looz, M., Rokas, A., 2020. Quartet-based computations of internode certainty provide robust measures of phylogenetic incongruence. Syst. Biol. 69, 308–324.

Zwick, P., 1967. Beschreibung der aquatischen Larve von *Neurorthus fallax* (Rambur) und Errichtung der neuen Planipennierfamilie Neurorthidae fam. nov. Gewässer und Abwässer 44/45:, 65–86.

# 3. Phylogenomics of the superfamily Dytiscoidea (Coleoptera: Adephaga) with an evaluation of phylogenetic conflict and systematic error

## 3.1. Introduction

Almost half of the ca. 13,000 beetle species with an aquatic lifestyle (Jäch and Balke, 2008) belong to the suborder Adephaga, which also contains more than 38,000 species of the terrestrial Carabidae and Trachypachidae. The aquatic (or semi-aquatic) adephagan families Amphizoidae, Dytiscidae, Gyrinidae, Haliplidae, Hygrobiidae, and Noteridae have traditionally been considered as monophyletic and collectively referred to as "Hydradephaga" (Crowson, 1960). The monophyly of "Hydradephaga" has not been corroborated in extensive phylogenetic analyses of morphological data or in recent phylogenomic investigations (e.g., Baca et al., 2017; Beutel, 1993; Beutel et al., 2008, 2006; Beutel and Haas, 1996; Beutel and Roughley, 1988; Dressler et al., 2011; Dressler and Beutel, 2010; S.Q. Zhang et al., 2018; but see López-López and Vogler, 2017). On the other hand, the monophyly of the superfamily Dytiscoidea (Amphizoidae, Aspidytidae, Dytiscidae, Hygrobiidae, Meruidae, and Noteridae) is well established (e.g., Baca et al., 2017; Beutel et al., 2013; Dressler et al., 2011; but see López-López and Vogler, 2017). Species of this superfamily can be encountered in virtually every kind of freshwater habitat, including springs, rivers, acidic swamps, lakes, and even in hypersaline or hygropetric habitats. Their widespread occurrence is primarily due to the astounding ecological versatility of species in the family Dytiscidae (Miller and Bergsten, 2016). Interestingly, the phylogenetic relationships within Dytiscoidea are still obscure, especially concerning the hypothesized monophyly of Aspidytidae and the phylogenetic affinities of its species to those of the families Amphizoidae and Hygrobiidae. In the present phylogenomic study, we investigate the above-outlined phylogenetic questions with the largest molecular

dataset compiled to date for studying phylogenetic relationships in this group of beetles.

Most species of Dytiscoidea are strictly aquatic, but two families with species inhabiting hygropetric habitats have recently been described. The species of these families occur in geographically disjunct regions. Meruidae, with the single species *Meru phyllisae* Spangler and Steiner, 2005, is known only from the Guiana Shield region of Venezuela (Spangler and Steiner, 2005). Aspidytidae contain two species, *Sinaspidytes wrasei* (Balke et al., 2003) from China (Balke et al., 2003; Toussaint et al., 2015) and *Aspidytes niobe* Ribera, Beutel, Balke, Vogler, 2002 from the Cape region of South Africa (Beutel et al., 2010; Ribera et al., 2002a). Phylogenetic analyses have placed these two families in the superfamily Dytiscoidea (Beutel et al., 2006; Ribera et al., 2002a), along with the Dytiscidae (diving beetles, 4489 species; Nilsson and Hájek, 2019), Noteridae (burrowing water beetles, 258 species; Nilsson, 2011), Hygrobiidae (squeak beetles, six species) and Amphizoidae (trout stream beetles, five species). The taxonomy of Dytiscoidea has been extensively studied, as have been its morphological and ecological adaptations (Balke and Hendrich, 2016; Miller and Bergsten, 2016) and the anatomy of adults and larvae (Belkaceme, 1991; Beutel, 1986a, 1986b, 1988, 1993; Dressler and Beutel, 2010). Moreover, species of the group are well documented in the fossil record and can be traced back to the Triassic (e.g. Beutel et al., 2013; Ponomarenko, 1993).

The phylogenetic relationships of dytiscoid beetles have been addressed in numerous studies investigating morphology, chemical gland compounds, fossil data, and DNA sequences (Alarie et al., 2011, 2004; Alarie and Bilton, 2005; Baca et al., 2017; Balke et al., 2008, 2005; Beutel et al., 2006, 2008, 2013; Beutel, 1993; Beutel

and Haas, 1996; Burmeister, 1976; Dettner, 1985; Kavanaugh, 1986; López-López and Vogler, 2017; McKenna et al., 2015; Ribera et al., 2002b; Toussaint et al., 2015). Analyses of these different data have not yielded congruent topologies (see Fig. 3.1 for selected hypotheses). The currently accepted view is that Meruidae + Noteridae represent the sister clade of the remaining four families of the superfamily Dytiscoidea (Fig. 3.1). However, the affinities of Amphizoidae, Aspidytidae, Dytiscidae, and Hygrobiidae remain unresolved. A clade consisting of Dytiscidae and Hygrobiidae is supported by some morphological features (Balke et al., 2005; Beutel et al., 2006; Dressler and Beutel, 2010), such as the presence of prothoracic glands (Beutel, 1986b, 1988; Forsyth, 1970) but molecular and total evidence analyses have yielded incongruent topologies (e.g. Baca et al., 2017; Balke et al., 2005; Ribera et al., 2002a; Toussaint et al., 2015).

A sister group relationship between Amphizoidae and Aspidytidae has been suggested in previous studies analyzing molecular data (Balke et al., 2005, 2008; Hawlitschek et al., 2012; Toussaint et al., 2015), but Toussaint et al. (2015) recovered paraphyletic Aspidytidae (in relation to Amphizoidae). Specifically, in a multigene analysis of nucleotide sequence data, and after excluding the highly saturated third codon positions, *A. niobe* was placed as the sister taxon of Amphizoidae (Fig. 3.1f). This new hypothesis contributed to the existing confusion on character evolution within Dytiscoidea (Balke et al., 2005; Beutel et al., 2006; Ribera et al., 2002a), because morphological characters of the adult beetles (antenna: configuration of scape and pedicel) suggest a monophyletic Aspidytidae, while morphological characters of the larvae of *S. wrasei* show considerable structural affinities with those of Amphizoidae (Toussaint et al., 2015).

Given the above-outlined uncertainties in the phylogenetic relationships of the families currently included in Dytiscoidea we (1) investigated whether Aspidytidae are monophyletic and (2) inferred the phylogenetic relationships among the families Amphizoidae, Aspidytidae, Dytiscidae, Hygrobiidae, and Noteridae based on an extensive transcriptomic dataset. In order to achieve these goals, we analyzed whole body transcriptomes of species of all major lineages of Dytiscoidea except Meruidae. We also investigated the effects of different potential sources of conflicting phylogenetic signal and phylogenomic incongruence when estimating phylogenetic relationships within Dytiscoidea, and evaluated the degree of confidence for alternative topologies using branch support tests and a data permutation approach.



**Fig. 3.1**: Overview of different phylogenetic hypotheses on family phylogenetic relationships among Dytiscoidea proposed in previous studies that had analyzed molecular and morphological data. (Note that Meruidae were not included in all studies. However, since their sister group relationship to Noteridae is generally considered undisputed, we consistently included them in the overview: "Meruidae+Noteridae"). (**a**) Balke et al. (2005) based on morphological data, (**b**) Baca et al. (2017) based on UCE data, (**c**) Beutel et al. (2006, 2013)

based on morphological data, (**d**) Ribera et al. (2002a) based on morphological and molecular data, (**e**) Balke et al. (2005, 2008) based on molecular data and Balke et al. (2005) based on morphological and molecular data, (**f**) Toussaint et al. (2015) based on molecular data and McKenna et al. (2015) based on molecular data with only *Aspidytes* included.

## 3.2. Materials and methods

### 3.2.1. Taxon sampling

We compiled a dataset consisting of *de novo*-sequenced transcriptomes and of previously published transcriptomes of Dytiscoidea (Table 3.1). The sampled species represent all extant families of Dytiscoidea except Meruidae (for which transcriptomic data were not available). As there is high confidence in the hypothesized sister group relationship between Meruidae and Noteridae (Baca et al., 2017; Balke et al., 2008; Beutel et al., 2006; Dressler et al., 2011; Toussaint et al., 2015), we do not deem the lack of the species *M. phyllisae* from our dataset as problematic for investigating the major relationships of Dytiscoidea (see Fig. 3.1). Representatives of Gyrinidae and Haliplidae were included as outgroups (Baca et al., 2017; Beutel et al., 2006, 2013; Beutel and Haas, 1996; Beutel and Roughley, 1988; Dressler et al., 2011; Dressler and Beutel, 2010).

The *de novo*-sequenced and assembled transcriptomes were screened for putative adaptor, vector and cross-contaminated sequences (see S1: Suppl. Text 1), and clean assemblies were subsequently submitted to the NCBI-TSA database (Table 3.1). For a detailed description of the procedures for specimen collection and preservation, RNA isolation, RNA library preparation, transcriptome sequencing, transcriptome assembly, cross-contamination screening and sequence submissions see the S1: Suppl.

Text 1. We used custom made Perl and Python scripts to calculate descriptive statistics for each transcriptome in our study (Table 3.1).

### 3.2.2. Orthology assignment and alignment refinement

We identified 3085 clusters of single-copy genes (COGs) that are non-homologous or out-paralogous among each other at the hierarchical level Endopterygota, based on a customized profile query in OrthoDB v. 9.1 (Zdobnov et al., 2017) (see S1: Suppl. Text 1). Our query was based on six endopterygote species (subsequently referred to as reference species) with well sequenced and annotated genomes (S2A: Suppl. Table 1). Each transcriptome was searched for transcripts orthologous to the sequences of a given COG (see Peters et al., 2017; Petersen et al., 2017). This search was performed with Orthograph v. 0.6.1 (Petersen et al., 2017). Orthologous sequences for each COG (including those of the reference species) were combined in two FASTA files: one containing sequences at the transcriptional level (i.e., nucleotides, nCOGs), the other containing sequences at the translational level (i.e., amino acids, aaCOGs). The resulting nCOGs and aaCOGs are deposited at MENDELEY DATA (see list of Supplementary materials).

Alignment of the amino-acid sequences in each aaCOG, was performed with MAFFT v. 7.309 (Katoh and Standley, 2013) using the algorithm L-INS-i. We screened the amino-acid multiple sequence alignments (MSAs) for potentially misaligned sequences and erroneously identified orthologs using the procedure outlined by Misof et al. (2014). We also adapted the alignment refinement procedure proposed by Misof et al. (2014). Amino-acid and nucleotide sequences that were still

identified as outliers after the alignment refinement procedure were removed from the MSAs.

Following the alignment refinement procedure, we removed all sequences of the reference species from the aligned aaCOGs and also discarded their corresponding nucleotide sequences. This resulted in FASTA files that comprised exclusively (aligned) amino-acid or (unaligned) nucleotide sequences of Dytiscoidea and of the outgroup families Gyrinidae and Haliplidae. Next, we discarded all COGs from the ortholog set containing transcripts from fewer than three species. After removing gap-only and ambiguous-only positions from the remaining 2,991 aaCOGs we generated codon-based nucleotide sequence alignments, with a modified version of the script Pal2nal.pl (Suyama et al., 2006) as described by Misof et al. (2014). The 2,991 aligned aaCOGs and the corresponding codon-based alignments are deposited at MENDELEY DATA (see list of Supplementary materials).

### 3.2.3. Concatenation-based and gene tree-based phylogenetic analyses of amino-acid sequence data

We generated eleven amino-acid supermatrices (Table 3.2, S3: Suppl. Fig. 1) and assessed the effects of different putative sources of topological incongruence on our concatenation-based phylogenetic inference, namely: (1) alignment masking (i.e., alignment column-filtering) of individual gene partitions when analyzed in a supermatrix context, (2) effects of data coverage and phylogenetic information content on the dytiscoid phylogenetic relationships, (3) taxonomic decisiveness of gene partitions with respect to a specific phylogenetic question, and (4) effects of compositionally heterogeneous genes in a supermatrix context. We modified the initial

supermatrix (supermatrix A, Table 3.2) by masking the effects of each of the above-mentioned factors one by one (e.g. by removing the randomly similar sections in each gene or removing partitions with low information content). This hierarchical masking strategy progressively resulted in supermatrices to be analyzed with fewer genes and fewer amino-acid alignment sites. We used each generated dataset (Table 3.2, S3: Suppl. Fig. 1) to infer the phylogeny of Dytiscoidea. The purpose of these analyses was to assess whether or not gradual masking of the initial supermatrix for any of the above factors affected the results of the phylogenetic inference. Amino-acid supermatrices A–K are deposited at MENDELEY DATA (see list of Supplementary materials).

*Masking of the individual amino-acid MSAs*

It has been suggested that current methods of alignment masking may lead to biased phylogenetic inferences because alignment columns are filtered too aggressively (Tan et al., 2015). To assess the effect of alignment masking on our results, we first concatenated the original MSAs of aaCOGs without applying alignment masking (supermatrix A). We then applied ALISCORE v. 1.2 (Kück et al., 2010; Misof and Misof, 2009) on each aaCOG separately with the options: -r $10^{27}$ (for the maximum number of pairwise sequence comparisons) and -e. The masked genes (aaCOGs) were then concatenated in a new masked supermatrix (supermatrix B). Concatenation of both masked and unmasked amino-acid MSAs was conducted with FASconCAT-G v. 1.02 (Kück and Longo, 2014).

**Table 3.1**: An overview of the newly sequenced and previously published transcriptomes that were analyzed in the present study. NCBI accession numbers and descriptive statistics to each transcriptome are provided. Species whose transcriptomes were analyzed are given in alphabetic order.

| Species name/Transcriptome | Family | TSA accesssion | BioSample accesion | Bioproject accession | Reference / Source | No. of contigs | After local Vec-Screen | After contam. check | Contigs published | Mean length | Median length | N50 length | Max. length |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| *Amphizoa insolens* LeConte, 1853 | Amphizoidae | GFUZ01000000 | SAMN07501457 | PRJNA398088 | NCBI-TSA | N/A | N/A | N/A | 23,404 | 1265 | 854 | 1858 | 17,558 |
| *Amphizoa lecontei* Matthews, 1872 | Amphizoidae | GFUH01000000 | SAMN07289768 | PRJNA392306 | this study | 53,433 | 53,331 | 53,298 | 53,272 | 869 | 467 | 1540 | 15,581 |
| *Aspidytes niobe* Ribera, Beutel, Balke, Vogler, 2002 | Aspidytidae | GFUO01000000 | SAMN07279561 | PRJNA391973 | this study | 22,688 | 22,683 | 22,269 | 22,272 | 1173 | 716 | 1996 | 9941 |
| *Batrachomatus nannup* (Watts, 1978) | Dytiscidae | GFUJ01000000 | SAMN07280954 | PRJNA392058 | this study | 43,890 | 43,601 | 43,554 | 43,521 | 741 | 446 | 1151 | 15,127 |
| *Cybister lateralimarginalis* (DeGeer, 1774) | Dytiscidae | GDLH01000000 | SAMN03799556 | PRJNA286512 | 1KITE, this study | 31,471 | 31,470 | 31,403 | 31,402 | 981 | 577 | 1586 | 47,239 |
| *Dineutus* sp. | Gyrinidae | GDNB01000000 | SAMN03799560 | PRJNA286516 | 1KITE, this study | 25,920 | 25,915 | 24,679 | 24,661 | 862 | 600 | 1281 | 11,252 |
| *Gyrinus marinus* Gyllenhal, 1808 | Gyrinidae | GAUY02000000 | SAMN02047132 | PRJNA219564 | 1KITE, Misof et al. (2014) | 23,637 | 23,637 | 23,510 | 23,491 | 866 | 535 | 1426 | 13,197 |
| *Haliplus fluviatilis* Aubé, 1836 | Haliplidae | GDMW01000000 | SAMN03799569 | PRJNA286525 | 1KITE, this study | 46,197 | 46,191 | 45,977 | 45,915 | 847 | 445 | 1504 | 34,051 |
| *Hygrobia hermanni* (Fabricius, 1775) | Hygrobiidae | GFUK01000000 | SAMN07297121 | PRJNA392382 | this study | 62,884 | 62,877 | 62,691 | 62,715 | 923 | 559 | 1430 | 19,834 |
| *Hygrobia nigra* (Clark, 1862) | Hygrobiidae | GFUN01000000 | SAMN07287246 | PRJNA392270 | this study | 28,837 | 28,835 | 28,561 | 28,569 | 918 | 567 | 1492 | 10,964 |
| *Liopterus haemorrhoidalis* (Fabricius, 1787) | Dytiscidae | GFUI01000000 | SAMN07280875 | PRJNA392045 | this study | 66,642 | 66,327 | 66,281 | 66,211 | 604 | 394 | 824 | 8663 |
| *Noterus clavicornis* (DeGeer, 1774) | Noteridae | GDNA01000000 | SAMN03799605 | PRJNA286561 | 1KITE, this study | 21,719 | 21,716 | 21,606 | 21,601 | 1046 | 639 | 1695 | 37,302 |
| *Sinaspidytes wrasei* (Balke, Ribera, Beutel, 2003) | Aspidytidae | GDNH01000000 | SAMN03799537 | PRJNA286492 | 1KITE, this study | 41,855 | 41,748 | 37,769 | 37,371 | 874 | 400 | 1725 | 25,916 |
| *Thermonectus intermedius* Crotch, 1873 | Dytiscidae | N/A | N/A | N/A | Boussau et al. (2014) | N/A | N/A | N/A | 15,833 | 1351 | 867 | 1938 | 38,615 |

*Increasing data coverage and phylogenetic information content*

We evaluated whether or not increasing the saturation (SV, the overall degree of data coverage with respect to gene presence or absence) and the phylogenetic information content (IC) of the supermatrix, as a function of data coverage and phylogenetic signal, had an effect on our tree reconstructions. IC and SV values were calculated with MARE v. 0.1.2-rc (MAtrix REduction) (Misof et al., 2013). We generated and assessed the following amino-acid supermatrices:

1. supermatrix C: selected optimal subset (SOS, default output supermatrix) of the software MARE when using supermatrix B as input;

2. supermatrix D: inferred from supermatrix B after removing those genes with IC=0;

3. supermatrix E: selected optimal subset (SOS, default output supermatrix) of the software MARE when using supermatrix D as input.

We also calculated the SV and the IC of every other amino-acid supermatrix (Table 3.2). In addition, we calculated the overall alignment completeness scores (Ca) of all supermatrices (Tables 3.2 and 3.3) with AliStat v. 1.6 (https://github.com/thomaskf/AliStat, see Misof et al., 2014). The overall completeness score provides a direct measure of the overall degree of missing data in each analyzed supermatrix. Moreover, we generated heatmaps of pairwise completeness scores for every amino-acid and nucleotide sequence supermatrix that we analyzed (S3: Suppl. Fig. 3–23).

**Table 3.2**: Detailed information and statistics of each generated amino-acid supermatrix analyzed in this study. The overall alignment completeness score of each matrix was calculated with the software AliStat. Matrix phylogenetic information content and saturation were calculated with the software MARE. The RCFV value was calculated with BaCoCa. Pairwise tests of symmetry for the Bowker's test were performed with SymTest. ($C_a$: overall alignment completeness score, SV: matrix saturation values, IC: matrix phylogenetic information content).

| Amino-acid matrix ID | No. of taxa | No. of amino-acid sites | No. of gene partitions | $C_a$ | SV | IC | Percentage of pairwise p-values < 0.05 for the Bowker's test | Optimization of partitioning scheme | No. of tree searches with unoptimized partitioning scheme | No. meta-partitions | No. of tree searches with optimized partitioning scheme | No. of bootstraps with unoptimized partitioning scheme | No. of tree searches with the PMSF model | No. of bootstraps with the PMSF CAT-like model | Information |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| A | 14 | 1,661,023 | 2,991 | 0.5976280 | 0.893 | 0.521 | 100.00 % | NO | 10 | - | - | 100 | - |  | Unmasked matrix |
| B | 14 | 1,384,486 | 2,991 | 0.6824300 | 0.891 | 0.523 | 100.00 % | NO | 10 | - | - | 100 | - |  | Masked genes of matrix A with ALISCORE |
| C | 14 | 955,158 | 1,901 | 0.6668550 | 0.921 | 0.650 | 96.70 % | NO | 10 | - | - | 100 | - |  | Default MARE matrix (SOS) of matrix B |
| D | 14 | 1,366,298 | 2,948 | 0.6888650 | 0.898 | 0.530 | 100.00 % | NO | 10 | - | - | 100 | 1 |  | Removed genes with IC=0 from matrix B. |
| E | 14 | 948,772 | 1,884 | 0.6654340 | 0.921 | 0.639 | 95.60 % | YES | 10 | 902 | 10 | 100 | 1 |  | Default MARE matrix (SOS) of matrix D. |
| F | 14 | 468,720 | 900 | 0.7548040 | 1.000 | 0.673 | 90.11 % | NO | 10 | - | - | 100 | - |  | Decisive 1: selected species with all genes from matrix E |
| G | 14 | 806,143 | 1,634 | 0.7016170 | 0.951 | 0.661 | 93.41 % | NO | 10 | - | - | 100 | - |  | Decisive 2: Aspidytidae both present and at least one species for each of the remaining families (filtered matrix E) |
| H | 14 | 211,275 | 416 | 0.8592440 | 1.000 | 0.660 | 73.63 % | YES | 10 | 170 | 10 | 100 | 1 |  | Removed genes with RCFV >= 0.1 from matrix F |
| I | 14 | 218,940 | 1 | 1.0000000 | N/A | N/A | 94.51 % | N/A (unpartitioned) | 10 | - | - | 100 | 1 |  | Selected sites with 100 % species coverage from matrix D |
| J | 14 | 391,961 | 814 | 0.7751530 | 0.927 | 0.639 | 84.62 % | NO | 10 | - | - | 100 | - |  | Removed genes with RCFV >= 0.1 from matrix E |
| K | 14 | 721,765 | 1,344 | 0.6862060 | 0.868 | 0.494 | 95.60 % | NO | 10 | - | - | 100 | - |  | Removed genes with RCFV >= 0.1 from matrix A |

*Controlling for data decisiveness*

We constructed two amino-acid sequence supermatrices to control for data decisiveness following the approach outlined by Dell'Ampio et al. (2014). Data decisiveness refers to the property of a partition to include data of every group of species that is relevant to address a specific phylogenetic question (e.g., the monophyly of Aspidytidae). We generated a subset of supermatrix E by including only those aaCOGs in which all 14 species were present (supermatrix F). An additional decisive dataset (supermatrix G) was constructed by including only those aaCOGs that included at least one representative of Amphizoidae, Dytiscidae, Gyrinidae, Haliplidae, Hygrobiidae, Noteridae, and both representatives of Aspidytidae (*A. niobe + S. wrasei*). These two amino-acid sequence datasets were considered decisive for addressing the inter-familiar relationships of Dytiscoidea and the monophyly of Aspidytidae.

*Controlling for among-species compositional heterogeneity*

Compositional heterogeneity among species in a dataset is often neglected as a source of systematic error in molecular phylogenetic studies (Jermiin et al., 2004; Nesnidal et al., 2010; Philippe and Roure, 2011; Romiguier et al., 2016; Whitfield and Kjer, 2008). We explicitly explored whether among-species compositional heterogeneity biased tree reconstructions. Compositionally heterogeneous aaCOGs were excluded from the decisive amino-acid dataset (supermatrix F) to generate a decisive and more compositionally homogeneous matrix (supermatrix H, S3: Suppl. Fig. 1). Among-species compositional heterogeneity was assessed for each partition separately, based on the partition-specific relative composition frequency variation

value (RCFV) (Zhong et al., 2011) calculated by BaCoCa v. 1.105 (Kück and Struck, 2014). We followed Fernandez et al. (2016) by considering compositional heterogeneity among species in a given aaCOG to be high when the overall RCFV value was greater than or equal to 0.1. We also filtered supermatrix A and supermatrix E using the same threshold (Table 3.3, supermatrices J and K) and compared results of tree reconstructions. Complementary to the RCFV approach, we used the software SymTest v. 2.0.47 (https://github.com/ottmi/symtest) to calculate the overall deviation from stationarity, reversibility, and homogeneity (SRH) (Jermiin et al., 2008) between the amino-acid (or nucleotide) sequences of the species in each generated supermatrix (see Misof et al., 2014 and S1: Suppl. Text 1). We generated heatmaps to visualize the pairwise deviations from SRH conditions in each generated supermatrix in our study (S1: Suppl. Text 1, S3: Suppl. Fig. 24–44).

*Maximum likelihood phylogenetic analyses of amino-acid sequence data*

For each of the amino-acid sequence supermatrices (A–K) ten independent partitioned tree searches were performed using IQ-TREE v. 1.5.5 (or later) (Nguyen et al., 2015) by specifying the aligned aaCOG boundaries. Model selection for each aaCOG was performed with ModelFinder (Kalyaanamoorthy et al., 2017), implemented in IQ-TREE. We considered the following amino-acid substitution models: DAYHOFF (Dayhoff et al., 1978), DCMUT (Kosiol and Goldman, 2005), JTT (Jones et al., 1992), JTTDCMUT (Kosiol and Goldman, 2005), LG (Le and Gascuel, 2008), LG4X (Le et al., 2012), and WAG (Whelan and Goldman, 2001) allowing all possible combinations of modeling rate heterogeneity among sites (options: -mrate E,I,G,I+G,R -gmedian -merit AICc). We used the edge-linked partitioned model for

tree reconstruction (option: -spp) allowing each gene to have its own rate but assuming a common topology and proportional branch lengths among all gene partitions (Chernomor et al., 2016). For each supermatrix the most appropriate model for each gene partition was selected during the first tree search (option -m MFP). The resulting NEXUS files of the first run were used as input for all remaining tree searches.

A common practice in phylogenomic analyses is to optimize the partitioning schemes and corresponding substitution models for the data within an algorithmic framework (Lanfear et al., 2012, 2014). Such optimizations of the partitioning schemes are time-consuming and could result in combining different genes in different meta-partition analyses due to the heuristic optimization procedures implemented in the existing software (Lanfear et al., 2014). This can lead to very different model assignments for different genes and therefore would add an additional uncontrollable effect when comparing different supermatrices. By defining the original masked gene boundaries for all supermatrices and by not optimizing the partitioning schemes we excluded the effects of differential model fit (due to the different composition of the inferred meta-partitions in each matrix) on the results of tree reconstructions. However, in order to avoid missing a unique topology of Dytiscoidea due to suboptimal model fit we optimized the partitioning scheme for a selection of amino-acid supermatrices. We selected the supermatrices H and E for this purpose, because they gave rise to different topologies when analyzing amino-acid sequence data. We used the relaxed clustering algorithm (rcluster) (Lanfear et al., 2014) and RaxML v. 8.2 (options: -raxml -rcluster-max 5000) (Stamatakis, 2014) in PartitionFinder v. 2.1.1 (Lanfear et al., 2017) to merge partitions according to the default weights under the AICc information criterion. We restricted the model search in PartitionFinder to the following amino-acid

substitution models: DAYHOFF+G, DAYHOFF+G+F, DCMUT+G, DCMUT+G+F, JTT+G, JTT+G+F, LG+G, LG+G+F, LG4X, WAG+G, and WAG+G+F. The inferred schemes and models for the corresponding meta-partitions were defined as input for the IQ-TREE tree searches (v. 1.5.5) again with the edge-linked model. Ten independent tree searches were performed with the optimized partitioning schemes of supermatrix E and H. The resulting NEXUS files with the optimized schemes of supermatrix E and of supermatrix H are deposited at MENDELEY DATA (see list of Supplementary materials). Statistical support of our inferred relationships was assessed based on the non-parametric bootstrap measure (Felsenstein, 1985) and the bootstrap by transfer (TBE) support measure (Lemoine et al., 2018). We calculated 100 non-parametric bootstrap replicates and TBE support using the unoptimized partitioning schemes of all the analyzed amino-acid datasets (Table 3.2). In addition, we calculated 100 nonparametric bootstrap replicates and TBE support for the optimized partitioning schemes of supermatrices E and H. Subsequently, we mapped the bootstrap support values on the maximum likelihood trees (i.e., trees with the best log-likelihood among all ten tree searches).

For the optimized partitioning schemes of the supermatrices E and supermatrix H we also performed one additional tree search with the options -bb 1000 -alrt 10000 -abayes to estimate different measures of branch support implemented in IQ-TREE v. 1.5.5: Ultrafast Bootstrap 1 (UFBoot1), SH-like aLRT, and aBayes respectively (Anisimova et al., 2011; Guindon et al., 2010; Minh et al., 2013). We also separately calculated branch support based on the updated version of Ultrafast Bootstrap in IQ-TREE v. 1.6.8 (UFBoot2, option: -bnni) with 1000 replicates (Hoang et al., 2017).

After verifying topological congruence to the maximum likelihood tree, we mapped the different branch support values on the maximum likelihood tree (Fig. 3.2).

For a selection of amino-acid supermatrices, we performed one additional tree search using IQ-TREE v. 1.5.5 (or later) by implementing the posterior-mean-site-frequency (PMSF) model (Wang et al., 2017), as a rapid approximation of the site-heterogeneous CAT-like mixture model (Quang et al., 2008) with 60 amino-acid profile categories and the exchange rates of the LG substitution matrix (option: -m LG+C60+G+F). We used the tree with the best log-likelihood that resulted from the analysis based on the partition model as a guide tree. The idea of applying this mixture model was to increase the biological realism of the modeled substitution processes, as it should be able to describe site-specific amino-acid preferences in the supermatrices. Moreover, proponents of the site-heterogeneous mixture models have recommended their use to alleviate systematic errors due to model violations (Lartillot et al., 2007). We calculated the non-parametric bootstrap measure (BS PMSF, Fig. 3.2a and 3.2b) when applying the PMSF model (LG+C60+G+F) with 100 replicates (Table 3.2).

*Coalescent-based phylogenetic analyses*

The supermatrix approach has been criticized for producing statistically inconsistent topologies as it fails to account for gene tree heterogeneity due to incomplete lineage sorting (ILS) (Kubatko and Degnan, 2007). However, research has shown that concatenation (even unpartitioned) can be more accurate than summary species tree methods under certain conditions (Bayzid and Warnow, 2013; Mirarab et al., 2016; Mirarab and Warnow, 2015; Xu and Yang, 2016) and that summary species tree methods can be sensitive to gene tree estimation errors or to low degree of

variation in the analyzed sets of loci (Bayzid and Warnow, 2013; Meiklejohn et al., 2016). In an attempt to explore the sensitivity of our phylogenetic results to the above-mentioned potentially biasing factors, we conducted coalescent species tree analyses with ASTRAL III v. 5.5.12 (Mirarab and Warnow, 2015; C. Zhang et al., 2018) as an alternative to the supermatrix approach. We expected that if both methods yield the same topologies for the datasets analyzed, any observed topological differences (between analyzed datasets) would unlikely be due to ILS, hybridization or due to biases resulting from gene tree estimation errors. We performed the coalescent approach on (1) a selected subset of COGs from supermatrix E and (2) the full set of COGs from supermatrix H. When analyzing supermatrix E, we discarded all COGs with fewer than 13 species and more than 20% ambiguous characters (X, -) to increase data coverage of the selected genes (Sayyari et al., 2017). When analyzing supermatrix H, we selected the full set of COGs to perform the species tree analysis, as this dataset had already a low proportion of missing data (Table 3.2, S3: Suppl. Fig. 10). Individual gene trees were constructed under the maximum likelihood optimality criterion in IQ-TREE v. 1.5.5. Model selection for each aaCOG was restricted to the amino-acid substitution matrices DCMUT, JTT, LG, and WAG under the AICc information criterion. We allowed a maximum of four free rate categories for modeling rate heterogeneity among sites in ModelFinder (option: -cmax 4). We calculated the branch lengths of the estimated species tree in coalescence units in ASTRAL with the option -q. We annotated the species tree with the option -t 2. This resulted in a tree labeled with quartet scores, total quartet support and local posterior probabilities (Sayyari and Mirarab, 2016). Quartet support values (q1, q2, q3) indicate the proportion of induced quartets in the gene trees that agree or disagree with a branch on the calculated species

tree. Each alternative value corresponds to the three possible topologies around each branch of interest. The local posterior probabilities are calculated based on the quartet support values (Sayyari and Mirarab, 2016). The first quartet support and local posterior probability for each branch (q1 and pp1 respectively) correspond to the topology that is depicted in the tree that resulted from the coalescent based species tree analysis.

### 3.2.4. Maximum likelihood phylogenetic analyses of nucleotide sequence data

We generated the codon-based nucleotide alignment of supermatrix D, by excluding partitions with IC=0 from supermatrix B (supermatrix nt.A, S3: Suppl. Fig. 2, Table 3.3). With this nucleotide supermatrix, we evaluated whether or not (1) there is congruence between amino-acid and nucleotide sequence-based trees, (2) excluding first and third codon positions had a topological effect in the resulting phylogeny of Dytiscoidea, (3) RY-recoding of the nucleotide matrix and subsequent tree reconstruction indicated that heterogeneous base composition is a confounding factor, (4) phylogenetic analyses by including compositionally heterogeneous nCOGs biased tree reconstructions and (5) relative evolutionary rates of COGs affected tree reconstructions. All generated nucleotide sequence supermatrices (Table 3.3, S3: Suppl. Fig. 2) are deposited at MENDELEY DATA (see list of Supplementary materials).

Saturation of nucleotide substitutions at third codon positions is a well-known problem when addressing deep phylogenetic relationships (Philippe et al., 2011; Xia et al., 2003) and was also relevant in a recent multigene phylogenetic study of the dytiscoid relationships (Toussaint et al., 2015). Additionally, nucleotide sequences with

highly heterogeneous GC content in the third codon positions may contribute to phylogenomic conflict (Romiguier et al., 2016). As a result, the authors of many studies have excluded saturated or compositionally heterogeneous sites prior to their phylogenetic analyses (e.g. Breinholt and Kawahara, 2013; Jarvis et al., 2014; Misof et al., 2014; Pauli et al., 2018; Peters et al., 2017). The second codon positions are arguably the most homogeneous sites among the codon triplets of a supermatrix (e.g. Misof et al., 2014; Timmermans et al., 2016) and should therefore deliver the least biased results. In order to dissect the influence of heterogeneous base composition or saturated substitutions on tree reconstructions, we compared the results of tree reconstructions when (1) including all codon positions of supermatrix nt.A for phylogenetic reconstruction, (2) including only the second codon positions and (3) recoding the nucleotide supermatrix nt.A into RY character states (R: Purines, Y: Pyrimidines). The expectation is that a recoded matrix should alleviate problems related to compositional heterogeneity and substitution saturation, at the cost of partially eliminating phylogenetic signal (Philippe and Roure, 2011).

We further explored the effect of masking (i.e., removing) the most compositionally heterogeneous genes (nCOGs) prior to the tree reconstructions (Table 3.3). In order to do so, we generated a decisive version of supermatrix nt.A by discarding those nCOGs with fewer than 14 taxa (S3: Suppl. Fig. 2). We did not perform any tree searches for this intermediate decisive dataset. Subsequently, two reduced versions of this decisive supermatrix were generated by excluding genes with RCFV value greater than 0.08 (supermatrix nt.A.homogeneous1, Table 3.3) and by excluding genes with RCFV value greater than 0.06 (supermatrix nt.A.homogeneous2, Table 3.3). In addition, because the evolutionary rates of individual genes are often

cited as an important predictor of their phylogenetic utility (Doyle et al., 2015; Klopfstein et al., 2017; Yang, 1998), we explored whether the relative evolutionary rates of the included sets of nCOGs biased tree reconstructions (S1: Suppl. Text 1, Table 3.3). Lastly, we tested whether removal of the species *S. wrasei* from supermatices nt.A and nt.A.homogeneous2 affected the phylogenetic placement of Hygrobiidae (Table 3.3). We decided to remove *S. wrasei*, because it is the species that was associated with the longest tree branches among the two species of Aspidytidae when analyzing codon-based nucleotide sequence data (Fig. 3.3).

Ten independent tree searches were performed for each generated nucleotide dataset with IQ-TREE v. 1.5.5 (or later). Tree searches and model selection in ModelFinder were based on an edge-linked partition model (options. -spp -gmedian -merit AICc), by considering the nCOG boundaries and the GTR substitution matrix (Tavaré, 1986), and by allowing all possible combinations for modeling among site rate variation. The RY recoded (in the form of binary data [0,1]) matrix was analyzed with an edge-linked partition model in IQ-TREE v. 1.6.8 (options: -spp -st BIN -m MFP -gmedian -merit AICc). For a selection of nucleotide supermatrices, we optimized the partitioning scheme in PartitionFinder v. 2.1.1 by restricting the model search to GTR and GTR+G with the options -raxml and -rcluster-max 5000 using the AICc information criterion. For this purpose, we selected the datasets with the lowest levels of among-species compositional heterogeneity (Table 3.3). The resulting combinations of partitions and models were used as input for IQ-TREE v. 1.5.5 for ten additional tree searches with the edge-linked model. Statistical branch support was estimated from 100 non-parametric bootstrap replicates, TBE support, 10,000 SH-like aLRT replicates, aBayes, 1000 UFBoot1 (IQ-TREE v. 1.5.5), and 1000 UFBoot2 (IQ-

TREE v. 1.6.8, -bnni) replicates on the datasets with the optimized partitioning schemes and on supermatrix nt.A. After verifying topological congruence to the maximum likelihood tree, we mapped these support values on the tree with the best log-likelihood among the trees that resulted from the ten maximum likelihood searches (Fig. 3.3, S3: Suppl. Fig. 69). We additionally calculated 100 non-parametric bootstrap replicates and TBE support for every other nucleotide sequence dataset (Table 3.3). The NEXUS files with the optimized schemes of the supermatrices nt.B and nt.A.homogeneous2, calculated with PartitionFinder, are deposited at MENDELEY DATA (see list of Supplementary materials).

### 3.2.5. Branch support tests with four-cluster likelihood-mapping and data permutations

We tested the statistical robustness of phylogenomic estimates of four selected phylogenetic hypotheses (S2B, S2C: Suppl. Tables 2 and 3) by means of the four-cluster likelihood-mapping approach (FcLM) on supermatrix E (Strimmer and von Haeseler, 1997). This approach considers the proportion of taxon quartets in a supermatrix that support each of the three alternative topologies around a specific branch of interest (for details, see also the supplementary material provided by Misof et al., 2014). The formulation of each hypothesis was based on the best tree topology inferred from phylogenetically analyzing supermatrix E (Fig. 3.2b). We assumed taxa within each group definition to be monophyletic. For each FcLM test (S2B, S2C: Suppl. Tables 2 and 3) we additionally permuted the original matrix in three ways as described by Misof et al. (2014) to evaluate (1) whether or not the quartet support for a certain hypothesis results from genuine phylogenetic signal, (2) whether or not it is

affected by confounding factors relating to compositional heterogeneity, (3) and whether or not the distribution of missing data affected the phylogenetic results (S1: Suppl. Text 1). The FcLM approach and the permutations for testing hypotheses 1 and 3 were also applied on different amino-acid and nucleotide supermatrices (see also Suppl. Text 1 and Sann et al., 2018 for a description of FcLM tests applied at the nucleotide sequence level) with the same taxon group definitions in an attempt to investigate the source of topological incongruence. For each phylogenetic hypothesis tested, we discarded partitions or meta-partitions (if an optimized scheme was calculated for the respective matrix) that were uninformative with respect to a specific taxon-group definition. For the original dataset we used the same models selected during the IQ-TREE tree search for the respective dataset with the option -spp. For the permuted matrices we used the models LG (for amino-acid alignments) and GTR (for the nucleotide alignments) and the option -q for the partition file. All FcLM analyses were conducted using IQ-TREE v. 1.5.5.

## 3.3. Results

### 3.3.1. Orthology assignment and dataset assembly

On average, 2689 transcripts per species (87% of 3085 COGs) passed the reciprocal best hit criterion (Min. = 2133, Max. = 2913) during the orthology assignment step. The dataset with the lowest number of assigned orthologs (2133) was the transcriptome of the diving beetle *Thermonectus intermedius*, while the transcriptome of the species *S. wrasei* was the dataset with the highest number of assigned orthologous transcripts (2913, Table 3.4). The average number of outlier sequences per species was 0.4% (i.e., a mean of 12 outliers per species across 2991

gene partitions). In total, 167 amino-acid (and corresponding nucleotide) sequences were removed after the alignment refinement step (S2D: Suppl. Table 4). The search for ambiguously aligned regions with ALISCORE resulted in the removal of a total number of 276,537 amino-acid sites from the original amino-acid sequence alignments of supermatrix A (and 829,611 sites from their corresponding codon-based nucleotide sequence alignments).



**Fig. 3.2:** Different phylogenetic hypotheses deduced from the analysis of amino-acid sequence data. (**a**) Phylogram with the best log-likelihood score on the optimized scheme of supermatrix H and (**b**) phylogram with the best log-likelihood score on the optimized scheme of supermatrix E. Branch support is denoted based on 100 non-parametric bootstrap replicates (BS), 100 non-parametric bootstraps based on the PMSF model (BS PMSF), 10,000 SH-like aLRT replicates (SH-aLRT), aBayes support, 1000 Ultrafast Bootstraps 1 (UFBoot1), 1000 Ultrafast

Bootstraps 2 (UFBoot2, -bnni), and 100 bootstraps by transfer (TBE). Both trees were rooted with Gyrinidae. Congruent and incongruent clades between the two trees (in terms of included terminal taxa) are illustrated in different colors. (**c**) Results of the FcLM analysis on the original data of supermatrix E for the phylogenetic hypothesis 1 (i.e., monophyly of Aspidytidae). (**d**) Results of the FcLM analysis on the original data of supermatrix E for the phylogenetic hypothesis 3 (i.e., Hygrobiidae are the sister group of Amphizoidae+Aspidytidae). Beetle photos: (1) *Sinaspidytes wrasei*, (2) *Noterus crassicornis*, (3) *Hygrobia hermanni*, (4) *Amphizoa lecontei*, (5) *Cybister lateralimarginalis* (photos and copyright: M. Balke).

### 3.3.2. Phylogenetic analyses of amino-acid sequence data

The different maximum likelihood searches for the same datasets resulted in congruent topologies (Fig. 3.2 and S3: Suppl. Fig. 45–59) irrespective of whether or not we optimized the partitioning scheme (for supermatrices E and H respectively). The phylogenetic analyses with the site-heterogeneous mixture models yielded topologies identical to those obtained when using partition models for the amino-acid datasets analyzed (S3: Suppl. Fig. 49, 51, 55, 57). All phylogenetic analyses inferred the monophyly Dytiscoidea as a whole and of each dytiscoid family, and supported a sister group relationship between Noteridae and all remaining families of Dytiscoidea. All the above relationships received high statistical support when analyzing amino-acid sequence data except for the monophyly of Aspidytidae when performing FcLM analysis on supermatrix E (see Section 3.3.4). Moreover, a clade comprising the families Amphizoidae and Aspidytidae was suggested in all maximum likelihood analyses of amino-acid sequence data and is fully supported by all branch support measures (Fig. 3.2a and 3.2b). FcLM analysis on both the original and the permuted data of supermatrix E indicate high support for a clade consisting of Amphizoidae and

Aspidytidae without detectable confounding signal (Section 3.3.4, Hypothesis 2, S2B: Suppl. Table 2).

The phylogenetic analyses of the amino-acid supermatrices which were not corrected for among-species compositional heterogeneity, suggested Hygrobiidae as the sister clade to Aspidytidae + Amphizoidae with strong statistical branch support. Analyses of these datasets suggested that the three families collectively form a clade sister to the diving beetles (e.g., Fig. 3.2b). The analysis of supermatrix H (RCFV-corrected version of supermatrix F) yielded a different arrangement with Hygrobiidae being placed as the sister group of (Amphizoidae + Aspidytidae) + Dytiscidae (Fig. 3.2a). Furthermore, the phylogenetic analysis of the supermatrices J and K (RCFV-corrected versions of supermatrices E and A respectively) also suggested the latter sister group relationship (S3: Suppl. Fig. 58–59). Non-parametric bootstrap support for the clade (Amphizoidae + Aspidytidae) + Dytiscidae is not very high (supermatrix H: 79%, Fig. 3.2a, see also S3: Suppl. Fig. 54, 58–59), but most measures such as BS PMSF, UFBoot1, aBayes, SH-aLRT and TBE strongly support this clade.

The coalescent-based species tree analyses with ASTRAL yielded topologies identical to those obtained from concatenation when analyzing supermatrices E and H (S3: Suppl. Fig. 71–72). Overall, the local posterior probabilities in favor of the monophyly of the dytiscoid lineages except Noteridae (i.e., Aspidytidae + Amphizoidae + Dytiscidae + Hygrobiidae), the monophyly of Aspidytidae, and the monophyly of Amphizoidae+Aspidytidae are high in both coalescent phylogenetic analyses. On the one hand, quartet support shows conflict among the selected gene trees of supermatrix E concerning the monophyly of Aspidytidae ($q1=0.44$; $q2=0.32$; $q3=0.22$) and the placement of Hygrobiidae as a sister group to Aspidytidae and

Amphizoidae (q1=0.37; q2=0.26; q3=0.36). On the other hand, the local posterior probabilities for the above relationships are high (0.99 and 0.90 respectively). A low quartet support for the monophyly of Aspidytidae is again observed when analyzing the gene trees of supermatrix H (q1=0.45; q2=0.32; q3=0.21), indicating conflict among the gene trees of this dataset for this relationship. A clade comprising Amphizoidae, Aspidytidae, and Dytiscidae (which resulted from the coalescent analysis of the genes in supermatrix H) received low quartet support (q1=0.37; q2=0.36; q3=0.26). This clade also received low support based on the local posterior probability value (0.73).

### 3.3.3. Phylogenetic analyses of nucleotide sequence data

In contrast to the analysis of the amino-acid sequence data, phylogenetic analysis of the codon-based nucleotide sequence data (supermatrix nt.A) yielded paraphyletic Aspidytidae, with *S. wrasei* placed as the sister taxon of Amphizoidae (Fig. 3.3b). However, after removal of the most compositionally heterogeneous genes, the phylogenetic analyses provided strong statistical branch support for the monophyly of Aspidytidae (Fig. 3.3a, S3: Suppl. Fig. 65–67). Analyzing exclusively second codon positions also provided strong support for the hypothesis of Aspidytidae representing a natural group (S3: Suppl. Fig. 60 and 69). The best tree from the analysis of the RY-recoded supermatrix supported the monophyly of Aspidytidae as well (S3: Suppl. Fig. 70). Some of the interfamiliar relationships recovered by the analysis of the recoded nucleotide sequence matrix are different than the relationships recovered from most of our analyses. The branch support values for those relationships are high but the internal branches of the tree are very short (S3: Suppl. Fig. 70). As expected, including

only the fastest evolving genes in the dataset delivered phylogenetic relationships (including paraphyletic Dytiscoidea) not seen in any of the other phylogenetic analyses (S3: Suppl. Fig. 62). In contrast, removing the ca. 25% or 75% of the fastest evolving genes did not result in topological alterations compared with the original results of the analysis of supermatrix nt.A (S3: Suppl. Fig. 61 and 63). Phylogenetic analyses of the concatenated codon-based nucleotide sequence dataset after removing outlier genes with respect to their relative evolutionary rate (S3: Suppl. Fig. 64), yielded the same topology as the analysis of the supermatrix composed of exclusively slowly evolving genes (S3: Suppl. Fig. 61).

Analysis of the nucleotide datasets did not corroborate the hypothesis of Hygrobiidae being the sister group to a clade comprising Aspidytidae, Dytiscidae and Amphizoidae, except when analyzing exclusively second codon positions. One additional difference between the trees derived from analyzing codon-based nucleotide sequence data and the tree based on the analysis of exclusively second codon positions is the placement of Amphizoidae as the sister group of Dytiscidae (S3: Suppl. Fig. 60 and 69). However, this placement is in conflict with the phylogenies inferred when analyzing amino-acid data and which suggested a sister group relationship of Amphizoidae and Aspidytidae (Fig. 3.2) with high support. The results of the FcLM analysis on the amino-acid supermatrix E (S2C: Suppl. Table 3) are also in support of a clade Amphizoidae+Aspidytidae without detectable confounding signal (see Section 3.3.4). Removal of the species *S. wrasei* from the selected codon-based datasets (nt.A and nt.A.homogeneous2) did not affect the phylogenetic placement of Hygrobiidae (S3: Suppl. Fig. 67–68). However, after removal of *S. wrasei* from the compositionally

homogeneous matrix the monophyly of (Amphizoidae + Aspidytidae) + Hygrobiidae is only weakly supported (S3: Suppl. Fig. 67).

### 3.3.4. Branch support tests with four-cluster likelihood mapping and data permutations

*Monophyly of Aspidytidae*

All trees based on the MSAs of amino-acid sequences recovered a monophyletic Aspidytidae. The FcLM analysis of the amino-acid sequence data did not, however, strongly support the monophyly of Aspidytidae (Fig. 3.2c: 55% of quartets support a monophyletic Aspidytidae when analyzing the original data of supermatrix E). The FcLM results when analyzing supermatrix E show some weaker signal for the placement of *A. niobe* as sister group to Amphizoidae (40% of quartets). Additionally, after eliminating phylogenetic signal in supermatrix E (permutation scheme I) putative confounding signal emerges supporting the monophyly of Aspidytidae (75% of quartets). This signal is reduced after having applied permutation scheme II on supermatrix E (40% of quartets), suggesting that it stems from non-stationary processes among species in supermatrix E (S2B: Suppl. Table 2). When the effect of among-species compositional heterogeneity is reduced in the original data (supermatrices H and K), the putative confounding signal supporting the monophyly of Aspidytidae decreases (25% and 20% of quartets, permutation scheme I, supermatrix H and K respectively) and the support for the monophyly of Aspidytidae when analyzing the original data increases (60% of quartets are in favor of the monophyly of Aspidytidae when analyzing the original data of supermatrices H and K).

**Table 3.3**: Detailed information and statistics of each generated nucleotide supermatrix analyzed in this study. The overall alignment completeness score of each matrix was calculated with AliStat. Pairwise tests of symmetry for the Bowker's test were performed with SymTest. Median p-values 0.00E+00 for the Bowker's test indicate very small numbers. ($C_a$: Overall alignment completeness score).
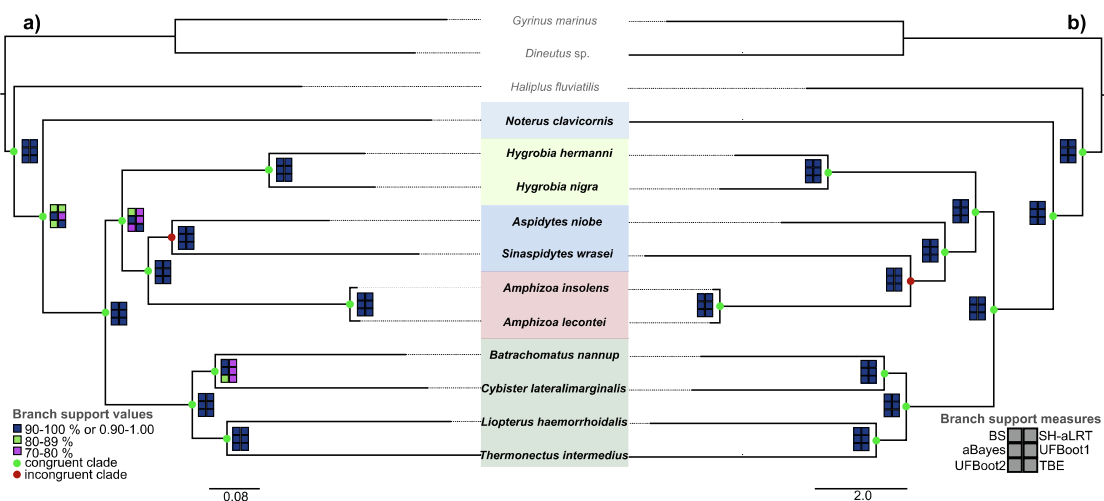
| Nucleotide dataset | No. of taxa | No. of nucleotide sites | No. of gene partitions | $C_a$ | Percentage of pairwise p-values < 0.05 for the Bowker's test | Median pairwise p-value for the Bowker's test | No. of tree searches with the unoptimized partitioning scheme | No. of bootstraps with the unoptimized partitioning scheme | Optimization of the partitioning scheme | No. of tree searches with the optimized partitioning scheme | No. of bootstraps with the optimized partitioning scheme | Information |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| supermatrix.nt.A | 14 | 4,098,894 | 2948 | 0.6889 | 98.90 % | 0.00E+00 | 10 | 100 | NO | - | | Codon-based nucleotide sequence alignment of -supermatrix C |
| supermatrix nt.B | 14 | 1,366,298 | 2948 | 0.6889 | 97.80 % | 3.20E-39 | 10 | 100 | YES | 10 | 100 | Second codon positions of supermatrix nt.A |
| supermatrix nt.A.recoded | 14 | 4,098,894 | 2948 | N/A | N/A | N/A | 10 | 100 | NO | - | | RY recoded matrix of -supermatrix nt.A |
| supermatrix nt.A.homogeneous1 | 14 | 617,355 | 498 | 0.8427 | 98.90 % | 0.00E+00 | 10 | 100 | NO | - | | Removed genes with RCFV > 0.08 from the decisive version of -supermatrix nt.A |
| supermatrix nt.A.homogeneous2 | 14 | 186,498 | 170 | 0.8849 | 98.90 % | 8.40E-75 | 10 | 100 | YES | 10 | 100 | Removed genes with RCFV > 0.06 from a decisive version of supermatrix nt.A |
| supermatrix nt.A.slow | 14 | 920,700 | 737 | 0.6074 | 98.90 % | 0.00E+00 | 10 | 100 | NO | - | | Removed genes with a relative rate > Q1 of sorted rates from -supermatrix nt.A |
| supermatrix nt.A.fast | 14 | 1,204,353 | 749 | 0.6623 | 100.00 % | 0.00E+00 | 10 | 100 | NO | - | | Removed genes with a relative rate < Q3 of sorted rates from -supermatrix nt.A |
| supermatrix nt.A.fast_removed | 14 | 2,913,135 | 2212 | 0.7002 | 100.00 % | 0.00E+00 | 10 | 100 | NO | - | | Removed genes with a relative rate > Q3 of sorted rates from -supermatrix nt.A |
| supermatrix nt.A.out_removed | 14 | 3,811,368 | 2804 | 0.7001 | 98.90 % | 0.00E+00 | 10 | 100 | NO | - | | Removed genes with outlier values of relative rates from -supermatrix nt.A |
| supermatrix.nt.A.sw | 13 | 4,092,338 | 2948 | 0.6805 | 98.72 % | 0.00E+00 | 10 | 100 | NO | - | | Removed species *Sinaspidytes wrasei* from supermatrix nt.A |
| supermatrix nt.A.homogeneous2.sw | 13 | 186,468 | 170 | 0.8810 | 98.72 % | 1.06E-48 | 10 | 100 | NO | - | | Removed species *Sinaspidytes wrasei* from supermatrix nt.A.homogeneous2 |

**Table 3.4**: Summarized statistics of the results of the transcript orthology assignment at the amino-acid sequence level. Species whose transcriptomes were analyzed are given in alphabetic order. The summary statistics were calculated with the helper scripts provided with the Orthograph package.

| Species name/Transcriptome | No. of orthologous hits | Proportion of COGs (%) | Total no. of amino acids | No. of X residues | No. of stop codons | N50 of protein lengths | Mean protein length | Median protein length | Maximum protein length | Minimum protein length |
|---|---|---|---|---|---|---|---|---|---|---|
| *Amphizoa insolens* LeConte, 1853 | 2820 | 91.41 % | 1,109,394 | 0 | 13 | 491 | 393 | 325 | 3633 | 30 |
| *Amphizoa lecontei* Matthews, 1872 | 2765 | 89.63 % | 984,227 | 0 | 39 | 446 | 355 | 304 | 2409 | 9 |
| *Aspidytes niobe* Ribera, Beutel, Balke, Vogler, 2002 | 2780 | 90.11 % | 1,077,674 | 20 | 26 | 485 | 387 | 328 | 2159 | 20 |
| *Batrachomatus nannup* (Watts, 1978) | 2561 | 83.01 % | 797,222 | 0 | 41 | 391 | 311 | 265 | 2142 | 6 |
| *Cybister lateralimarginalis* (DeGeer, 1774) | 2680 | 86.87 % | 1,084,064 | 16 | 21 | 508 | 404 | 332 | 6510 | 10 |
| *Dineutus* sp. | 2642 | 85.64 % | 781,715 | 72 | 11 | 362 | 295 | 259 | 2168 | 15 |
| *Gyrinus marinus* Gyllenhal, *1808* | 2571 | 83.34 % | 830,399 | 12 | 16 | 395 | 322 | 291 | 1478 | 13 |
| *Haliplus fluviatilis* Aubé, 1836 | 2891 | 93.71 % | 1,171,464 | 88 | 33 | 502 | 405 | 337 | 2924 | 17 |
| *Hygrobia hermanni* (Fabricius, 1775) | 2903 | 94.10 % | 1,249,213 | 17 | 40 | 541 | 430 | 351 | 3455 | 12 |
| *Hygrobia nigra* (Clark, 1862) | 2662 | 86.29 % | 950,213 | 13 | 32 | 444 | 356 | 309 | 1977 | 9 |
| *Liopterus haemorrhoidalis* (Fabricius, 1787) | 2450 | 79.42 % | 698,178 | 0 | 48 | 351 | 284 | 246 | 2249 | 13 |
| *Noterus clavicornis* (DeGeer, 1774) | 2868 | 92.97 % | 1,128,976 | 6 | 38 | 485 | 393 | 329 | 6482 | 6 |
| *Sinaspidytes wrasei* (Balke, Ribera, Beutel, 2003) | 2913 | 94.42 % | 1,187,784 | 51 | 28 | 515 | 407 | 340 | 3305 | 8 |
| *Thermonectus intermedius* Crotch, 1873 | 2133 | 69.14 % | 897,627 | 0 | 6 | 524 | 420 | 340 | 6828 | 6 |

Maximum likelihood phylogenetic analysis of the supermatrix nt.A strongly supports the sister group relationship between *S. wrasei* and Amphizoidae, as indicated by all applied branch support measures (Fig. 3.3b). This arrangement also received relatively high quartet support from the FcLM analysis on the original data of supermatrix nt.A (70% of quartets, S2C: Suppl. Table 3). There is however strong putatively confounding phylogenetic signal in favor of this hypothesis after applying permutation scheme I on supermatrix nt.A (70% of quartets). This signal is greatly reduced in permutation number II of the same matrix (20% of quartets), suggesting that it stems from non-stationary processes among species in the supermatrix nt.A. The total number of different quartets that are informative with respect to the monophyly of Aspidytidae is low (20 quartets, S2B: Suppl. Table 2) due to the low number of species in our dataset.



**Fig. 3.3:** Comparison of phylogenetic hypotheses resulted from the analysis of the codon-based nucleotide sequence data. Congruent and incongruent clades between the two trees (in terms of included terminal taxa) are illustrated in different colors. (**a**) Phylogram with the best log-likelihood score on the optimized scheme of supermatrix nt.A.homogeneous2. (**b**) Phylogram with the best log-likelihood score on the unoptimized partitioning scheme of

supermatrix nt.A. Branch support is denoted based on 100 non-parametric bootstrap replicates (BS), 10,000 SH-like aLRT replicates (SH-aLRT), aBayes support, 1000 Ultrafast Bootstraps 1 (UFBoot1), 1000 Ultrafast Bootstraps 2 (UFBoot2, -bnni), and 100 bootstraps by transfer (TBE). Both trees were rooted with Gyrinidae.

*Phylogenetic relationships of the dytiscoid families*

In all our tree reconstructions, Noteridae were inferred as the sister taxon of all remaining Dytiscoidea (e.g., Fig. 3.2a, 3.2b, 3.3a, 3.3b). This phylogenetic placement received strong support from most applied statistics, and is also supported by the FcLM and data permutation tests on supermatrix E (100% of quartets support a clade of Dytiscidae + Hygrobiidae + Amphizoidae + Aspidytidae as the sister group of Noteridae, S2B: Suppl. Table 2, Hypothesis 4). In addition, a clade of Aspidytidae + Amphizoidae is fully supported by all analyses based on the amino-acid and nucleotide sequences, except for the analyses of the second codon positions (S3: Suppl. Fig. 60 and 69). We observed a strong signal in favor of Amphizoidae+Aspidytidae when analyzing the original data of supermatrix E (95.3% of quartets support Amphizoidae+Aspidytidae, S2B: Suppl. Table 2), and no detectable confounding signal for this arrangement after applying permutation scheme I on the same amino-acid dataset (39.1% of quartets support Amphizoidae+Aspidytidae when eliminating phylogenetic signal in supermatrix E).

The position of Hygrobiidae with respect to Amphizoidae, Aspidytidae and Dytiscidae differs between the trees that were inferred at the amino-acid sequence level when allowing for different degrees of compositional heterogeneity among species in the dataset (e.g., Fig. 3.2). The two prevailing phylogenetic hypotheses that were inferred from analyzing amino-acid sequence data (Fig. 3.2a and 3.2b) received almost

equally high support in the FcLM analyses of the different amino-acid and nucleotide data matrices with no detectable confounding factors (Fig. 3.2d, S2B, S2C: Suppl. Tables 2 and 3). This result indicates the substantial phylogenetic conflict among the analyzed quartets for this particular phylogenetic question. Again, the total number of quartets for investigating the phylogenetic hypothesis number 3 was not very high (128 quartets) due to taxon sampling limitations in our dataset.

## 3.4. Discussion

### 3.4.1. The phylogeny of the dytiscoid families and the monophyly of Aspidytidae

Previous analyses based on either morphological or molecular data were unable to deliver congruent reconstructions of dytiscoid phylogenetic relationships (e.g., Baca et al., 2017; Balke et al., 2005, 2008 Beutel et al., 2008, 2013; Toussaint et al., 2015). We addressed these phylogenetic problems with an unprecedented amount of phylogenomic data representing all dytiscoid families except Meruidae. Results of our phylogenomic analyses are consistent with the hypothesis of Noteridae (plus most likely Meruidae) being the sister group of a clade comprising the families Amphizoidae, Aspidytidae, Dytiscidae, and Hygrobiidae (Baca et al., 2017; Beutel et al., 2008; Dressler et al., 2011; McKenna et al., 2015). The monophyly of the latter clade received strong statistical support in all of our analyses. The phylogenetic relationships within this clade, however, are not robustly resolved and resolution depends on the phylogenetic approach and dataset. Nevertheless, our analyses demonstrate that selecting the datasets that violate model assumptions the least support a sister group relationship between Hygrobiidae and a clade comprising Amphizoidae, Aspidytidae, and Dytiscidae. The monophyly of the latter three families is also

suggested by an unusual morphological apomorphy, a pair of large and sclerotized epipharyngeal sensilla (Dressler and Beutel, 2010). A clade comprising the squeak beetles and the diving beetles (Hygrobiidae + Dytiscidae), as suggested by some studies based on the analysis of morphological characters (e.g., Alarie and Bilton, 2005; Beutel et al., 2013; Beutel and Roughley, 1988; Dressler et al., 2011) was not recovered in any of our analyses. This suggests that prothoracic glands (Forsyth, 1970) have evolved independently in the two families.

All analyses of amino-acid sequence data and nucleotide sequence data with reduced levels of among-species compositional heterogeneity suggest monophyletic Aspidytidae. This result is congruent with the analysis of the morphological characters of the adults of Aspidytidae (Balke et al., 2003). Moreover, we received high branch support and high FcLM support for a clade consisting of Amphizoidae and Aspidytidae in all analyses of amino-acid sequence data, and this phylogenetic relationship is also supported by the analysis of codon-based nucleotide sequence data. On the other hand, the analysis of second codon positions suggest a sister group relationship of Amphizoidae and Dytiscidae. The cause of this incongruent result is unclear, but may be due to insufficient or conflicting signal for this relationship in the second codon positions. Overall, we consider a sister group relationship of Amphizoidae and monophyletic Aspidytidae as the most plausible scenario suggested by our data.

The disjunct geographical distribution of Amphizoidae, Aspidytidae and Hygrobiidae in combination with the extensive molecular divergence among the three families, and between the two aspidytid species in particular, suggests that these groups represent old and relictual lineages. In this aspect, we corroborate the results put forth by Toussaint et al. (2015) and Hawlitschek et al. (2012), who came to similar

conclusions, but these conclusions were based on phylogenetic results from only a few molecular loci. Thus, our results provide a base line for future phylogenomic analyses of dytiscoid relationships and help to identify the most pressing open questions. Additionally, we want to emphasize that the disjunct, relict and micro-endemic distribution of Aspidytidae demands appropriate actions to conserve their habitats and future existence.

The instability of the phylogenetic placement of Hygrobiidae among the different datasets analyzed deserves special attention. The lack of resolution in phylogenetics is often attributed to biological phenomena of ancient rapid cladogenesis (Whitfield and Kjer, 2008). Signatures of such processes when analyzing genome-scale data are illustrated by either low levels of phylogenetic signal or highly conflicting phylogenetic signal (Suh, 2016; Whitfield and Kjer, 2008). Our FcLM results as well as the coalescent analyses showed substantial levels of phylogenomic conflict for the interrelationships of the dytiscoid families Amphizoidae, Aspidytidae and Hygrobiidae. The large molecular divergence observed between these families and within Aspidytidae, together with their disjunct geographical distributions and the high levels of gene tree conflict for the interfamiliar relationships observed here, are indications that these lineages may have originated via rapid cladogenesis. On the other hand, such ancient rapid speciation events can be difficult to distinguish from other causes related to data quality and conflict in the analyzed datasets (Whitfield and Kjer, 2008) and this hypothesis should be further tested using molecular dating and diversification analyses.

The lack of phylogenetic resolution can be the result of deficient taxon sampling (Nabhan and Sarkar, 2012). We acknowledge the sensitivity of phylogenetic reconstructions to taxon sampling, yet we consider our dataset as the most

comprehensive genome-scale dataset to date in terms of the number of included species within the small families Amphizoidae, Aspidytidae and Hygrobiidae. Furthermore, we acknowledge that the statistical power of the FcLM approach is highly dependent on the number of sampled species. Increasing the available genomic data, especially within the species-rich Dytiscidae and Noteridae, will inevitably boost the statistical power of the FcLM analyses and further facilitate addressing the persisting phylogenetic uncertainties. Lastly, the analysis of other kind of data such as whole genome sequences, and genomic meta-characters can provide additional or complementary evidence to decipher the evolutionary history of Dytiscoidea (Niehuis et al., 2012).

### 3.4.2. Model violations bias the reconstruction of the phylogeny of Dytiscoidea

We pointed out that model violations are one very likely source of the observed phylogenetic discrepancies among the different datasets that we analyzed. This is not an unknown phenomenon, as violations of model assumptions, uneven distribution of data coverage, data-type effects, or unnoticed cross-contamination are some of the factors that can strongly bias the results of tree reconstructions (Borowiec et al., 2019; Feuda et al., 2017; Jeffroy et al., 2006; Jermiin et al., 2004; Nesnidal et al., 2013; Philippe et al., 2011; Reddy et al., 2017; Whitfield and Kjer, 2008). In the presented analyses of the dytiscoid relationships we are able to show that masking the genes with the highest levels of among-species compositional heterogeneity altered the topologies of the inferred phylogenetic trees. This was the case irrespective of whether we analyzed amino-acid sequence data or nucleotide sequence data. We deduce from this

that scientists should seek to take measures against violations of model assumptions in order to more accurately infer the real evolutionary history of the taxa of interest.

At the amino-acid sequence level, we reconstructed phylogenetic relationships of Dytiscoidea based on three supermatrices for which the most compositionally heterogeneous genes had been removed (supermatrices H, J, and K). All of these reconstructions yielded congruent topologies, with respect to the interrelationships of the dytiscoid families, which differed from the topologies that resulted from the analyses of the compositionally heterogeneous amino-acid sequence datasets. The effects of among-species compositional heterogeneity at the amino-acid sequence level is further corroborated by our FcLM tests. Although Aspidytidae are recovered as a monophylum when analyzing amino-acid sequence data, there is detectable confounding signal supporting this monophyly in the compositionally heterogeneous supermatrix E. This putatively confounding signal most likely stems from compositional heterogeneity among species in the alignment because it is reduced when analyzing the datasets with reduced levels of among-species compositional heterogeneity. Furthermore, despite the fact that phylogenetic analysis of both the compositionally homogeneous and the compositionally heterogeneous amino-acid datasets yielded monophyletic Aspidytidae, the compositionally homogeneous supermatrices showed slightly increased phylogenetic signal supporting the monophyly of Aspidytidae. We conclude from these observations that gene partitions with high degrees of among-species compositional heterogeneity biased some of our phylogenetic analyses and are one very likely source of incongruence between tree topologies inferred from analyzing amino-acid sequence data.

Summary coalescent phylogenetic analyses (Mirarab and Warnow, 2015) suggested topologies identical to those obtained when applying a concatenation approach. The observation that both approaches resulted in the same topology irrespective of what dataset we analyzed makes us confident that the incongruence between topologies of different datasets are not due to high levels of incomplete lineage sorting or ancient introgression. This observation further suggests that the applied summary species tree method is sensitive to the same compositional bias as the supermatrix approach.

Our results showed that reducing the degree of missing data and indecisive gene partitions in the amino-acid supermatrices did not affect the topology of the reconstructed dytiscoid phylogeny. The analysis of the amino-acid sequence supermatrix with 100% data coverage across all species delivered the same topology as the analyses of the non-homogeneous datasets, further supporting the idea that non-random distribution of missing data unlikely accounts for the observed topological differences. Additionally the use of site-heterogeneous amino-acid mixture models in a maximum likelihood framework yielded identical topologies compared with the analysis based on site-homogeneous partition models. The overall information content of the supermatrices (Misof et al., 2013) could not be related to the topological incongruence.

It has been argued that alignment masking might be detrimental to reliable phylogenetic reconstructions (Tan et al., 2015). Tan et al. (2015) argue that alignment masking eliminates too much phylogenetic signal and therefore reduces the resolution of single-gene phylogenetic inferences. We found no evidence that alignment masking

affected the topology of the dytiscoid phylogeny in the analyses of concatenated and masked aaCOGs.

The analysis of the nucleotide sequence data revealed that first and third codon positions are heterogeneous in their base composition, because their inclusion results in a major deviation from SRH conditions. Congruently, the Bowker's pairwise symmetry tests corroborate previous hypotheses that the smallest deviations from SRH conditions are consistently observed in datasets composed solely of second codon positions. Reducing among-species compositional heterogeneity, by recoding the nucleotide sequence data or by removing compositionally heterogeneous genes, restored the monophyly of the cliff water beetles, congruent with tree reconstructions based on the amino-acid sequence datasets. These results indicate that the paraphyly of Aspidytidae as it was found by Toussaint et al. (2015) could also be an artifact resulting from compositional biases in the underlying dataset. Additional evidence for the effect of compositional bias on the analysis of the nucleotide sequence data comes from the results of the FcLM. The FcLM results on supermatrix nt.A suggest that the paraphyletic Aspidytidae stems from non-stationary processes among species in the analyzed dataset, as the signal in favor of this relationship is greatly reduced when applying permutation scheme II. The FcLM results of the nucleotide matrix after reducing among-species compositional heterogeneity shows that there is weak signal supporting the original results (40%) but there are no detectable confounding effects observed for this arrangement. Taken together these results suggest that the observed paraphyly Aspidytidae obtained when analyzing supermatrix nt.A probably stems from systematic bias owing to among-species compositional heterogeneity in first and third codon positions.

We compared the resolution of three distinct sets of genes relative to their evolutionary rate and found that except for the set of genes with the highest relative evolutionary rates, the selection of gene sets did not influence the results. In the extreme case of analyzing a set of the ca. 25% of the fastest evolving genes in our supermatrix, we recovered many unexpected relationships, which in turn suggests that including only fast evolving genes results in erroneous phylogenetic estimates of the dytiscoid relationships. Analyses based on the 25% of the most slowly evolving genes yielded results congruent with those obtained when analyzing all genes (i.e., those of supermatrix nt.A). We also find that after extending the phylogenetic analysis to the 75% of the slowest evolving genes (i.e., by removing only the 25% of the fastest evolving genes), the relationships recovered are the same as when analyzing supermatrix nt.A, including the paraphyly of Aspidytidae. Hence, we hypothesize that the paraphyly of Aspidytidae, obtained when analyzing the nucleotide sequence data of supermatrix nt.A, is very likely not driven by the confounding effects of genes with very high evolutionary rates.

## 3.5. Conclusions

Our extensive phylogenomic analyses resolve some outstanding issues in adephagan beetle phylogeny, as well as pointing to some problems which apply to phylogenomic approaches more generally. We present evidence that the cliff water beetles (Aspidytidae) constitute a monophylum despite their highly disjunct geographical distribution and large molecular divergence. In addition, our analyses suggest that Aspidytidae are the closest relatives of Amphizoidae. The close affinity of Amphizoidae and Aspidytidae is supported by most of our phylogenetic analyses and

by FcLM tests of amino-acid sequence data. Our study could not provide conclusive evidence for some of the interfamiliar relationships of Dytiscoidea, yet we show that excluding genomic regions with high among-species compositional heterogeneity yields different topologies for our transcriptomic dataset. After accounting for most potential tree confounding factors, we consider a sister group relationship between Hygrobiidae and a clade comprising Amphizoidae, Aspidytidae, and Dytiscidae to most likely represent the evolutionary relationships. Overall, we demonstrated in our study how confounding parameters can lead to misleading results. Our study also highlights the importance of interpreting, integrating and summarizing across different datasets and tree-inference approaches for drawing major phylogenetic conclusions. It is obvious that incongruence due to model violations, uneven distribution of missing data, unequal evolutionary rates, as well as conflicting phylogenetic signal among gene trees will prevail in primarily sequence-based phylogenomic analyses, and measures need to be taken against violations of model assumptions. An alternative or complementary route would be the comparative analyses of genomic meta-characters such as the position of introns, the evolution of gene families, or the structure of genes. The tremendous advances in sequencing technologies are currently opening a window into these fields of research (Niehuis et al., 2012).

## 3.6. References

Alarie, Y., Beutel, R.G., Watts, H.S., 2004. Larval morphology of three species of Hygrobiidae (Coleoptera: Adephaga: Dytiscoidea) with phylogenetic considerations. Eur. J. Entomol. 101, 293–311.

Alarie, Y., Bilton, D.T., 2005. Larval morphology of Aspidytidae (Coleoptera: Adephaga) and its phylogenetic implications. Ann. Entomol. Soc. Am. 98, 417–430..

Alarie, Y., Short, A.E.Z., Garcia, M., Joly, L., 2011. Larval morphology of Meruidae (Coleoptera: Adephaga) and its phylogenetic implications. Ann. Entomol. Soc. Am. 104, 25–36.

Anisimova, M., Gil, M., Dufayard, J.F., Dessimoz, C., Gascuel, O., 2011. Survey of branch support methods demonstrates accuracy, power, and robustness of fast likelihood-based approximation schemes. Syst. Biol. 60, 685–699.

Baca, S.M., Alexander, A., Gustafson, G.T., Short, A.E.Z., 2017. Ultraconserved elements show utility in phylogenetic inference of Adephaga (Coleoptera) and suggest paraphyly of 'Hydradephega'. Syst. Entomol. 42, 786–795.

Balke, M., Hendrich, L., 2016. 7.6 Dytiscidae Leach. In: Beutel, R.G., Leschen, R.A.B. (Eds.), Handbook of Zoology. Volume IV Arthropoda: Insecta Part 38. Coleoptera, Beetles. Morphology and Systematics (Archostemata, Adephaga, Myxophaga, Polyphaga Partim), vol. 1. W. DeGruyter, Berlin, pp. 116–147.

Balke, M., Ribera, I., Beutel, R., Viloria, A., Garcia, M., Vogler, A.P., 2008. Systematic placement of the recently discovered beetle family Meruidae (Coleoptera: Dytiscoidea) based on molecular data. Zool. Scr. 37, 647–650.

Balke, M., Ribera, I., Beutel, R.G., 2005. The systematic position of Aspidytidae, the diversification of Dytiscoidea (Coleoptera, Adephaga) and the phylogenetic signal of third codon positions. J. Zool. Syst. Evol. Res. 43, 223–242.

Balke, M., Ribera, I., Beutel, R.G., 2003. ASPIDYTIDAE: On the discovery of a new beetle family: detailed morphological analysis, description of a second species, and key to fossil and extant adephagan families (Coleoptera). In: Jäch, M.A., Ji, L. (Eds.), Water Beetles of China. Zoologisch-Botanische Gesellschaft & Wiener Coleopterologenverein, Wien, pp. 53–66.

Bayzid, M.S., Warnow, T., 2013. Naive binning improves phylogenomic analyses. Bioinformatics 29, 2277–2284.

Belkaceme, T., 1991. Skelet und Muskulatur des Kopfes und Thorax von *Noterus laevis* Sturm. Ein Beitrag zur Morphologie und Phylogenie der Noteridae (Coleoptera: Adephaga). Stuttgarter Beiträge zur Naturkunde Ser. A 462, 1–94.

Beutel, R.G., 1993. Phylogenetic analysis of Adephaga (Coleoptera) based on characters of the larval head. Syst. Entomol. 18, 127–147.

Beutel, R.G., 1988. Studies of the metathorax of the trout-stream beetle, *Amphizoa lecontei* Matthews (Coleoptera:Amphizoidae): contribution towards clarification of the systematic position of Amphizoidae. Int. J. Insect Morphol. Embryol. 17, 63–81.

Beutel, R.G., 1986a. Skelet und Muskulatur des Kopfes der Larve von *Haliplus lineatocollis* Mrsh. (Coleoptera). Stutt. Beitr. Naturkd. 390, 1–15.

Beutel, R.G., 1986b. Skelet und Muskulatur des Kopfes und Thorax von *Hygrobia tarda* (Herbst). Ein Beitrag zur Klärung der phylogenetischen Beziehungen der Hydradephaga (Insecta: Coleoptera). Stutt. Beitr. Naturkd. 388, 1–54.

Beutel, R.G., Balke, M., Ribera, I., 2010. 3.1. Aspidytidae Ribera, Beutel, Balke and Vogler, 2002. In: Leschen, R.A.B., Beutel, R.G., Lawrence, J.F. (Eds.), Handbook of Zoology, Arthropoda: Insecta. Coleoptera, Beetles. Morphology and Systematics (Elateroidea, Bostrichiformia, Cucujiformia Partim), vol. 2. W. DeGruyter, Berlin, pp. 21–28.

Beutel, R.G., Balke, M., Steiner, W.E., 2006. The systematic position of Meruidae (Coleoptera, Adephaga) and the phylogeny of the smaller aquatic adephagan beetle families. Cladistics 22, 102–131.

Beutel, R.G., Haas, A., 1996. Phylogenetic analysis of larval and adult characters of Adephaga (Coleoptera) using cladistic computer programs. Entomol. Scand. 27, 197–205.

Beutel, R.G., Ribera, I., Bininda-Emonds, O.R.P., 2008. A genus-level supertree of Adephaga (Coleoptera). Org. Divers. Evol. 7, 255–269.

Beutel, R.G., Roughley, R.E., 1988. On the systematic position of the family Gyrinidae (Coleoptera: Adephaga). J. Zool. Syst. Evol. Res. 26, 380–400.

Beutel, R.G., Wang, B., Tan, J.J., Ge, S.Q., Ren, D., Yang, X.K., 2013. On the phylogeny and evolution of Mesozoic and extant lineages of Adephaga (Coleoptera, Insecta). Cladistics 29, 147–165.

Borowiec, M.L., Rabeling, C., Brady, S.G., Fisher, B.L., Schultz, T.R., Ward, P.S., 2019. Compositional heterogeneity and outgroup choice influence the internal phylogeny of the ants. Mol. Phylogenet. Evol. 134, 111–121.

Boussau, B., Walton, Z., Delgado, J.A., Collantes, F., Beani, L., Stewart, I.J., Cameron, S.A., Whitfield, J.B., Johnston, J.S., Holland, P.W.H., Bachtrog, D., Kathirithamby, J., Huelsenbeck, J.P., 2014. Strepsiptera, phylogenomics and the long branch attraction problem. PLoS One 9, e107709.

Breinholt, J.W., Kawahara, A.Y., 2013. Phylotranscriptomics: Saturated third codon positions radically influence the estimation of trees based on next-gen data. Genome Biol. Evol. 5, 2082–2092.

Burmeister, E.G., 1976. Der Ovipositor der Hydradephaga (Coleoptera) und seine phylogenetische Bedeutung unter besonderer Berücksichtigung der Dytiscidae. Zoomorphologie 85, 165–257.

Chernomor, O., Von Haeseler, A., Minh, B.Q., 2016. Terrace aware data structure for phylogenomic inference from supermatrices. Syst. Biol. 65, 997–1008.

Crowson, R.A., 1960. The phylogeny of Coleoptera. Annu. Rev. Entomol. 5, 111–134.

Dayhoff, M.O., Schwartz, R.M., Orcutt, B.C., 1978. A model of evolutionary change in proteins. Atlas Protein Seq. Struct. 5, 345–351.

Dell'Ampio, E., Meusemann, K., Szucsich, N.U., Peters, R.S., Meyer, B., Borner, J., Petersen, M., Aberer, A.J., Stamatakis, A., Walzl, M.G., Minh, B.Q., von Haeseler, A., Ebersberger, I., Pass, G., Misof, B., 2014. Decisive data sets in phylogenomics: lessons from studies on the phylogenetic relationships of primarily wingless insects. Mol. Biol. Evol. 31, 239–249.

Dettner, K., 1985. Ecological and phylogenetic significance of defensive compounds from pygidial glands of Hydradephaga (Coleoptera). Proc. Acad. Nat. Sci. Philadelphia 137, 156–171.

Doyle, V.P., Young, R.E., Naylor, G.J.P., Brown, J.M., 2015. Can we identify genes with increased phylogenetic reliability? Syst. Biol. 64, 824–837.

Dressler, C., Beutel, R.G., 2010. The morphology and evolution of the adult head of Adephaga (Insecta: Coleoptera). Arthropod Syst. Phylogeny 68, 239–287.

Dressler, C., Ge, S.Q., Beutel, R.G., 2011. Is *Meru* a specialized noterid (Coleoptera, Adephaga)? Syst. Entomol. 36, 705–712.

Felsenstein, J., 1985. Confidence limits on phylogenies: an approach using the bootstrap. Evolution 39, 783–791.

Fernandez, R., Edgecombe, G.D., Giribet, G., 2016. Exploring phylogenetic relationships within Myriapoda and the effects of matrix composition and occupancy on phylogenomic reconstruction. Syst. Biol. 65, 871–889.

Feuda, R., Dohrmann, M., Pett, W., Philippe, H., Rota-Stabelli, O., Lartillot, N., Wörheide, G., Pisani, D., 2017. Improved modeling of compositional heterogeneity supports sponges as sister to all other animals. Curr. Biol. 27, 3864–3870.

Forsyth, D.J., 1970. The structure of the defence glands of the Cicindelidae, Amphizoidae, and Hygrobiidae (Insecta: Coleoptera). J. Zool. 160, 51–69.

Guindon, S., Dufayard, J.F., Lefort, V., Anisimova, M., Hordijk, W., Gascuel, O., 2010. New algorithms and methods to estimate maximum-likelihood phylogenies: assessing the performance of PhyML 3.0. Syst. Biol. 59, 307–321.

Hawlitschek, O., Hendrich, L., Balke, M., 2012. Molecular phylogeny of the squeak beetles, a family with disjunct Palearctic-Australian range. Mol. Phylogenet. Evol. 62, 550–554.

Hoang, D.T., Chernomor, O., von Haeseler, A., Minh, B.Q., Vinh, L.S., 2017. UFBoot2: improving the ultrafast bootstrap approximation. Mol. Biol. Evol. 35, 518–522.

Jäch, M.A., Balke, M., 2008. Global diversity of water beetles (Coleoptera) in freshwater. Hydrobiologia 595, 419–442.

Jarvis, E.D., Mirarab, S., Aberer, A.J., Li, B., Houde, P., Li, C., Ho, S.Y.W., Faircloth, B.C., Nabholz, B., Howard, J.T., Suh, A., Weber, C.C., da Fonseca, R.R., Li, J., Zhang, F., Li, H., Zhou, L., Narula, N., Liu, L., Ganapathy, G., Boussau, B., Bayzid, M.S., Zavidovych, V., Subramanian, S., Gabaldón, T., Capella-Gutiérrez, S., Huerta-Cepas, J., Rekepalli, B., Munch, K., Schierup, M., Lindow, B., Warren, W.C., Ray, D., Green, R.E., Bruford, M.W., Zhan, X., Dixon, A., Li, S., Li, N., Huang, Y., Derryberry, E.P., Bertelsen, M.F., Sheldon, F.H., Brumfield, R.T., Mello, C.V., Lovell, P.V., Wirthlin, M., Schneider, M.P.C., Prosdocimi, F., Samaniego, J.A., Velazquez, A.M.V., Alfaro-Núñez, A., Campos, P.F., Petersen, B., Sicheritz-Ponten, T., Pas, A., Bailey, T., Scofield, P., Bunce, M., Lambert, D.M., Zhou, Q., Perelman, P., Driskell, A.C., Shapiro, B., Xiong, Z., Zeng, Y., Liu, S., Li, Z., Liu, B., Wu, K., Xiao, J., Yinqi, X., Zheng, Q., Zhang, Y., Yang, H.,

Wang, J., Smeds, L., Rheindt, F.E., Braun, M., Fjeldsa, J., Orlando, L., Barker, F.K., Jønsson, K.A., Johnson, W., Koepfli, K.-P., O'Brien, S., Haussler, D., Ryder, O.A., Rahbek, C., Willerslev, E., Graves, G.R., Glenn, T.C., McCormack, J., Burt, D., Ellegren, H., Alström, P., Edwards, S.V., Stamatakis, A., Mindell, D.P., Cracraft, J., Braun, E.L., Warnow, T., Jun, W., Gilbert, M.T.P., Zhang, G., 2014. Whole-genome analyses resolve early branches in the tree of life of modern birds. Science 346, 1320–1331.

Jeffroy, O., Brinkmann, H., Delsuc, F., Philippe, H., 2006. Phylogenomics: the beginning of incongruence? Trends Genet. 22, 225–231.

Jermiin, L.S., Ho, S.Y.W., Ababneh, F., Robinson, J., Larkum, A.W.D., 2004. The biasing effect of compositional heterogeneity on phylogenetic estimates may be underestimated. Syst. Biol. 53, 638–643.

Jermiin, L.S., Jayaswal, V., Ababneh, F., Robinson, J., 2008. Phylogenetic model evaluation. In: Keith, J. (Ed.), Bioinformatics, data, sequences analysis and evolution, vol. I. Humana Press, Totowa, pp. 331–363.

Jones, D.T., Taylor, W.R., Thornton, J.M., 1992. The rapid generation of mutation data matrices from protein sequences. Comput. Appl. Biosci. 8, 275–282.

Kalyaanamoorthy, S., Minh, B.Q., Wong, T.K.F., von Haeseler, A., Jermiin, L.S., 2017. ModelFinder: fast model selection for accurate phylogenetic estimates. Nat. Methods 14, 587–589.

Katoh, K., Standley, D.M., 2013. MAFFT multiple sequence alignment software version 7: improvements in performance and usability. Mol. Biol. Evol. 30, 772–780.

Kavanaugh, D.H., 1986. A systematic review of Amphizoid beetles (Amphizoidae: Coleoptera) and their phylogenetic relationships to other Adephaga. Proc. Calif. Acad. Sci. 44, 67–109.

Klopfstein, S., Massingham, T., Goldman, N., 2017. More on the best evolutionary rate for phylogenetic analysis. Syst. Biol. 66, 769–785.

Kosiol, C., Goldman, N., 2005. Different versions of the dayhoff rate matrix. Mol. Biol. Evol. 22, 193–199.

Kubatko, L.S., Degnan, J.H., 2007. Inconsistency of phylogenetic estimates from concatenated data under coalescence. Syst. Biol. 56, 17–24.

Kück, P., Longo, G.C., 2014. FASconCAT-G: extensive functions for multiple sequence alignment preparations concerning phylogenetic studies. Front. Zool. 11, 81.

Kück, P., Meusemann, K., Dambach, J., Thormann, B., von Reumont, B.M., Wägele, J.W., Misof, B., 2010. Parametric and non-parametric masking of randomness in sequence alignments can be improved and leads to better resolved trees. Front. Zool. 7, 10.

Kück, P., Struck, T.H., 2014. BaCoCa – a heuristic software tool for the parallel assessment of sequence biases in hundreds of gene and taxon partitions. Mol. Phylogenet. Evol. 70, 94–98.

Lanfear, R., Calcott, B., Ho, S.Y.W., Guindon, S., 2012. PartitionFinder: combined selection of partitioning schemes and substitution models for phylogenetic analyses. Mol. Biol. Evol. 29, 1695–1701.

Lanfear, R., Calcott, B., Kainer, D., Mayer, C., Stamatakis, A., 2014. Selecting optimal partitioning schemes for phylogenomic datasets. BMC Evol. Biol. 14, 82.

Lanfear, R., Frandsen, P.B., Wright, A.M., Senfeld, T., Calcott, B., 2017. Partitionfinder 2: new methods for selecting partitioned models of evolution for molecular and morphological phylogenetic analyses. Mol. Biol. Evol. 34, 772–773.

Lartillot, N., Brinkmann, H., Philippe, H., 2007. Suppression of long-branch attraction artefacts in the animal phylogeny using a site-heterogeneous model. BMC Evol. Biol. 7, S4.

Le, S.Q., Dang, C.C., Gascuel, O., 2012. Modeling protein evolution with several amino acid replacement matrices depending on site rates. Mol. Biol. Evol. 29, 2921–2936.

Le, S.Q., Gascuel, O., 2008. An improved general amino acid replacement matrix. Mol. Biol. Evol. 25, 1307–1320.

Lemoine, F., Domelevo Entfellner, J.B., Wilkinson, E., Correia, D., Dávila Felipe, M., De Oliveira, T., Gascuel, O., 2018. Renewing Felsenstein's phylogenetic bootstrap in the era of big data. Nature 556, 452–456.

López-López, A., Vogler, A.P., 2017. The mitogenome phylogeny of Adephaga (Coleoptera). Mol. Phylogenet. Evol. 114, 166–174.

McKenna, D.D., Wild, A.L., Kanda, K., Bellamy, C.L., Beutel, R.G., Caterino, M.S., Farnum, C.W., Hawks, D.C., Ivie, M.A., Jameson, M.L., Leschen, R.A.B., Marvaldi, A.E., Mchugh, J.V., Newton, A.F., Robertson, J.A., Thayer, M.K., Whiting, M.F., Lawrence, J.F., Ślipiński, A., Maddison, D.R., Farrell, B.D., 2015.

The beetle tree of life reveals that Coleoptera survived end-Permian mass extinction to diversify during the Cretaceous terrestrial revolution. Syst. Entomol. 40, 835–880.

Meiklejohn, K.A., Faircloth, B.C., Glenn, T.C., Kimball, R.T., Braun, E.L., 2016. Analysis of a rapid evolutionary radiation using ultraconserved elements: evidence for a bias in some multispecies coalescent methods. Syst. Biol. 65, 612–627.

Miller, K.B., Bergsten, J., 2016. Diving beetles of the world. Systematics and biology of the Dytiscidae. Johns Hopkins University Press, Baltimore.

Minh, B.Q., Nguyen, M.A.T., Von Haeseler, A., 2013. Ultrafast approximation for phylogenetic bootstrap. Mol. Biol. Evol. 30, 1188–1195.

Mirarab, S., Bayzid, M.S., Warnow, T., 2016. Evaluating summary methods for multilocus species tree estimation in the presence of incomplete lineage sorting. Syst. Biol. 65, 366–380.

Mirarab, S., Warnow, T., 2015. ASTRAL-II: coalescent-based species tree estimation with many hundreds of taxa and thousands of genes. Bioinformatics 31, i44–i52.

Misof, B., Liu, S., Meusemann, K., Peters, R.S., Donath, A., Mayer, C., Frandsen, P.B., Ware, J., Flouri, T., Beutel, R.G., Niehuis, O., Petersen, M., Izquierdo-Carrasco, F., Wappler, T., Rust, J., Aberer, A.J., Aspöck, U., Aspöck, H., Bartel, D., Blanke, A., Berger, S., Bohm, A., Buckley, T.R., Calcott, B., Chen, J., Friedrich, F., Fukui, M., Fujita, M., Greve, C., Grobe, P., Gu, S., Huang, Y., Jermiin, L.S., Kawahara, A.Y., Krogmann, L., Kubiak, M., Lanfear, R., Letsch, H., Li, Y., Li, Z., Li, J., Lu, H., Machida, R., Mashimo, Y., Kapli, P., McKenna, D.D., Meng, G., Nakagaki, Y.,

Navarrete-Heredia, J.L., Ott, M., Ou, Y., Pass, G., Podsiadlowski, L., Pohl, H., von Reumont, B.M., Schutte, K., Sekiya, K., Shimizu, S., Slipinski, A., Stamatakis, A., Song, W., Su, X., Szucsich, N.U., Tan, M., Tan, X., Tang, M., Tang, J., Timelthaler, G., Tomizuka, S., Trautwein, M., Tong, X., Uchifune, T., Walzl, M.G., Wiegmann, B.M., Wilbrandt, J., Wipfler, B., Wong, T.K.F., Wu, Q., Wu, G., Xie, Y., Yang, S., Yang, Q., Yeates, D.K., Yoshizawa, K., Zhang, Q., Zhang, R., Zhang, W., Zhang, Y., Zhao, J., Zhou, C., Zhou, L., Ziesmann, T., Zou, S., Xu, X., Yang, H., Wang, J., Kjer, K.M., Zhou, X., 2014. Phylogenomics resolves the timing and pattern of insect evolution. Science 346, 763–767.

Misof, B., Meyer, B., von Reumont, B.M., Kück, P., Misof, K., Meusemann, K., 2013. Selecting informative subsets of sparse supermatrices increases the chance to find correct trees. BMC Bioinformatics 14, 348.

Misof, B., Misof, K., 2009. A Monte Carlo approach successfully identifies randomness in multiple sequence alignments: a more objective means of data exclusion. Syst. Biol. 58, 21–34.

Nabhan, A.R., Sarkar, I.N., 2012. The impact of taxon sampling on phylogenetic inference: a review of two decades of controversy. Brief. Bioinform. 13, 122–134.

Nesnidal, M.P., Helmkampf, M., Bruchhaus, I., Hausdorf, B., 2010. Compositional heterogeneity and phylogenomic inference of metazoan relationships. Mol. Biol. Evol. 27, 2095–2104.

Nesnidal, M.P., Helmkampf, M., Meyer, A., Witek, A., Bruchhaus, I., Ebersberger, I., Hankeln, T., Lieb, B., Struck, T.H., Hausdorf, B., 2013. New phylogenomic data support the monophyly of Lophophorata and an Ectoproct-Phoronid clade and

indicate that Polyzoa and Kryptrochozoa are caused by systematic bias. BMC Evol. Biol. 13, 253.

Nguyen, L.T., Schmidt, H.A., Von Haeseler, A., Minh, B.Q., 2015. IQ-TREE: a fast and effective stochastic algorithm for estimating maximum-likelihood phylogenies. Mol. Biol. Evol. 32, 268–274.

Niehuis, O., Hartig, G., Grath, S., Pohl, H., Lehmann, J., Tafer, H., Donath, A., Krauss, V., Eisenhardt, C., Hertel, J., Petersen, M., Mayer, C., Meusemann, K., Peters, R.S., Stadler, P.F., Beutel, R.G., Bornberg-Bauer, E., McKenna, D.D., Misof, B., 2012. Genomic and morphological evidence converge to resolve the enigma of Strepsiptera. Curr. Biol. 22, 1309–1313.

Nilsson, A.N., 2011. A world catalogue of the family Noteridae, or the burrowing water beetles (Coleoptera, Adephaga). Version 16.VIII.2011. Distributed as a PDF file via Internet. Available from:<http://www.waterbeetles.eu>(accessed 30 June 2018) [WWW Document].

Nilsson, A.N., Hájek, J., 2019. A world catalogue of the family Dytiscidae, or the diving beetles (Coleoptera, Adephaga). Version 1.I.2019. Distributed as a PDF file via Internet. Available from:<http://www.waterbeetles.eu>(accessed 07 February 2019) [WWW Document].

Pauli, T., Burt, T.O., Meusemann, K., Bayless, K., Donath, A., Podsiadlowski, L., Mayer, C., Kozlov, A., Vasilikopoulos, A., Liu, S., Zhou, X., Yeates, D., Misof, B., Peters, R.S., Mengual, X., 2018. New data, same story: phylogenomics does not support Syrphoidea (Diptera: Syrphidae, Pipunculidae). Syst. Entomol. 43, 447–459.

Peters, R.S., Krogmann, L., Mayer, C., Donath, A., Gunkel, S., Meusemann, K., Kozlov, A., Podsiadlowski, L., Petersen, M., Lanfear, R., Diez, P.A., Heraty, J., Kjer, K.M., Klopfstein, S., Meier, R., Polidori, C., Schmitt, T., Liu, S., Zhou, X., Wappler, T., Rust, J., Misof, B., Niehuis, O., 2017. Evolutionary history of the Hymenoptera. Curr. Biol. 27, 1013–1018.

Petersen, M., Meusemann, K., Donath, A., Dowling, D., Liu, S., Peters, R.S., Podsiadlowski, L., Vasilikopoulos, A., Zhou, X., Misof, B., Niehuis, O., 2017. Orthograph: a versatile tool for mapping coding nucleotide sequences to clusters of orthologous genes. BMC Bioinformatics 18, 111.

Philippe, H., Brinkmann, H., Lavrov, D.V., Littlewood, D.T.J., Manuel, M., Wörheide, G., Baurain, D., 2011. Resolving difficult phylogenetic questions: why more sequences are not enough. PLoS Biol. 9, e1000602.

Philippe, H., Roure, B., 2011. Difficult phylogenetic questions: more data, maybe; better methods, certainly. BMC Biol. 9, 91.

Ponomarenko, A.G., 1993. Two new species of Mesozoic dytiscoid beetles from Asia. Paleont. J. 27, 182–191.

Quang, L.S., Gascuel, O., Lartillot, N., 2008. Empirical profile mixture models for phylogenetic reconstruction. Bioinformatics 24, 2317–2323.

Reddy, S., Kimball, R.T., Pandey, A., Hosner, P.A., Braun, M.J., Hackett, S.J., Han, K.-L., Harshman, J., Huddleston, C.J., Kingston, S., Marks, B.D., Miglia, K.J., Moore, W.S., Sheldon, F.H., Witt, C.C., Yuri, T., Braun, E.L., 2017. Why do phylogenomic data sets yield conflicting trees? Data type influences the avian tree of life more than taxon sampling. Syst. Biol. 66, 857–879.

Ribera, I., Beutel, R.G., Balke, M., Vogler, A., 2002a. Discovery of Aspidytidae, a new family of aquatic Coleoptera. Proc. R. Soc. B Biol. Sci. 269, 2351–2356.

Ribera, I., Hogan, J.R., Vogler, A.P., 2002b. Phylogeny of hydradephagan water beetles inferred from 18S rRNA sequences. Mol. Phylogenet. Evol. 23, 43–62.

Romiguier, J., Cameron, S.A., Woodard, S.H., Fischman, B.J., Keller, L., Praz, C.J., 2016. Phylogenomics controlling for base compositional bias reveals a single origin of eusociality in corbiculate bees. Mol. Biol. Evol. 33, 670–678.

Sann, M., Niehuis, O., Peters, R.S., Mayer, C., Kozlov, A., Podsiadlowski, L., Bank, S., Meusemann, K., Misof, B., Bleidorn, C., Ohl, M., 2018. Phylogenomic analysis of Apoidea sheds new light on the sister group of bees. BMC Evol. Biol. 18, 71.

Sayyari, E., Mirarab, S., 2016. Fast coalescent-based computation of local branch support from quartet frequencies. Mol. Biol. Evol. 33, 1654–1668.

Sayyari, E., Whitfield, J.B., Mirarab, S., 2017. Fragmentary gene sequences negatively impact gene tree and species tree reconstruction. Mol. Biol. Evol. 34, 3279–3291.

Spangler, P.J., Steiner, W.E., 2005. A new aquatic beetle family, Meruidae, from Venezuela (Coleoptera: Adephaga). Syst. Entomol. 30, 339–357.

Stamatakis, A., 2014. RAxML version 8: a tool for phylogenetic analysis and post-analysis of large phylogenies. Bioinformatics 30, 1312–1313.

Strimmer, K., von Haeseler, A., 1997. Likelihood-mapping: a simple method to visualize phylogenetic content of a sequence alignment. Proc. Natl. Acad. Sci. U. S. A. 94, 6815–6819.

Suh, A., 2016. The phylogenomic forest of bird trees contains a hard polytomy at the root of Neoaves. Zool. Scr. 45, 50–62.

Suyama, M., Torrents, D., Bork, P., 2006. PAL2NAL: robust conversion of protein sequence alignments into the corresponding codon alignments. Nucleic Acids Res. 34, 609–612.

Tan, G., Muffato, M., Ledergerber, C., Herrero, J., Goldman, N., Gil, M., Dessimoz, C., 2015. Current methods for automated filtering of multiple sequence alignments frequently worsen single-gene phylogenetic inference. Syst. Biol. 64, 778–791.

Tavaré, S., 1986. Some probabilistic and statistical problems in the analysis of DNA sequences. Lect. Math. Life Sci. 17, 57–86.

Timmermans, M.J.T.N., Barton, C., Haran, J., Ahrens, D., Culverwell, C.L., Ollikainen, A., Dodsworth, S., Foster, P.G., Bocak, L., Vogler, A.P., 2016. Family-level sampling of mitochondrial genomes in Coleoptera: Compositional heterogeneity and phylogenetics. Genome Biol. Evol. 8, 161–175.

Toussaint, E.F.A., Beutel, R.G., Morinière, J., Jia, F., Xu, S., Michat, M.C., Zhou, X., Bilton, D.T., Ribera, I., Hájek, J., Balke, M., 2015. Molecular phylogeny of the highly disjunct cliff water beetles from South Africa and China (Coleoptera: Aspidytidae). Zool. J. Linn. Soc. 176, 537–546.

Wang, H.C., Minh, B.Q., Susko, E., Roger, A.J., 2017. Modeling site heterogeneity with posterior mean site frequency profiles accelerates accurate phylogenomic estimation. Syst. Biol. 67, 216–235.

Whelan, S., Goldman, N., 2001. A general empirical model of protein evolution derived from multiple protein families using a maximum-likelihood approach. Mol. Biol. Evol. 18, 691–699.

Whitfield, J.B., Kjer, K.M., 2008. Ancient rapid radiations of insects: challenges for phylogenetic analysis. Annu. Rev. Entomol. 53, 449–472.

Xia, X., Xie, Z., Salemi, M., Chen, L., Wang, Y., 2003. An index of substitution saturation and its application. Mol. Phylogenet. Evol. 26, 1–7.

Xu, B., Yang, Z., 2016. Challenges in species tree estimation under the multispecies coalescent model. Genetics 204, 1353–1368.

Yang, Z., 1998. On the best evolutionary rate for phylogenetic analysis. Syst. Biol. 47, 125–133.

Zdobnov, E.M., Tegenfeldt, F., Kuznetsov, D., Waterhouse, R.M., Simao, F.A., Ioannidis, P., Seppey, M., Loetscher, A., Kriventseva, E.V., 2017. OrthoDB v9.1: cataloging evolutionary and functional annotations for animal, fungal, plant, archaeal, bacterial and viral orthologs. Nucleic Acids Res. 45, D744–D749.

Zhang, C., Rabiee, M., Sayyari, E., Mirarab, S., 2018. ASTRAL-III: polynomial time species tree reconstruction from partially resolved gene trees. BMC Bioinformatics 19, 153.

Zhang, S.Q., Che, L.-H., Li, Y., Liang, Dan, Pang, H., Ślipiński, A., Zhang, P., 2018. Evolutionary history of Coleoptera revealed by extensive sampling of genes and species. Nat. Commun. 9, 205.

Zhong, M., Hansen, B., Nesnidal, M., Golombek, A., Halanych, K.M., Struck, T.H., 2011. Detecting the symplesiomorphy trap: a multigene phylogenetic analysis of terebelliform annelids. BMC Evol. Biol. 11, 369.

# 4. Phylogenomic analyses of taxon-rich datasets consolidate the evolution of Adephaga (Coleoptera) and highlight biases due to model misspecification and excessive data trimming

This chapter is intended for publication in Systematic Entomology.

List of authors:

Alexandros Vasilikopoulos, Michael Balke, Oliver Niehuis, Sandra Kukowka, James M. Pflug, Sebastian Martin, Karen Meusemann, Lars Hendrich, Alexander Donath, Christoph Mayer, David R. Maddison, Rolf G. Beutel, Bernhard Misof

Authors' contributions:

AV, MB, ON, RGB and BM conceived and designed the study. AV, MB and RGB wrote the manuscript draft with AV taking the lead. MB, ON, RGB, DRM and BM contributed to funding acquisition. MB and LH collected and provided insect specimens. AV designed the set of DNA-hybridization baits *in silico*. AV and SK performed all molecular laboratory experiments for hybrid enrichment. JMP and DRM provided new transcriptomic data. AV and JMP performed assembly and cross-contamination checks for transcriptomes. SM, JMP and AD performed NCBI sequence submissions. AV performed assembly, contamination checks and further processing of hybrid-enrichment data. AV and KM performed phylogenetic analyses. CM provided bioinformatic methods for outlier sequence detection.

## 4.1. Introduction

Beetles (Coleoptera) are the most speciose group of animals and their phylogeny has been the focus of attention for many decades (e.g., Beutel et al., 2020, 2019a; Crowson, 1960; Hunt et al., 2007; Lawrence et al., 2011; Lawrence and Newton, 1982; McKenna et al., 2019). Polyphaga is the largest beetle suborder with predominantly phytophagous species whereas Adephaga, which mostly includes predatory species, is the second largest beetle suborder with more than 45,000 species assigned into 11 families (Beutel et al., 2020; Duran and Gough, 2020). The family-level phylogenetic relationships of Adephaga have been extensively debated but scientists are now reaching a consensus on the most likely scenario of their evolution (Beutel et al., 2020; Gustafson et al., 2020; McKenna et al., 2019). Despite this, open questions remain, such as the phylogenetic relationships of terrestrial families, some relationships within Dytiscoidea and the intra-familial relationships within the species-rich families Carabidae, Cicindelidae, Dytiscidae and Gyrinidae (Beutel et al., 2020; Gustafson et al., 2020; Michat et al., 2017; Vasilikopoulos et al., 2019). In addition, previous analyses of familial relationships of Adephaga have suggested that the results of previous studies might be artifacts due to phylogenetic errors (Cai et al., 2020). In this study, we address these unresolved issues by combining newly generated exon-capture sequence data with transcriptomes to infer the phylogeny of Adephaga based on extensive sampling of species.

The majority of species diversity in Adephaga belong to the terrestrial families Carabidae (ground beetles, >35,000 extant species), Cicindelidae (tiger beetles, >2400 spp.) and Trachypachidae (6 spp.) (Beutel et al., 2020; Duran and Gough, 2020;

Lorenz, 2020). The terrestrial families of Adephaga have been collectively referred to as "Geadephaga" (Crowson, 1960). The monophyly of Geadephaga has been disputed in the past based on analyses of morphological characters (e.g., Burmeister 1976; Beutel and Roughley 1988), but most recent morphological and molecular analyses suggest a single origin of the terrestrial families (Beutel et al., 2020; Gustafson et al., 2020; Maddison et al., 2009; McKenna et al., 2019). Despite this, the phylogenetic relationships among Carabidae, Cicindelidae and Trachypachidae remain obscure, mainly because different phylogenomic analyses have produced different topologies for the relationships of these groups. Specifically, phylotranscriptomic analyses have concluded that Trachypachidae are sister to a clade of Cicindelidae + Carabidae (McKenna et al., 2019). In contrast, analyses of mitochondrial genomes suggested a weakly supported clade of Trachypachidae + Cicindelidae as sister to Carabidae (López-López and Vogler, 2017), while analyses of ultraconserved elements (UCEs) suggested a clade of Cicindelidae + (Trachypachidae + Carabidae) (Gustafson et al. 2020). It should be noted, however, that the taxon sampling of previous phylogenomic studies was not sufficient to test the monophyly of Carabidae and Cicindelidae and to robustly infer the phylogenetic position of the small family Trachypachidae (e.g., Zhang et al. 2018b; McKenna et al. 2019; Gough et al. 2020; Gustafson et al. 2020). In addition, the results of some molecular analyses do not agree with results of morphological studies that suggest Trachypachidae as sister to Carabidae + Cicindelidae (e.g., Beutel et al. 2020). Therefore a re-evaluation of the relationships of Geadephaga with a careful examination of potential sources of systematic error and increased species sampling is needed.

The species of the remaining eight families of Adephaga (Amphizoidae, Aspidytidae, Dytiscidae, Haliplidae, Hygrobiidae, Meruidae, Noteridae, Gyrinidae) occur primarily in aquatic or semi-aquatic habitats (Jäch and Balke, 2008; Short, 2018). Most species of Gyrinidae, Dytiscidae, Hygrobiidae and Noteridae are strictly aquatic, while members of Amphizoidae, Aspidytidae and Meruidae occur in hygropetric or semi-aquatic habitats (Balke et al., 2003; Kavanaugh, 1986; Spangler and Steiner, 2005; Vasilikopoulos et al., 2019). Crowson (1960) suggested that all these groups constitute a monophylum to which he referred to as "Hydradephaga". Only a few molecular phylogenetic studies have supported monophyletic "Hydradephaga" (López-López and Vogler, 2017; McKenna et al., 2015; Shull et al., 2001), whereas its monophyly has been refuted in more comprehensive studies based on analyses of morphological characters and phylogenomics (e.g., Beutel and Roughley 1988; Baca et al. 2017; Gustafson et al. 2019; McKenna et al. 2019; Beutel et al. 2020). More specifically, the placement of Gyrinidae as sister to all other Adephaga is currently a well-accepted scenario (e.g., Baca et al., 2017; Beutel et al., 2020; Gustafson et al., 2020 but see Freitas et al., 2020). In addition, most analyses suggest a sister group relationship of Haliplidae to the superfamily Dytiscoidea (which includes Amphizoidae, Aspidytidae, Dytiscidae, Hygrobiidae, Meruidae, Noteridae) and a clade Meruidae + Noteridae as sister to all remaining families of Dytiscoidea (Baca et al., 2017; Beutel et al., 2006; Gustafson et al., 2020; Vasilikopoulos et al., 2019). Despite this, the phylogenetic position of the family Hygrobiidae (squeak beetles) within Dytiscoidea remains contentious (Baca et al., 2017; Cai et al., 2020; Gustafson et al., 2020; Toussaint et al., 2016; Vasilikopoulos et al., 2020, 2019). Model misspecification, compositional biases, incomplete lineage sorting (ILS), fast evolving alignment sites

and deficient taxon sampling are among the factors that have been proposed to affect the internal phylogeny of Dytiscoidea including the monophyly of the relictual family Aspidytidae (Cai et al., 2020; Gustafson et al., 2020; Vasilikopoulos et al., 2019). Consequently, a thorough assessment of the phylogenetic relationships within Dytiscoidea in the light of increased taxon sampling of genomic data is pending.

In the last decade, a plethora of hybrid enrichment (or sequence capture) approaches for phylogenomics have been developed (Bragg et al., 2016; Faircloth et al., 2012; Lemmon et al., 2012; Mayer et al., 2016). The ultraconserved element (UCE) approach has been proven useful for inferring phylogenetic relationships both at deep and shallow timescales (Faircloth et al., 2012) and is the only sequence-capture approach that has been applied to infer the phylogeny of Adephaga to date (Baca et al., 2017; Gustafson et al., 2020). There is also an extensive set of available bioinformatic toolkits for processing of UCE data in a consistent and efficient way (Faircloth, 2017, 2016). However, there are several reasons why scientists might want to apply exon-capture or transcriptomic approaches in addition to- or independently of the UCE approach. Firstly, orthology predictions for UCEs are based on the core ultraconserved regions but the analyzed flanking regions might not necessarily be homologous to each other (Bank et al., 2017; Li et al., 2013). Secondly, the extension of the selected UCEs beyond the ultraconserved regions is based on arbitrary length criteria that differ across different experiments or taxonomic clades (e.g., Faircloth, 2017). Thirdly, UCE data can be analyzed only at the nucleotide sequence level because there is usually no information on whether they overlap with coding regions (Bank et al., 2017). Additionally, individual UCE loci may not harbor sufficient information to infer reliable locus-specific phylogenetic trees (Meiklejohn et al., 2016). Lastly, cross-

validation of the results of analyses based on different types of data constitutes the basis for substantiating the conclusions of different studies (Vasilikopoulos et al., 2020). Therefore the exon-capture approach can provide complementary or independent evidence for testing the validity of previously suggested phylogenetic hypotheses of Adephaga.

Ideally, scientists would like to have universal sets of DNA-hybridization baits that capture a large number of orthologous genes across a wide range of species (Glenn and Faircloth, 2016). However, previous research suggests that exon-capture approaches are effective for investigating taxonomic clades characterized by small to moderate levels of molecular divergence (Bi et al., 2012; Bragg et al., 2016; Mayer et al., 2016). The advantage of the exon-capture approach is that the target regions are well defined genomic units for which orthology assignment is more straightforward and facilitates their integration with other types of protein-coding data such as transcriptomes. In addition, protein-coding exons usually undergo some purifying selection on protein structure and this in turn makes multiple sequence alignment (MSA) of these regions more straightforward than UCEs due to the development of accurate translation-based alignment algorithms (Karin et al., 2020). Despite this, the success of the exon-capture approach depends on the availability of transcriptomic or genomic data from closely related species, which can be used as basis for designing baits, and the degree of molecular divergence within the clade of interest (Mayer et al., 2016). Therefore, it has been put forward that the UCE approach should be preferred over exon-capture at deep phylogenetic time-scales because UCEs are more conserved across highly divergent species (Bragg et al., 2016). However, if transcriptomic resources are available for a broad set of species within the clade of interest, they can

be used for testing the applicability of exon-specific DNA-hybridization baits for at deeper phylogenetic scales. Recently developed bioinformatic approaches are able to automatically detect suitable regions for bait design in aligned DNA sequence data, including exonic alignments, by minimizing overall bait-to-target distances (Mayer et al., 2016), therefore offering a promising solution to the problem of designing probes that have broad phylogenetic applicability (Lemmon and Lemmon, 2013). Additionally, transcriptomic and genomic resources for adephagan beetles have increased considerably in the last years (Gustafson et al., 2019; McKenna et al., 2019; Vasilikopoulos et al., 2019). These resources combined with the recently developed bioinformatic approaches make it now possible to test the applicability and efficiency of the exon-capture approach for deep-level phylogenetics in Adephaga.

In this study, we develop a novel set of DNA-hybridization baits specifically tailored to capture hundreds of single-copy genes across adephagan lineages and generate new exon-capture data to infer the phylogeny of Adephaga. We test the efficiency of this new bait set for locus recovery in a large number of specimens from different families of Adephaga and we combine the newly generated exon-capture data with transcriptomes to generate the most taxon-rich phylogenomic dataset for adephagan beetles presented to date. In order to avoid biased estimates of phylogeny of Adephaga we take measures to minimize phylogenetic artifacts by employing realistic evolutionary models and by reducing potentially biasing factors in the data using data-filtering strategies that select conserved alignment sites. We evaluate the effects of model misspecification and excessive data trimming both on the results of phylogenetic tree reconstructions and on quartet-based analyses of phylogenetic signal in an attempt to acquire a more detailed view of resolution, conflict and bias in the backbone

phylogeny of Adephaga. Additionally, we explore whether or not incongruence between concatenation and summary coalescent analyses can possibly be explained by gene-tree errors and we suggest possible strategies for selecting informative genes that minimize these errors and therefore reduce incongruence. Lastly, we discuss our results in the context of the morphological evolution of Adephaga.
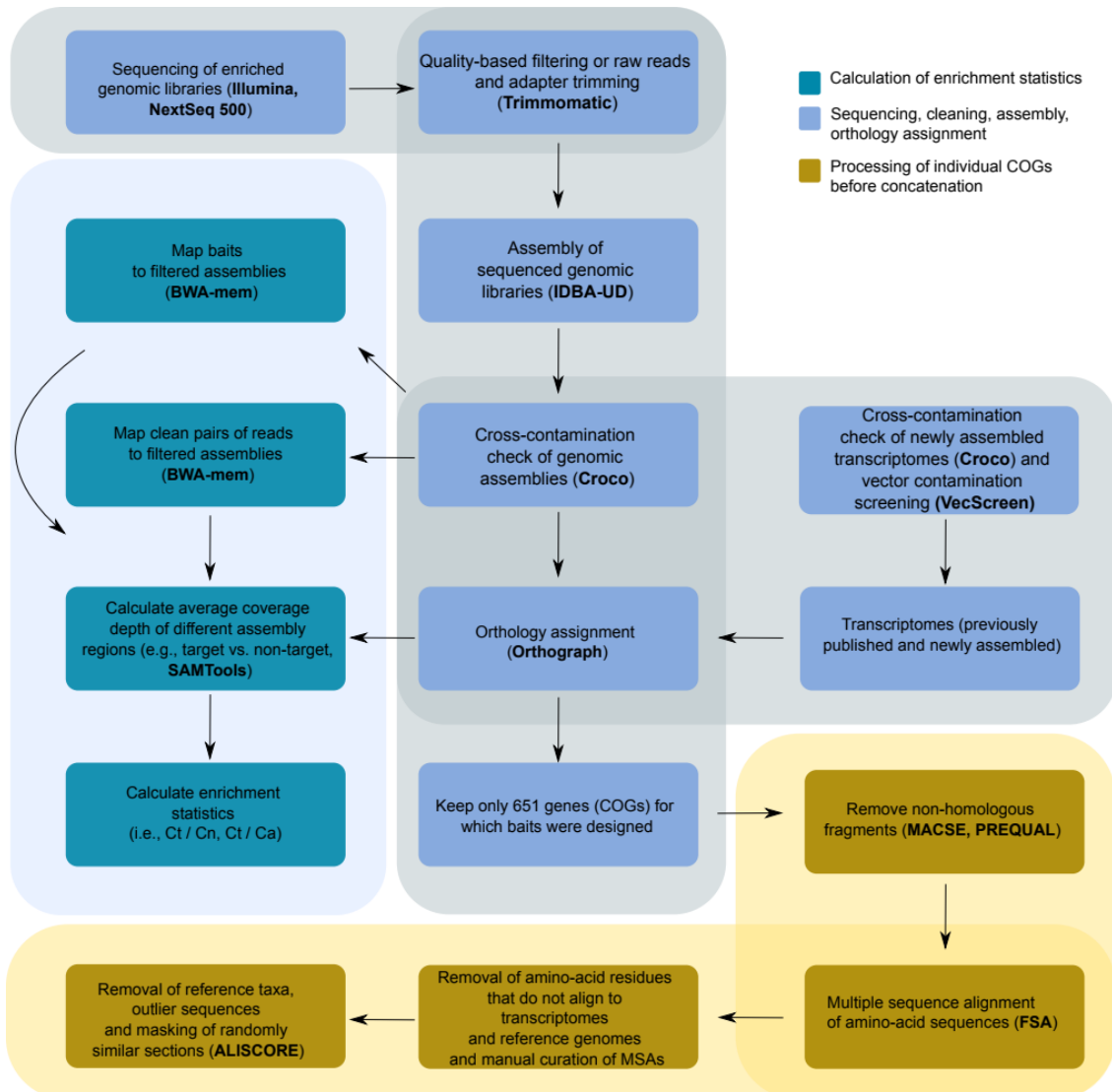
## 4.2. Materials and methods

### 4.2.1. Taxon sampling

We combined 38 transcriptomes from 23 species of Adephaga and 15 outgroup species (S4: Table S1) with newly generated exon-capture sequence data from 95 species of Adephaga (S4: Table S2). In total, our initial taxon sampling comprised data from 118 species of Adephaga representing all families except the monotypic Meruidae and 21 outgroups (two Hymenoptera, three Mecopterida, two Strepsiptera, four Neuropterida, two Myxophaga, two Archostemata, six Polyphaga). The initial taxon sampling includes the six reference species of the ortholog set (see below).

### 4.2.2. Inference of bait sequences for hybrid enrichment of protein-coding exons

We used 24 transcriptomes of Adephaga as a basis to build codon-based nucleotide multiple sequence alignments (MSAs) of orthologous genes and search for MSA regions that are suitable for bait design within Adephaga. The transcriptome of *Metrius contractus* was only used for bait design and was not included in phylogenetic reconstructions (S4: Table S1). First, we used a custom ortholog gene set consisting of 3085 clusters of orthologous and single-copy genes (COGs) at the hierarchical level Holometabola (Vasilikopoulos et al., 2019) to assign orthologous transcripts from each

transcriptome to each COG. Orthology assignment of transcripts to each COG was performed with Orthograph v. 0.6.1 (Petersen et al., 2017). Subsequently, we followed procedures for amino-acid multiple sequence alignment, alignment refinement, outlier sequence removal and removal of reference taxa before generating codon-based nucleotide MSAs (see supplementary information of Misof et al., 2014a for details on these procedures). We then used Baitfisher v. 1.2.7 (Mayer et al., 2016) to screen the codon-based MSAs for regions that are appropriate for bait design within the Adephaga clade (see S5: Supplementary Text 1). We conducted seven different tiling designs experiments, corresponding to different lengths of bait regions, bait offsets, and total number of baits in order to capture as many promising exons as possible while accounting for variable exon length, possibly large amount of missing data or hypervariable regions in some parts of the gene alignments (S4: Table S3). In order to exclude baits targeting multiple genomic regions in adephagan genomes, we filtered the resulting baits (separately for each each tiling design experiment) by blasting them against a draft genome assembly of the beetle *Bembidion* sp. nr. *transversale* (Gustafson et al., 2019, see S5: Supplementary Text 1 for options). We then selected only one bait region per exon in each tiling design experiment: the one that required the minimum amount of baits (Mayer et al., 2016). Subsequently, for those exons that were captured in multiple tiling-design experiments only the longest bait regions among experiments were considered. The last task was accomplished by adding the bait regions from the different experiments (non-redundantly for exons, from longer to shorter bait regions) to a combined file with the baits until the maximum size of ~5.99Mbp of DNA was reached (i.e., max. size of bait sequences for the DNA target enrichment kit that was used: SureSelectXT2 Target Enrichment System, Agilent

Technologies). The last task was performed with custom Perl scripts. In total, we inferred 49,787 120bp-long bait sequences for targeting 923 protein-coding exons from 651 protein-coding genes.



**Fig. 4.1:** Summarized workflow of the steps that were used to sequence, clean, assemble and combine the hybrid capture sequence data with transcriptomes to generate individual COGs. A short workflow for calculating the hybrid-enrichment statistics is also provided.

### 4.2.3. Tissue preservation, total genomic DNA extraction, next-generation sequencing (NGS) library preparation and hybrid enrichment

Most specimens used for enrichment of target genomic DNA (gDNA) were freshly collected and preserved in 96% ethanol but we also used a few dry pinned museum specimens (see S4: Table S4). Total genomic DNA (gDNA) was extracted from 96 specimens of adephagan species (S4: Table S2) using the DNeasy Blood & Tissue Kit (Qiagen, Hilden, Germany) and eluted in 100 µl nuclease-free water. Whenever available voucher material has been deposited at Zoological State Collections, Munich, Germany (tissue or extracted DNA, S4: Table S4). Quality and quantity of the extracted gDNA was assessed with a Fragment Analyzer (Agilent Technologies Inc., Santa Clara, U.S.A.) and a Quantus Fluorometer (Promega, Fitchburg, Wisconsin, U.S.A.). Whenever sufficient amount of extracted DNA was available, we used 100 ng of DNA diluted in 10 µl for fragmentation before library preparation, otherwise less than 100ng were used. First, gDNA was sheared into fragments of 150–400 bp using a Bioruptor Pico sonication device (Diagenode s.a., Seraing, Belgium). Multiple shearing steps were performed for each sample until at least ~90% of fragments was within the desired length threshold. The quality and quantity of the fragmented gDNA was assessed with a Fragment Analyzer at the end of each shearing step. For library preparation, we followed the SureSelectXT2 Target Enrichment System Protocol for Illumina Paired-End Multiplexed Sequencing (Version E1 published in June 2015 by Agilent Technologies Inc.) with some minor modifications (see Bank et al., 2017). Specifically, in the library preparation steps "End Repair", "A-tailing", we reduced the reaction volume specified in Agilent's protocol (pages 43–49 for 100 ng DNA samples) by 50% as described by Bank et al. (2017).

Subsequently, adapter ligation was performed with the NEBNext Quick Ligation Module and the adapters from the NEBNext Multiplex Oligos for Illumina (Dual Index Set1) kit. NGS library PCR was then performed with the NEBNext Multiplex Oligos for Illumina and the NEBNext Q5 HotStart HiFi PCR Master Mix, to dual-index the libraries. Cycles of the NGS library PCR were adjusted as follows (due to the concentration measurements after "A-tailing"): 98 °C for 30 sec., followed by 8–10 cycles of 98 °C for 10 sec. and 65 °C for 75 sec., followed by 5 min. at 65 °C followed by 4 °C until the samples were removed from the thermocycler. Subsequently, all steps of the target DNA enrichment followed the protocol given by Bank et al. (2017) with modifications adjusted to the number of library pools and volume concentrations in our study (see S5: Supplementary Text 1).

### 4.2.4. Sequencing and assembly of the enriched genomic libraries

The enriched genomic libraries for the 95 samples of Adephaga were paired-end sequenced (150bp) on a single flow cell of an Illumina NextSeq 500 sequencer (Illumina Inc., San Diego, CA, U.S.A., Fig. 4.1). Sequenced raw reads per genomic library were trimmed to remove Illumina adapter sequences and low quality reads with Trimmomatic v. 0.38 (Bolger et al., 2014, see S5: Supplementary Text 1 for options). Only full pairs of trimmed reads were used for *de novo* assembly of the enriched genomic libraries (S4: Table S2). *De novo* assembly of each genomic library was performed with the software IDBA-UD v. 1.1.3 (see S5: Supplementary Text 1, Fig. 4.1) that is optimized to assemble genomic data with highly unequal coverage depth (Peng et al., 2012).

**Fig. 4.2**: Box-plots of Ct / Cn ratios inferred separately for each family of Adephaga. The plots were calculated by pooling the ratios for species of the same family into the same box-plot.


### 4.2.5. Calculation of hybrid-enrichment statistics

We calculated the ratio of average per base coverage depth of target regions (Ct) divided by the average coverage depth of the non-target regions (Ct / Cn, S4: Table S2, Fig. 4.2) as an approximate measure of the enrichment success for each genomic library in our analyses. To identify the target regions, we first identified bait-binding regions in each assembled genomic library by mapping the bait sequences to the clean assembly files (i.e., after putative cross-contaminated contigs had been removed) using the

software BWA-mem v. 0.7.17 (Li and Durbin, 2009). Subsequently, we separately mapped the trimmed reads to the assemblies with the same version of BWA-mem. A summarized file with the coverage depth of each assembly position was generated with SAMtools v. 1.7 (Li et al., 2009). We used a custom Python script and the IDs of the contigs that contained orthologous sequence (contigs assigned to any of the 651 target COGs, see below) to calculate the average coverage depth of the bait-binding regions but only on those contigs that contained orthologous sequence (i.e., target regions, Ct, Fig. 4.1). We subsequently calculated the average coverage depth of all remaining regions in the assembly for each genomic library (i.e., non-target regions, Cn). Lastly, we calculated the average coverage depth of the whole assembly for each assembled genomic library (Ca). Any positions with zero coverage were excluded from the above calculations to avoid the inflation of enrichment statistics. We considered the statistics: Ct / Cn and Ct / Ca as approximate measures of the enrichment success for each of the 95 genomic libraries (S4: Table S2, Fig. 4.1). We generated box-plots of these statistics separately for each adephagan family and performed pairwise Mann-Whitney-Wilcoxon tests between families in order to assess whether or not the values for different families were drawn from the same underlying distribution. The pairwise statistical tests were performed in R v. 3.6.3 (S4: Table S5) (R Core Team, 2020).

**Table 4.1**: Summarized statistics and description for each generated and analyzed amino-acid supermatrix (see S6: Fig. S1). Saturation statistics of each supermatrix (adjusted $R^2$ and slope) based on the patristic and *p*-distances are also presented. Saturation of each supermatrix was also measured with the average pairwise lambda score (see text). P.I.: parsimony informative, $C_a$ : Overall alignment completeness scores, IC: information content (MARE), *p* dist: observed pairwise distances, N/A: Not applicable, SHETU: site-heterogeneous unpartitioned, SHOMU: site-homogeneous unpartitioned, SHOMP: site-homogeneous partitioned.[1]Note: analyzed under the Bayesian site-heterogeneous model CAT+GTR+G4 (BSHETU).

| Amino-acid supermatrix ID | No. of species | No. of alignment sites | P.I. sites | Percent. (%) of P.I. sites | Average pairwise λ score | Adjusted R² (SHETU) | Slope (SHETU) | Adjusted R² (SHOMU) | Slope (SHOMU) | Adjusted R² (SHOMP) | $C_a$ |
|---|---|---|---|---|---|---|---|---|---|---|---|
| A | 136 | 200,017 | 104,221 | 52.1% | 0.163 | - | - | - | - | - | 0.504 |
| B | 136 | 49,468 | 21,917 | 44.3% | 0.118 | 0.425 | 0.126 | 0.486 | 0.213 | 0.479 | 0.831 |
| C | 136 | 55,521 | 26,220 | 47.2% | 0.135 | 0.369 | 0.111 | 0.403 | 0.182 | 0.405 | 0.790 |
| D[1] | 136 | 49,797 | 21,401 | 43.0% | 0.116 | 0.451 | 0.133 | 0.512 | 0.226 | N/A | 0.846 |
| D - recoded[1] | 136 | 49,797 | 12,699 | 25.5% | 0.069 | - | - | - | - | - | 0.846 |
| E | 136 | 50,614 | 21,773 | 43.0% | 0.116 | 0.454 | 0.133 | 0.515 | 0.227 | N/A | 0.846 |
| F[1] | 136 | 36,511 | 14,143 | 38.7% | 0.095 | 0.510 | 0.155 | 0.569 | 0.256 | N/A | 0.882 |
| G[1] | 120 | 36,511 | 10,879 | 29.8% | 0.079 | 0.396 | 0.230 | 0.393 | 0.272 | N/A | 0.880 |
| H[1] | 100 | 36,511 | 9658 | 26.5% | 0.074 | 0.570 | 0.247 | 0.575 | 0.306 | N/A | 0.892 |
| I[1] | 136 | 29,361 | 11,711 | 39.9% | 0.104 | 0.418 | 0.135 | 0.480 | 0.225 | N/A | 0.857 |
| J[1] | 136 | 23,442 | 7684 | 32.8% | 0.069 | 0.556 | 0.177 | 0.642 | 0.299 | N/A | 0.911 |

(*Table continues on the next page*)

**Table 4.1 (con.)**: Summarized statistics and description for each generated and analyzed amino-acid supermatrix (see S6: Fig. S1). Saturation statistics of each supermatrix (adjusted $R^2$ and slope) based on the patristic and $p$-distances are also presented. Saturation of each supermatrix was also measured with the average pairwise lambda score (see text). P.I.: parsimony informative, $C_a$ : Overall alignment completeness scores, IC: phylogenetic information content (MARE), $p$ dist: observed pairwise distances, N/A: Not applicable, SHETU: site-heterogeneous unpartitioned, SHOMU: site-homogeneous unpartitioned, SHOMP: site-homogeneous partitioned.

| Amino-acid supermatrix ID | Average $p$-dist | Median pairwise p-value to the Bowker's test | Median pairwise p-value to the Stuart's test | IC | Percent. (%) of pairwise p-values < 0.05. Bowker's test | Percent. (%) of pairwise p-values < 0.05. Stuart's test | Description |
|---|---|---|---|---|---|---|---|
| A | 0.154 | 2.14E-02 | 7.38E-05 | 0.672 | 58.92% | 82.94% | Concatenated supermatrix of masked genes with ALISCORE after partitions with IC=0 had been removed |
| B | 0.111 | 1.07E-01 | 1.15E-02 | 0.620 | 37.44% | 64.07% | Trimmed each gene partition of supermatrix A with BMGE, BLOSUM62, h=0.4, keep only genes with length >= 50 amino-acid sites |
| C | 0.127 | 9.46E-02 | 6.73E-03 | 0.599 | 40.10% | 68.27% | Trimmed each partition of supermatrix A with BMGE, BLOSUM62, h=0.5, keep only genes with length >= 80 amino-acid sites and <= 30% missing data |
| D[1] | 0.109 | 1.26E-01 | 1.19E-02 | N/A | 34.69% | 64.11% | Removed genes that fail symmetry tests (IQ-TREE) from supermatrix A. Subsequently, trimmed resulting supermatrix with BMGE (h=0.5, BLOSUM62) |
| D - recoded[1] | 0.052 | 2.16E-01 | - | N/A | 24.67% | - | Dayhoff-6 recoded version of supermatrix D |
| E | 0.109 | 1.22E-01 | 1.14E-02 | N/A | 35.02% | 64.19% | Trimmed supermatrix A with BMGE, BLOSUM62, h=0.5 |
| F[1] | 0.089 | 1.99E-01 | 4.15E-02 | N/A | 24.98% | 51.94% | Trimmed supermatrix A with BMGE, BLOSUM62, h=0.4 |
| G[1] | 0.074 | 2.27E-01 | 6.99E-02 | N/A | 20.35% | 45.27% | Removed distantly related outgroup species from supermatrix F (see supplementary information) |
| H[1] | 0.070 | 2.34E-01 | 8.53E-02 | N/A | 18.85% | 41.92% | Removed fast evolving ingroup species (20 ingroup species with highest LB scores) from supermatrix G (see supplementary information) |
| I[1] | 0.098 | 1.75E-01 | 4.51E-02 | N/A | 25.59% | 50.94% | Remove 50% of genes with the highest RCFV value from matrix A. Trim resulting supermatrix with BMGE, BLOSUM62, h=0.5 |
| J[1] | 0.065 | 2.96E-01 | 1.51E-01 | N/A | 13.97% | 35.21% | Trimmed supermatrix A with BMGE, BLOSUM62, h=0.3 |

### 4.2.6. Cross-contamination checks and orthology assignment

Putative cross-contaminated sequences or sequences of ambiguous origin within the assembled genomic sequence-capture data were identified with the software package CroCo v. 1.1 (Simion et al., 2018). CroCo is primarily designed to screen RNA-seq data but can also potentially identify cross-contaminants from genomic data based on the assumption that the coverage of the contaminated differs between the source library of contamination and the contaminated library respectively (see Simion et al., 2018 and manual of CroCo for details and also Mayer et al., 2016 for a similar approach). We considered contigs that were 99% similar over a fragment of 200 nucleotides as suspicious for cross-contamination (--tool K and otherwise default options). Contigs that were identified as putative contaminants as well as those of ambiguous origin were deleted from the assemblies before downstream analyses (see S4: Table S6 and S5: Supplementary Text 1 for cross-contamination checks applied for some of the transcriptomes).

Orthology assignment of genomic fragments to each of the COGs of the ortholog set was performed with Orthograph v. 0.6.3 (Petersen et al., 2017). From the 3085 COGs of the ortholog set, we conservatively chose to analyze only the 651 COGs for which we had originally designed baits (S4: Tables S1, S2). Orthograph-reporter was run with the "protein2dna" exonerate model for all hybrid capture data (S4: Table S2), whereas the default "protein2genome" model was used for all transcriptomes in the dataset (S4: Table S1, see S5: Supplementary Text 1 for additional options).

**Table 4.2**: Detailed results of the four-cluster likelihood mapping analyses for the two examined phylogenetic hypotheses. Results (i.e., percentages) are shown only for the fully resolved quartets (i.e., quartets falling within the corner areas of the triangular Vonoroi diagrams, see Strimmer and von Haeseler, 1997). Amp.: Amphizoidae, Asp.: Aspidytidae, Hyg.: Hygrobiidae, Dyt.: Dytiscidae, Rem.: Remaining species, Cici.: Cicindelidae, Cara.: Carabidae, Tr.: Trachypachidae, SHETU: site-heterogeneous unpartitioned, SHOMU: site-homogeneous unpartitioned.

| | SHETU Model (original data) | | | | SHETU Model (permuted data) | | | | SHOMU Model (original data) | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | Given topology supermatrix D | Alternative topology 1 supermatrix D | Alternative topology 2 supermatrix D | Total resolved quartets (%) | Given topology supermatrix D | Alternative topology 1 supermatrix D | Alternative topology 2 supermatrix D | Total resolved quartets (%) | Given topology supermatrix D | Alternative topology 1 supermatrix D | Alternative topology 2 supermatrix D | Total resolved quartets (%) |
| **Hypo1 (25,296 quartets)** | (Hyg.+Amp.+ Asp.), (Dyt.+ Rem.) | (Hyg.+ Rem.), (Dyt.+ Amp.+Asp.) | (Hyg. + Dyt..), (Rem. + Amp.+Asp.) | | (Hyg.+ Amp. +Asp.), (Dyt.+ Rem.) | (Hyg.+ Rem.), (Dyt.+Amp. +Asp.) | (Hyg. + Dyt..), (Rem.+Amp. +Asp.) | | (Hyg.+ Amp.+ Asp.), (Dyt.+ Rem.) | (Hyg.+ Rem.), (Dyt.+ Amp.+ Asp.) | (Hyg. + Dyt..), (Rem.+ Amp. +Asp.) | |
| Supermatrix D | 59.80% | 28.80% | 8.10% | 96.70% | 10.70% | 43.40% | 43.60% | 97.70% | 65.30% | 27.30% | 6.30% | 98.90% |
| Supermatrix E | 58.80% | 29.60% | 8.30% | 96.70% | 8.80% | 39.40% | 48.00% | 96.20% | 64.40% | 27.90% | 6.60% | 98.90% |
| Supermatrix F | 53.70% | 29.00% | 11.70% | 94.40% | 36.60% | 18.40% | 34.50% | 89.50% | 61.50% | 28.20% | 8.60% | 98.30% |
| Supermatrix J | 44.00% | 36.00% | 10.60% | 90.60% | 3.30% | 3.90% | 2.40% | 9.60% | 51.00% | 36.60% | 9.10% | 96.70% |
| | | | | | | | | | | | | |
| **Hypo2 (30,912 quartets)** | (Cici.+ Cara.), (Tr. + Rem.) | (Tr. + Cici.), (Cara.+Rem.) | (Tr.+ Cara.), (Cici.+ Rem.) | | (Cici.+ Cara.), (Tr. + Rem.) | (Tr. + Cici.), (Cara.+ Rem.) | (Tr. + Cara.), (Cici.+ Rem.) | | (Cici.+ Cara.), (Tr. + Rem.) | (Tr. + Cici.), (Cara. + Rem.) | (Tr. + Cara.), (Cici.+ Rem.) | |
| Supermatrix D | 76.50% | 5.90% | 12.60% | 95.00% | 32.50% | 35.10% | 28.60% | 96.20% | 67.20% | 5.70% | 24.80% | 97.70% |
| Supermatrix E | 75.90% | 6.20% | 12.80% | 94.90% | 29.10% | 42.60% | 24.20% | 95.90% | 66.80% | 5.80% | 25.10% | 97.70% |
| Supermatrix F | 68.10% | 7.30% | 16.60% | 92.00% | 28.80% | 38.00% | 22.00% | 88.80% | 60.90% | 6.10% | 29.90% | 96.90% |
| Supermatrix J | 50.20% | 13.80% | 21.90% | 85.90% | 2.50% | 11.20% | 3.50% | 17.20% | 48.30% | 11.40% | 35.00% | 94.70% |

## 4.2.7. Data filtering, multiple sequence alignment, outlier-sequence removal and masking of randomly similar sections

The output of Orthograph could still possibly contain non-exonic residues due to random extension of open reading frames beyond the protein-coding regions (Bank et al., 2017). Therefore we followed additional procedures for filtering sequences within each COG. Specifically, we used the software MACSE v. 2.03 (Ranwez et al., 2018) (option: -trimNonHomologous) to remove long individual sequence fragments that shared no homology with other sequences in each COG, such as those of possibly unidentified intronic fragments (Ranwez et al., 2018). The software PREQUAL v. 1.02 was subsequently used to remove shorter non-homologous fragments such as those resulting from assembly artifacts or annotation errors (default parameters) (Whelan et al., 2018). These filtering steps were applied at the nucleotide sequence level and the resulted COGs (aaCOGs and nCOGs) were used for further downstream filtering. We used the software FSA v. 1.15.9 (option –fast) to infer amino-acid MSAs for each filtered aaCOG (Bradley et al., 2009). We selected the software FSA because it shows higher accuracy (i.e., lower false-positive alignment rate) than other MSA software and tends to leave non-homologous amino-acid residues unaligned (Bradley et al., 2009). By aligning the amino-acid sequences with FSA we greatly reduced the possibility of aligning non-homologous fragments to each other. Subsequently, we filtered the aligned aaCOGs so that amino-acid residues from hybrid enrichment data that did not align to amino-acid residues of at least one reference species (i.e., official gene set) and at least one transcriptome were masked with an "X". Transcriptomic amino-acid residues that did not align to the protein-coding sequences of at least one reference taxon were also masked with an "X". As a last quality check we manually curated all aligned aaCOGs

to mask putative non-homologous amino-acid fragments. We used these filtered amino-acid alignments as a blueprint to generate corresponding codon-based nucleotide alignments with a modified version of PAL2NAL (Suyama et al., 2006) as described by Misof et al. (2014a). A custom python script was then used to mask all corresponding codons of the previously masked amino-acids with "NNN". We performed additional identification and removal of individual outlier sequences in each aligned aaCOG, based on BLOSUM62 expected distances among taxa (see Dietz et al., 2019 and S4: Supplementary Text 1). We subsequently removed all sequences of the reference taxa, except for the sequences of the two hymenopteran species (*Harpegnathos saltator*, *Nasonia vitripennis*) and those of *Tribolium castaneum* that we included as outgroups. Lastly, alignment sections of random similarity within each aaCOG were identified with ALISCORE v. 1.2 (Kück et al., 2010; Misof and Misof, 2009) and were subsequently removed with ALICUT v. 2.31 (https://github.com/PatrickKueck/AliCUT, last access 16.06.2020) both at the amino-acid and the nucleotide sequence levels. The filtered and aligned aaCOGs were finally concatenated into a supermatrix with FASconCAT-G v. 1. 04 (Kück and Longo, 2014).

**4.2.8. Supermatrix evaluation and optimization for phylogenetic analyses**

We opted for an informative subset of the above-described amino-acid supermatrix by using the software MARE v. 0.1.2rc and by removing partitions with an information content of zero (IC = 0) (Misof et al., 2013). After careful visual inspection of the resulted supermatrix (supermatrix A, Table 4.1) we observed that it still contained hypervariable alignment blocks. In addition, supermatrix A contained a large proportion of missing data (~50%, Table 4.1), which can bias phylogenetic

reconstructions if missing characters are not randomly distributed (Lemmon et al., 2009; Misof et al., 2014b). Additionally, supermatrix A showed evidence for deviation from the assumption of stationarity, reversibility and homogeneity (SRH) as measured with the Bowker's and Stuart's tests of symmetry in Symtest v. 2.0.47 (Bowker, 1948; Stuart, 1955) (see Misof et al., 2014a and Table 4.1). Therefore, we chose to filter supermatrix A by applying strategies designed to select conserved alignment sites and reduce the degree of missing data and the potential effects of model violations in phylogenetic reconstructions (e.g., Laumer et al., 2019; Misof et al., 2001; Sharma et al., 2014). First, we identified and removed individual gene partitions within that deviate from model assumptions using the -symtest option in IQ-TREE v. 2.0.4 (Minh et al., 2020; Naser-Khdour et al., 2019). The resulting filtered amino-acid supermatrix was then trimmed with the software BMGE v. 1.12 (h = 0.5, amino-acid replacement matrix: BLOSUM62) to remove hypervariable alignment sites (resulting in supermatrix D). We selected the software BMGE for removing hypervariable sites because it selects informative sites by inferring biologically realistic variability for each column of the alignment (Cai et al., 2020; Criscuolo and Gribaldo, 2010). We also generated five additional and independent amino-acid supermatrices by directly trimming supermatrix A or the partitions of supermatrix A with BMGE in order to examine the effects of progressively more aggressive filtering on the phylogenetic results (see Table 4.1). Additional matrices were generated by using three degrees of stringency (h = 0.5, h = 0.4 and h = 0.3, see Table 4.1 and S6: Fig. S1).

Among-species compositional heterogeneity is a potential source of systematic error that is frequently associated with fast evolving sites (Kocot et al., 2017; Misof et al., 2001). In order to reduce the sensitivity of our phylogenetic analyses to

compositional heterogeneity among species, we generated and analyzed a Dayhoff6-recoded version of supermatrix D. As an alternative approach to reduce among-species compositional heterogeneity in the data, another independent supermatrix was generated for the same purpose by keeping only the 50% of genes with the lowest degree of among-species compositional heterogeneity (RCFV values calculated with BaCoCa v. 1.105, Kück and Struck, 2014). The 322 compositionally homogeneous genes were then concatenated into a new supermatrix which was subsequently trimmed with BMGE (h = 0.5, BLOSUM62) to remove hypervariable sites (supermatrix I, Table 4.1). We also tested whether the removal of distantly related outgroup species or the removal of long-branched ingroup taxa (based on long-branch scores, LB, see S5: Supplementary Text 1) affected the phylogenetic relationships.

We also tested whether the removal of distantly related outgroup species affected the phylogenetic relationships as has been previously suggested for other taxonomic groups (Philippe et al., 2009; Pisani et al., 2015). Therefore, we generated one additional matrix by removing distantly related outgroup species from supermatrix F (i.e., supermatrix G). In addition, we tested whether removal of long-branched ingroup species affected phylogenetic reconstructions by removing the 20 species of the ingroup with the highest long-branch scores from supermatrix G (LB scores, see Supplementary Text 1). Species-specific LB scores were calculated TreSpEx v. 1.1 (Struck, 2014).

We performed a large number of statistical tests on each generated supermatrix in order to evaluate its suitability for phylogenetic reconstruction (Table 4.1). First, we inferred substitution saturation plots for most analyzed supermatrices (Table 4.1, Misof et al. 2001; Nosenko et al. 2013) by calculating pairwise amino-acid $p$-distances and

pairwise patristic distances. Pairwise patristic and *p*-distances were calculated with TreSpEx v. 1.1 (Struck, 2014) by providing the best maximum-likelihood (ML) trees and their corresponding amino-acid supermatrices. Substitution saturation plots were then inferred in R v. 3.6.3 (R Core Team, 2020). We also inferred an alternative measure of substitution saturation that is independent on the patristic distances and the evolutionary model; the average lambda score for each supermatrix (i.e., $\lambda$, ranging from 0.0 to 1.0) that was recently introduced for pairs of aligned sequenced data (higher values indicate higher degree of saturation, Jermiin and Misof, 2020). All pairwise $\lambda$ scores in each supermatrix were calculated with the software SatuRation v. 1.0 (Jermiin and Misof, 2020). For each filtered supermatrix as well as for the original supermatrix A we also measured the overall deviation from SRH conditions with the software SymTest v. 2.0.47 (current version available at https://github.com/ottmi/symtest, last access 20.04.2020, see also Misof et al., 2014a) and by applying the Bowker's and Stuart's tests of symmetry (Table 4.1). Additionally, we calculated the overall completeness scores of the analyzed supermatrices and generated heatmaps of pairwise completeness scores with AliStat v. 1.11 (Wong et al., 2020) (Table 4.1). Lastly, we screened each generated supermatrix for taxa with heterogeneous sequence divergence by generating heatmaps of pairwise mean similarity scores with ALIGROOVE v. 1.06 (Kück et al., 2014).

**Table 4.3**: Branch support (% ultrafast bootstrap and posterior probabilities) for specific well-established phylogenetic clades of Adephaga and their outgroups (based on morphology and other molecular phylogenetic studies) depending on the dataset that was analyzed. Results are shown only for supermatrices that are comparable because they resulted from direct trimming of supermatrix A. Supermatrix D resulted from trimming a slightly different version of supermatrix A from which only 12 genes had been removed. N.M.: Not monophyletic. SHOMU: site-homogeneous unpartitioned, SHETU: site-heterogeneous unpartitioned, PMSF: posterior mean site frequency profile model, BSHETU: Bayesian site-heterogeneous CAT+GTR+G4 model.

| | Adephaga excluding Gyrinidae | | | | Coleoptera | | | | Haliplidae + Dytiscoidea | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | SHOMU | SHETU | PMSF | BSHETU | SHOMU | SHETU | PMSF | BSHETU | SHOMU | SHETU | PMSF | BSHETU |
| **Supermatrix D*** | 99 | 100 | 100 | 0.99 | 84 | 99 | 98 | 1.00 | 99 | 100 | 100 | 1.00 |
| **Supermatrix E** | 98 | 99 | 99 | - | 80 | 97 | 99 | - | 98 | 100 | 100 | - |
| **Supermatrix F*** | 98 | 93 | 95 | 0.94 | N.M. | 94 | N.M. | 1.00 | 98 | 99 | 99 | 1.00 |
| **Supermatrix J** | 98 | 87 | 89 | 0.78 | 74 | 53 | 91 | N.M. | 98 | 92 | 98 | 0.99 |

**Table 4.3** (continued, see legend above)

| | Coleopterida | | | | Aspidytidae | | | |
|---|---|---|---|---|---|---|---|---|
| | SHOMU | SHETU | PMSF | BSHETU | SHOMU | SHETU | PMSF | BSHETU |
| **Supermatrix D*** | 100 | 1.00 | 100 | 1.00 | 100 | 100 | 100 | 1.00 |
| **Supermatrix E** | 100 | 1.00 | 100 | - | 100 | 100 | 100 | - |
| **Supermatrix F*** | 100 | 1.00 | 100 | 1.00 | 97 | 81 | 82 | 0.88 |
| **Supermatrix J** | N.M. | 1.00 | 100 | 0.99 | 100 | 99 | 99 | 1.00 |

*Note: The BSHETU analyses of supermatrix D did not reach convergence while BSHETU analyses of supermatrix F have reached the acceptable value of maxdiff. = 0.3.

**4.2.9. Concatenation-based phylogenetic analyses of amino-acid supermatrices**

Modeling site-specific propensities of amino-acid frequencies has been shown to be more important than modeling partition-wise heterotachy in concatenation-based phylogenomic analyses (Feuda et al., 2017; Wang et al., 2019). In order to account for site-specific amino-acid preferences in the datasets, we analyzed most amino-acid supermatrices under the site-heterogeneous model CAT+GTR+G4 (Bayesian site-heterogeneous model, BSHETU) using the software Phylobayes MPI v. 1.8 (Table 4.1) (Lartillot et al., 2013). Two independent MCMC chains were run for each dataset until more than 20,000 samples were collected or until convergence (maxdiff < 0.3, Table S7). We assessed convergence of the runs on the tree space as well as for the summary variables of the model with bpcomp and tracecomp respectively (see manual of Phylobayes, Lartillot et al. 2013).

We also analyzed the amino-acid supermatrices using a maximum likelihood approach (ML) with IQ-TREE v. 1.6.12 (Nguyen et al., 2015). We first selected the best-fitting substitution models in ModelFinder based on the AICc criterion on the unpartitioned matrices (S4: Table S9, Akaike, 1974; Kalyaanamoorthy et al., 2017). In order to test the relative fit of site-heterogeneous versus site-homogeneous models in a ML framework, we also included empirical site-heterogeneous mixture models in our model-selection procedure (Quang et al., 2008). In total, more than 270 models were tested on each of supermatrices B–J (unpartitioned data) except for the recoded dataset which was only analyzed with the BSHETU model in Phylobayes. For the partitioned supermatrices (B, C, Table 4.1), we also calculated an optimal partitioning scheme using an edge-linked partition model using the same version of IQ-TREE (see S5: Supplementary Text 1) (Chernomor et al., 2016; Lanfear et al., 2014). For these

supermatrices we assessed the relative model fit of site-homogeneous unpartitioned (SHOMU), site-homogeneous partitioned (SHOMP) and site-heterogeneous unpartitioned (SHETU) models by using a fixed neighbor-joining tree (Table S8, see S5: Supplementary Text 1). Phylogenetic tree inference was performed for each matrix under the SHOMU, SHETU and SHOMP models (where applicable) in order to explore the extent to which using a suboptimal model affected phylogenetic reconstructions (S4: Tables S8–S9). Lastly, for each supermatrix we also performed ML analyses using the posterior mean-site frequency profile (PMSF) approximation to the site-heterogeneous models (S5: Supplementary Text 1) (Wang et al., 2018). Statistical branch support of the inferred relationships in all concatenation-based ML analyses was estimated based on 2000 ultrafast bootstrap (UFB) replicates (Hoang et al., 2018). Lastly, we calculated pairwise normalized RF distances among the inferred trees under the same model (SHOMU, SHETU or PMSF) for amino-acid datasets with full taxon sampling using ETE v. 3.1.1 (Huerta-Cepas et al., 2016).

### 4.2.10. Phylogenetic analyses of nucleotide sequence data

To assess the stability of phylogenetic results among analyses of different types of data we also generated and analyzed four supermatrices at the nucleotide sequence level (S4: Table S10). Analyses of these supermatrices was performed with the same version of IQ-TREE and by selecting best-fitting SHOMP and SHOMU models (see S5: Supplementary Text 1). We also inferred phylogenetic relationships using a model that accounts for heterotachy among sequences but has only been extensively tested for application in nucleotide sequence analyses (see S5: Supplementary Text 1, Table S10, Crotty et al., 2020).

**4.2.11. Estimating alternative and confounding signals in supermatrices via four-cluster likelihood mapping and data permutations**

We applied the four-cluster likelihood mapping approach (FcLM) to assess the robustness of phylogenetic results, and to measure the strength of alternative phylogenetic signals with respect to specific phylogenetic hypotheses that resulted from the analyses of supermatrix D (Fig. 4.3) (Strimmer and von Haeseler, 1997). The hypotheses that we tested were the following: a) Hygrobiidae are sister to a clade of Amphizoidae + Aspidytidae (hypothesis 1) and b) Cicindelidae are the sister group of Carabidae (hypothesis 2). FcLM analyses were performed on different amino-acid supermatrices that were trimmed with different degrees of stringency and were based on both SHETU and SHOMU models, in an attempt to assess whether model misspecification affected the phylogenetic signal in favor of specific hypotheses (Table 4.2). In addition, FcLM analyses under the better-fitting SHETU models were performed with permutations of data (i.e., randomization of phylogenetic signal, permutation no. I in Misof et al., 2014a), in order to assess whether or not the FcLM support for a particular inferred relationship under the SHETU models resulted from misleading signal (Table 4.2) (Misof et al., 2014a).

**4.2.12. Summary coalescent phylogenetic analyses (SCA)**

To explore the sensitivity of our concatenation-based analyses to the putative effects of ILS we conducted summary coalescent phylogenetic analyses with ASTRAL III v. 5.7.3 (C. Zhang et al., 2018). Due to the fact that SCA are prone to gene-tree estimation errors (Mirarab et al., 2016; Sayyari et al., 2017) we took steps to reduce these effects on the results of our analyses. Alignment trimming methods have been

shown to be detrimental in phylogenetic inference of gene trees (Tan et al., 2015), and therefore we selected the unmasked amino-acid alignments for these analyses (before trimming with ALISCORE, Fig. 4.1, S6: Fig. S1). However, in order to reduce the negative effects of fragmentary sequences (Sayyari et al., 2017) as they are common for sequence capture data (Hosner et al., 2016), we 1) removed sites with more than or equal to 50% ambiguous characters, and then 2) removed sequences for which more than 75% of sequence length contained ambiguous characters. Finally, we kept only genes that had a length of at least 150 amino acids and less than 50% total missing data. The filtering tasks were performed with custom PERL scripts. In total, 348 filtered gene alignments were used for SCA. Gene trees were inferred after selecting the best-fitting models (same set of models that were tested for SHOMP and SHOMU analyses of amino-acid supermatrices) with the same version of IQ-TREE (see S5: Supplementary Text 1 for details). Branch support of individual gene trees was calculated based on 10,000 SH-aLRT replicates (Guindon et al., 2010). SCA were then conducted with ASTRAL after collapsing weakly supported branches (< 50% SH-aLRT support) with ETE v. 3.1.1.

Because SCA resulted in different topologies from the concatenation analyses we explored whether or not selecting genes with the highest levels of phylogenetic information resulted in higher congruence with the concatenation-based analyses. Potential phylogenetic information of each of the 348 filtered genes was assessed based on three criteria: (a) average SH-aLRT branch support of inferred gene trees (SH), (b) percentage of fully resolved quartets by likelihood mapping (see S5: Supplementary Text 1, LM, Strimmer and von Haeseler, 1997), and (c) number of parsimony informative sites per gene (PI). Subsets of genes with the highest scores were then

obtained for downstream analyses (i.e., with values larger than the median for criteria a and b and larger or equal to the median for criterion c, Fig. 4.4c-4.4f). Subsequently, SCA were repeated for all selected subsets of genes as well as for the overlaps of genes that were selected by different approaches (i.e., see Fig. 4.4f and 4.4g, again after collapsing weakly supported branches with lower then 50% support). In order to evaluate gene-tree support for competing hypotheses and to assess whether or not gene-tree error might have contributed the conflicting phylogenetic results and low branch support values in the SCA, we performed gene-tree discordance analyses (GTD) with DiscoVista v. 1.0 (Sayyari et al., 2018). GTD was separately performed for each subset of gene trees (i.e., full set of gene trees and for the three selected subsets with optimal phylogenetic information) using a branch support threshold of 70% for clades to be considered strongly accepted or rejected (Sayyari et al., 2018). GTD was calculated for predefined hypotheses but also for clades that are well established based on previous analyses of molecular and morphological data (e.g., monophylettic Coleoptera, monophyletic Adephaga, a clade Amphizoidae + Aspidytidae, see Fig. 4.4b). We postulated that the concatenation-based tree under the better fitting SHETU model (Fig. 4.3) provides a good approximation of the true familial phylogenetic relationships of Adephaga, because it is highly congruent with morphology-based phylogenies and latest molecular phylogenetic analyses of the group (Baca et al., 2017; Beutel et al., 2020; Gustafson et al., 2020). Based on that premise, we calculated normalized Robinson-Foulds (RF) distances (Robinson and Foulds, 1981) between this tree and the different species-trees that resulted from the SCA analyses under various gene-subsampling strategies in order to assess which gene subsampling strategy results in higher topological congruence with the concatenation-based tree. RF distances between

the concatenation-based tree (Fig. 4.3) and the SCA trees were calculated with ETE v. 3.1.1 and visualized in R v. 3.6.3.

## 4.3. Results

### 4.3.1. Sequencing, assembly, cross-contamination check and orthology assignment for the hybrid-capture data

On average, we generated 1,670,754 pairs of sequence reads per genomic library (S4: Table S2). Quality and adapter trimming resulted in the removal of 239,210 paired reads from each sequenced genomic library on average. After assembly and cross-contamination removal, each of the clean assemblies contained 28,923 contigs on average. The summarized results of the orthology assignment for the sequence-capture data show that more than half of the 651 genes of the bait set were identified in the species of each family of Adephaga (S6: Fig. S3, median values: Cicindelidae = 523, Carabidae = 547.5, Dytiscidae = 532, Gyrinidae = 497, Haliplidae = 596, Noteridae = 549.5). On average, 534 genes where identified in the orthology assignment step in each genomic assembly (median = 542, , max. = 642, min. = 177, Table S2). Results of the orthology assignment for the transcriptomes are separately presented in S4: Table S1 (no. of orthologous transcripts: mean = 640.5, median = 650, max. = 651, min. = 533).

### 4.3.2. Statistics of the hybrid enrichment

The results show that the overall Ct / Cn ratio is much higher than one for the majority of species which in turn suggests that the enrichment of the target regions was successful for the majority of species in our dataset (S4: Table S2, Ct / Cn median

values: Carabidae = 22.163, Noteridae = 39.414, Haliplidae = 50.231, Gyrinidae = 5.682, Cicindelidae = 11.312, Dytiscidae = 9.976, Fig. 4.2). The same applies for the Ct / Ca ratio (S4: Table S2, S6: Fig. S2). However, the calculated statistics showed that the enrichment was potentially more successful for some adephagan families than others (Fig. 4.2, S6: Fig. S2). For example, Noteridae and Haliplidae have the highest overall Ct / Cn scores that are statistically significantly higher than values for Gyrinidae, Dytiscidae and Cicindelidae (Fig. 4.2, Table S5). The calculated enrichment statistics for Carabidae suggest that the enrichment was potentially more successful for this family than for the species in the families Cicindelidae and Dytiscidae, although not statistically different from the species of Gyrinidae, Haliplidae and Noteridae (Fig. 4.2, Table S5).

### 4.3.3. Family-level phylogenetic relationships of Adephaga

Most concatenation-based analyses delivered a congruent picture on the evolution of adephagan beetles irrespective of the data type used or the model applied (e.g., Fig. 4.3). Specifically, a clade of Archostemata + Myxophaga as sister to Adephaga was recovered in all analyses under the best-fitting SHETU models in a ML framework (Fig. 4.3, S6: Fig. S4–S11), in most BSHETU analyses of amino-acid data (e.g., S6: Fig. S12–S18) but also in the analyses of nucleotide sequence data under site-homogeneous models and models that account for heterotachy (S6: Fig. S19–S23). The family Gyrinidae was inferred as sister to all other Adephaga in all concatenation-based analyses under full taxon sampling except for the unconverged BSHETU analyses of the Dayhoff6-recoded supermatrix D (S6: Fig. S13). Interestingly, removal of distantly related outgroups from supermatrix F (i.e., supermatrix G) without also removing long-

branched ingroup taxa (i.e., supermatrix H), resulted in the equivocal placement of Gyrinidae (i.e., polytomy) under the BSHETU model, whereas SHOMU, SHOMP, SHETU and PMSF models consistently recovered Gyrinidae as sister to all other Adephaga for these datasets (S6: Fig. S4–S43). Geadephaga were consistently inferred as monophyletic and sister to Haliplidae + Dytiscoidea under concatenated analyses of different models and data types (Fig. 4.3, S6: Fig. S4–S43). Within Dytiscoidea, the family Noteridae was inferred as sister to all other dytiscoid families and Amphizoidae was inferred as sister to monophyletic Aspidytidae in all concatenation-based analyses of amino-acids and nucleotides (e.g., Fig. 4.3, S6: Fig. S19–S23). Within Geadephaga, the monophyly of the tiger beetles (Cicindelidae) and their placement as sister to monophyletic ground beetles (Carabidae) was inferred in all analyses under the more complex site-heterogeneous models (BSHETU, SHETU, PMSF, Table 4.4, S6: Fig. S4– S18, S33–S41) and was also supported by analyses of nucleotide sequence data (Fig. S19–S23). In contrast, Trachypachidae was inferred as sister to Carabidae only in the analyses of supermatrix J and only under conditions of model misspecification (i.e., SHOMU and PMSF models) yet with no strong statistical branch support (Table 4.4, S6: Fig. S32, S41).

Concerning the inferred position of the family Hygrobiidae all ML analyses under the best-fitting SHETU models supported a clade of Hygrobiidae + (Amphizoidae + Aspidytidae) and most of them with strong branch support (e.g., Fig. 4.3). The branch support of this clade under SHETU models was lower when more stringent trimming criteria were applied but the inference of this clade remained robust to the selection of dataset when a SHETU model was applied (Table 4.4). On the other hand, analyses under the SHOMU models were inconsistent regarding this hypothesis

(Table 4.4). Specifically, SHOMU analyses of the most stringently trimmed supermatrix under full taxon sampling (supermatrix J) supported a clade Dytiscidae + (Amphizoidae + Aspidytidae) as sister to Hygrobiidae but not with strong statistical branch support (Table 4.4). In general, progressive trimming with more stringent criteria resulted in shift from a strongly or moderately supported Hygrobiidae + (Amphizoidae + Aspidytidae) clade (supermatrix D and E) to a poorly supported Dytiscidae + (Amphizoidae + Aspidytidae) clade (supermatrix F and J) but only in conditions of model misspecification (SHOMU models). This pattern is also observed under BSHETU model although only for the most stringently trimmed suppermatrix (supermatrix J, Table 4.4). Phylogenetic analyses with PMSF approximation to the SHETU model (using a SHOMU-based guide tree) restored the monophyly of Hygrobiidae + (Amphizoidae + Aspidytidae) for most supermatrices (Table 4.4, except supermatrix C) suggesting that the clade Dytiscidae + (Amphizoidae + Aspidytidae) inferred under SHOMU models is likely an artifact due to model misspecification. Overall, a clade that includes Dytiscidae, Amphizoidae and Aspidytidae is either inferred under conditions of site-homogeneous models or is never strongly supported (Table 4.4, Fig. S17, S18, S23, S34, S42).

Additional support for a clade Hygrobiidae + (Amphizoidae + Aspidytidae) comes from the results after removing distant outgroups and long-branched ingroup taxa from supermatrix F. Specifically, removing distantly related outgroups did not result in strong support for this clade (93%, Table 4.4) under the SHETU model but when long-branched ingroup taxa were also subsequently removed, the support for the above-mentioned clade increased under the same model (98%). Additionally, the topology flipped from the clade Dytiscidae + (Amphizoidae + Aspidytidae) to the clade

Hygrobiidae + (Amphizoidae + Aspidytidae) clade under the SHOMU and BSHETU models when long-branched ingroup species were also removed (although not with strong support under the BSHETU, Table 4.4). This suggests that removal of distant outgroups without also accounting for branch-length heterogeneity of the ingroup might result in erroneous topology even when a site-heterogeneous model is used. Phylogenetic analyses of the Dayhoff6-recoded matrix D recovered unexpected and poorly supported clades with respect to the internal phylogeny of Dytiscoidea and more generally Adephaga (e.g., Gyrinidae + Geadephaga and Amphizoidae + Dytiscidae with low support, S6: Fig. S13). Although the BSHETU analyses of the recoded matrix failed to reach robust convergence statistics (S6: Fig. S13, S4: Table S7, maxdiff = 0.49, more than 29,000 samples per MCMC chain), these observations suggest that amino-acid data-recoding might be detrimental in those cases that excessive alignment trimming and data filtering has been applied before recoding of the data.

### 4.3.4. Phylogeny of Carabidae, Cicindelidae and Dytiscidae

Analyses of amino-acid and nucleotide supermatrices in a concatenation framework resulted in the monophyly of all subfamilies of diving beetles in the family Dytiscidae (e.g., Fig. 4.3). However, phylogenetic relationships among constituent subfamilies were unstable and not consistently resolved in all analyses except for a few cases. For instance, the subfamily Hydrodytinae was always inferred as sister to Hydroporinae with strong support (e.g., Fig. 4.3). The subfamilies Lancetinae and Coptotominae were always inferred as sister groups (e.g., Fig. 4.3, S6: Fig. S4–S18). In addition, all concatenation-based analyses resulted in a clade that includes all subfamilies of Dytiscidae excluding Lancetinae, Coptotominae and Laccophilinae (Fig.

4.3, S6: Fig. S4–S43). Specifically, most analyses with the best-fitting SHETU models recovered Lancetinae + Coptotominae as sister to Laccophilinae + remaining Dytiscidae (e.g., Fig. 4.3, S6: Fig. S4–S11). In addition, most concatenation-based analyses of amino acids suggested the placement of Copelatinae as sister to a clade Matinae + (Hydrodytinae + Hydroporinae) (e.g., Fig. 4.3, Fig. S4–S18, S24–S43). Lastly, all concatenation-based analyses recovered Colymbetinae as sister to Agabinae and Cybistrinae as sister to Dytiscinae with strong statistical branch support (Fig. 4.3, S6: Fig. S4–S43).

Concerning the phylogeny of Cicindelidae, the tribe Manticorini was inferred as sister to all other subfamilies of Cicindelidae in concatenation-based analyses (Fig. 4.3, S6: Fig. S4–S43). Although paraphyletic Manticorini was inferred in a few instances, this result was likely an artifact due to the extremely high degree of missing data for the species *Manticora latipennis* (S4: Table S2). The tribe Megacephalini was placed as sister to all remaining Cicindelidae except Manticorini, while the tribe Collyridini was inferred as sister to a clade that included Cicindelini and Oxycheilinini (Fig. 4.3, S6: Fig. S4–S18, S24–S43). The internal phylogeny of the megadiverse Carabidae remained largely unstable across analyses of different supermatrices and models (e.g., S6: Fig. S4–S43). Despite this, some relationships were robustly inferred. For example, the subfamily Trechinae was always inferred as sister to Brachininae + monophyletic Harpalinae, whereas the subfamilies Paussinae, Rhysodinae and Siagoninae were placed in a monophyletic group close to the base of the tree of Carabidae in analyses of amino-acid supermatrices (Fig. 4.3, S6: Fig. S4–S18, S24–S43). Lastly, the subfamily Carabinae was inferred as sister to Nebriinae in most concatenation-based phylogenetic analyses of amino-acid sequence data (Fig. 4.3, S6: Fig. S4–S18, S24–S43).

**Fig. 4.3**: Phylogenetic relationships of Adephaga as they resulted from the analysis of supermatrix D under the JTT+C60+F+R8 site-heterogeneous model (i.e., SHETU model).

**Fig. 4.3** (*caption continued from previous page*): Circles on tree nodes indicate branch support based on 2000 ultrafast bootstraps. All beetle photos by M. Balke.

### 4.3.5. Comparison of different schemes of evolutionary modeling and predictability of substitution saturation among different modeling schemes

In total 277 models were tested on each unpartitioned amino-acid supermatrix with ModelFinder. The results show that SHETU models significantly outperformed the best SHOMU models for all supermatrices in an unpartitioned context (S6: Table S9). All the best-fitting SHETU models included 60 categories of fixed empirical amino-acid frequencies (i.e., C60 site-heterogeneous models) suggesting that the most complex SHETU models fitted the data better even for the most stringently trimmed supermatrices (e.g., supermatrices F and J, S4: Table S9). Comparison of the optimal partitioning schemes (SHOMP) for supermatrices B and C with the complex SHETU models showed that site-heterogeneous models (SHETU) fitted these datasets better than both partitioned and unpartitioned site-homogeneous models (SHOMP and SHOMU, S4: Table S8, S9). Based on the observation that SHETU models fit the data better, the saturation statistics showed that using a site-homogeneous model (SHOMP or SHOMU) resulted in underestimation of the amount of substitution saturation in the amino-acid supermatrices when a measure that is dependent on patristic distances was used (i.e., adjusted $R^2$, Table 4.1).

**Table 4.4**: Branch support statistics (% ultrafast bootstrap support and posterior probabilities) for the two most controversial clades of Adephaga under different models in all analyzed supermatrices. In those cases that the clade Hygrobiidae + (Amphizoidae + Aspidytidae) was not inferred, a clade Dytiscidae + (Amphizoidae + Aspidytidae) was inferred instead (as sister to Hygrobiidae) with low branch support (i.e., lower than 95 ultrafast bootstrap support or lower than 0.95 posterior probability). In all cases that a clade Cicindelidae + Carabidae was not inferred, a clade Trachypachidae + Carabidae was recovered instead (as sister to Cicindelidae) with low branch support (i.e., lower than 95 ultrafast bootstrap support). N.I.: not inferred, N.A.: Not applicable, SHETU: site-heterogeneous unpartitioned, SHOMP:site-homogeneous partitioned, SHOMU: site-homogeneous unpartitioned, BSHETU: Bayesian CAT+GTR model.

| Dataset | No. of species | Clade | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|
| | | Hygrobiidae + (Amphizoidae + Aspidytidae) | | | | | Carabidae + Cicindelidae | | | | |
| | | SHOMU | SHOMP | SHETU | PMSF | BSHETU | SHOMU | SHOMP | SHETU | PMSF | BSHETU |
| Supermatrix B | 136 | 70 | N.I. | 96 | 98 | - | 100 | 100 | 100 | 100 | - |
| Supermatrix C | 136 | N.I. | 77 | 96 | N.I. | - | 100 | 100 | 100 | 100 | - |
| Supermatrix D* | 136 | 96 | N.A. | 100 | 100 | 1.00 | 100 | N.A. | 100 | 100 | 1.00 |
| Supermatrix E | 136 | 94 | N.A. | 100 | 100 | - | 100 | N.A. | 100 | 100 | - |
| Supermatrix F* | 136 | N.I. | N.A. | 96 | 90 | 0.99 | 100 | N.A. | 100 | 100 | 1.00 |
| Supermatrix G | 120 | N.I. | N.A. | 93 | 87 | N.I. | 100 | N.A. | 100 | 100 | 1.00 |
| Supermatrix H | 100 | 95 | N.A. | 98 | 99 | 0.83 | 100 | N.A. | 100 | 100 | 1.00 |
| Supermatrix I | 136 | N.I. | N.A. | 96 | 94 | N.I. | 100 | N.A. | 100 | 100 | 1.00 |
| Supermatrix J | 136 | N.I. | N.A. | 92 | 86 | N.I. | N.I. | N.A. | 61 | N.I. | 0.88 |

*Note: The BSHETU analyses of supermatrix D did not reach convergence.

**Fig. 4.4**: Results of summary coalescent (SCA) and gene-tree discordance (GTD) analyses. (**a**) Summarized phylogram that resulted from the SCA analyses with all genes (n = 348), (**b**) GTD analyses showing proportion of gene trees that support or reject different relationships of Adephaga and outgroups when the full set of gene trees was used, (**c**) distribution of the percentage of resolved quartets among different genes in our dataset (n = 348), (**d**) distribution of the average SH-aLRT branch support of inferred gene trees among different genes in our

dataset (n = 348), (**e**) distribution of the number of parsimony informative sites among different genes in our dataset (n = 348), (**f**) Venn diagram showing number of genes selected based on each criterion applied (LM, SH, PI) and number of overlapping genes among selected subsets (i.e., LM+SH = 104, LM+PI = 87, SH+PI = 130), (**g**) normalized RF distances of the trees inferred under different subset of genes in Fig. 4.4f to the concatenation-based species tree in Fig. 4.3. Note: dashed lines in  histograms of Fig. 4.4c–4.4e indicate median values.

### 4.3.6. Stability of inferred relationships of Adephaga across analyses with different evolutionary models

We calculated all pairwise normalized RF distances among trees inferred under the same model (SHOMU, SHETU or PMSF) for those amino-acid datasets with full taxon sampling (seven trees per model, supermatrices B, C, D, E, F, I, J, Fig. 4.5). We did this to assess whether or not topological distances between inferred trees differ when using different evolutionary models. Although, median RF distances of inferred trees did not significantly differ between PMSF and SHOMU models ($p$-value = 0.237, Mann-Whitney-Wilcoxon test with continuity correction) or between PMSF and SHETU models ($p$-value = 0.136, Mann-Whitney-Wilcoxon test with continuity correction), RF distances of inferred trees were lower in analyses of SHETU models when compared to the SHOMU models ($p$-value = 0.013, Mann-Whitney-Wilcoxon test with continuity correction, Fig. 4.5). This result is congruent with the consistent inference of the clade Hygrobiidae + (Amphizoidae + Aspidytidae) under SHETU models that was instead not consistently inferred under the SHOMU models, and constitutes further evidence that full site-heterogeneous empirical mixture models (C60, ML-based) result in greater stability of the inferred relationships than the less complex SHOMU models (Table 4.4, Fig. 4.5).

**Fig. 4.5**: Box-plots of all pairwise Robinson-Foulds distances among trees that were inferred from different amino-acid supermatrices under the same type of model (normalized distances, only maximum likelihood analyses). Note: we only included distances among trees that were inferred with full taxon sampling (i.e., supermatrices: B, C, D, E, F, I, J). SHETU: site-heterogeneous unpartitioned model, PMSF: posterior mean site frequency profile model, SHOMU: site-homogeneous unpartitioned model.

### 4.3.7. Effects of removing hypervariable sites, distantly related outgroups and long-branched taxa on the statistical properties of amino-acid supermatrices

On the one hand, removal of hypervariable sites had a positive impact on the statistical properties of amino-acid supermatrices in terms of eliminating potential confounding factors (Table 4.1). In particular, trimming the supermatrices with BMGE

resulted in reduction of total and pairwise missing data (Table 4.1 and S6: Fig. S44–S53) and reduced deviation from SRH conditions as indicated by the reduced percentage of pairwise comparisons that failed the corresponding symmetry tests in the analyzed supermatrices (Table 4.1, Bowker's test: 35.02%, 24.98%, 13.97% failed tests in supermatrices E, F and J respectively, see also S6: Fig. S54–S64). More specifically, progressive removal of hypervariable sites resulted in progressively fewer failed pairwise symmetry tests, as well as in progressively increasing completeness of the supermatrices ($C_a$ scores: 0.846, 0.882 and 0.911 for supermatrices E, F and J respectively, Table 4.1, S6: Fig. S44–S53 ). Supermatrices D and E did not significantly differ when comparing their statistical properties because only 12 genes from supermatrix A failed the symmetry tests in IQ-TREE and had therefore been removed before trimming (Table 4.1). Pairwise alignment similarity scores of taxa and indices for substitution saturation also improved with BMGE trimming (Table 4.1, Fig. S65–S94, supermatrices D, E, F and J), suggesting that progressively removing hypervariable sites results in progressively less saturated amino-acid supermatrices (supermatrices D, E, F and J). The average λ scores within each supermatrix also showed that progressive removal of hypervariable sites resulted in supermatrices with less decay of potential historical signal (i.e., lower average λ scores, supermatrices D, E, F and J in Table 4.1). On the other hand, progressive and more aggressive trimming of hypervariable sites resulted in progressive reduction of total parsimony informative sites and reduced percentage of parsimony informative sites (from 43.00% in supermatrix E to 32.80% in supermatrix J, Table 4.1). In a similar fashion, Dayhoff6-recoding resulted in removal of 40.66% of parsimony informative sites from supermatrix D (Table 4.1).

Removal of distantly related outgroups from supermatrix F, resulted in a less saturated supermatrix according to average λ score, whereas the regression of *p-* and patristic distances under the SHOMU and SHETU models showed reduced adjusted $R^2$ value (i.e., suggesting higher saturation) compared to the dataset before removing distantly related outgroups (i.e., supermatrix F). This in turn suggests that the different saturation statistics (adjusted $R^2$ and average λ) do not always provide congruent estimates on which datasets are the most saturated (Table 4.1). Moreover, comparisons of saturation statistics among datasets and models showed that conventional statistics of substitution saturation ($R^2$ and slope of regression) are highly dependent on the model (e.g., supermatrices F, G, H and I, Table 4.1). Despite this, removal of distantly related outgroups from supermatrix F resulted in reduced proportion of failed pairwise symmetry tests (Bowker's test: 24.98%, 20.35, 18.85% failed tests in supermatrices F, G, H respectively). Removal of long-branched ingroup taxa resulted in further decrease in potential deviations from SRH conditions and also in further reduction in the degree of saturation (Bowker's test: 24.98%, 20.35%, 18.85% failed tests, λ scores: 0.095, 0.079, 0.074 in supermatrices F, G, H respectively).

### 4.3.8. Effects of removing hypervariable sites on the branch support statistics for well-established adephagan relationships

We examined how removing hypervariable sites with BMGE using different degrees of stringency affected phylogenetic branch support for well-established relationships of Adephaga and their outgroups. A clade that includes all adephagan families except Gyrinidae was strongly supported when using a moderate trimming strategy (supermatrices D, E and I, see Table 4.3) but UFB support for this relationship

decreased with more aggressive trimming of the data under the SHETU and PMSF models (SHETU: 93% and 87% support in supermatrices F and J respectively, Table 4.3). This pattern is also observed under the complex BSHETU model (0.94 and 0.78 posterior probability in supermatrices F and J respectively, Table 4.3), whereas analyses under a misspecified model (SHOMU) were the still strong support for this relationship (98% in supermatrix J). A similar pattern is observed for the monophyly of a clade Haliplidae + Dytiscoidea which is inferred under all models but receives lower support in the analyses of supermatrices that were trimmed more aggressively (99% and 92% UFB support in supermatrices F and J under the SHETU model respectively, Table 4.3). In addition excessive trimming of the supermatrix A resulted in very low support for the monophyly of Coleoptera under the better fitting SHETU model and even resulted in non-monophyletic Coleoptera in cases of model misspesification (Table 4.3, supermatrix F). The monophyly of the family Aspidytidae is also less well-supported in the analyses of supermatrices that were produced by very stringent trimming (supermatrices F and J, 81% and 99% respectively under the SHETU model, Table 4.3). Lastly, progressive trimming of the supermatrices resulted in reduction of the overall proportion of clades that are well supported under the better-fitting SHETU models (total proportion of branches with > 95% UFB support, Fig. 4.6).

**Fig. 4.6**: Percentage of branches with support lower than 100% (red bars) and lower than 95% (blue bars) in the the phylogenetic trees inferred under the SHETU models using amino-acid supermatrices that were trimmed with different degrees of stringency (i.e., BLOSUM62 and h = 0.5, h = 0.4 or h = 0.3). Note: we only included supermatrices that are comparable because they resulted from direct trimming of supermatrix A. Supermatrix D resulted from trimming a slightly different version of supermatrix A from which only 12 genes had been removed.

## 4.3.9. Measuring alternative and confounding signals using a combination of four-cluster likelihood mapping and data permutations

Overall, more aggressive trimming of hypervariable sites (i.e., h = 0.4, h = 0.3) resulted in a reduction of the total number of resolved quartets for the two tested hypotheses under SHOMU models and even more profoundly for the better-fitting SHETU models (Table 4.2). More specifically for  hypothesis 2 less than 90% of the total number of quartets were fully resolved after applying the most stringent trimming regime under the SHETU models (85.90% in the analyses of supermatrix J, Table 4.2).

Concerning the position of Hygrobiidae (hypothesis 1) there was moderate to strong support for a clade Hygrobiidae + (Amphizoidae + Aspidytidae) in the analyses of moderately trimmed matrices (D and E, 59.80% and 58.80%) without detectable confounding signal (see FcLM of permuted data, Table 4.2). However, the signal in favor of this clade was reduced in the supermatrices that were trimmed more aggressively and a shift in phylogenetic support for the other two alternatives was observed (Table 4.2). The same pattern was observed in hypothesis 2 in which a clade Carabidae + Cicindelidae was strongly supported in the FcLM analyses of supermatrices D and E (76.50% and 75.90% respectively), whereas there was reduction in support for this hypothesis when more stringent criteria were applied (68.10% and 50.20% respectively). The absence of detectable confounding signal supporting the original results of tree reconstructions in the moderately trimmed matrices (D and E, permuted data) suggests that the shift in support from the original strongly supported hypotheses to the other two alternatives is likely not due to removal of potentially confounding signal when trimming the data but likely due to removal of genuine phylogenetic signal. This is likely the case for both hypotheses (Table 4.2).

### 4.3.10. Summary coalescent phylogenetic analyses (SCA)

The SCA from the analyses of all genes produced topologies that were mostly congruent with concatenation-based analyses concerning the familial relationships of Adephaga and outgroups with some exceptions (Fig. 4a). Despite this, SCA resulted in weakly supported well-established clades (e.g., monophyly of Coleoptera and Dytiscoidea, Fig. 4a). The clade Archostemata + Myxophaga was disrupted in most SCA, that instead resulted in a clade Archostemata + Adephaga but with low branch

support (Fig. 4.4a, S6: Fig. S95–102). In addition, SCA did not recover Gyrinidae as sister to all other Adephaga but the exact phylogenetic position of the family differed based on the subset of genes that was analyzed (Fig. S3: S95–102, Fig. 4.4a). Despite this, SCA with either all genes included or with an LM-based optimal subset of genes did not reject the monophyly of the families of Adephaga excluding Gyrinidae, because the branch length of the inferred clade Gyrinidae + Geadephaga was estimated to zero resulting in a polytomy (Fig. 4.4a, S6: Fig. S95, S96). Other conflicts between concatenation-based analyses and SCA is the inference of a clade Trachypachidae + Carabidae with low branch support in the latter (Fig. 4.4a, S6: Fig. S95–S102), and also the differences in inferred relationships within the adephagan families that are generally poorly supported in all SCA (S6: Fig. S95–S102). GTD analyses showed that the vast majority of the best gene trees strongly reject all well-established hypotheses of Adephaga and outgroups as well as the monophyly of some individual families (e.g., monophyly of Coleoptera, monophyly of Haliplidae + Dytiscoidea, monophyly of the families Carabidae and Dytiscidae). In addition, when considering only the most controversial familial relationships of Adephaga (e.g., position of Gyrinidae, Hygrobiidae and Trachypachidae) the vast majority of gene trees reject the different alternative topologies suggesting that gene-tree error is prominent in the data. Additionally, the distribution of potential phylogenetic signal among the sampled genes by any applied criterion (LM, PI, SH) shows that many of the genes were highly uninformative (Fig. 4.4c, 4.4d, 4.4e, median LM = 58.14, median SH = 67.97, median PI = 119) and therefore unlikely to have produced correct gene trees.

We tested whether or not choosing genes with higher potential phylogenetic information resulted in higher topological congruence with the concatenation-based

species-tree and whether or not GTD for well-established hypotheses of Adephaga was reduced when applying gene subsampling strategies. Overall, the different SCA delivered different topologies suggesting that the applied summary coalescent method is extremely sensitive to the set of gene trees used (Fig. 4.4f–4.4g, S6: S95–S102). Selecting subsets of gene trees based on the phylogenetic information of genes (here measured with LM, PI and SH criteria, Fig. 4.4f) resulted in higher topological congruence with the concatenation-based tree (Fig. 4.3) than the analyses utilizing all gene trees (Fig. 4.4g). Nevertheless, SCA of the SH- and PI-based gene subsets failed to recover some familial relationships of Adephaga (Fig. 4.4a, such as the placement of Gyrinidae as sister to all other Adephaga and the monophyly of Dytiscoidea, Fig. 4.4g, S6: Fig. S97, S98). Overall, subsampling genes based on the LM criterion (percentage of resolved quartets) resulted in the lower RF distances to the concatenation-based tree (Fig. 4.3) than subsampling based on the PI and SH criteria (Fig. 4.4g). SCA of the LM-subset of genes resulted in familial relationships identical to the SCA with all genes (Fig. 4.4a, S6: Fig. S95–S96) but with higher overall topological congruence to the concatenation-based tree (i.e., lower RF distance, Fig. 4.4g). Despite this, GTD analyses on the LM-, SH- and PI-selected subsets of genes showed that gene-tree error was still very prominent even for the analyses of selected gene subsets as the majority of gene trees strongly rejected all tested phylogenetic hypotheses similarly to the GTD analysis performed for the full set of gene trees (Fig. 4.4b, S6: Fig. S103–S105).

It is also noteworthy, that different criteria of potential phylogenetic informativeness produced different predictions on which genes are the most informative, with SH- and PI-based selected gene subsets showing a greater overlap of selected genes (Fig. 4.4f). We performed SCA analyses of overlapping subsets of

selected genes among different subsampling approaches. When overlapping sets of genes between filtered subsets were analyzed, the subsets LM+PI and LM+SH (87 and 104 genes respectively) resulted in higher topological congruence to the concatenation-based tree than the SCA of the PI+SH subset (130 genes) despite the lower number of genes analyzed (Fig. 4.4f). These observations provide further evidence that subsampling genes based on LM may be superior to subsampling based on the other two criteria.

## 4.4. Discussion

### 4.4.1. A novel and universally applicable set of DNA-hybridization baits for evolutionary genomic studies of Adephaga

We tested the applicability of the exon-capture approach for locus recovery in a wide range of species in the clade Adephaga. Our orthology assignment results show that the newly designed set of baits can be used to capture the majority of target loci in different species of Adephaga. Our calculated enrichment statistics confirm this result, as they suggest that the coverage of the target regions is generally higher than the coverage of non-target regions which is an indication that the recovery of the target loci was not due to random sequencing, rather due to successful enrichment of the target loci in the species of interest. It should however, be noted that the calculated hybrid-enrichment statistics could have been potentially inflated due to inability of the assembler to include regions of low coverage, resulting in low coverage regions being potentially underrepresented in the assembly relative to high coverage regions. Despite this, we used a genomic assembler that is potentially robust to uneven coverage depth among different genomic regions (Peng et al., 2012). In addition, potential off-target

binding of baits is expected to reduce the actual differences in coverage, therefore balancing out the potential inflation of the calculated enrichment statistics. Hence, our calculated statistics do not provide an exact quantification of the target enrichment in each species, but they rather constitute an approximate comparison of coverage between target and non-target regions in the assemblies, that was used here as a proxy for evaluating target DNA enrichment success.

Despite the success of the hybrid enrichment in all families of Adephaga, the statistics show that the baits may be more successful for enriching the target loci in some families than others. Specifically, the statistics were higher for species in the families Noteridae, Haliplidae, and Carabidae than in the families Dytiscidae, Cicindelidae and Gyrinidae. The observed differences are difficult to evaluate and interpret and could be due to technical factors, such as specimen quality and processing of the samples of species in some families, but also due to biological factors such as the smaller evolutionary distances among the species of these families for the genes analyzed here. It should also be noted that the taxon sampling in some families (e.g., Gyrinidae) was not large enough to provide conclusive evidence on the relative success of target DNA enrichment in these families and therefore our results should be further corroborated in future studies with increased species sampling. In summary, our results show that our newly designed bait set is a valuable resource for future phylogenomic and potentially other evolutionary genomic (such as population genomic) studies of Adephaga. Additionally, given that the Adephaga clade is very old (i.e., the last common ancestor of Adephaga is more than 200 million years old, McKenna et al., 2019), our results show that when available transcriptomic resources are sampled

broadly within the clade of interest, they can be utilized to successfully infer exon-specific baits that are applicable for phylogenetics at deep evolutionary timescales.

**4.4.2. Consolidation of the evolutionary tree of Adephaga by using a combination of transcriptomes and exon-capture sequence data**

Familial relationships of Adephaga in our concatenation-based analyses are generally highly congruent with the most recent phylogenomic studies of Adephaga that are based on analyses of UCEs and transcriptomes (Baca et al., 2017; Gustafson et al., 2020; McKenna et al., 2019). This constitutes further evidence for the phylogenetic utility of our baits at deep evolutionary timescales and helps to further consolidate the phylogenetic pattern of Adephaga. In addition, the results of the analyses of concatenated data largely confirm the pattern of morphological evolution outlined by Beutel et al. (2020). The first split into the highly specialized surface swimming Gyrinidae and the remaining families of Adephaga, suggested for the first time by Beutel and Roughley (1988), is well supported by transformations of larval and adults features and consolidates the paraphyly of "Hydradephaga" as previously suggested by recent concatenated analyses of UCEs (Baca et al., 2017; Gustafson et al., 2020). Recent SCA analyses that resulted in monophyletic "Hydradephaga" do not have any plausibility from a morphological standpoint and did not provide strong clade support for this hypothesis (Freitas et al., 2020). The sampling of Gyrinidae was limited in our study as it did not include *Spanglerogyrus* and *Heterogyrus*, the sister group of the remaining family and of the large subfamily Gyrininae respectively (Beutel et al., 2019b, 2017; Gustafson et al., 2017). A clade comprising Orectochilini and Dineutini, as suggested previously based on morphological data (e.g., Beutel et al., 2006; Beutel

and Roughley, 1993), was confirmed in our concatenated analyses of amino-acids under the best models.

All of our analyses consolidate the monophyly of Geadephaga that is mainly supported by the presence of a specific protibial antenna cleaner and a dense antennal pubescence (Beutel et al., 2006). The clade composed of monophyletic Carabidae and Cicindelidae is well supported by morphological apomorphies, notably by various larval features (Beutel et al., 2020). This result is in agreement with analyses of transcriptomes (McKenna et al., 2019) but not with analyses of UCE data that suggested a clade Trachypachidae + Carabidae (Gustafson et al., 2020). Here we find, that a clade Trachypachidae + Carabidae is only inferred in cases of excessive alignment trimming and only under conditions of suboptimal models or in SCA analyses but it is never strongly supported. The sister group relationship between Trachypachidae and the clade Cicindelidae + Carabidae indicates that a broad procoxal process and broad prothoracic postcoxal bridge are autapomorphies of tiger beetles, in addition to numerous derived features of the highly specialized ambush predating larvae. However, the interpretation of the prothoracic features remains somewhat ambiguous, as the same supposedly derived conditions occur in the wood-associated Rhysodinae (Beutel, 1992a). Evolutionary changes in larvae and adults of Carabidae have been outlined in several studies (Beutel, 1992b, 1992a; Dressler and Beutel, 2010). However, robust molecular phylogeny with a dense taxon sampling of Carabidae is required for a solid reconstruction of the character evolution in this megadiverse lineage. Concerning the phylogeny of Cicindelidae, our inferred tribal relationships are mostly congruent to previous phylogenetic hypotheses of the family, with Manticorini placed as sister to all other tribes (Duran and Gough, 2020; Gough et al., 2020, 2019).

The placement of Haliplidae as sister to Dytiscoidea is in agreement with recent large-scale phylogenomic and morphological studies (Beutel et al., 2020; Gustafson et al., 2020; McKenna et al., 2019), and with Beutel et al. (2013), a study that included extant and extinct lineages of Adephaga. Morphological arguments supporting this clade are sparse, but an important implication is a that the common ancestor invaded the aquatic environment for a second time after Gyrinidae. Aquatic habits in the groundplan of Adephaga would be equally parsimonious, but given the very different adaptations of larvae in these families this appears unlikely. The phylogenetic pattern recovered within Haliplidae is consistent with (Beutel and Ruhnau, 1990), with *Peltodytes* placed as sister to the rest of the family, and *Brychius* as sister to all remaining Haliplidae except *Peltodytes*.

Dytiscoidea are characterized by many well-defined morphological synapomorphies (see Beutel et al., 2020) and our analyses corroborate previous morphological and molecular analyses in that sense (Gustafson et al., 2020; McKenna et al., 2019). The sister group relationship between a clade Noteridae + Meruidae (note: Meruidae was not included in the present study) and the remaining Dytiscoidea is robust (Balke et al., 2008; Beutel et al., 2006), supported for instance by elongate caudal tentorial arms and an entire series of ventral pharyngeal dilators and is corroborated by our results and by other recent phylogenomic studies (Baca et al., 2017; Gustafson et al., 2020; McKenna et al., 2019; Vasilikopoulos et al., 2019). A placement of *Notomicrus* as sister to all other Noteridae, as inferred in our analyses, is compatible with earlier results based on morphology (Belkaceme, 1991; Beutel and Roughley, 1987). However, the taxon sampling of Noteridae in the present study is not sufficient for a reconstruction of the character evolution in Noteridae.

Resolving the phylogenetic relationships of the small families of Dytiscoidea (i.e, Amphizoidae, Aspidytidae, Hygrobiidae) has proven an extremely difficult task (e.g., Cai et al., 2020; Vasilikopoulos et al., 2020, 2019). What likely impedes the reconstruction of phylogenetic relationships of Dytiscoidea based on morphology is a large taxonomic gap caused by the extinction of †Colymbothetidae, †Liadytidae, †Parahygrobiidae, and all subfamilies of †Coptoclavidae, which are likely a polyphyletic assemblage (Beutel et al., 2013). The notoriously difficult placement of the small relict families Amphizoidae, Aspidytidae and Hygrobiidae (e.g., Vasilikopoulos et al., 2019) may also be due to extinction events in these ancient groups, resulting in a drastically reduced extant diversity and relict geographical distributions and potentially also to the accumulation of multiple substitutions along the phylogenetic branches that make modeling of evolutionary processes particularly difficult. Concerning the placement of Hygrobiidae, our SCA analyses agree with most of our concatenation-based analyses under the best models and therefore incongruencies due to ILS do not seem likely. Specifically, our analyses consolidate the monophyly of Aspidytidae and their sister group relationship to Amphizoidae (Cai et al., 2020; Gustafson et al., 2020; Vasilikopoulos et al., 2019), while most concatenated and all SCA analyses suggest Hygrobiidae as sister to Amphizoidae + Aspidytidae in agreement with analyses of UCE data (Gustafson et al., 2020) and with analyses based on a few molecular markers (McKenna et al., 2015; Toussaint et al., 2016). Despite this, the clade comprising Amphizoidae, Aspidytidae and Hygrobiidae is not supported by any solid morphological evidence so far. It implies that the reduction of the duplicatures of the metacoxal plates occurred independently in *Hygrobia* and Dytiscidae, and also the independent acquisition of prothoracic defensive glands

(Beutel et al., 2020). It also implies that the absence of the unusual structures in *Hygrobia* (Beutel, 1986) is due to secondary loss. Despite this, Forsyth (1970) pointed out that prothoracic glands might have evolved independently in members of Dytiscidae and Hygrobiidae. The presence of large and sclerotized epipharyngeal sensorial lobes is a shared derived feature of Dytiscidae, Aspidytidae and Amphizoidae (Dressler and Beutel, 2010). A clade Dytiscidae + (Aspidytidae + Amphizoidae) was previously inferred in analyses of specific subsets of transcriptomic data (Vasilikopoulos et al., 2019) or under conditions of incomplete taxon sampling in analyses of UCE data (Baca et al., 2017; Gustafson et al., 2020). In the present study, a clade Dytiscidae + (Aspidytidae + Amphizoidae) is only inferred under conditions of model misspecification or in some analyses of stringently trimmed supermatrices under the complex BSHETU model but is never strongly supported.

Overall, our results confirm those of Vasilikopoulos et al. (2019) in that removing sites that deviate from the model assumptions (i.e., compositionally heterogeneous genes and hypervariable sites) result in a shift from a strongly or moderately supported Hygrobiidae + (Amphizoidae + Aspidytidae) clade to a less well-supported Dytiscidae + (Aspidytidae + Amphizoidae) clade. This change in topology is here observed only for the less-fitting site-homogeneous models (SHOMU). In contrast to previous phylotranscriptomic analyses (Vasilikopoulos et al., 2019), most phylogenetic reconstructions under the PMSF model resulted in a clade Hygrobiidae + (Amphizoidae + Aspidytidae) irrespective of the trimming threshold and topology of the guide trees. The same applies for most analyses under the better-fitting SHETU models, but also BSHETU models in analyses of moderately trimmed supermatrices under full taxon sampling.

The shift in phylogenetic signal between these two hypotheses is evident in quartet analyses (i.e., FcLM) of both SHETU and SHOMU models. We postulate that this shift in FcLM support under the SHETU model, when more aggressive trimming is applied, is likely due to elimination of useful phylogenetic information because no biasing factors in favor of Hygrobiidae + (Amphizoidae + Aspidytidae) were detected in the quartet analyses of permuted data. Given these observations, we suggest that a clade Dytiscidae + (Aspidytidae + Amphizoidae) seems less likely and probably stems from model misspecification or excessive data removal.

Within the Dytiscidae, we recover all subfamilies as monophyletic and also consolidate the major phylogenetic patterns within these subfamiliar units as suggested by adult and larval morphology (Michat et al., 2017; Miller, 2001), or by Sanger sequencing data combined with morphology (Désamoré et al., 2018; Miller and Bergsten, 2014a). Copelatinae and *Hydrotrupes* were previously assumed be early branches in the tree of Dytiscidae based on their mandibles with open mesal grooves (Beutel, 1994). As in other recent molecular analyses (Miller and Bergsten, 2014b; Ribera et al., 2008), their separate position in the inferred trees implies that larval mandibular sucking channels are an apomorphy of Dytscidae and were secondarily lost in these two groups. This is a possible scenario of character reversal, very likely linked to shifts in larval feeding behavior. Furthermore, our analyses establish Coptotominae as the sister group of Lancetinae. This hypothesis was supported by an evaluation of larval morphology (Michat et al., 2017), who placed a clade Coptotominae + Lancetinae as the sister to Dytiscinae + Cybistrinae. In most analyses under the best models we instead recovered Coptotominae + Lancetinae as the sister to all remaining diving beetles, as opposed to other studies in which Matinae were identified as the

sister to all other Dytiscidae (Désamoré et al., 2018; Miller, 2001). The placement of Matinae as sister to all other Dytiscidae is morphologically established based on their female genital structure (Miller, 2001). Here, we retrieved Matinae nested within the family and as the sister taxon of Hydrodytinae + Hydroporinae. In terms of morphological characters, this could possibly imply a reversal from closer metacoxal lines to more widely separated coxal liones, and also the reversal in the case of the of the separated bursa copulatrix and vagina (Miller, 2001). These observations corroborate Nilsson's (1989) claim, that "The dytiscid phylogeny will most probably be very difficult to reconstruct, because of the widespread convergent evolution." (of morphological characters). This scenario was meticulously discussed in detail by Michat et al. (2017) based on a dataset of 303 larval characters. Overall, previous morphological analyses of Dytiscidae recovered the same major clades as in our study, but identified widespread character homoplasy and ambiguity along the backbone nodes of the tree. The generally unstable backbone topology of Dytiscidae revealed in our study and the inconclusive results of previous morphological studies are clear indications that more careful examination of phylogenomic data, and also a careful re-assessment of the evolution of morphological characters in Dytiscidae is needed.

## 4.4.3. Excessive trimming of supermatrices results in reduced resolution of phylogenetic relationships of Adephaga

Our sensitivity analyses to remove hypervariable sites with different degrees of stringency show that there is a clear trade-off between removing sites that potentially violate the model assumptions and removing sites that contain phylogenetic information. The negative effects of excessive alignment trimming on the phylogenetic

reconstructions have been demonstrated before using different trimming algorithms (Portik and Wiens, 2020; Talavera and Castresana, 2007; Tan et al., 2015). However, the authors of these studies examined loss of phylogenetic information on single genes or loci and not at the supermatrix level. Our FcLM analyses for the examined phylogenetic hypotheses of Adephaga show that the number of resolved quartets and parsimony informative sites decreases in the supermatrices that are trimmed with high degree of stringency. In general, we confirm recent analyses that showed that BMGE trimming of hypervariable sites with very stringent thresholds (e.g., h=0.3) results in reduced phylogenetic accuracy (Steenwyk et al., 2020). Moreover, the overall branch support of our inferred trees under the SHETU model also decreases in the analyses of the stringently trimmed datasets. Lastly, when looking at the phylogenetic branch support for specific hypotheses of Adephaga and outgroups, it is obvious that very stringent BMGE trimming results in poor support or even non-monophyly of some well-established insect clades, such as the clades Coleoptera and Haliplidae + Dytiscoidea. Low and conflicting branch support is a well documented phenomenon for shorter multiple sequence alignments that is due to stochastic error (Delsuc et al., 2005; Phillips et al., 2004) but such phenomena have been mostly observed in phylogenies inferred based on a few or single loci (e.g., Gontcharov et al., 2004). In our study, even the most stringently trimmed supermatrices are long enough (i.e., D and J, > 20,000 amino-acid sites) to be considered genomic-scale datasets, yet the proportion of well-supported clades in their inferred trees is drastically reduced in comparison to less stringently trimmed datasets. This observation further suggests that a balance between removing bias and phylogenetic information should be pursued in phylogenomic analyses.

**4.4.4. Site-heterogeneous models outperform site-homogeneous models and are more robust to the selection of dataset**

Models that account for site-specific amino-acid propensities in the supermatrices by incorporating heterogeneity in the amino-acid equilibrium frequencies among sites (i.e., site-heterogeneous models, SHETU, BSHETU) have been shown to provide a better fir to the data than site-homogeneous models (partitioned or unpartitioned) (Feuda et al., 2017). Our analyses confirm these results although our model selection procedure was not performed in a Bayesian framework to include the most complex site-heterogeneous models (BSHETU) (Lartillot and Philippe, 2004). Despite this, recent research shows that even when the number of equillibrium frequency categories is fixed (e.g., C60 models), the models can potentially describe heterogeneous processes in the data as well as the unconstrained CAT mixture models (Li et al., 2020). An interesting and novel (to our knowledge) outcome of our study is that C60 site-heterogeneous models result in more stable phylogenetic relationships than unpartitioned site-homogeneous models. Specifically, we observed that irrespective of the inferred phylogenetic position of Hygrobiidae under SHOMU model, analyses under the SHETU model (and most analyses under the PMSF model) resulted in a clade Hygrobiidae + (Amphizoidae + Aspidytidae). In addition, comparison of the pairwise RF distances of inferred trees among different models suggests that SHETU models result in more stable phylogenetic relationships of Adephaga and show that analyses under SHETU model are potentially less affected by the trimming or gene selection regimes. Due to computational limitations we were not able to test this hypothesis for the CAT+GTR model as not all analyses reached robust convergence statistics and also we were not able to perform BSHETU analyses for all

datasets. Nevertheless, we suggest that SHETU models may help to reduce incongruence among analyses of different subsets of amino-acid supermatrices. Lastly, we corroborate previous claims that site-homogeneous models underestimate substitution saturation (e.g., Lozano-Fernandez et al., 2019) for a wide selection of amino-acid datasets and trimming regimes. This observation implies that saturation indices that are calculated based on patristic distances are highly dependent on the evolutionary model. Therefore we suggest the employment of alternative substitution saturation measures that are independent of model-based patristic distances of the phylogenetic trees (e.g., Jermiin and Misof, 2020).

## 4.4.5. Gene-tree discordance analyses combined with locus subsampling strategies highlight excessive gene-tree errors in the data

Gene-tree discordance analyses on the complete set of loci but also on the selected subsets of loci suggest that our inferred gene trees are characterized by widespread gene-tree errors. Specifically, the vast majority of gene trees strongly rejected any given well-known clade in Adephaga or in their outgroup but also any alternative phylogenetic hypotheses for the controversial clades of Adephaga. Further indirect evidence for the extent of gene-tree errors in our dataset is provided by observing the distribution of putative phylogenetic information among the inferred gene trees. Many of the inferred gene trees are characterized by low percentage of resolved quartets or very low average branch support and low number of parsimony informative sites, all of which are factors that can result in biases in gene tree estimation. It is frequently assumed that gene-tree discordances are mainly due to biological factors in the data such as ILS (e.g., Cloutier et al., 2019; Linkem et al., 2016). Despite this, we

consider unlikely that ILS has affected all possible deep nodes in the phylogeny of Adephaga and their outgroups and therefore suggest that the observed GTD patterns are probably due to gene-tree errors. This is more apparent when considering that our GTD analyses mostly show strongly rejected alternative phylogenetic hypotheses, rather than strongly supported discordance concerning different phylogenetic hypotheses. Our results confirm the views of other authors who suggest that the biasing effects of biological gene-tree discordance is possible but nevertheless less important than other biasing factors such as model misspecification and gene-tree errors at deep evolutionary timescales (Bryant and Hahn., 2020; Gatesy and Springer, 2014). Although there is no direct evidence from our analyses that the errors affect specific branches of our inferred species tree, our observations suggest that our results of the different SCA analyses cannot be trusted with confidence. This is further corroborated from comparing the distances of the best concatenation-based tree to the trees inferred with SCA using different subsets of genes. These comparisons show that the SCA method is highly sensitive to the set of input gene trees. It is however, encouraging that the SCA could still recover many well-established relationships of Adephaga when all genes are sampled (e.g., Haliplidae + Dytiscoidea, Dytiscoidea, Geadephaga) although some with low support.

It should be noted here that the inability of the SCA to infer congruent results to the concatenation-based tree or strongly supported results might also be related to the small number of genes in the analyzed gene subsets. Specifically, we observed that species trees inferred using the four smallest subsets of genes had the highest topological distance from the concatenation tree. This is in agreement with recent evidence that the ASTRAL method can more accurately infer species trees when

thousands of loci are sampled (Tilic et al., 2020). Furthermore, the potential of increasing the accuracy of summary coalescent analyses by applying empirical site-heterogeneous models (e.g., C10, Quang et al., 2008) for the inference of individual gene trees has to be explored. Despite this, our results show that selecting informative genes based on the likelihood-mapping criterion may be a superior approach to selecting genes based on the number of parsimony informative sites or the average branch support values when scientists want to reduce incongruence to the concatenation-based species tree. This result is in agreement with previous research that suggests likelihood mapping may a good *a priori* estimator of phylogenetic informativeness (Klopfstein et al., 2017).

## 4.5. Conclusions

We provide a novel set of exon-specific DNA-hybridization baits shows great promise in recovering orthologous loci for phylogenomic investigations in different families of Adephaga. Using an extensive sampling of species by combining hybrid-capture data and transcriptomes and we are able to consolidate the phylogenetic relationships of the major groups of Adephaga such as the sister group relationship of Gyrinidae to all other families, a clade Haliplidae + Dytiscoidea, and the sister group relationship Trachypachidae to a clade Carabidae + Cicindelidae. Furthermore, our extensive analyses under different trimming strategies and models shed light on the evolution of the families in Dytiscoidea and show that when moderate trimming and a well-fitting site-heterogeneous model is used, Hygrobiidae is recovered as sister to Amphizoidae + Aspidytidae. Excessive removal of hypervariable sites using stringent trimming strategies should be avoided as it can lead to reduction in phylogenetic signal

and reduced resolution of phylogenetic relationships, as we observed here for the phylogeny of Adephaga. Site-heterogeneous models always fit the data better but most interestingly our results show that analyses with C60 site-heterogeneous models result in increased stability of inferred phylogenetic relationships of Dytiscoidea and Adephaga in general. Therefore, incongruence between analyses of different subsets of amino-acid supermatrices may be ameliorated with the use of C60 site-heterogeneous models. Moreover, our analyses of a carefully curated set of genes suggest that gene-tree errors are prominent in the data and possibly responsible for poorly supported or incongruent species trees in SCA analyses or for incongruence between concatenation and SCA. Hence, our results show that scientists should take measures to eliminate or minimize gene-tree errors before attributing gene-tree discordance and phylogenomic incongruence to other factors (e.g., ILS). As we have shown, a promising solution for reducing incongruence between coalescent-based and concatenation-based analyses is to select informative genes based on the likelihood mapping criterion.

## 4.6. References

Akaike, H., 1974. A new look at the statistical model identification. IEEE Trans. Automat. Contr. 19, 716–723.

Baca, S.M., Alexander, A., Gustafson, G.T., Short, A.E.Z., 2017. Ultraconserved elements show utility in phylogenetic inference of Adephaga (Coleoptera) and suggest paraphyly of 'Hydradephaga'. Syst. Entomol. 42, 786–795.

Balke, M., Ribera, I., Beutel, R., Viloria, A., Garcia, M., Vogler, A.P., 2008. Systematic placement of the recently discovered beetle family Meruidae (Coleoptera: Dytiscoidea) based on molecular data. Zool. Scr. 37, 647–650.

Balke, M., Ribera, I., Beutel, R.G., 2003. ASPIDYTIDAE: on the discovery of a new beetle family: detailed morphological analysis, description of a second species, and key to fossil and extant adephagan families (Coleoptera), in: Jach, M.A., Ji, L. (Eds.), Water Beetles of China. Zoologisch-Botanische Gesellschaft & Wiener Coleopterologenverein, pp. 53–66.

Bank, S., Sann, M., Mayer, C., Meusemann, K., Donath, A., Podsiadlowski, L., Kozlov, A., Petersen, M., Krogmann, L., Meier, R., Rosa, P., Schmitt, T., Wurdack, M., Liu, S., Zhou, X., Misof, B., Peters, R.S., Niehuis, O., 2017. Transcriptome and target DNA enrichment sequence data provide new insights into the phylogeny of vespid wasps (Hymenoptera: Aculeata: Vespidae). Mol. Phylogenet. Evol. 116, 213–226.

Belkaceme, T., 1991. Skelet und Muskulatur des Kopfes und Thorax von *Noterus laevis* Sturm. Ein Beitrag zur Morphologie und Phylogenie der Noteridae (Coleoptera: Adephaga). Stuttgarter Beiträge zur Naturkunde, Ser. A. 462, 1–94.

Beutel, R.G., 1994. On the systematic position of *Hydrotrupes palpalis* Sharp (Coleoptera: Dytiscidae). Aquat. Insects 16, 157–164.

Beutel, R.G., 1992a. Larval head structures of *Omoglymmius hamatus* and their implications for the relationships of Rhysodidae (Coleoptera: Adephaga). Insect Syst. Evol. 23, 169–184.

Beutel, R.G., 1992b. Phylogenetic analysis of thoracic structures of Carabidae (Coleoptera: Adephaga). J. Zool. Syst. Evol. Res. 30, 53–74.

Beutel, R.G., 1986. Skelet und Muskulatur des Kopfes und Thorax von *Hygrobia tarda* (Herbst). Ein Beitrag zur Klärung der phylogenetischen Beziehungen der Hydradephaga (Insecta: Coleoptera). Stutt. Beitr. Naturkd. 388, 1–54.

Beutel, R.G., Balke, M., Steiner, W.E., 2006. The systematic position of Meruidae (Coleoptera, Adephaga) and the phylogeny of the smaller aquatic adephagan beetle families. Cladistics 22, 102–131.

Beutel, R.G., Pohl, H., Yan, E. V., Anton, E., Liu, S.-P., Ślipiński, A., McKenna, D., Friedrich, F., 2019a. The phylogeny of Coleopterida (Hexapoda) – morphological characters and molecular phylogenies. Syst. Entomol. 44, 75–102.

Beutel, R.G., Ribera, I., Fikáček, M., Vasilikopoulos, A., Misof, B., Balke, M., 2020. The morphological evolution of the Adephaga (Coleoptera). Syst. Entomol. 45, 378–395.

Beutel, R.G., Roughley, R.E., 1993. Phylogenetic analysis of Gyrinidae based on characters of the larval head (Coleoptera: Adephaga). Insect Syst. Evol. 24, 459–468.

Beutel, R.G., Roughley, R.E., 1988. On the systematic position of the family Gyrinidae (Coleoptera: Adephaga). J. Zool. Syst. Evol. Res. 26, 380–400.

Beutel, R.G., Roughley, R.E., 1987. On the systematic position of the genus *Notomicrus* Sharp (Hydradephaga, Coleoptera) . Can. J. Zool. 65, 1898–1905.

Beutel, R.G., Ruhnau, S., 1990. Phylogenetic analysis of the genera of Haliplidae (Coleoptera) based on characters of adults. Aquat. Insects 12, 1–17.

Beutel, R.G., Wang, B., Tan, J.J., Ge, S.Q., Ren, D., Yang, X.K., 2013. On the phylogeny and evolution of Mesozoic and extant lineages of Adephaga (Coleoptera, Insecta). Cladistics 29, 147–165.

Beutel, R.G., Yan, E., Richter, A., Büsse, S., Miller, K.B., Yavorskaya, M., Wipfler, B., 2017. The head of Heterogyrus milloti (Coleoptera: Gyrinidae) and its phylogenetic implications. Arthropod Syst. Phylogeny 75, 261–280.

Beutel, R.G., Yan, E., Yavorskaya, M., Büsse, S., Gorb, S.N., Wipfler, B., 2019b. On the thoracic anatomy of the Madagascan *Heterogyrus milloti* and the phylogeny of Gyrinidae (Coleoptera). Syst. Entomol. 44, 336–360.

Bi, K., Vanderpool, D., Singhal, S., Linderoth, T., Moritz, C., Good, J.M., 2012. Transcriptome-based exon capture enables highly cost-effective comparative genomic data collection at moderate evolutionary scales. BMC Genomics 13, 403.

Bolger, A.M., Lohse, M., Usadel, B., 2014. Trimmomatic: a flexible trimmer for Illumina sequence data. Bioinformatics 30, 2114–2120.

Bowker, A.H., 1948. A test for symmetry in contingency tables. J. Am. Stat. Assoc. 43, 572–574.

Bradley, R.K., Roberts, A., Smoot, M., Juvekar, S., Do, J., Dewey, C., Holmes, I., Pachter, L., 2009. Fast statistical alignment. PLoS Comput. Biol. 5, e1000392.

Bragg, J.G., Potter, S., Bi, K., Moritz, C., 2016. Exon capture phylogenomics: efficacy across scales of divergence. Mol. Ecol. Resour. 16, 1059–1068.

Bryant, D., Hahn., M.W., 2020. The concatenation question., in: Scornavacca, C., Delsuc, F., Galtier, N. (Eds.), Phylogenetics in the genomic era. No commercial publisher | Authors' open access book, pp. 3.4:1–3.4:23.

Burmeister, E.-G., 1976. Der Ovipositor der Hydradephaga (Coleoptera) und seine phylogenetische Bedeutung unter besonderer Berücksichtigung der Dytiscidae. Zoomorphologie 85, 165–257.

Cai, C., Tihelka, E., Pisani, D., Donoghue, P.C.J., 2020. Data curation and modeling of compositional heterogeneity in insect phylogenomics: A case study of the phylogeny of Dytiscoidea (Coleoptera: Adephaga). Mol. Phylogenet. Evol. 147, 106782.

Chernomor, O., von Haeseler, A., Minh, B.Q., 2016. Terrace aware data structure for phylogenomic inference from supermatrices. Syst. Biol. 65, 997–1008.

Cloutier, A., Sackton, T.B., Grayson, P., Clamp, M., Baker, A.J., Edwards, S. V., 2019. Whole-genome analyses resolve the phylogeny of flightless birds (Palaeognathae) in the presence of an empirical anomaly zone. Syst. Biol. 68, 937–955.

Criscuolo, A., Gribaldo, S., 2010. BMGE (Block Mapping and Gathering with Entropy): a new software for selection of phylogenetic informative regions from multiple sequence alignments. BMC Evol. Biol. 10, 210.

Crotty, S.M., Minh, B.Q., Bean, N.G., Holland, B.R., Tuke, J., Jermiin, L.S., von Haeseler, A., 2020. GHOST: recovering historical signal from heterotachously evolved sequence alignments. Syst. Biol. 69, 249–264.

Crowson, R.A., 1960. The phylogeny of Coleoptera. Annu. Rev. Entomol. 5, 111–134.

Delsuc, F., Brinkmann, H., Philippe, H., 2005. Phylogenomics and the reconstruction of the tree of life. Nat. Rev. Genet. 6, 361–375.

Désamoré, A., Laenen, B., Miller, K.B., Bergsten, J., 2018. Early burst in body size evolution is uncoupled from species diversification in diving beetles (Dytiscidae). Mol. Ecol. 27, 979–993.

Dietz, L., Dömel, J.S., Leese, F., Mahon, A.R., Mayer, C., 2019. Phylogenomics of the longitarsal Colossendeidae: The evolutionary history of an Antarctic sea spider radiation. Mol. Phylogenet. Evol. 136, 206–214.

Dressler, C., Beutel, R.G., 2010. The morphology and evolution of the adult head of Adephaga (Insecta: Coleoptera). Arthropod Syst. Phylogeny 68, 239–287.

Duran, D.P.., Gough, H.M., 2020. Validation of tiger beetles as distinct family (Coleoptera: Cicindelidae), review and reclassification of tribal relationships. Syst. Entomol. 45, 723–729.

Faircloth, B.C., 2017. Identifying conserved genomic elements and designing universal bait sets to enrich them. Methods Ecol. Evol. 8, 1103–1112.

Faircloth, B.C., 2016. PHYLUCE is a software package for the analysis of conserved genomic loci. Bioinformatics 32, 786–788.

Faircloth, B.C., McCormack, J.E., Crawford, N.G., Harvey, M.G., Brumfield, R.T., Glenn, T.C., 2012. Ultraconserved elements anchor thousands of genetic markers spanning multiple evolutionary timescales. Syst. Biol. 61, 717–726.

Feuda, R., Dohrmann, M., Pett, W., Philippe, H., Rota-Stabelli, O., Lartillot, N., Wörheide, G., Pisani, D., 2017. Improved Modeling of Compositional Heterogeneity Supports Sponges as Sister to All Other Animals. Curr. Biol. 27, 3864–3870.

Forsyth, D.J., 1970. The structure of the defence glands of the Cicindelidae, Amphizoidae, and Hygrobiidae (Insecta: Coleoptera). J. Zool. 160, 51–69.

Freitas, F. V, Branstetter, M.G., Griswold, T., Almeida, E.A.B., 2020. Partitioned gene-tree analyses and gene-based topology testing help resolve incongruence in a phylogenomic study of host-specialist bees (Apidae: Eucerinae). Mol. Biol. Evol. doi: 10.1093/molbev/msaa277.

Gatesy, J., Springer, M.S., 2014. Phylogenetic analysis at deep timescales: Unreliable gene trees, bypassed hidden support, and the coalescence/concatalescence conundrum. Mol. Phylogenet. Evol. 80, 231–266.

Glenn, T.C., Faircloth, B.C., 2016. Capturing Darwin's dream. Mol. Ecol. Resour. 16, 1051–1058.

Gontcharov, A.A., Marin, B., Melkonian, M., 2004. Are combined analyses better than single gene phylogenies? A case study using SSU rDNA and rbcL sequence comparisons in the Zygnematophyceae (Streptophyta). Mol. Biol. Evol. 21, 612–624.

Gough, H.M., Allen, J.M., Toussaint, E.F.A., Storer, C.G., Kawahara, A.Y., 2020. Transcriptomics illuminate the phylogenetic backbone of tiger beetles. Biol. J. Linn. Soc. 129, 740–751.

Gough, H.M., Duran, D.P., Kawahara, A.Y., Toussaint, E.F.A., 2019. A comprehensive molecular phylogeny of tiger beetles (Coleoptera, Carabidae, Cicindelinae). Syst. Entomol. 44, 305–321.

Guindon, S., Dufayard, J.F., Lefort, V., Anisimova, M., Hordijk, W., Gascuel, O., 2010. New algorithms and methods to estimate maximum-likelihood phylogenies: assessing the performance of PhyML 3.0. Syst. Biol. 59, 307–321.

Gustafson, G.T., Alexander, A., Sproul, J.S., Pflug, J.M., Maddison, D.R., Short, A.E.Z., 2019. Ultraconserved element (UCE) probe set design: base genome and initial design parameters critical for optimization. Ecol. Evol. 9, 6933–6948.

Gustafson, G.T., Baca, S.M., Alexander, A.M., Short, A.E.Z., 2020. Phylogenomic analysis of the beetle suborder Adephaga with comparison of tailored and generalized ultraconserved element probe performance. Syst. Entomol. 45, 552–570.

Gustafson, G.T., Bergsten, J., Ranarilalatiana, T., Randriamihaja, J.H., Miller, K.B., 2017. The morphology and behavior of the endemic Malagasy whirligig beetle *Heterogyrus milloti* Legros, 1953 (Coleoptera: Gyrinidae: Heterogyrinae). Coleopt. Bull. 71, 315–328.

Hoang, D.T., Chernomor, O., von Haeseler, A., Minh, B.Q., Le, S.V., 2018. UFBoot2: Improving the Ultrafast Bootstrap Approximation. Mol. Biol. Evol. 35, 518–522.

Hosner, P.A., Faircloth, B.C., Glenn, T.C., Braun, E.L., Kimball, R.T., 2016. Avoiding missing data biases in phylogenomic inference: an empirical study in the landfowl (Aves: Galliformes). Mol. Biol. Evol. 33, 1110–1125.

Huerta-Cepas, J., Serra, F., Bork, P., 2016. ETE 3: Reconstruction, analysis, and visualization of phylogenomic data. Mol. Biol. Evol. 33, 1635–1638.

Hunt, T., Bergsten, J., Levkanicova, Z., Papadopoulou, A., John, O.S., Wild, R., Hammond, P.M., Ahrens, D., Balke, M., Caterino, M.S., Gómez-zurita, J., Ribera, I., Barraclough, T.G., Bocakova, M., Bocak, L., Vogler, A.P., 2007. A comprehensive phylogeny of beetles reveals the evolutionary origins of a superradiation. Science. 318, 1913–1916.

Jäch, M.A., Balke, M., 2008. Global diversity of water beetles (Coleoptera) in freshwater. Hydrobiologia 595, 419–442.

Jermiin, L.S., Misof, B., 2020. Measuring historical and compositional signals in phylogenetic data. BioRxiv. doi: 10.1101/2020.01.03.894097.

Kalyaanamoorthy, S., Minh, B.Q., Wong, T.K.F., von Haeseler, A., Jermiin, L.S., 2017. ModelFinder: fast model selection for accurate phylogenetic estimates. Nat. Methods 14, 587–589.

Karin, B.R., Gamble, T., Jackman, T.R., Vidal, N., 2020. Optimizing phylogenomics with rapidly evolving long exons: comparison with anchored hybrid enrichment and ultraconserved elements. Mol. Biol. Evol. 37, 904–922.

Kavanaugh, D.H., 1986. A systematic review of Amphizoid beetles (Amphizoidae: Coleoptera) and their phylogenetic relationships to other Adephaga. Proc. Calif. Acad. Sci. 44, 67–109.

Klopfstein, S., Massingham, T., Goldman, N., 2017. More on the best evolutionary rate for phylogenetic analysis. Syst. Biol. 66, 769–785.

Kocot, K.M., Struck, T.H., Merkel, J., Waits, D.S., Todt, C., Brannock, P.M., Weese, D.A., Cannon, J.T., Moroz, L.L., Lieb, B., Halanych, K.M., 2017. Phylogenomics of Lophotrochozoa with consideration of systematic error. Syst. Biol. 66, 256–282.

Kück, P., Longo, G.C., 2014. FASconCAT-G: extensive functions for multiple sequence alignment preparations concerning phylogenetic studies. Front. Zool. 11, 81.

Kück, P., Meid, S.A., Groß, C., Wägele, J.W., Misof, B., 2014. AliGROOVE-- visualization of heterogeneous sequence divergence within multiple sequence alignments and detection of inflated branch support. BMC Bioinformatics 15, 294.

Kück, P., Meusemann, K., Dambach, J., Thormann, B., von Reumont, B.M., Wägele, J.W., Misof, B., 2010. Parametric and non-parametric masking of randomness in sequence alignments can be improved and leads to better resolved trees. Front. Zool. 7, 10.

Kück, P., Struck, T.H., 2014. BaCoCa - A heuristic software tool for the parallel assessment of sequence biases in hundreds of gene and taxon partitions. Mol. Phylogenet. Evol. 70, 94–98.

Lanfear, R., Calcott, B., Kainer, D., Mayer, C., Stamatakis, A., 2014. Selecting optimal partitioning schemes for phylogenomic datasets. BMC Evol. Biol. 14, 82.

Lartillot, N., Philippe, H., 2004. A Bayesian mixture model for across-site heterogeneities in the amino-acid replacement process. Mol. Biol. Evol. 21, 1095–1109.

Lartillot, N., Rodrigue, N., Stubbs, D., Richer, J., 2013. PhyloBayes MPI : phylogenetic reconstruction with infinite mixtures of profiles in a parallel environment. Syst. Biol. 62, 611–615.

Laumer, C.E., Fernández, R., Lemer, S., Combosch, D., Kocot, K.M., Riesgo, A., Andrade, S.C.S., Sterrer, W., Sørensen, M. V., Giribet, G., 2019. Revisiting metazoan phylogeny with genomic sampling of all phyla. Proc. R. Soc. B Biol. Sci. 286, 20190831.

Lawrence, J.F., Newton, A.F., 1982. Evolution and classification of beetles. Annu. Rev. Ecol. Syst. 13, 261–290.

Lawrence, J.F., Ślipiński, A., Seago, A.E., Thayer, M.K., Newton, A.F., Marvaldi, A.E., 2011. Phylogeny of the Coleoptera based on morphological characters of adults and larvae. Ann. Zool. 61, 1–217.

Lemmon, A.R., Brown, J.M., Stanger-Hall, K., Lemmon, E.M., 2009. The effect of ambiguous data on phylogenetic estimates obtained by maximum likelihood and bayesian inference. Syst. Biol. 58, 130–145.

Lemmon, A.R., Emme, S.A., Lemmon, E.M., 2012. Anchored hybrid enrichment for massively high-throughput phylogenomics. Syst. Biol. 61, 727–744.

Lemmon, E.M., Lemmon, A.R., 2013. High-throughput genomic data in systematics and phylogenetics. Annu. Rev. Ecol. Evol. Syst 44, 99–121.

Li, C., Hofreiter, M., Straube, N., Corrigan, S., Naylor, G.J.P., 2013. Capturing protein-coding genes across highly divergent species. Biotechniques 54, 321–326.

Li, H., Durbin, R., 2009. Fast and accurate short read alignment with Burrows–Wheeler transform. Bioinformatics 25, 1754–1760.

Li, H., Handsaker, B., Wysoker, A., Fennell, T., Ruan, J., Homer, N., Marth, G., Abecasis, G., Durbin, R., 2009. The sequence alignment/map format and SAMtools. Bioinformatics 25, 2078–2079.

Li, Y., Shen, X.-X., Evans, B., Dunn, C.W., Rokas, A., 2020. Rooting the animal tree of life. BioRxiv. doi: 10.1101/2020.10.27.357798

Linkem, C.W., Minin, V.N., Leaché, A.D., 2016. Detecting the anomaly zone in species trees and evidence for a misleading signal in higher-level skink phylogeny (Squamata: Scincidae). Syst. Biol. 65, 465–477.

López-López, A., Vogler, A.P., 2017. The mitogenome phylogeny of Adephaga (Coleoptera). Mol. Phylogenet. Evol. 114, 166–174.

Lorenz, W., 2020. CarabCat: Global database of ground beetles (version Oct 2017) [WWW Document]. Species 2000 ITIS Cat. Life, 2020-12-01. URL https://www.catalogueoflife.org/ (accessed 12.11.20).

Lozano-Fernandez, J., Tanner, A.R., Giacomelli, M., Carton, R., Vinther, J., Edgecombe, G.D., Pisani, D., 2019. Increasing species sampling in chelicerate genomic-scale datasets provides support for monophyly of Acari and Arachnida. Nat. Commun. 10, 2295.

Maddison, D.R., Moore, W., Baker, M.D., Ellis, T.M., Ober, K.A., Cannone, J.J., Gutell, R.R., 2009. Monophyly of terrestrial adephagan beetles as indicated by three nuclear genes (Coleoptera: Carabidae and Trachypachidae). Zool. Scr. 38, 43–62.

Mayer, C., Sann, M., Donath, A., Meixner, M., Podsiadlowski, L., Peters, R.S., Petersen, M., Meusemann, K., Liere, K., Wägele, J.-W., Misof, B., Bleidorn, C., Ohl, M., Niehuis, O., 2016. BaitFisher: A software package for multi-species target DNA enrichment probe design. Mol. Biol. Evol. 33, 1875–1886.

McKenna, D.D., Shin, S., Ahrens, D., Balke, M., Beza-Beza, C., Clarke, D.J., Donath, A., Escalona, H.E., Friedrich, F., Letsch, H., Liu, S., Maddison, D., Mayer, C., Misof, B., Murin, P.J., Niehuis, O., Peters, R.S., Podsiadlowski, L., Pohl, H., Scully, E.D., Yan, E. V, Zhou, X., Ślipiński, A., Beutel, R.G., 2019. The evolution and genomic basis of beetle diversity. Proc. Natl. Acad. Sci. U. S. A. 116, 24729–24737.

McKenna, D.D., Wild, A.L., Kanda, K., Bellamy, C.L., Beutel, R.G., Caterino, M.S., Farnum, C.W., Hawks, D.C., Ivie, M.A., Jameson, M.L., Leschen, R.A.B., Marvaldi, A.E., Mchugh, J. V, Newton, A.F., Robertson, J.A., Thayer, M.K., Whiting, M.F., Lawrence, J.F., Ślipiński, A., Maddison, D.R., Farrell, B.D., 2015. The beetle tree of life reveals that Coleoptera survived end-Permian mass extinction to diversify during the Cretaceous terrestrial revolution. Syst. Entomol. 40, 835–880.

Meiklejohn, K.A., Faircloth, B.C., Glenn, T.C., Kimball, R.T., Braun, E.L., 2016. Analysis of a rapid evolutionary radiation using ultraconserved elements: evidence for a bias in some multispecies coalescent methods. Syst. Biol. 65, 612–627.

Michat, M.C., Alarie, Y., Miller, K.B., 2017. Higher-level phylogeny of diving beetles (Coleoptera: Dytiscidae) based on larval characters. Syst. Entomol. 42, 734–767.

Miller, K.B., 2001. On the phylogeny of the Dytiscidae (insecta: Coleoptera) with emphasis on the morphology of the female reproductive system. Insect Syst. Evol. 32, 45–92.

Miller, K.B., Bergsten, J., 2014. The phylogeny and classification of predaceous diving beetles. In: Yee, D. (Ed.), Ecology, systematics, and the natural history of predaceous diving beetles (Coleoptera: Dytiscidae). Springer, Dordrecht, pp. 49–172.

Minh, B.Q., Schmidt, H.A., Chernomor, O., Schrempf, D., Woodhams, M.D., von Haeseler, A., Lanfear, R., 2020. IQ-TREE 2: new models and efficient methods for phylogenetic inference in the genomic era. Mol. Biol. Evol. 37, 1530–1534.

Mirarab, S., Bayzid, M.S., Warnow, T., 2016. Evaluating summary methods for multilocus species tree estimation in the presence of incomplete lineage sorting. Syst. Biol. 65, 366–380.

Misof, B., Liu, S., Meusemann, K., Peters, R.S., Donath, A., Mayer, C., Frandsen, P.B., Ware, J., Flouri, T., Beutel, R.G., Niehuis, O., Petersen, M., Izquierdo-Carrasco, F., Wappler, T., Rust, J., Aberer, A.J., Aspöck, U., Aspöck, H., Bartel, D., Blanke, A., Berger, S., Bohm, A., Buckley, T.R., Calcott, B., Chen, J., Friedrich, F., Fukui, M., Fujita, M., Greve, C., Grobe, P., Gu, S., Huang, Y., Jermiin, L.S., Kawahara, A.Y., Krogmann, L., Kubiak, M., Lanfear, R., Letsch, H., Li, Y., Li, Z., Li, J., Lu, H., Machida, R., Mashimo, Y., Kapli, P., McKenna, D.D., Meng, G., Nakagaki, Y., Navarrete-Heredia, J.L., Ott, M., Ou, Y., Pass, G., Podsiadlowski, L., Pohl, H., von Reumont, B.M., Schutte, K., Sekiya, K., Shimizu, S., Slipinski, A., Stamatakis, A., Song, W., Su, X., Szucsich, N.U., Tan, M., Tan, X., Tang, M., Tang, J.,

Timelthaler, G., Tomizuka, S., Trautwein, M., Tong, X., Uchifune, T., Walzl, M.G., Wiegmann, B.M., Wilbrandt, J., Wipfler, B., Wong, T.K.F., Wu, Q., Wu, G., Xie, Y., Yang, S., Yang, Q., Yeates, D.K., Yoshizawa, K., Zhang, Q., Zhang, R., Zhang, W., Zhang, Y., Zhao, J., Zhou, C., Zhou, L., Ziesmann, T., Zou, S., Xu, X., Yang, H., Wang, J., Kjer, K.M., Zhou, X., 2014a. Phylogenomics resolves the timing and pattern of insect evolution. Science. 346, 763–767.

Misof, B., Meusemann, K., von Reumont, B.M., Kück, P., Prohaska, S.J., Stadler, P.F., 2014b. A priori assessment of data quality in molecular phylogenetics. Algorithms Mol. Biol. 9, 22.

Misof, B., Meyer, B., von Reumont, B.M., Kück, P., Misof, K., Meusemann, K., 2013. Selecting informative subsets of sparse supermatrices increases the chance to find correct trees. BMC Bioinformatics 14, 348.

Misof, B., Misof, K., 2009. A Monte Carlo approach successfully identifies randomness in multiple sequence alignments: a more objective means of data exclusion. Syst. Biol. 58, 21–34.

Misof, B., Rickert, A.M., Buckley, T.R., Fleck, G., Sauer, K.P., 2001. Phylogenetic signal and its decay in mitochondrial SSU and LSU rRNA gene fragments of Anisoptera. Mol. Biol. Evol. 18, 27–37.

Naser-Khdour, S., Minh, B.Q., Zhang, W., Stone, E.A., Lanfear, R., Bryant, D., 2019. The prevalence and impact of model violations in phylogenetic analysis. Genome Biol. Evol. 11, 3341–3352.

Nguyen, L.-T., Schmidt, H.A., Von Haeseler, A., Minh, B.Q., 2015. IQ-TREE: a fast and effective stochastic algorithm for estimating maximum-likelihood phylogenies. Mol. Biol. Evol. 32, 268–274.

Nilsson, A.N., 1989. On the genus *Agabetes* Crotch (Coleoptera, Dytiscidae), with a new species from Iran. Ann. Entomol. Fenn. 55, 35–40.

Nosenko, T., Schreiber, F., Adamska, M., Adamski, M., Eitel, M., Hammel, J., Maldonado, M., Müller, W.E.G., Nickel, M., Schierwater, B., Vacelet, J., Wiens, M., Wörheide, G., 2013. Deep metazoan phylogeny: when different genes tell different stories. Mol. Phylogenet. Evol. 67, 223–233.

Peng, Y., Leung, H.C.M., Yiu, S.M., Chin, F.Y.L., 2012. IDBA-UD: A *de novo* assembler for single-cell and metagenomic sequencing data with highly uneven depth. Bioinformatics 28, 1420–1428.

Petersen, M., Meusemann, K., Donath, A., Dowling, D., Liu, S., Peters, R.S., Podsiadlowski, L., Vasilikopoulos, A., Zhou, X., Misof, B., Niehuis, O., 2017. Orthograph: a versatile tool for mapping coding nucleotide sequences to clusters of orthologous genes. BMC Bioinformatics 18, 111.

Philippe, H., Derelle, R., Lopez, P., Pick, K., Borchiellini, C., Boury-Esnault, N., Vacelet, J., Renard, E., Houliston, E., Quéinnec, E., Da Silva, C., Wincker, P., Le Guyader, H., Leys, S., Jackson, D.J., Schreiber, F., Erpenbeck, D., Morgenstern, B., Wörheide, G., Manuel, M., 2009. Phylogenomics revives traditional views on deep animal relationships. Curr. Biol. 19, 706–712.

Phillips, M.J., Delsuc, F., Penny, D., 2004. Genome-scale phylogeny and the detection of systematic biases. Mol. Biol. Evol. 21, 1455–1458.

Pisani, D., Pett, W., Dohrmann, M., Feuda, R., Rota-Stabelli, O., Philippe, H., Lartillot, N., Wörheide, G., 2015. Genomic data do not support comb jellies as the sister group to all other animals. Proc. Natl. Acad. Sci. U.S.A. 112, 15402–15407.

Portik, D.M., Wiens, J.J., 2020. Do alignment and trimming methods matter for phylogenomic (UCE) analyses? Syst. Biol. doi: 10.1093/sysbio/syaa064

Quang, L.S., Gascuel, O., Lartillot, N., 2008. Empirical profile mixture models for phylogenetic reconstruction. Bioinformatics 24, 2317–2323.

R Core Team, 2020. R: A language and environment for statistical computing.

Ranwez, V., Douzery, E.J.P., Cambon, C., Chantret, N., Delsuc, F., 2018. MACSE v2: toolkit for the alignment of coding sequences accounting for frameshifts and stop codons. Mol. Biol. Evol. 35, 2582–2584.

Ribera, I., Vogler, A.P., Balke, M., 2008. Phylogeny and diversification of diving beetles (Coleoptera: Dytiscidae). Cladistics 24, 563–590.

Robinson, D.F., Foulds, L.R., 1981. Comparison of phylogenetic trees. Math. Biosci. 53, 131–147.

Sayyari, E., Whitfield, J.B., Mirarab, S., 2018. DiscoVista: interpretable visualizations of gene tree discordance. Mol. Phylogenet. Evol. 122, 110–115.

Sayyari, E., Whitfield, J.B., Mirarab, S., 2017. Fragmentary gene sequences negatively impact gene tree and species tree reconstruction. Mol. Biol. Evol. 34, 3279–3291.

Sharma, P.P., Kaluziak, S.T., Pérez-Porro, A.R., González, V.L., Hormiga, G., Wheeler, W.C., Giribet, G., 2014. Phylogenomic interrogation of Arachnida reveals systemic conflicts in phylogenetic signal. Mol. Biol. Evol. 31, 2963–2984.

Short, A.E.Z., 2018. Systematics of aquatic beetles (Coleoptera): current state and future directions. Syst. Entomol. 43, 1–18.

Shull, V.L., Vogler, A.P., Baker, M.D., Maddison, D.R., Hammond, P.M., 2001. Sequence alignment of 18S ribosomal RNA and the basal relationships of adephagan beetles: evidence for monophyly of aquatic families and the placement of Trachypachidae Syst. Biol. 50, 945–969.

Simion, P., Belkhir, K., François, C., Veyssier, J., Rink, J.C., Manuel, M., Philippe, H., Telford, M.J., 2018. A software tool 'CroCo' detects pervasive cross-species contamination in next generation sequencing data. BMC Biol. 16, 28.

Spangler, P.J., Steiner, W.E., 2005. A new aquatic beetle family, Meruidae, from Venezuela (Coleoptera: Adephaga). Syst. Entomol. 30, 339–357.

Steenwyk, J.L., Buida III, T.J., Li, Y., Shen, X.-X., Rokas, A., 2020. ClipKIT: a multiple sequence alignment trimming software for accurate phylogenomic inference. PLOS Biol. 18, e3001007.

Strimmer, K., von Haeseler, A., 1997. Likelihood-mapping: a simple method to visualize phylogenetic content of a sequence alignment. Proc. Natl. Acad. Sci. U. S. A. 94, 6815–6819.

Struck, T.H., 2014. TreSpEx-Detection of misleading signal in phylogenetic reconstructions based on tree information. Evol. Bioinforma. 10, 51–67.

Stuart, A., 1955. A test for homogeneity of the marginal distributions in a two-way classification. Biometrika 42, 412–416.

Suyama, M., Torrents, D., Bork, P., 2006. PAL2NAL: robust conversion of protein sequence alignments into the corresponding codon alignments. Nucleic Acids Res. 34, 609–612.

Talavera, G., Castresana, J., 2007. Improvement of phylogenies after removing divergent and ambiguously aligned blocks from protein sequence alignments. Syst. Biol. 56, 564–577.

Tan, G., Muffato, M., Ledergerber, C., Herrero, J., Goldman, N., Gil, M., Dessimoz, C., 2015. Current methods for automated filtering of multiple sequence alignments frequently worsen single-gene phylogenetic inference. Syst. Biol. 64, 778–791.

Tilic, E., Sayyari, E., Stiller, J., Mirarab, S., Rouse, G.W., 2020. More is needed—thousands of loci are required to elucidate the relationships of the 'flowers of the sea' (Sabellida, Annelida). Mol. Phylogenet. Evol. 151, 106892.

Toussaint, E.F.A., Beutel, R.G., Morinière, J., Jia, F., Xu, S., Michat, M.C., Zhou, X., Bilton, D.T., Ribera, I., Hájek, J., Balke, M., 2016. Molecular phylogeny of the highly disjunct cliff water beetles from South Africa and China (Coleoptera: Aspidytidae). Zool. J. Linn. Soc. 176, 537–546.

Vasilikopoulos, A., Balke, M., Beutel, R.G., Donath, A., Podsiadlowski, L., Pflug, J.M., Waterhouse, R.M., Meusemann, K., Peters, R.S., Escalona, H.E., Mayer, C., Liu, S., Hendrich, L., Alarie, Y., Bilton, D.T., Jia, F., Zhou, X., Maddison, D.R., Niehuis, O., Misof, B., 2019. Phylogenomics of the superfamily Dytiscoidea (Coleoptera: Adephaga) with an evaluation of phylogenetic conflict and systematic error. Mol. Phylogenet. Evol. 135, 270–285.

Vasilikopoulos, A., Gustafson, G.T., Balke, M., Niehuis, O., Beutel, R.G., Misof, B., 2020. Resolving the phylogenetic position of Hygrobiidae (Coleoptera: Adephaga) requires objective statistical tests and exhaustive phylogenetic methodology: a response to Cai et al. (2020). Mol. Phylogenet. Evol. 106923.

Wang, H.-C., Minh, B.Q., Susko, E., Roger, A.J., 2018. Modeling site heterogeneity with posterior mean site frequency profiles accelerates accurate phylogenomic estimation. Syst. Biol. 67, 216–235.

Wang, H.-C., Susko, E., Roger, A.J., 2019. The relative importance of modeling site pattern heterogeneity versus partition-wise heterotachy in phylogenomic inference. Syst. Biol. 68, 1003–1019.

Whelan, S., Irisarri, I., Burki, F., 2018. PREQUAL: detecting non-homologous characters in sets of unaligned homologous sequences. Bioinformatics 34, 3929–3930.

Wong, T.K.F., Kalyaanamoorthy, S. Meusemann, K., Yeates, D.K., Misof, B., Jermiin, L.S., 2020. A minimum reporting standard for multiple sequence alignments. NAR Genomics Bioinforma. 2. doi: 10.1093/nargab/lqaa024

Zhang, C., Rabiee, M., Sayyari, E., Mirarab, S., 2018. ASTRAL-III: Polynomial time species tree reconstruction from partially resolved gene trees. BMC Bioinformatics 19, 153.

Zhang, S.-Q., Che, L.-H., Li, Y., Dan Liang, Pang, H., Ślipiński, A., Zhang, P., 2018. Evolutionary history of Coleoptera revealed by extensive sampling of genes and species. Nat. Commun. 9, 205.

# 5. General discussion

In the previous chapters I focused on answering fundamental questions in the evolutionary history of two groups of Holometabola, the beetle suborder Adephaga and the superorder Neuropterida using analyses of NGS data. The aims of my dissertation were several, including: the inference of new ortholog sets for phylogenetic reconstructions in these groups, inferring the phylogeny of both groups using new and carefully curated NGS data, reconciliation of morphological with molecular phylogenies, assessing robustness of phylogenetic results with alternative measures of support and the inference of a new set of DNA-hybridization baits that can be useful in phylogenetic analyses of Adephaga. My results show that the most difficult task was to reconcile my phylogenomic results with morphological phylogenies. The inferred phylogenetic trees provide substantial progress in this direction, for example by resolving the relationships of Geadephaga, the phylogenetic position of Gyrinidae (chapter 4) and by consolidating the inter-ordinal relationships of Neuropterida (chapter 2). However, I find that certain results based on my sequenced-based analyses that are moderately or strongly supported in molecular phylogenetic analyses cannot be reconciled with results of morphological studies (e.g., phylogenetic position of Coniopterygidae in chapter 2 and Hygrobiidae in chapter 4). I have shown, however, that with a thorough investigation of phylogenomic data by and assessing phylogenetic support with alternative measures and by using analyses under different evolutionary models and data types it is possible to consolidate or exclude specific phylogenetic hypotheses and to identify phylogenetic relationships that are difficult to resolve. The herewith identified difficult relationships correspond in some instances to those that contradict morphology-based hypotheses. In this chapter, I discuss the most important phylogenetic results of my study in the historical context of the phylogeny of Adephaga

and Neuropterida and provide general directions for investigating unresolved questions in the phylogeny of these insect groups. Subsequently, I provide a general discussion of the approaches that hold great promise for overcoming the limitations of methods and data in phylogenomics and for evaluating inferred evolutionary hypotheses. Lastly, I briefly discuss the potential of the two genome-reduction approaches that I used here for future phylogenomic studies.

## 5.1. Phylogeny of Adephaga and Neuropterida – current status, controversies and future challenges

### 5.1.1. Phylogeny of Neuropterida

Almost two decades ago, Aspöck (2002) provided a comprehensive review on the *status quo* of phylogenetic relationships of Neuropterida that were at the time based mostly on analyses of morphological characters. Molecular phylogenetic analyses before that study had focused on the ordinal relationships of holometabolous insects but did not specifically address the familial relationships within Neuropterida and Neuroptera (e.g., Whiting et al., 1997). Ten years after the work of Aspöck (2002), another comprehensive review was published with a similar purpose (Aspöck et al., 2012), but this time the first molecular phylogenetic analyses specifically designed to address the familial relationships within Neuropterida and Neuroptera were also considered (Haring and Aspöck, 2004; Winterton et al., 2010). Almost ten years after the last comprehensive review of the phylogeny of Neuropterida, it seems that views on specific aspects on the phylogeny of the group have become more stable partially due to the stable results based on analyses of genome-scale data (e.g., Misof et al., 2014; Vasilikopoulos et al., 2020b; Wang et al., 2017; Winterton et al., 2018). In particular, it

seems that morphological and molecular hypotheses have since been reconciled to a large extent concerning: 1) the phylogenetic placement of Neuropterida within Holometabola, 2) the monophyly of the orders of Neuropterida and 3) the relationships among the orders of Neuropterida (Aspöck et al., 2001; Aspöck and Aspöck, 2008; Boussau et al., 2014; Misof et al., 2014; Peters et al., 2014; Vasilikopoulos et al., 2020b; Wang et al., 2017; Winterton et al., 2018). This consensus suggests that specific hypotheses in the phylogeny of Neuropterida and Holometabola such as a sister group relationship of Megaloptera and Raphidioptera (e.g., Kristensen, 1981) and the potential sister group relationships of Neuropterida + Coleoptera as sister to all other Holometabola (Kristensen, 1999)   should be abandoned in the light of this new evidence.

The accumulation of molecular sequence data for a larger taxon sampling within the orders of Neuropterida in the beginning of the 21[st] century allowed the exploration of relationships among the families within the orders (Haring and Aspöck, 2004; Winterton, 2003; Winterton et al., 2010). Despite the reconciliation of molecular and morphological phylogenies concerning the origin and interordinal relationships of Neuropterida in the last decade, genomic-scale datasets that were designed to address familial relationships of Neuroptera, including the dataset analyzed in the present thesis, have brought up new challenges in the phylogeny of the group. In particular, incongruence between phylogenomic and morphological phylogenies is now prominent and a consensus concerning neuropteran relationships among morphologists and molecular systematists is a challenging task.

Several examples in the phylogeny of Neuroptera highlight the disagreement between molecular and morphological studies. For example, Coniopterygidae is

unequivocally supported as a sister group to all other Neuroptera in all phylogenomic analyses as well as in my analyses using different data types data subsets and is also supported by analyses of alternaive mesures of branch support (quartet-based, chapter 2, Vasilikopoulos et al., 2020b). A potential close affinity of Chrysopidae + Hemerobiidae based on morphology (Aspöck and Aspöck, 2008) has been disputed in previous phylogenomic analyses (Winterton et al., 2018) and is also disputed based on the results presented here. The monophyly of Myrmeleontiformia (Aspöck and Aspöck, 2008; Badano et al., 2018) has also been disputed in previous molecular analyses (Winterton et al., 2018) and my results of codon-based sequence data also cast doubt on the monophyly this group (Vasilikopoulos et al., 2020b). Despite these incongruencies between molecular and morphological phylogenies my analyses show that some clades that contradict morphological hypotheses either a) do not receive strong support from quartet analyses or b) are sensitive to data-type effects (e.g., relationships of Chrysopidae, Hemerobiidae and monophyly of Myrmeleontiformia). At the same time some relationships that contradict morphological hypotheses (e.g., Coniopterygidae versus Nevrorthidae as sister to all other Neuroptera) are recovered in analyses of different subsets of the data and also receive strong support from quartet-based analyses but also from conventional branch support measures. All these observations suggest that the inferred relationships that receive strong support from several independent phylogenomic analyses (e.g., a clade Megaloptera + Neuroptera and a clade of all Neuroptera except Coniopterygidae) should be considered for now as the most likely scenario of neuropterid evolution but at the same time it is obvious that a top priority for future studies should be: a) the re-examinationa and re-evaluation of the morphological characters supporting Nevrorthidae as sister to all other Neuroptera (and

other characters supporting conflicting relationships) and b) the thorough assessment of phylogenomic results using better new methods and models.

A has been shown in chapter 4 and has been discussed in several studies the use of appropriate models of molecular evolution is critical for the inference of accurate phylogenetic relationships (Feuda et al., 2017; Philippe et al., 2009; Pisani et al., 2015). The two largest phylogenomic analyses of Neuropterida, the one presented here and the one by Winterton et al. (2018), did not include analyses under the CAT+GTR model that has been shown to fit data better than partitioned models and alleviate long-branch attraction artifacts (e.g., Lartillot et al., 2007). Since this model cannot be applied on very large matrices (such as the matrices presented in the present thesis) the development of objective criteria to subsample loci or sites for analyses should also be a priority (see next subsection and chapter 4). Despite this, the PMSF approximation to the site-heterogeneous C60 model that was used here (Wang et al., 2018) did not result in topological differences in comparison to my analyses of the partitioned models. Nevertheless, the PMSF model is only an approximation to the empirical site-heterogeneous mixture models (C10–C60) (Quang et al., 2008) and therefore analyses using these full empirical mixture models or a more complex CAT+GTR model may provide new insights in future studies.

The sampling of transcriptomic data for families that were missing in my analyses (i.e., Psychopsidae and Rhachiberothidae) should also be a priority for future studies as they might facilitate answering of open phylogenetic questions such as the monophyly of Myrmeleontiformia. Divergence time analyses presented here suggest an older origin of Neuropterida than previous studies (Misof et al., 2014; Tong et al., 2015)and these results should also be continuously evaluated on the basis of new

discovered fossils. Using a more comprehensive sampling of non-neuropterid taxa than previous studies (Wang et al., 2017; Winterton et al., 2018), I have shown the aquatic lifestyle of the larvae in some groups of Neuropterida was likely acquired secondarily from the primarily terrestrial lifestyle of the larvae of Neuropterida and Neuroptera and this result might be further corroborated using more comprehensive sampling of neuropterid and non-neuroptrid taxa, for example in the context of evolution of larval ecologies of the entire clade of Insecta. Needless to say, of particular interest for future studies should be the relative contribution of ILS and other tree-heterogeneous processes, such as ancient introgression, in shaping ancient divergences of insects in general but specifically also concerning those ancient splits in Neuropterida that seem to be separated by relatively short internal branches (e.g., Mantispoidea + Chrysopidae and Hamerobiidae + Ithonidae + Myrmeleontiformia). The identification of such processes in ancient phylogenetic clades constitutes a big challenge since these biological processes might be combined with systematic or stochastic errors (e.g., gene tree errors) making the relative contributions of different biasing factors difficult to detect (Kapli et al., 2020; Whitfield and Kjer, 2008).

### 5.1.2. Phylogeny of Adephaga

In contrast to the phylogeny of Neuropterida, the familial relationships of the beetle subgroup Adephaga has lately reached a stable consensus between molecular and morphological analyses (Beutel et al., 2020). The familial relationships of Adephaga have not significantly changed based on analyses of morphology in the last 30 years (Beutel et al., 2020; Beutel and Haas, 1996; Beutel and Roughley, 1988; Dressler and Beutel, 2010). Molecular studies based on a few genes produced some conflicting

results (Maddison et al., 2009; Ribera et al., 2002; Shull et al., 2001), but the accumulation of genome-scale data has largely reconciled morphological and molecular phylogenies of familial relationships of Adephaga (see chapter 4 and Baca et al., 2017; Gustafson et al., 2020). The major incongruence between molecular and morphological phylogenies at the moment concerns the position of Adephaga within Coleoptera (Beutel et al., 2019). My analyses corroborate this debate between morphological and molecular studies. Specifically, in most of my analyses of both site-homogeneous and site-heterogenous models a clade Myxophaga + Archostemata is inferred as sister to Adephaga while all three suborders together form the sister group of Polyphaga. This result is identical to large scale phylogenomic analyses of Coleoptera (McKenna et al., 2019; Zhang et al., 2018) but in conflict with morphology-based phylogenies that suggest Archostemata as sister to all other suborders (Beutel et al., 2019). Thus, future molecular studies should put more focus on potential systematic errors affecting the inferred relationships of beetle suborders based on phylogenomics but morphologists should also potential re-evaluation of several morphological features (e.g., Beutel et al., 2019).

Here, I take one step further in reconciling morphological and molecular phylogenies of Adephaga by showing that Aspidytidae is unequivocally inferred as monophyletic when steps are taken to reduce potentially biasing factors in the data (chapters 3 and 4, Vasilikopoulos et al., 2019). A sister group relationship of Aspidytidae and Amphizoidae is also hypothesized (chapter 3) in congruence with morphology-based phylogenies (Beutel et al., 2020) and later was corroborated based on analyses of UCEs (Gustafson et al., 2020) and also based on my latest analyses of combined exon-capture and transcriptome data (chapter 4). Most importantly,

Trachypachidae is robustly inferred as sister to Carabidae + Cicindelidae which is in agreement with morphology-based phylogenies but in contrast to the latest analyses of UCEs (Gustafson et al., 2020). In addition, my results corroborate previous hypotheses originally based on morphological analyses that the aquatic groups of Adephaga are not a monophylum and that there have been probably more than one transitions from terrestrial to aquatic lifestyle within the suborder (Beutel and Roughley, 1988).

Similarly to the results for Neuropterida, I find that difficult phylogenetic questions in the phylogenetic backbone of Adephaga do not receive unequivocal support from quartet analyses and might be sensitive to the model or dataset used. Specifically, the phylogenetic placement of Hygrobiidae is unstable between analyses of different models and datasets and as expected concatenation-based (FcLM) and gene-tree-based quartet analyses suggest that the inferred relationships of the family is not strongly supported based in analyses of transcriptomes (chapter 2). Despite this, with increased taxon sampling the position of Hygrobiidae is more stable among analyses of different datasets when a better-fitting site-heterogeneous model is used (chapter 4) and quartet-based support in favor of Hygrobiidae (Amphizoidae + Aspidytidae) is increased (chapter 4). Using a larger taxon sampling than Cai et al. (2020) and more thorough investigation of several data and taxon properties (e.g., removal of distant outgroup taxa and long-branched ingroup taxa, chapter 4), I find that application of a site-heterogeneous model does not support the morphology-based hypothesis Hygrobiidae + Dytiscidae in any analysis. Additionally, by performing analyses under both best-fitting and less-fitting evolutionary models in chapter 4, I find that Dytiscidae + (Amphizoidae + Aspidytidae) is likely an artifact due to model misspecification. Thus, although my analyses in chapter 4 are in agreement with UCE

studies and presently support Hygrobiidae (Amphizoidae + Aspidytidae), this part of the tree of Adephaga definitely deserves more attention in future studies both from a morphological and molecular perspective to address the existing incongruence.

The potential of the exon capture-approach and of my new set of baits for investigating other shallower relationships of Adephaga has to be further explored in future studies, for example by attempting to resolve the relationships within the megadiverse Carabidae and the Dytiscidae. In particular, the exploitation of the rich museum collections of Adephaga using the herewith newly presented set of baits has to be further pursued including analyses that are focused but not limited to a phylogenetic scope (e.g., biodiversity monitoring and population genomic studies). Furthermore, a comprehensive molecular dating analysis based on genomic data that is focused specifically on Adephaga is currently pending. The rich fossil record of Adephaga (Beutel et al., 2013) as well as the available phylogenomic data (e.g., Gustafson et al., 2020 and my data in chapter 4) can be utilized in future studies to infer a robust temporal framework of adephagan diversification. Lastly, formal analyses of ancestral character state reconstructions concerning the lifestyles of the adults and larvae should also be a topic for future studies. An underappreciated but essential issue is the lack of complete and annotated genomes for the groups of Adephaga and Neuropterida to date (e.g., McKenna, 2018). The sequencing, assembly and annotation of new complete genomes for both groups can facilitate the acquisition of molecular sequence and data and potentially provide new sets of markers for phylogenetic analyses of these groups.

## 5.2. Prospects for overcoming the limitations of methods and data in future phylogenomic studies

In the context of the phylogenetic inference of Adephaga and Neuropterida, my analyses show that there exist some difficult to resolve phylogenetic relationships in these groups. The inability to provide conclusive answers to some phylogenetic questions or to conclusively reconcile results between morphological analyses and phylogenomics can be due to limitations of the available molecular data (i.e., NGS data based on genome-reduction methods) or limitations of the analytical methods for inferring phylogenies from molecular sequence data (e.g., tree reconstruction methods and models of sequence evolution). Similar problems have previously been documented in numerous other phylogenomic studies that have attempted to address deep phylogenetic questions of insects and other animals (Feuda et al., 2017; Kocot et al., 2017; Reddy et al., 2017; Szucsich et al., 2020; Whelan et al., 2015). In this section, I will provide a summary of the most promising approaches that can be utilized in future molecular systematic studies in order to decipher difficult phylogenetic questions in the backbone phylogeny of Adephaga and Neuropterida and of other insect groups.

### 5.2.1. Development of better statistical models and methods to accommodate heterogeneous processes in the data

A straightforward solution to the problem of insufficient evolutionary modeling is the development of better statistical models that accommodate different heterogeneous substitution processes across sites and species of the alignment (Simion et al., 2020). Additionally, further research needs to be conducted on deciphering the potential biasing effect of less well-studied heterogeneous processes in the data, such as

the heterogeneity of the substitution process over time called heteropecilly (Roure and Philippe, 2011), and accordingly towards the development of suitable models to accommodate such under-studied heterogeneous processes. In the past, models that accommodate possible deviation of the data from stationarity and homogeneity have been developed (Blanquart and Lartillot, 2006; Jayaswal et al., 2011) and also models that accommodate non-stationarity combined with modeling of among-site heterogeneity of amino-acid propensities (Blanquart and Lartillot, 2008). However, due to their complexity these models are not scalable to large phylogenomic supermatrices. Moreover, phylogenetic tree-inference methods that are robust to various processes of genealogical heterogeneity across genes have to be developed and tested (e.g., Davidson et al., 2015; Zhang et al., 2020). This is because heterogeneous processes that generated the data are not known *a priori* (Vasilikopoulos et al., 2020a) and therefore scientists should ideally use methods that are robust to different sources of genealogical heterogeneity (Smith and Hahn, 2020).

An important aspect of applying new and better models and methods includes understanding of the limitations and pitfalls of each method (Bryant and Hahn, 2020). At the moment no method is good for all purposes, and until further progress is made, scientists might have to select appropriate methods based on the specific questions asked (Bryant and Hahn, 2020). This means that there is no panacea for analyzing the phylogenetic relationships of different groups species with genome-scale data. For instance, some site-heterogeneous mixture models might be more robust than site-homogeneous models in cases of gene-tree heterogeneity due to ILS (Wang et al., 2019), but these models still assume a common topology across all sites of the supermatrix, an assumption that may be unrealistic for biological data (Bryant and

Hahn, 2020). Moreover, the accuracy of coalescent methods for inferring ancient divergences where the molecular clock is seriously violated also needs to be more thoroughly investigated (Kapli et al., 2020). Lastly, in the context of summary coalescent analyses, scientists might want to take into consideration problems arising from stochastic and systematic error and from poor data quality (e.g., alignments, contamination or paralogous sequences) in the inference of gene trees before attributing gene-tree discordance to biological factors (Simion et al., 2020). This is because inference of gene trees is particularly difficult for ancient divergences with short internal branches (Salichos and Rokas, 2013), a pattern also observed here for some of the relationships of the families of Adephaga and Neuropterida. Summary coalescent and concatenation-based methods have been shown to be highly sensitive to the signal from a few outlier genetic loci (Gatesy et al., 2019; Shen et al., 2017; Walker et al., 2018) and these problems might be overcome in future studies by screening phylogenomic data for errors (e.g., De Vienne et al., 2012; Mai and Mirarab, 2018), or by performing phylogenomic subsampling analyses (Edwards, 2016; Simmons et al., 2016).

Selecting models and methods that are scalable to the size of the dataset should also be an important consideration in future phylogenomic studies. Due to the fact the most complex models of sequence evolution are not applicable to very large phylogenomic supermatrices, many authors have employed jacknifing approaches for analyzing phylogenomic data under such complex models (e.g., Delsuc et al., 2008; Simion et al., 2017). Others have reduced the size of datasets by randomly subsampling an arbitrary number of sites (Kapli and Telford, 2020) or by removing hypervariable sites (Cai et al., 2020; Lozano-Fernandez et al., 2019). These solutions are useful for

the purpose of analyzing very large supermatrices, or for selecting subsets of supermatrices based on specific properties of data (e.g., Cai et al., 2020; Delsuc et al., 2008), but there is no guarantee that the result from different random subsamples of the data, for example, will be more accurate than the phylogenetic result from the full dataset. In addition, the approach of repeating the phylogenetic analyses for different subsets of data is not well founded from a statistical point of view because each subset needs to be analyzed independently and this might require large amounts of computational resources (Simion et al., 2020). Therefore, the development of complex models that are scalable to large phylogenomic datasets should be given priority in the next years. Recently developed methods allow for the application of empirical site-heterogeneous models with a large number of amino-acid categories to large datasets and these methods are definitely worth exploring in future phylogenomic studies of Adephaga and Neuropterida (Schrempf et al., 2020).

An important last step in the analyses, even under the most complex models, is the application of tests of goodness-of-fit or test of absolute fit. These tests assess the ability of the model to accurately describe heterogeneous processes in the data after phylogenetic reconstruction has been performed (Shepherd and Klaere, 2019). Specifically, the development of better and more complex phylogenetic models that are designed to accommodate heterogeneous processes in the data does not mean that these models should be applied blindly in future phylogenetic analyses. A first step would be to perform model selection (relative test of fit) (Sullivan and Joyce, 2005) and also assess whether or not the assumptions of the selected model are violated (Jermiin et al., 2020). The step of model selection based on objective statistical criteria is critical in order to avoid over-parameterization (Sullivan and Joyce, 2005) and is widely accepted

as an important first step before phylogeny reconstruction (Buckley and Cunningham, 2002; Kalyaanamoorthy et al., 2017; Lanfear et al., 2012; Sullivan and Joyce, 2005). Assessment of model violations are also important, but the best models from the pool of available models might fail to adequately describe the underlying data even when the model assumptions are met (Jermiin et al., 2020; Shepherd and Klaere, 2019). Statistical methods to assess the model adequacy (i.e., goodness-of-fit or absolute fit) *a posteriori* are available in a ML (Duchêne et al., 2018; Jermiin et al., 2020) and BI framework (Bollback, 2002; Lartillot et al., 2013). In addition, statistical tests exist to assess the absolute fit of the multi-species coalescent model to the data (Reid et al., 2014). These tests might help scientists in future phylogenomic studies of Adephaga and Neuropterida and of other insect groups to assess if the applied model provides a good description of the processes in their data. If tests of model adequacy show that a less complex model is able to describe the heterogeneity of the data as well as a more complex model, then there is no good argument for using a more complex model for phylogenetic analysis (e.g., Li et al., 2020).

## 5.2.2. Better understanding of the data properties to select optimal sets of loci for analyses

An alternative or complementary approach to developing better phylogenetic models is to select high-quality data for phylogenetic analysis. Assuming that errors in alignment quality, orthology detection and cross-contaminations have been eliminated, one might wish to have an objective criterion of selecting genes or alignment sites for answering a specific phylogenetic question. This is because using all available genes genes might be 1) unrealistic (as described above) or 2) suboptimal depending on

specific criteria (e.g., deviation from model assumptions, phylogenetic information content, Jermiin et al., 2020; Klopfstein et al., 2017; Misof et al., 2013).

The question of how to best select data for phylogenetic inference is very old with extensive debate among scientists over the years (Dell'Ampio et al., 2014; Doyle et al., 2015; Evangelista et al., 2020; Goldman, 1998; Klopfstein et al., 2017; Misof et al., 2013; Mongiardino Koch and Thompson, 2020; Yang, 1998). One promising idea would be to select genes or alignment sites with reduced deviation from model assumptions (Simion et al., 2020). It should be noted, however, that removing biased sites comes at the cost of removing phylogenetic information (Mongiardino Koch and Thompson, 2020; Vasilikopoulos et al., 2020a). This is usually overlooked in phylogenomic studies because it is assumed that the benefits of removing bias will outweigh the limitations of removing phylogenetic information. In that sense, other authors have suggested using genes with high phylogenetic information for resolving difficult ancient divergences (Salichos and Rokas, 2013).

Several studies have examined other properties of genes as a measure of their phylogenetic utility to resolve ancient divergences, including among-species compositional homogeneity (Nesnidal et al., 2010), clock-likeness (Doyle et al., 2015), evolutionary rate (Doyle et al., 2015; Klopfstein et al., 2017; Yang, 1998), overall branch length heterogeneity (Kocot et al., 2017), length of the alignment and number of variable sites (Shen et al., 2016). Kocot et al. (2017) found that specific properties of genes can not be treated independently from other properties. For example, among-species compositional heterogeneity seems to be correlated with evolutionary rate which is not surprising when considering that fast evolving genes might develop differences in amino-acid compositions across different lineages over time due to

substitutional biases (Kocot et al., 2017). Additionally, Mongiardino Koch and Thompson (2020) found that systematic bias (measured with several different indices) is positively correlated with phylogenetic signal of genes (measured with average bootstrap support values as a proxy) and negatively correlated with gene-tree error, meaning that genes with low systematic bias also contain low phylogenetic information and vice versa. This suggests that there is no starightforward way to select optimized subsets of genes by increasing phylogenetic signal and reducing systematic bias at the same time (Mongiardino Koch and Thompson, 2020). My results of the phylogenetic inference of Adephaga corroborate this finding as they suggest that removing a lot of potentially biasing alignment sites results in reduced resolution of relationships of Adephaga most likely due to the removal of useful phylogenetic information. These results together with the results of previous studies show that biologists should be more careful when selecting genes or alignment regions for phylogenetic inferences by taking into account multiple properties of these regions. For example, maintaining a balance between removing bias and removing phylogenetic signal can be pursued in future studies. In addition, more research is needed to shed light on the putative dependency of different gene properties and to identify the factors that contribute to the most of heterogeneity in the datasets (e.g., Kocot et al., 2017). Another topic that has received little attention is the use of universal statistical criteria for selecting optimal numbers of loci according to specific gene properties as most scientists use arbitrary thresholds for reducing the size of their datasets (but see Klopfstein et al., 2017 for suggestions on universally optimal values of evolutionary rates). A promising solution to these problems would be to calculate different statistical properties of genes, and subsequently employ multivariate statistical analyses to select groups of genes that

when put together show optimized property values in comparison to the full set of genes or in comparison to random subsamples of genes (e.g., Mongiardino Koch and Thompson, 2020). Other authors have suggested that not all nodes in a phylogeny can be resolved by the same subsets of genes and therefore have suggested the selection of question-specific genes to reduce incongruence (Chen et al., 2015). Selecting question-specific genes or optimized subsets of genes requires the generation of a larger pool of available genes to select from. Therefore, the sequencing of larger portions of genomes is a prerequisite for both approaches and should be further pursued in future studies. Accordingly, sequencing more genomes of the species under investigation, or larger proportions of genomes, should also be a big priority for resolving recalcitrant nodes in the tree of life of Adephaga and Neuropterida and of other contentious phylogenetic relationships of insects (e.g., Meusemann et al., 2020).

### 5.2.3. Adding more data for more species

As mentioned before, one solution to increase the accuracy of phylogenetic analyses would be to add more data for more species. Taxon sampling is an important aspect of phylogenetic analyses and increasing the number of species in the analyses is expected to improve phylogenetic inference (Zwickl and Hillis, 2002). In a similar fashion, adding more genes increases phylogenetic accuracy (Rokas and Carroll, 2005) and provides a larger pool of available loci for subsampling (see previous subsection). Sequencing more genomes for more species has therefore the potential to improve phylogenetic analyses and is facilitated by the dramatic reduction of sequencing costs (Bleidorn, 2017). Sequencing larger portions of genomes for the species in question can also provide insights into the evolutionary processes of different genomic regions (e.g.,

coding versus non-coding) and provide the basis for comparison of phylogenies of these regions (Reddy et al., 2017). It is therefore logical to expect that a second revolution in phylogenomics will take place once full genomes will be available for the majority of species under investigation (Bravo et al., 2019).

Complementary to sequencing whole genomes, there exist other approaches to increase the number of available single-copy orthologs for phylogenetic analyses. One way to do this would be to add single-copy orthologs with missing data in the analyses (Smith and Hahn, 2020). The negative impact of missing data on the phylogenetic inference has been demonstrated in some previous studies (Dell'Ampio et al., 2014; Lemmon et al., 2009; Roure et al., 2013; Sayyari et al., 2017), but there is no strong consensus among scientists concerning these effects and how they specifically affect phylogenetic inference (Smith and Hahn, 2020; Wiens and Morrill, 2011). For example, it has been pointed out that highly incomplete taxa might increase accuracy of phylogenetic analyses (Wiens, 2005). In addition, it is possible that when we examine a particular phylogenetic question (i.e., node of the tree), genes might have to be decisive with respect to this particular node of the tree but missing data for shallower nodes can be tolerated (e.g., Dell'Ampio et al., 2014). Other studies have shown that the most problematic sequences for summary coalescent methods are partial DNA sequences rather than completely missing sequences (Hosner et al., 2016; Sayyari et al., 2017) whereas concatenation may not be negatively affected by increasing sparseness of the supermatrix (Hosner et al., 2016).

Another way of increasing the number of single-copy orthologs in the analyses would be to include genes with lineage-specific duplications (Smith and Hahn, 2020). For example, a gene might be characterized by lineage-specific duplications in some

terminal leaves of the tree but this should not affect the inference of the phylogenetic backbone of the tree as long as one copy is sampled per species for this COG (Smith and Hahn, 2020). This approach for increasing number of available data for phylogeny reconstruction has hardly ever been applied in phylogenomics and deserves more attention (Smith and Hahn, 2020).

Lastly, alternative approaches might allow the phylogenetic inference of species trees from multi-copy gene families (Emms and Kelly, 2018; Smith and Hahn, 2020; Zhang et al., 2020). By using these approaches, the amount of genes in a phylogentic analysis is drastically increased (Emms and Kelly, 2018) and therefore they are worth exploring in future phylogenomic studies. These approaches are particularly important when considering that it is hard to find orthologs that are single-copy and present in all species in the analysis (e.g., Thomas et al., 2020). This problem is especially apparent when the number of analyzed genomes is large and the divergence of the species under investigation is very old (e.g., Thomas et al., 2020).

## 5.2.4. Combination of phylogenomic data with morphological and paleontological data

The genomic revolution has undoubtedly changed the landscape of molecular phylogenetics by providing an enormous amount of data to answer long-standing phylogenetic questions (e.g., Dunn et al., 2008; Misof et al., 2014; Wickett et al., 2014). Due to these advances in NGS sequencing it has been postulated that the amount of molecular data will swamp out inferences from morphological characters when molecular data matrices are more than an order of magnitude larger than morphological data matrices (Giribet, 2010; Wortley and Scotland, 2006). This is one of the reasons

that systematists will often ignore the analyses of morphological data in the context of large scale phylogenomic analyses (Neumann et al., 2020). However, it has been recently suggested that morphological characters can drastically affect topological inferences from phylogenomic data when analyzed in combination, even when the upweighting of morphological data is minimal and the molecular data matrices are orders of magnitute larger than morphological matrices (Neumann et al., 2020). In addition, it has been shown that molecular sequence data can alter the phylogenetic placement of fossil taxa in combined analyses, thereby offering a solution to the placement fossils that are problematic to assign based on morphology (Reeder et al., 2015; Wiens et al., 2010). Most importantly, it has also been postulated that combined analyses of morphological and phylogenomic data can reduce incongruence from analyses of different types of data (Mongiardino Koch and Thompson, 2020). These observations suggest that morphological data matrices have still an important role to play for inferring species trees in the era of phylogenomics.

Some obvious impediments in the analyses of combined datasets are: 1) the potential lack of overlap of species between different types of data, 2) the selection of an objective weighting scheme for concatenated and combined data (Schierwater et al., 2016) and 3) the disagreement concerning the appropriate method to analyze the combined data (i.e., concatenation or supertree methods) (Bininda-Emonds, 2004; de Queiroz and Gatesy, 2007). Despite these problems, the majority of insect systematists have abandoned supertree methods due to early criticism (Bininda-Emonds, 2014; Gatesy et al., 2004, 2002) and instead have chosen to analyze combined data matrices with the supermatrix approach (e.g., Wahlberg et al., 2005; Wiegmann et al., 2002; Winterton et al., 2010, 2001). Concerning the the problem of non-overlapping taxa, if

the purpose of the study is to investigate the higher-level phylogeny (e.g., familial relationships) of the target group and morphological data matrices are composed of characters specific to resolve these higher-level relationships, then the generation of composite phylogenetic terminals (i.e., by sampling one molecular sequence per family or clade and assigning it to a different species of the same family) should not be a problem for inferring higher-level relationships (e.g., Mongiardino Koch and Thompson, 2020; Wiegmann et al., 2002). The problem of subjective weighting scheme can also be seen as an opportunity (in exploratory analyses) to inform ourselves on the relative power of morphological data to affect inferences in combined analyses (Neumann et al., 2020). All things considered, these observations suggest that combined analyses of morphological, paleontological and phylogenomic data might be worth exploring in future phylogenomic studies as they might bring new insights into the phylogeny of species under investigation .

Concerning the phylogeny of Neuropterida, one previously combined phylogenetic analysis of a few molecular markers and morphological characters of extant species (Winterton et al., 2010) resulted in relatively congruent results with contemporary phylogenomic analyses of Neuropterida (Wang et al., 2017; Winterton et al., 2018). Another study that combined morphological data from extinct and extant taxa and molecular sequence data from a few genetic markers (Yang et al., 2012) resulted in less congruent results with contemporary phylogenomic studies. It is, however, possible that the accuracy of results of those studies was hampered by the paucity of available molecular sequence data. Combined analyses of morphological and molecular sequence data for the adephagan superfamily Dytiscoidea have previously resulted in the placement of the family Hygrobiidae as sister to Dytiscidae +

(Amphizoidae + Aspidytidae) (Balke et al., 2005) in accordance to the phylotranscriptomic analyses of Dytiscoidea (see chapter 3, Vasilikopoulos et al., 2019), but in contrast to most analyses of the phylogenomic dataset of Adephaga (chapter 4). It should be noted, however, that the previously performed combined phylogenetic analyses of Dytiscoidea were based on equal weighting of molecular and morphological characters. Apart from this, they may also have been biased by the small number of molecular sequence data (Balke et al., 2005), similarly to previous combined analyses of Neuropterida.

The combination of large phylogenomic data with morphological or paleontological data is still pending for Adephaga and Neuropterida, and the ability of morphological data to affect the inference of specific phylogenetic hypotheses in combined analyses is therefore unclear. It is also unclear whether combined phylogenetic analyses can reconcile the results from analyses of these different types of data (e.g., Mongiardino Koch and Thompson, 2020). Furthermore, the question of how phylogenomic data might affect the phylogenetic placement of neuropterid and adephagan fossil taxa, especially in case of those are not conclusively placed based on morphology, requires further investigation. The available phylogenomic (e.g., Gustafson et al., 2020; Vasilikopoulos et al., 2020b; Winterton et al., 2018) and morphological matrices (e.g., Beutel et al., 2020; Winterton et al., 2010; Yang et al., 2012), combined with the rich fossil records of Adephaga and Neuropterida (e.g., Beutel et al., 2013; Yang et al., 2012) can be utilized in future studies to answer such questions and shed light on the evolutionary history of these groups.

**5.2.5. Beyond the limits of sequence-based phylogenomics: analyses of genomic metacharacters**

A possible solution to resolving difficult ancient divergences, such as those of Adephaga and Neuropterida, is the use of genomic metacharacters (or rare genomic changes) such as: retroelement insertions, the structure of genes, gene adjacency and synteny mapping, gene duplications, gene losses, gene fusions and the order of genes along the genome (Bleidorn, 2017; Boore and Fuerstenberg, 2008; Drillon et al., 2020; Krauss et al., 2008; Rokas and Holland, 2000; Schierwater et al., 2016). Some authors have referred to these types characters as the morphology of the genome or "molecular morphology" (Schierwater et al., 2016). For instance, phylogenetic reconstructions based on gene content is an old idea (Huson and Steel, 2004; Snel et al., 1999) that is based on comparing the proportion of genes shared by different genomes as a measure of phylogenetic relatedness (e.g., Snel et al., 1999). However, attempts to reconstruct phylogenetic relationships of Adephaga and Neuropterida based on gene content have not been made due to the lack of completely sequenced and annotated nuclear genomes in these groups.

Another promising approach that has been lately applied to resolve ancient divergences of birds and mammals utilizes information from insertions of low-homoplasy retroelements (Cloutier et al., 2019; Hallström et al., 2011; Springer et al., 2020; Suh et al., 2015). This approach is particularly intriguing especially considering that many retroelement insertions are characterized by low homoplasy in comparison to DNA sequences (Hallström et al., 2011). Insertions within intronic regions might be particularly useful due to their potentially being easier to orthologize based on adjacent exon information (Bleidorn, 2017). Interestingly, Springer et al. (2020) developed an

approach that uses low-homoplasy retroelements insertions in a multispecies coalescent framework to infer species trees. This approach carries high potential for resolving old rapid radiations, when considering that the pattern of inheritance of these insertions due to ILS might resemble the pattern of persistence of ancestral polymorphisms of gene data (Suh et al., 2015). Despite their promise in resolving old divergences, retroelement insertions are probably not useful for resolving very ancient divergences, (i.e., older than 50 million years) (Bleidorn, 2017) because homologous insertions might be harder to detect when mutations have accumulate over longer periods of time (Bleidorn, 2017; Shedlock and Okada, 2000). Therefore the retroelement approach is unlikely to be useful for inferring the familial relationships of Adephaga and Neuropterida or other deep splits in the phylogeny of insects. However, it is likely that they may be useful for inferring shallower phylogenetic relationships within the families of these groups once complete genomes become available.

Mitochondrial genomic rearrangements constitute another type of genomic metacharacter that has been successfully applied for inferring some relationships of within insects (e.g., Tyagi et al., 2020) but also for deciphering deeper arthropod relationships (Boore et al., 1998). Mitochondrial rearrangements have been previously studied in the context of the phylogenetic inference of Neuropterida (Wang et al., 2017; Zhao et al., 2013) and have been useful in identifying one rearrangement that is synanapomorphic for all families of Neuroptera except Coniopterygidae, Nevrorthidae, Osmylidae and Sisyridae (Wang et al., 2017). On the other hand, no rearrangement that contains useful information has been detected in the mitochondrial genomes of Adephaga so far (López-López and Vogler, 2017). The advancements in NGS technologies might facilitate the acquisition of complete adephagan and neuropterid

nuclear genomes in the near future and therefore enable the use of other genomic metacharacters to infer their phylogeny.

An example of genomic metacharacters that are based on nuclear genomic data and that have been applied successfully to resolve the relationships of other insect groups are near intron pairs (Krauss et al., 2008; Niehuis et al., 2012). Near intron pairs have also been useful for inferring other ancient metazoan relationships (Lehmann et al., 2013). The approach of near intron pairs was developed to overcome the limitation of homoplasy in patterns of intron gain and loss and is based on the idea that very short exons are rarely found in nature (Bleidorn, 2017). Therefore intronic sequences separated by 50 bp or less cannot have co-existed and probably represent different character states (Krauss et al., 2008). Once genomic data become available, the study of the structure of genes will be a promising solution to deciphering difficult questions in the backbone phylogeny of Adephaga and Neuropterida.

One last example of genomic metacharacters that hold great promise for phylogenetic reconstructions of difficult phylogenetic relationships are genomic rearrangements along the chromosomes and in particular the utilization of synteny breakpoints as markers for phylogenetic analysis (Drillon et al., 2020). Recently, methods and software has been developed for inferring phylogenetic relationships of species by using information from synteny conservation along the chromosomes and do not require whole-genome alignments (Drillon et al., 2020, 2014). These methods have been applied to successfully infer the phylogeny of vertebrates, a clade that is older than Neuropterida and much older than Adephaga (Irisarri et al., 2017; McKenna et al., 2019; Vasilikopoulos et al., 2020b). The method is based on all pairwise comparisons of syntenic blocks in a dataset, therefore taking into consideration differences among

distantly and closely related species, and uses a bottom-up approach (distance-like) to infer the phylogeny of the species (Drillon et al., 2020). Because the approach requires assembled and annotated genomes from the species under investigation, this approach holds great promise for deciphering contentious relationships in Adephaga and Neuropterida only when complete genomes are available for these groups.

## 5.3. Evaluation of inferred evolutionary hypotheses

### 5.3.1. Alternative measures of support to assess confidence in specific hypotheses

Even when scientists are confident that the phylogenetic methods used are appropriate for the analyses of their data, it is imperative to find objective ways to assess an evolutionary hypothesis or to assign a certain degree of confidence in a particular phylogenetic result. More specifically, assigning statistical support for specific phylogenetic branches helps to evaluate the degree of confidence for these branches and is one highly desirable yet very complex task (Kumar et al., 2012; Minh et al., 2020). Nowadays, phylogenomicists use the classical non-parametric bootstraping in a ML context (Felsenstein, 1985) or the Bayesian posterior probabilities in a BI context (Rannala and Yang, 1996) to calculate statistical support for specific branches on a phylogenetic tree. With increasing amounts of data in molecular systematics it became clear that these measures alone are not sufficient to assess credibility and measure phylogenetic support in favor of specific hypotheses (Minh et al., 2020; Pease et al., 2018).

In particular concerning the resampling approaches for assessing branch support (e.g., non-parametric bootstrap), these approaches are meant to be used as an approximation of data from a larger ideal population (Pease et al., 2018). Whole

genomes constitute the entire set of genomic data for an organism and therefore are not part of a larger population of data. This in turn suggests that non-parametric bootstrapping is inappropriate for analyses of whole genome data in a phylogenetic context (Pease et al., 2018). In essence, the non-parametric bootstrap assesses the potential lack of repeatability of analyses due to sampling effects but does not constitute a measure of phylogenetic accuracy (Felsenstein, 1985; Soltis and Soltis, 2003). Another reason that conventional branch support measures such as bootstrap and posterior probabilities are insufficient for analyses of phylogenomic data is that in many instances they show high support for mutually exclusive clades between different studies or analyzed datasets (Jarvis et al., 2014; Prum et al., 2015). I have shown, for example, that is is the case for some relationships in Neuroptera, one concerning the putative sister group of the clade Ithonidae + Myrmeleontiformia (i.e., Hemerobiidae vs. Chrysopidae, Ultrafast bootstrap support and SH-aLRT support). The same pattern is observed in the analyses of Winterton et al. (2018) concerning this particular phylogenetic relationship using Bayesian posterior probabilities.

The first realization when coming across mutually exclusive branches that are strongly supported is that one of them or both of those branches have to be wrong because they cannot be both true at the same time. When issues of outlier genes and other errors of data quality have been eliminated, it is valid to assume that one of the two or both analyses were performed under the wrong evolutionary model (e.g., Reddy et al., 2017). Firstly, this is because given sufficient amount of data (and therefore sufficient phylogenetic signal), analyses of different data types at the transcriptional and the translational sequence level (i.e., amino acid vs. nucleotides) should result in the same topology. Secondly, when looking at different genomic regions, if those are

sampled randomly along the genome, there is no valid reason to suggest that UCEs, for example, have a different evolutionary history than protein-coding regions. Bootstrap and ultrafast bootstrap support is known to produce high support for incorrect topologies due to the use of an incorrect model (Hoang et al., 2018; Huang et al., 2020) and the same applies for Bayesian posterior probabilities (Suzuki et al., 2002; Yang, 2014). Given the sensitivity of some branch support measures to model misspecification, the observations of strongly supported incongruent clades in empirical studies and the expectation that model misspecification is more prominent with increasing amount of data (due to potential unknown heterogeneous processes), it becomes clear how accurate assessment of branch support with some existing measures (e.g., bootstrap, posterior probability, ultrafast bootstrap) is likely obstructed by model misspecification. This also highlights the need for the development of branch support measures that are robust to model violations (Hoang et al., 2018; Sayyari and Mirarab, 2016).

Another reason that conventional branch support measures (e.g., bootstrapping, ultrafast bootstrapping, SH-aLRT, posterior probabilities) should not be applied in isolation from other measures is because they do not allow the exploration of particular biological properties or potential errors in the data (e.g., Minh et al., 2020; Pease et al., 2018). For example, other measures of support such as gene concordance factors are interesting tools to investigate levels of gene-tree discordance but also to detect potential gene-tree errors, similarly to my gene-tree discordance analyses for Adephaga in chapter 4. Site-concordance factors, as implemented in available software (Minh et al., 2020), are also interesting measures that are useful for evaluating effects of ILS in phylogenetic reconstructions but also potential biasing effects of a few outlier genes in

ML and BI analyses (Minh et al., 2020). Both measures are potentially useful in cases that concatenated analyses might provide strongly supported but incorrect clades due to ILS (Minh et al., 2020). However, it is not entirely clear yet how site-concordance factors are affected by confounding factors such as unequal evolutionary rates or compositional heterogeneity among taxa in the alignments. For example, assuming that two of the four subsets of taxa around the focal branch are characterized by the presence of fast evolutionary rates in all of their species, it is possible that quartet-based calculation of site-concordance factor might be biased for quartets of species due to the presence of fast evolving unrelated branches (Felsenstein, 1978).

Other measures of support sample informative quartet trees from the dataset to assess the robustness of specific phylogenetic branches. They do this by first extracting quartet trees informative for the focal branch and comparing the likelihoods of all three quartet topologies in these quartet trees (Pease et al., 2018; Strimmer and von Haeseler, 1997). This can be done using likelihoods of quartet trees from concatenated data (Pease et al., 2018; Zhou et al., 2020) or likelihoods of quartet trees from individual partitions or genetic loci (Pease et al., 2018). The four-cluster likelihood mapping (FcLM) approach, for example, allows the evaluation of particular nodes in an inferred tree by looking at the signal emerging from informative quartets of species around the focal branch (Strimmer and von Haeseler, 1997). A similar approach developed by (Pease et al., 2018) is intended to automatically detect support for secondary evolutionary histories in the branches of the entire tree and separate cases of conflict from cases of low support. Therefore such approaches can be used to detect cases of genomic introgression as skewed (i.e., uneven) support for alternative topologies are typical for such processes (Pease et al., 2018). Another similar approach allows for

similar quartet-based calculations based on a per-calculated set of gene trees and is therefore sensitive to the accuracy of the input gene trees (Zhou et al., 2020).

Despite the promise of these quartet-based approaches for estimating incongruence and support for specific phylogenetic relationships and their complementary nature to bootstrap support values and posterior probabilities (Zhou et al., 2020) they have certain limitations. Firstly, they are dependent on the models applied and therefore are not guaranteed to provide accurate estimates, especially considering that phylogenetic inference from quartets of species might not be as accurate as the inferences based on the full taxon sampling. This potential bias could be more prominent with unequal rates of evolutionary change in the data (Hendy and Penny, 1989). Secondly, it is difficult to distinguish between conflicts due to model misspecification from conflicts related to true biological processes that might generate similar conflict patterns. In such cases a permutation approach, as described by Misof et al. (2014) and that was also applied here, may help to separate between cases of true conflict from cases of conflict due to systematic bias. Lastly, another disadvantage of these approaches is that the data from quartets of taxa are analyzed under the same evolutionary model as the dataset under full-taxon sampling (Pease et al., 2018; Strimmer and von Haeseler, 1997) an assumption that may be unrealistic.

In my thesis, I have used multiple branch support measures to investigate the evolutionary relationships of Adephaga and Neuropterida. Despite their shortcomings, I have demonstrated that quartet-based measures of phylogenomic incongruence (i.e., FcLM, and gene-tree-based quartet scores) can be useful to detect clades with inflated branch support or more generally clades that are difficult to resolve and unstable between different analyses of datasets and methods. These combined approaches have

been applied here for the first time in the phylogenetic investigation of Adephaga and Neuropterida and in the future may be combined with other measures such as concordance factors and gene-wise likelihood scores (Minh et al., 2020; Shen et al., 2017) in order to assess potential phylogenetic artifacts and to gain insight into various evolutionary processes in the history of these groups.

### 5.3.2. Simulation-based studies for evaluating alternative phylogenetic hypotheses

Except for alternative measures of phylogenetic support, simulations can be used to assess the reliability of alternative evolutionary hypotheses. Simulated data based on two tree hypotheses can be generated under the best-fitting models that provide the best explanation of the data and subsequently explore whether tree reconstructions of simulated data using the same model or a less-fitting model result in different phylogenetic results (e.g., Kapli and Telford, 2020). This approach is essentially resembles a parametric bootstrapping approach but tree reconstructions are conducted under both true and wrong models of evolution to assess the potential effects of systematic errors in the results concerning a particular phylogenetic relationship. Specifically, the approach has been applied lately to assess which of the two most prevailing hypotheses in the early evolution of Metazoa (i.e., Ctenophora or Porifera sister) is likely to be the result of systematic error (Kapli and Telford, 2020). The investigation of relationships within the neuropteran superfamily Osmyloidea could, for example, benefit from this investigative simulation-based approach, by simulating and analyzing data using the two mutually exclusive clades Nevrorthidae + Sisyridae and Nevrorthidae + Osmylidae. The available phylogenomic data for these groups (Vasilikopoulos et al., 2020b; Wang et al., 2017; Winterton et al., 2018) and the

published phylogenetic hypotheses could serve as starting points for in combined approaches that include simulations and empirical data under site-homogeneous and site-heterogeneous models in future studies. The same approaches can also be applied to investigate the relationships of Dytiscoidea and more specifically the placement of the family Hygrobiidae.

### 5.3.3. Congruency tests to evaluate alternative phylogenetic hypotheses

In a previous subsection of this chapter it was mentioned that morphological characters are frequently analyzed in isolation from molecular sequence data and that combined analyses might bridge the gap between the different analyses. One problem of these combined analyses is that the modeling framework that is frequently more appropriate for analyzing molecular sequence data (Drummond and Rambaut, 2007; Lartillot and Philippe, 2004) is not suited for the analyses of morphological data matrices. This problem sometimes prompts for separate analyses of these two types of data. Even if genomic data are analyzed separately from morphological data, congruency tests based on morphology are important. It is recommended, therefore, that the results from phylogenomic analyses are compared to results from morphological characters because this is the strongest complementary evidence currently available for the groups under investigation (e.g., Pisani et al., 2007). However, this is expected to change with larger proportion of genomes becoming available for more species. For example, it will be interesting to conduct congruency tests at a larger genomic scale by comparing sequence-based inferences from different regions of the genome (e.g., non-coding UCEs versus protein-coding exons) with evolutionary analyses of genomic metacharacters. All these different types of molecular

data constitute independent evidence for assessing the reliability and plausibility of different evolutionary hypotheses and their results can also be compared with results of morphological analyses.

## 5.4. The future of genome-reduction sequencing strategies in phylogenomics

In my thesis, I have used two types of genome-reduction sequence data for inferring the phylogenetic relationships of Adephaga and Neuropterida: a) transcriptomes (chapters 2, 3 and 4) and b) hybrid-enrichment data of protein-coding exons (chapter 4). Transcriptome-based approaches for the purpose of inferring insect and arthropod evolutionary relationships were extensively applied in the first years of next-generation sequencing (e.g., Dunn et al., 2008; Meusemann et al., 2010; Peters et al., 2014; von Reumont et al., 2012) and hybrid-enrichment approaches have also been developed and extensively used since then (Bank et al., 2017; Faircloth et al., 2015; Sann et al., 2018; Young et al., 2016). The main advantage of these approaches in comparison to whole-genome sequencing in molecular systematic studies is their lower sequencing costs (Jones and Good, 2016). However, with continuously decreasing sequencing costs, sequencing of entire genomes is eventually going to become so cheap that the benefits of sequencing whole genomes versus sequencing selected genomic regions will outweigh the drawback of difference in the cost of the two approaches. In this subsection, I submit that these two genome-reduction approaches, and especially hybrid enrichment, will continue to be important and complementary tools for future molecular systematic studies and I discuss a few cases in which they might be useful

even when whole-genome sequencing becomes the standard data-collection strategy in phylogenomics.

First and foremost, hybrid-enrichment approaches might still be useful for phylogenetic research in species with a small body size that are rare or generally difficult to sample. Such cases might refer to species with limited geographical distributions, protected species, or even species that are extinct (Delsuc et al., 2018; Thomsen et al., 2009). This is because only a few or a single individual of some of these species might be available and given a potentially small body size, whole-genome sequencing might not be an achievable goal from such specimens. Furthermore, hybrid-enrichment approaches constitute the golden standard for capturing historical DNA from old museum specimens (Jones and Good, 2016; McCormack et al., 2016) in which cases whole-genome sequencing might be unrealistic. Similarly to ancient DNA, historical DNA from museum samples can be characterized by increased levels of contamination (e.g., from bacterial sources, Jones and Good, 2016). Therefore, hybrid-enrichment approaches might be particularly useful in cases of extensive contamination of historical samples if they are tailored to only capture the target regions in a specific taxonomic clade (Jones and Good, 2016). Another obvious advantage of these approaches is when the genomes of the species in question are very large and therefore difficult to assemble (e.g., Verlinden et al., 2020). Lastly, even when whole-genome sequencing becomes the standard procedure for investigating ancient insect divergences, it will long before high quality genome assemblies are generated for several species of the same genus and family.

An additional application of hybrid enrichment data in phylogenetics concerns the utilization of these data to extract phylogenetic information from genomic

metacharacters. The potential utility of hybrid-enrichment data for this type of analysis can already be explored in future studies before whole-genome sequencing becomes a common practice. For instance, it is interesting to consider that UCEs are generally found in regions of the genome with an increased degree of synteny (McCole et al., 2018). Because of this it could be theoretically possible to study the order of genomic elements close the ultraconserved regions to identify potential gene adjacencies and the presence of potentially informative synteny breakpoints in different species. However, this would require that the lab protocols are aimed at generating long sequenced reads using third-generation sequencing (Amarasinghe et al., 2020; Eid et al., 2009) in order to capture genomic information much further upstream and downstream from the UCE regions. In the case of the exon capture approach, when the baits are designed to capture neighboring exons of the same gene, it would in theory be possible to identify cases of intron loss in different species by identifying cases of adjacent exons stitched together in the assembled contigs. Segregation of the exons in different contigs or their separation by sequences that lack open reading frames could indicate the presence of an intron separating the two exons in other species. Despite this, the separation of exons into two separate contigs is not necessarily due to the presence of an intron but could also be due to assembly artifacts. In the cases of adjacent exons, the combination of hybrid enrichment with long-read sequencing (Amarasinghe et al., 2020) could help reduce these artifacts and potentially also allow the study of near intron pairs, particularly in cases of exons separated by long introns, similarly to analyses based on whole-genome data (Niehuis et al., 2012). Therefore the combination of hybrid enrichment and long-read sequencing can facilitate the use of potentially informative genomic metacharacters before entire genomes become available. One last

consideration for scientists with regard to this approach concerns the missing data might be prevalent for many species in hybrid-enrichment studies (Hosner et al., 2016). For example, whole exons might be captured in some species, partial exons in other species or the same regions may not be captured at all in other species due to divergent DNA sequences but also due to poor sample quality. Therefore such an approach would potentially be more beneficial when DNA is of good quality especially with regard to long-read sequencing approaches.

In contrast to hybrid enrichment, transcriptomics is not useful for collecting historical or ancient DNA from museum samples because the approach requires that RNA is extracted from fresh tissue (Bleidorn, 2017). Overall, the amount of effort needed is much lower in transcriptome sequencing than in whole-genome sequencing and the outcome is much more predictable and safe with respect to recovering loci for phylogenomics. One reason for this is for example that genome assembly and annotation is a much more time-consuming and complex task than transcriptome assembly and requires a great deal of bioinformatics expertise and manual work (Allen et al., 2017; Johnson, 2019; Richards, 2018; Wilbrandt et al., 2019). Transcriptome sequencing is almost always needed for building gene models from assembled genomes which allows selecting appropriate groups of genes for phylogenomics. However, specific software has lately been developed that assembles loci for phylogenomics from genomic raw sequenced reads without the need for a genome assembly and annotation (Allen et al., 2018). Because of this, transcriptome sequencing will not offer any specific advantages over the whole-genome sequencing approach but will still be a part of molecular systematic toolkit when: a) the species under investigation have large genomes that are not easy to sequence and assemble (although it is possible to sample

loci for phylogenomics without genome assembly and annotation as described above),

b) scientists want to use transcriptomes as a basis to design baits for exon capture (e.g.,

Bi et al., 2012) or c) transcriptomic evidence is needed to annotate the assembled

genomes and identify genes for phylogenomics in a new clade of interest. The last two

approaches do not directly use transcriptomic evidence for phylogenetic inference.

Transcriptome sequencing might also be advantageous over whole-genome sequencing

when only one or very few fresh individuals are available from a rare species sampled

in the field. This is because the small size of the animal might impede the extraction of

sufficient amounts of total genomic DNA for whole-genome sequencing. Despite this,

transcriptomes are not better than hybrid-enrichment approaches for recovering loci for

phylogenomics from such small and rare samples. Lastly, transcriptome sequencing

might be useful in cases in which scientists would like to examine additional aspects of

the sampled loci for phylogenomics such as the their expression levels in different

organisms or between different sexes.


## 5.5. Concluding remarks

In the present thesis several new insights into the phylogeny and evolution of

Adephaga and Neuropterida are presented based on analyses of data obtained with two

different genome-reduction approaches (transcriptomes and hybrid-enrichment data). In

addition, many useful methodological insights and prospects for future phylogenomic

studies have emerged. First, the comprehensive analyses presented here provide a

consolidation of most phylogenetic relationships of the Neuropterida based on analyses

of genome-scale data. For example, a sister group relationships of Megaloptera and

Neuroptera and the sister group relationship of Coniopterygidae to all other Neuroptera

is strongly supported using the largest data matrix analyzed to date. A new timeline of evolution of the major lineages of Neuropterida is established and a new hypothesis on the evolution of larval ecologies is presented that suggests more than one transition from terrestrial to aquatic lifestyle of the larvae within Neuropterida.

Secondly, the inferred familial phylogenetic relationships of Adephaga using a combination of exon-capture sequence data and transcriptomes reconcile results from previous molecular and morphological phylogenies to a large extent, and they provide a solid framework for future evolutionary and comparative studies in adephagan beetles. For example, Gyrinidae is inferred as sister to all other Adephaga in agreement with previous morphology-based phylogenies that suggested the paraphyly of "Hydradephaga" and therefore two independent transitions from terrestrial to aquatic lifestyle within Adephaga. Furthermore, Trachypachidae is restored as the sister group of Carabidae + Cicindelidae in agreement with morphological studies. Most importantly, the exon-capture approach was successful in recovering the target loci in divergent lineages of Adephaga and the phylogenetic analyses of the combined dataset showed that the captured regions carry useful phylogenetic signal to answer questions both at deep and shallow timescales (based on comparisons of results with morphology-based phylogenies). Based on these observations, I suggest that the presented set of DNA-hybridization baits shows great promise for future phylogenomic and potentially other evolutionary genomic or ecological studies in Adephaga..

The use of congruency tests with morphology-based hypotheses is one widely used approach to assess the validity of an inferred hypotheses. In my thesis, I have shown that reconciliation of results from molecular and morphological data is largely possible for 1) the familial relationships of Adephaga, 2) the ordinal relationships of

Neuropterida, and 3) the inferred position of Neuropterida within Holometabola. However, such a reconciliation is not possible for some familial relationships of Neuropterida (e.g., position of Coniopterygidae and Dilaridae), one familial phylogenetic relationship within Adephaga (i.e., position of Hygrobiidae) as well as the relationships of the suborders of Coleoptera. The inferred relationships from molecular analyses that receive strong support from several independent analyses of molecular data (i.e., position of Coniopterygidae) or those that are supported under the use of the best-fitting evolutionary models (e.g., Hygrobiidae as sister to Amphizoidae + Aspidytidae) should, for now, be considered as the most likely scenario for the evolution of these groups based on phylogenomics.

Conflicting results between morphology and phylogenomics, however, should be the focus of future morphological and molecular studies in order to identify the source of these incongruencies and possibly reconcile results. This is imperative for all conflicting relationships whether or not strongly supported in molecular studies. To that end, several promising strategies exist or will soon exist for overcoming the limitations of existing methods and data in future phylogenomic studies. Examples of such strategies are the development of better models of sequence evolution, the combined analyses of morphological and molecular data and the exploitation of potential phylogenetic information of genomic metacharacters. In particular, whole-genome sequencing has the potential to drastically increase the available data for analysis in Adephaga, Neuropterida and other insect groups with undersampled genomes. Nevertheless, genome-reduction approaches, such as hybrid enrichment and trancsriptomics, will continue to have a complementary role in molecular systematics even when whole genomes are routinely used for inferring species phylogenies.

Another promising approach for assessing the robustness of phylogenetic estimates and therefore for evaluating inferred hypotheses is to compare conventional measures of branch support with quartet-based measures of phylogenomic incongruence. I have demonstrated the usefulness of this approach for excluding specific evolutionary hypotheses and for detecting difficult phylogenetic relationships, such as those described above. I have also shown that performing analyses under both best-fitting and less-fitting models can be used for assessing the reliability of specific evolutionary hypotheses. In addition, I have shown that difficult phylogenetic questions in the backbone phylogeny of Adephaga and Neuropterida, and potentially also in phylogenomic analyses of other groups, are possible to identify using an integrative and comparative analysis of results of phylogenomic data, for example, by using analyses under different data types. Interestingly, some of these difficult phylogenetic problems correspond to inferred clades that contradict morphology-based hypotheses. This observation highlights the complementary nature of the different approaches for evaluating evolutionary hypotheses.

## 5.6. References

Allen, J.M., Boyd, B., Nguyen, N.-P., Vachaspati, P., Warnow, T., Huang, D.I., Grady, P.G.S., Bell, K.C., Cronk, Q.C.B., Mugisha, L., Pittendrigh, B.R., Leonardi, M.S., Reed, D.L., Johnson, K.P., 2017. Phylogenomics from Whole Genome Sequences Using aTRAM. Syst. Biol. 66, 786–798.

Allen, J.M., LaFrance, R., Folk, R.A., Johnson, K.P., Guralnick, R.P., 2018. aTRAM 2.0: An Improved, Flexible Locus Assembler for NGS Data. Evol Bioinform Online 14, 1–4.

Amarasinghe, S.L., Su, S., Dong, X., Zappia, L., Ritchie, M.E., Gouil, Q., 2020. Opportunities and challenges in long-read sequencing data analysis. Genome Biol. 21, 30.

Aspöck, U., 2002. Phylogeny of the Neuropterida (Insecta: Holometabola). Zool. Scr. 31, 51–55.

Aspöck, U., Aspöck, H., 2008. Phylogenetic relevance of the genital sclerites of Neuropterida (Insecta: Holometabola). Syst. Entomol. 33, 97–127.

Aspöck, U., Haring, E., Aspöck, H., 2012. The phylogeny of the Neuropterida: long lasting and current controversies and challenges (Insecta: Endopterygota). Arthropod Syst. Phylogeny 70, 119–129.

Aspöck, U., Plant, J.D., Nemeschkal, H.L., 2001. Cladistic analysis of Neuroptera and their systematic position within Neuropterida (Insecta: Holometabola: Neuropterida: Neuroptera). Syst. Entomol. 26, 73–86.

Baca, S.M., Alexander, A., Gustafson, G.T., Short, A.E.Z., 2017. Ultraconserved elements show utility in phylogenetic inference of Adephaga (Coleoptera) and suggest paraphyly of 'Hydradephaga'. Syst. Entomol. 42, 786–795.

Badano, D., Engel, M.S., Basso, A., Wang, B., Cerretti, P., 2018. Diverse Cretaceous larvae reveal the evolutionary and behavioural history of antlions and lacewings. Nat. Commun. 9, 3257.

Balke, M., Ribera, I., Beutel, R.G., 2005. The systematic position of Aspidytidae, the diversification of Dytiscoidea (Coleoptera, Adephaga) and the phylogenetic signal of third codon positions. J. Zool. Syst. Evol. Res. 43, 223–242.

Bank, S., Sann, M., Mayer, C., Meusemann, K., Donath, A., Podsiadlowski, L., Kozlov, A., Petersen, M., Krogmann, L., Meier, R., Rosa, P., Schmitt, T., Wurdack, M., Liu, S., Zhou, X., Misof, B., Peters, R.S., Niehuis, O., 2017. Transcriptome and target DNA enrichment sequence data provide new insights into the phylogeny of vespid wasps (Hymenoptera: Aculeata: Vespidae). Mol. Phylogenet. Evol. 116, 213–226.

Beutel, R.G., Haas, A., 1996. Phylogenetic analysis of larval and adult characters of Adephaga (Coleoptera) using cladistic computer programs. Insect Syst. Evol. 27, 197–205.

Beutel, R.G., Pohl, H., Yan, E. V., Anton, E., Liu, S.-P., Ślipiński, A., McKenna, D., Friedrich, F., 2019. The phylogeny of Coleopterida (Hexapoda) – morphological characters and molecular phylogenies. Syst. Entomol. 44, 75–102.

Beutel, R.G., Ribera, I., Fikáček, M., Vasilikopoulos, A., Misof, B., Balke, M., 2020. The morphological evolution of the Adephaga (Coleoptera). Syst. Entomol. 45, 378–395.

Beutel, R.G., Roughley, R.E., 1988. On the systematic position of the family Gyrinidae (Coleoptera: Adephaga). J. Zool. Syst. Evol. Res. 26, 380–400.

Beutel, R.G., Wang, B., Tan, J.J., Ge, S.Q., Ren, D., Yang, X.K., 2013. On the phylogeny and evolution of Mesozoic and extant lineages of Adephaga (Coleoptera, Insecta). Cladistics 29, 147–165.

Bi, K., Vanderpool, D., Singhal, S., Linderoth, T., Moritz, C., Good, J.M., 2012. Transcriptome-based exon capture enables highly cost-effective comparative genomic data collection at moderate evolutionary scales. BMC Genomics 13, 403.

Bininda-Emonds, O.R.P., 2014. An introduction to supertree construction (and partitioned phylogenetic analyses) with a view toward the distinction between gene trees and species trees, in: Garamszegi, L.Z. (Ed.), Modern phylogenetic comparative methods and their application in evolutionary biology. Springer-Verlag Berlin Heidelberg, Berlin, Heidelberg, pp. 49–76.

Bininda-Emonds, O.R.P., 2004. The evolution of supertrees. Trends Ecol. Evol. 19, 315–322.

Blanquart, S., Lartillot, N., 2008. A site- and time-heterogeneous model of amino acid replacement. Mol. Biol. Evol. 25, 842–858.

Blanquart, S., Lartillot, N., 2006. A Bayesian compound stochastic process for modeling nonstationary and nonhomogeneous sequence evolution. Mol. Biol. Evol. 23, 2058–2071.

Bleidorn, C., 2017. Phylogenomics - an introduction, 1st ed. Springer International Publishing.

Bollback, J.P., 2002. Bayesian model adequacy and choice in phylogenetics. Mol. Biol. Evol. 19, 1171–1180.

Boore, J.L., Fuerstenberg, S.I., 2008. Beyond linear sequence comparisons: the use of genome-level characters for phylogenetic reconstruction. Philos Trans R Soc L. B Biol Sci. 363, 1445–1451.

Boore, J.L., Lavrov, D. V, Brown, W.M., 1998. Gene translocation links insects and crustaceans. Nature 392, 667–668.

Boussau, B., Walton, Z., Delgado, J.A., Collantes, F., Beani, L., Stewart, I.J., Cameron, S.A., Whitfield, J.B., Johnston, J.S., Holland, P.W.H., Bachtrog, D., Kathirithamby, J., Huelsenbeck, J.P., 2014. Strepsiptera, phylogenomics and the long branch attraction problem. PLoS One 9, e107709.

Bravo, G.A., Antonelli, A., Bacon, C.D., Bartoszek, K., Blom, M.P.K., Huynh, S., Jones, G., Knowles, L.L., Lamichhaney, S., Marcussen, T., Morlon, H., Nakhleh, L.K., Oxelman, B., Pfeil, B., Schliep, A., Wahlberg, N., Werneck, F.P., Wiedenhoeft, J., Willows-Munro, S., Edwards, S. V, 2019. Embracing heterogeneity: coalescing the Tree of Life and the future of phylogenomics. PeerJ 7, e6399.

Bryant, D., Hahn, M.W., 2020. The concatenation question, in: Scornavacca, C., Delsuc, F., Galtier, N. (Eds.), Phylogenetics in the genomic era. No commercial publisher | Authors' open access book, pp. 3.4:1–3.4:23.

Buckley, T.R., Cunningham, C.W., 2002. The effects of nucleotide substitution model assumptions on estimates of nonparametric bootstrap support. Mol. Biol. Evol. 19, 394–405.

Cai, C., Tihelka, E., Pisani, D., Donoghue, P.C.J., 2020. Data curation and modeling of compositional heterogeneity in insect phylogenomics: a case study of the phylogeny of Dytiscoidea (Coleoptera: Adephaga). Mol. Phylogenet. Evol. 147, 106782.

Chen, M.Y., Liang, D., Zhang, P., 2015. Selecting question-specific genes to reduce incongruence in phylogenomics: a case study of jawed vertebrate backbone phylogeny. Syst. Biol. 64, 1104–1120.

Cloutier, A., Sackton, T.B., Grayson, P., Clamp, M., Baker, A.J., Edwards, S. V., 2019. Whole-genome analyses resolve the phylogeny of flightless birds (Palaeognathae) in the presence of an empirical anomaly zone. Syst. Biol. 68, 937–955.

Davidson, R., Vachaspati, P., Mirarab, S., Warnow, T., 2015. Phylogenomic species tree estimation in the presence of incomplete lineage sorting and horizontal gene transfer. BMC Genomics 16, S1.

de Queiroz, A., Gatesy, J., 2007. The supermatrix approach to systematics. Trends Ecol. Evol. 22, 34–41.

De Vienne, D.M., Ollier, S., Aguileta, G., 2012. Phylo-MCOA: a fast and efficient method to detect outlier genes and species in phylogenomics using multiple co-inertia analysis. Mol. Biol. Evol. 29, 1587–1598.

Dell'Ampio, E., Meusemann, K., Szucsich, N.U., Peters, R.S., Meyer, B., Borner, J., Petersen, M., Aberer, A.J., Stamatakis, A., Walzl, M.G., Minh, B.Q., von Haeseler,

A., Ebersberger, I., Pass, G., Misof, B., 2014. Decisive data sets in phylogenomics: lessons from studies on the phylogenetic relationships of primarily wingless insects. Mol. Biol. Evol. 31, 239–249.

Delsuc, F., Kuch, M., Gibb, G.C., Hughes, J., Szpak, P., Southon, J., Enk, J., Duggan, A.T., Poinar, H.N., 2018. Resolving the phylogenetic position of Darwin's extinct ground sloth (*Mylodon darwinii*) using mitogenomic and nuclear exon data. Proc. R. Soc. B Biol. Sci. 285, 20180214.

Delsuc, F., Tsagkogeorga, G., Lartillot, N., Philippe, H., 2008. Additional molecular support for the new chordate phylogeny. Genesis 46, 592–604.

Doyle, V.P., Young, R.E., Naylor, G.J.P., Brown, J.M., 2015. Can we identify genes with increased phylogenetic reliability? Syst. Biol. 64, 824–837.

Dressler, C., Beutel, R.G., 2010. The morphology and evolution of the adult head of Adephaga (Insecta: Coleoptera). Arthropod Syst. Phylogeny 68, 239–287.

Drillon, G., Carbone, A., Fischer, G., 2014. SynChro: a fast and easy tool to reconstruct and visualize synteny blocks along eukaryotic chromosomes. PLoS One 9, e92621.

Drillon, G., Champeimont, R., Oteri, F., Fischer, G., Carbone, A., 2020. Phylogenetic reconstruction based on synteny block and gene adjacencies. Mol. Biol. Evol. 37, 2747–2762.

Drummond, A.J., Rambaut, A., 2007. BEAST: Bayesian evolutionary analysis by sampling trees. BMC Evol. Biol. 7, 214.

Duchêne, D.A., Duchêne, S., Ho, S.Y.W., 2018. PhyloMAd: efficient assessment of phylogenomic model adequacy. Bioinformatics 34, 2300–2301.

Dunn, C.W., Hejnol, A., Matus, D.Q., Pang, K., Browne, W.E., Smith, S. a, Seaver, E., Rouse, G.W., Obst, M., Edgecombe, G.D., Sørensen, M. V, Haddock, S.H.D., Schmidt-Rhaesa, A., Okusu, A., Kristensen, R.M., Wheeler, W.C., Martindale, M.Q., Giribet, G., 2008. Broad phylogenomic sampling improves resolution of the animal tree of life. Nature 452, 745–749.

Edwards, S. V, 2016. Phylogenomic subsampling: a brief review. Zool. Scr. 45, 63–74.

Eid, J., Fehr, A., Gray, J., Luong, K., Lyle, J., Otto, G., Peluso, P., Rank, D., Baybayan, P., Bettman, B., Bibillo, A., Bjornson, K., Chaudhuri, B., Christians, F., Cicero, R., Clark, S., Dalal, R., DeWinter, A., Dixon, J., Foquet, M., Gaertner, A., Hardenbol, P., Heiner, C., Hester, K., Holden, D., Kearns, G., Kong, X., Kuse, R., Lacroix, Y., Lin, S., Lundquist, P., Ma, C., Marks, P., Maxham, M., Murphy, D., Park, I., Pham, T., Phillips, M., Roy, J., Sebra, R., Shen, G., Sorenson, J., Tomaney, A., Travers, K., Trulson, M., Vieceli, J., Wegener, J., Wu, D., Yang, A., Zaccarin, D., Zhao, P., Zhong, F., Korlach, J., Turner, S., 2009. Real-time DNA sequencing from single polymerase molecules. Science. 323, 133–138.

Emms, D.M., Kelly, S., 2018. STAG: species tree inference from all genes. BioRxiv. doi: 10.1101/267914.

Evangelista, D., Simon, S., Wilson, M.M., Kawahara, A.Y., Kohli, M.K., Ware, J.L., Wipfler, B., Béthoux, O., Grandcolas, P., Legendre, F., 2020. Assessing support for Blaberoidea phylogeny suggests optimal locus quality. Syst. Entomol. doi: 10.1111/syen.12454

Faircloth, B.C., Branstetter, M.G., White, N.D., Brady, S.G., 2015. Target enrichment of ultraconserved elements from arthropods provides a genomic perspective on relationships among Hymenoptera. Mol. Ecol. Resour. 15, 489–501.

Felsenstein, J., 1985. Confidence limits on phylogenies: an approach using the bootstrap. Evolution 39, 783–791.

Felsenstein, J., 1978. Cases in which parsimony or compatibility methods will be positively misleading. Syst. Biol. 27, 401–410.

Feuda, R., Dohrmann, M., Pett, W., Philippe, H., Rota-Stabelli, O., Lartillot, N., Wörheide, G., Pisani, D., 2017. Improved Modeling of Compositional Heterogeneity Supports Sponges as Sister to All Other Animals. Curr. Biol. 27, 3864–3870.

Gatesy, J., Baker, R.H., Hayashi, C., 2004. Inconsistencies in arguments for the supertree approach: supermatrices versus supertrees of Crocodylia. Syst. Biol. 53, 342–355.

Gatesy, J., Matthee, C., DeSalle, R., Hayashi, C., 2002. Resolution of a supertree/supermatrix paradox. Syst. Biol. 51, 652–664.

Gatesy, J., Sloan, D.B., Warren, J.M., Baker, R.H., Simmons, M.P., Springer, M.S., 2019. Partitioned coalescence support reveals biases in species-tree methods and detects gene trees that determine phylogenomic conflicts. Mol. Phylogenet. Evol. 139, 106539.

Giribet, G., 2010. A new dimension in combining data? The use of morphology and phylogenomic data in metazoan systematics. Acta Zool. 91, 11–19.

Goldman, N., 1998. Phylogenetic information and experimental design in molecular systematics. Proceedings. Biol. Sci. 265, 1779–1786.

Gustafson, G.T., Baca, S.M., Alexander, A.M., Short, A.E.Z., 2020. Phylogenomic analysis of the beetle suborder Adephaga with comparison of tailored and generalized ultraconserved element probe performance. Syst. Entomol. 45, 552–570.

Hallström, B.M., Schneider, A., Zoller, S., Janke, A., 2011. A genomic approach to examine the complex evolution of laurasiatherian mammals. PLoS One 6, e28199.

Haring, E., Aspöck, U., 2004. Phylogeny of the Neuropterida: a first molecular approach. Syst. Entomol. 29, 415–430.

Hendy, M.D., Penny, D., 1989. A framework for the quantitative study of evolutionary trees. Syst. Zool. 38, 297–309.

Hoang, D.T., Chernomor, O., von Haeseler, A., Minh, B.Q., Le, S.V., 2018. UFBoot2: improving the ultrafast bootstrap approximation. Mol. Biol. Evol. 35, 518–522.

Hosner, P.A., Faircloth, B.C., Glenn, T.C., Braun, E.L., Kimball, R.T., 2016. Avoiding missing data biases in phylogenomic inference: an empirical study in the landfowl (Aves: Galliformes). Mol. Biol. Evol. 33, 1110–1125.

Huang, J., Liu, Y., Zhu, T., Yang, Z., 2020. The asymptotic behavior of bootstrap support values in molecular phylogenetics. Syst. Biol. doi: 10.1093/sysbio/syaa100

Huson, D.H., Steel, M., 2004. Phylogenetic trees based on gene content. Bioinformatics 20, 2044–2049.

Irisarri, I., Baurain, D., Brinkmann, H., Delsuc, F., Sire, J.-Y., Kupfer, A., Petersen, J., Jarek, M., Meyer, A., Vences, M., Philippe, H., 2017. Phylotranscriptomic consolidation of the jawed vertebrate timetree. Nat. Ecol. Evol. 1, 1370–1378.

Jarvis, E.D., Mirarab, S., Aberer, A.J., Li, B., Houde, P., Li, C., Ho, S.Y.W., Faircloth, B.C., Nabholz, B., Howard, J.T., Suh, A., Weber, C.C., da Fonseca, R.R., Li, J., Zhang, F., Li, H., Zhou, L., Narula, N., Liu, L., Ganapathy, G., Boussau, B., Bayzid, M.S., Zavidovych, V., Subramanian, S., Gabaldón, T., Capella-Gutiérrez, S., Huerta-Cepas, J., Rekepalli, B., Munch, K., Schierup, M., Lindow, B., Warren, W.C., Ray, D., Green, R.E., Bruford, M.W., Zhan, X., Dixon, A., Li, S., Li, N., Huang, Y., Derryberry, E.P., Bertelsen, M.F., Sheldon, F.H., Brumfield, R.T., Mello, C. V, Lovell, P. V, Wirthlin, M., Schneider, M.P.C., Prosdocimi, F., Samaniego, J.A., Velazquez, A.M.V., Alfaro-Núñez, A., Campos, P.F., Petersen, B., Sicheritz-Ponten, T., Pas, A., Bailey, T., Scofield, P., Bunce, M., Lambert, D.M., Zhou, Q., Perelman, P., Driskell, A.C., Shapiro, B., Xiong, Z., Zeng, Y., Liu, S., Li, Z., Liu, B., Wu, K., Xiao, J., Yinqi, X., Zheng, Q., Zhang, Y., Yang, H., Wang, J., Smeds, L., Rheindt, F.E., Braun, M., Fjeldsa, J., Orlando, L., Barker, F.K., Jønsson, K.A., Johnson, W., Koepfli, K.-P., O'Brien, S., Haussler, D., Ryder, O.A., Rahbek, C., Willerslev, E., Graves, G.R., Glenn, T.C., McCormack, J., Burt, D., Ellegren, H., Alström, P., Edwards, S. V, Stamatakis, A., Mindell, D.P., Cracraft, J., Braun, E.L., Warnow, T., Jun, W., Gilbert, M.T.P., Zhang, G., 2014. Whole-genome analyses resolve early branches in the tree of life of modern birds. Science. 346, 1320–1331.

Jayaswal, V., Jermiin, L.S., Poladian, L., Robinson, J., 2011. Two stationary nonhomogeneous Markov models of nucleotide sequence evolution. Syst. Biol. 60, 74–86.

Jermiin, L.S., Catullo, R.A., Holland, B.R., 2020. A new phylogenetic protocol: dealing with model misspecification and confirmation bias in molecular phylogenetics. NAR Genomics Bioinforma. 2. doi: 10.1093/nargab/lqaa041

Johnson, K.P., 2019. Putting the genome in insect phylogenomics. Curr. Opin. Insect Sci. 36, 111–117.

Jones, M.R., Good, J.M., 2016. Targeted capture in evolutionary and ecological genomics. Mol. Ecol. 25, 185–202.

Kalyaanamoorthy, S., Minh, B.Q., Wong, T.K.F., von Haeseler, A., Jermiin, L.S., 2017. ModelFinder: fast model selection for accurate phylogenetic estimates. Nat. Methods 14, 587–589.

Kapli, P., Telford, M.J., 2020. Topology dependent asymmetry in systematic errors affects phylogenetic placement of Ctenophora and Xenacoelomorpha. Sci. Adv. 6, eabc5162.

Kapli, P., Yang, Z., Telford, M.J., 2020. Phylogenetic tree building in the genomic age. Nat. Rev. Genet. 21, 428–444.

Klopfstein, S., Massingham, T., Goldman, N., 2017. More on the best evolutionary rate for phylogenetic analysis. Syst. Biol. 66, 769–785.

Kocot, K.M., Struck, T.H., Merkel, J., Waits, D.S., Todt, C., Brannock, P.M., Weese, D.A., Cannon, J.T., Moroz, L.L., Lieb, B., Halanych, K.M., 2017. Phylogenomics of Lophotrochozoa with consideration of systematic error. Syst. Biol. 66, 256–282.

Krauss, V., Thümmler, C., Georgi, F., Lehmann, J., Stadler, P.F., Eisenhardt, C., 2008. Near intron positions are reliable phylogenetic markers: an application to holometabolous insects. Mol. Biol. Evol. 25, 821–830.

Kristensen, N.P., 1999. Phylogeny of endopterygote insects, the most successful lineage of living organisms. Eur. J. Entomol. 96, 237–253.

Kristensen, N.P., 1981. Phylogeny of insect orders. Annu. Rev. Entomol. 26, 135–157.

Kumar, S., Filipski, A.J., Battistuzzi, F.U., Kosakovsky Pond, S.L., Tamura, K., 2012. Statistics and truth in phylogenomics. Mol. Biol. Evol. 29, 457–472.

Lanfear, R., Calcott, B., Ho, S.Y.W., Guindon, S., 2012. PartitionFinder: combined selection of partitioning schemes and substitution models for phylogenetic analyses. Mol. Biol. Evol. 29, 1695–1701.

Lartillot, N., Brinkmann, H., Philippe, H., 2007. Suppression of long-branch attraction artefacts in the animal phylogeny using a site-heterogeneous model. BMC Evol. Biol. 7, S4.

Lartillot, N., Philippe, H., 2004. A Bayesian mixture model for across-site heterogeneities in the amino-acid replacement process. Mol. Biol. Evol. 21, 1095–1109.

Lartillot, N., Rodrigue, N., Stubbs, D., Richer, J., 2013. PhyloBayes MPI : Phylogenetic reconstruction with infinite mixtures of profiles in a parallel environment. Syst. Biol. 62, 611–615.

Lehmann, J., Stadler, P.F., Krauss, V., 2013. Near intron pairs and the metazoan tree. Mol. Phylogenet. Evol. 66, 811–823.

Lemmon, A.R., Brown, J.M., Stanger-Hall, K., Lemmon, E.M., 2009. The effect of ambiguous data on phylogenetic estimates obtained by maximum likelihood and Bayesian inference. Syst. Biol. 58, 130–145.

Li, Y., Shen, X.-X., Evans, B., Dunn, C.W., Rokas, A., 2020. Rooting the animal tree of life. BioRxiv. doi: 10.1101/2020.10.27.357798

López-López, A., Vogler, A.P., 2017. The mitogenome phylogeny of Adephaga (Coleoptera). Mol. Phylogenet. Evol. 114, 166–174.

Lozano-Fernandez, J., Tanner, A.R., Giacomelli, M., Carton, R., Vinther, J., Edgecombe, G.D., Pisani, D., 2019. Increasing species sampling in chelicerate genomic-scale datasets provides support for monophyly of Acari and Arachnida. Nat. Commun. 10, 2295.

Maddison, D.R., Moore, W., Baker, M.D., Ellis, T.M., Ober, K.A., Cannone, J.J., Gutell, R.R., 2009. Monophyly of terrestrial adephagan beetles as indicated by three nuclear genes (Coleoptera: Carabidae and Trachypachidae). Zool. Scr. 38, 43–62.

Mai, U., Mirarab, S., 2018. TreeShrink: fast and accurate detection of outlier long branches in collections of phylogenetic trees. BMC Genomics 19, 272.

McCole, R.B., Erceg, J., Saylor, W., Wu, C., 2018. Ultraconserved elements occupy specific arenas of three-dimensional mammalian genome organization. Cell Rep. 24, 479–488.

McCormack, J.E., Tsai, W.L.E., Faircloth, B.C., 2016. Sequence capture of ultraconserved elements from bird museum specimens. Mol. Ecol. Resour. 16, 1189–1203.

McKenna, D.D., 2018. Beetle genomes in the 21st century: prospects, progress and priorities. Curr. Opin. Insect Sci. 25, 76–82.

McKenna, D.D., Shin, S., Ahrens, D., Balke, M., Beza-Beza, C., Clarke, D.J., Donath, A., Escalona, H.E., Friedrich, F., Letsch, H., Liu, S., Maddison, D., Mayer, C., Misof, B., Murin, P.J., Niehuis, O., Peters, R.S., Podsiadlowski, L., Pohl, H., Scully, E.D., Yan, E. V, Zhou, X., Ślipiński, A., Beutel, R.G., 2019. The evolution and genomic basis of beetle diversity. Proc. Natl. Acad. Sci. U. S. A. 116, 24729–24737.

Meusemann, K., Trautwein, M., Friedrich, F., Beutel, R.G., Wiegmann, B.M., Donath, A., Podsiadlowski, L., Petersen, M., Niehuis, O., Mayer, C., Bayless, K.M., Shin, S., Liu, S., Hlinka, O., Minh, B.Q., Kozlov, A., Morel, B., Peters, R.S., Bartel, D., Grove, S., Zhou, X., Misof, B., Yeates, D.K., 2020. Are fleas highly modified Mecoptera? Phylogenomic resolution of Antliophora (Insecta: Holometabola). BioRxiv. doi: 10.1101/2020.11.19.390666

Meusemann, K., von Reumont, B.M., Simon, S., Roeding, F., Strauss, S., Kück, P., Ebersberger, I., Walzl, M., Pass, G., Breuers, S., Achter, V., von Haeseler, A.,

Burmester, T., Hadrys, H., Wägele, J.W., Misof, B., 2010. A phylogenomic approach to resolve the arthropod tree of life. Mol. Biol. Evol. 27, 2451–2464.

Minh, B.Q., Hahn, M.W., Lanfear, R., 2020. New methods to calculate concordance factors for phylogenomic datasets. Mol. Biol. Evol. 37, 2727–2733.

Misof, B., Liu, S., Meusemann, K., Peters, R.S., Donath, A., Mayer, C., Frandsen, P.B., Ware, J., Flouri, T., Beutel, R.G., Niehuis, O., Petersen, M., Izquierdo-Carrasco, F., Wappler, T., Rust, J., Aberer, A.J., Aspöck, U., Aspöck, H., Bartel, D., Blanke, A., Berger, S., Bohm, A., Buckley, T.R., Calcott, B., Chen, J., Friedrich, F., Fukui, M., Fujita, M., Greve, C., Grobe, P., Gu, S., Huang, Y., Jermiin, L.S., Kawahara, A.Y., Krogmann, L., Kubiak, M., Lanfear, R., Letsch, H., Li, Y., Li, Z., Li, J., Lu, H., Machida, R., Mashimo, Y., Kapli, P., McKenna, D.D., Meng, G., Nakagaki, Y., Navarrete-Heredia, J.L., Ott, M., Ou, Y., Pass, G., Podsiadlowski, L., Pohl, H., von Reumont, B.M., Schutte, K., Sekiya, K., Shimizu, S., Slipinski, A., Stamatakis, A., Song, W., Su, X., Szucsich, N.U., Tan, M., Tan, X., Tang, M., Tang, J., Timelthaler, G., Tomizuka, S., Trautwein, M., Tong, X., Uchifune, T., Walzl, M.G., Wiegmann, B.M., Wilbrandt, J., Wipfler, B., Wong, T.K.F., Wu, Q., Wu, G., Xie, Y., Yang, S., Yang, Q., Yeates, D.K., Yoshizawa, K., Zhang, Q., Zhang, R., Zhang, W., Zhang, Y., Zhao, J., Zhou, C., Zhou, L., Ziesmann, T., Zou, S., Xu, X., Yang, H., Wang, J., Kjer, K.M., Zhou, X., 2014. Phylogenomics resolves the timing and pattern of insect evolution. Science. 346, 763–767.

Misof, B., Meyer, B., von Reumont, B.M., Kück, P., Misof, K., Meusemann, K., 2013. Selecting informative subsets of sparse supermatrices increases the chance to find correct trees. BMC Bioinformatics 14, 348.

Mongiardino Koch, N., Thompson, J.R., 2020. A total-evidence dated phylogeny of Echinoidea combining phylogenomic and paleontological data. Syst. Biol. doi: 10.1093/sysbio/syaa069

Nesnidal, M.P., Helmkampf, M., Bruchhaus, I., Hausdorf, B., 2010. Compositional heterogeneity and phylogenomic inference of metazoan relationships. Mol. Biol. Evol. 27, 2095–2104.

Neumann, J.S., Desalle, R., Narechania, A., Schierwater, B., Tessler, M., 2020. Morphological characters can strongly influence early animal relationships inferred from phylogenomic data sets. Syst. Biol. doi: 10.1093/sysbio/syaa038

Niehuis, O., Hartig, G., Grath, S., Pohl, H., Lehmann, J., Tafer, H., Donath, A., Krauss, V., Eisenhardt, C., Hertel, J., Petersen, M., Mayer, C., Meusemann, K., Peters, R.S., Stadler, P.F., Beutel, R.G., Bornberg-Bauer, E., McKenna, D.D., Misof, B., 2012. Genomic and morphological evidence converge to resolve the enigma of Strepsiptera. Curr. Biol. 22, 1309–1313.

Pease, J.B., Brown, J.W., Walker, J.F., Hinchliff, C.E., Smith, S.A., 2018. Quartet sampling distinguishes lack of support from conflicting support in the green plant tree of life. Am. J. Bot. 105, 385–403.

Peters, R.S., Meusemann, K., Petersen, M., Mayer, C., Wilbrandt, J., Ziesmann, T., Donath, A., Kjer, K.M., Aspöck, U., Aspöck, H., Aberer, A., Stamatakis, A., Friedrich, F., Hünefeld, F., Niehuis, O., Beutel, R.G., Misof, B., 2014. The evolutionary history of holometabolous insects inferred from transcriptome-based phylogeny and comprehensive morphological data. BMC Evol. Biol. 14, 52.

Philippe, H., Derelle, R., Lopez, P., Pick, K., Borchiellini, C., Boury-Esnault, N., Vacelet, J., Renard, E., Houliston, E., Quéinnec, E., Da Silva, C., Wincker, P., Le Guyader, H., Leys, S., Jackson, D.J., Schreiber, F., Erpenbeck, D., Morgenstern, B., Wörheide, G., Manuel, M., 2009. Phylogenomics revives traditional views on deep animal relationships. Curr. Biol. 19, 706–712.

Pisani, D., Benton, M.J., Wilkinson, M., 2007. Congruence of morphological and molecular phylogenies. Acta Biotheor. 55, 269–281.

Pisani, D., Pett, W., Dohrmann, M., Feuda, R., Rota-Stabelli, O., Philippe, H., Lartillot, N., Wörheide, G., 2015. Genomic data do not support comb jellies as the sister group to all other animals. Proc. Natl. Acad. Sci. 112, 15402–15407.

Prum, R.O., Berv, J.S., Dornburg, A., Field, D.J., Townsend, J.P., Moriarty Lemmon, E., Lemmon, A.R., 2015. A comprehensive phylogeny of birds (Aves) using targeted next-generation DNA sequencing. Nature 526, 569–573.

Quang, L.S., Gascuel, O., Lartillot, N., 2008. Empirical profile mixture models for phylogenetic reconstruction. Bioinformatics 24, 2317–2323.

Rannala, B., Yang, Z., 1996. Probability distribution of molecular evolutionary trees: A new method of phylogenetic inference. J. Mol. Evol. 43, 304–311.

Reddy, S., Kimball, R.T., Pandey, A., Hosner, P.A., Braun, M.J., Hackett, S.J., Han, K.-L., Harshman, J., Huddleston, C.J., Kingston, S., Marks, B.D., Miglia, K.J., Moore, W.S., Sheldon, F.H., Witt, C.C., Yuri, T., Braun, E.L., 2017. Why do phylogenomic data sets yield conflicting trees? Data type influences the avian Tree of Life more than taxon sampling. Syst. Biol. 66, 857–879.

Reeder, T.W., Townsend, T.M., Mulcahy, D.G., Noonan, B.P., Wood, P.L., Sites, J.W., Wiens, J.J., 2015. Integrated analyses resolve conflicts over squamate reptile phylogeny and reveal unexpected placements for fossil taxa. PLoS One 10, e0118199.

Reid, N.M., Hird, S.M., Brown, J.M., Pelletier, T.A., McVay, J.D., Satler, J.D., Carstens, B.C., 2014. Poor fit to the multispecies coalescent is widely detectable in empirical data. Syst. Biol. 63, 322–333.

Ribera, I., Hogan, J.R., Vogler, A.P., 2002. Phylogeny of hydradephagan water beetles inferred from 18S rRNA sequences. Mol. Phylogenet. Evol. 23, 43–62.

Richards, S., 2018. Full disclosure: genome assembly is still hard. PLOS Biol. 16, e2005894.

Rokas, A., Carroll, S.B., 2005. More genes or more taxa? The relative contribution of gene number and taxon number to phylogenetic accuracy. Mol. Biol. Evol. 22, 1337–1344.

Rokas, A., Holland, P.W.H., 2000. Rare genomic changes as a tool for phylogenetics. Trends Ecol. Evol. 15, 454–459.

Roure, B., Baurain, D., Philippe, H., 2013. Impact of missing data on phylogenies inferred from empirical phylogenomic data sets. Mol. Biol. Evol. 30, 197–214.

Roure, B., Philippe, H., 2011. Site-specific time heterogeneity of the substitution process and its impact on phylogenetic inference. BMC Evol. Biol. 11, 17.

Salichos, L., Rokas, A., 2013. Inferring ancient divergences requires genes with strong phylogenetic signals. Nature 497, 327–331.

Sann, M., Niehuis, O., Peters, R.S., Mayer, C., Kozlov, A., Podsiadlowski, L., Bank, S., Meusemann, K., Misof, B., Bleidorn, C., Ohl, M., 2018. Phylogenomic analysis of Apoidea sheds new light on the sister group of bees. BMC Evol. Biol. 18, 71.

Sayyari, E., Mirarab, S., 2016. Fast coalescent-based computation of local branch support from quartet frequencies. Mol. Biol. Evol. 33, 1654–1668.

Sayyari, E., Whitfield, J.B., Mirarab, S., 2017. Fragmentary gene sequences negatively impact gene tree and species tree reconstruction. Mol. Biol. Evol. 34, 3279–3291.

Schierwater, B., Holland, P.W.H., Miller, D.J., Stadler, P.F., Wiegmann, B.M., Wörheide, G., Wray, G.A., DeSalle, R., 2016. Never ending analysis of a century old evolutionary debate: "unringing" the Urmetazoon bell. Front. Ecol. Evol. 4, 5.

Schrempf, D., Lartillot, N., Szöllősi, G., 2020. Scalable empirical mixture models that account for across-site compositional heterogeneity. Mol. Biol. Evol. 37, 3616–3631.

Shedlock, A.M., Okada, N., 2000. SINE insertions: powerful tools for molecular systematics. BioEssays 22, 148–160.

Shen, X.-X., Hittinger, C.T., Rokas, A., 2017. Contentious relationships in phylogenomic studies can be driven by a handful of genes. Nat. Ecol. Evol. 1, 0126.

Shen, X.-X., Salichos, L., Rokas, A., 2016. A genome-scale investigation of how sequence, function, and tree-based gene properties influence phylogenetic inference. Genome Biol. Evol. 8, 2565–2580.

Shepherd, D.A., Klaere, S., 2019. How well does your phylogenetic model fit your data? Syst. Biol. 68, 157–167.

Shull, V.L., Vogler, A.P., Baker, M.D., Maddison, D.R., Hammond, P.M., 2001. Sequence alignment of 18S ribosomal RNA and the basal relationships of adephagan beetles: evidence for monophyly of aquatic families and the placement of Trachypachidae. Syst. Biol. 50, 945–969.

Simion, P., Delsuc, F., Philippe, H., 2020. To What Extent Current Limits of Phylogenomics Can Be Overcome?, in: Scornavacca, C., Delsuc, F., Galtier, N. (Eds.), Phylogenetics in the Genomic Era. No commercial publisher | Authors' open access book, pp. 2.1:1–2.1:34.

Simion, P., Philippe, H., Baurain, D., Jager, M., Richter, D.J., Di Franco, A., Roure, B., Satoh, N., Quéinnec, É., Ereskovsky, A., Lapébie, P., Corre, E., Delsuc, F., King, N., Wörheide, G., Manuel, M., 2017. A large and consistent phylogenomic dataset supports sponges as the sister group to all other animals. Curr. Biol. 27, 958–967.

Simmons, M.P., Sloan, D.B., Gatesy, J., 2016. The effects of subsampling gene trees on coalescent methods applied to ancient divergences. Mol. Phylogenet. Evol. 97, 76–89.

Smith, M.L., Hahn, M.W., 2020. New approaches for inferring phylogenies in the presence of paralogs. Trends Genet. doi: 10.1016/j.tig.2020.08.012

Snel, B., Bork, P., Huynen, M.A., 1999. Genome phylogeny based on gene content. Nat. Genet. 21, 108–110.

Soltis, P.S., Soltis, D.E., 2003. Applying the bootstrap in phylogeny reconstruction. Stat. Sci. 18, 256–267.

Springer, M.S., Molloy, E.K., Sloan, D.B., Simmons, M.P., Gatesy, J., 2020. ILS-aware analysis of low-homoplasy retroelement insertions: inference of species trees and introgression using quartets. J. Hered. 111, 147–168.

Strimmer, K., von Haeseler, A., 1997. Likelihood-mapping: a simple method to visualize phylogenetic content of a sequence alignment. Proc. Natl. Acad. Sci. U. S. A. 94, 6815–6819.

Suh, A., Smeds, L., Ellegren, H., 2015. The dynamics of incomplete lineage sorting across the ancient adaptive radiation of neoavian birds. PLOS Biol. 13, e1002224.

Sullivan, J., Joyce, P., 2005. Model selection in phylogenetics. Annu. Rev. Ecol. Evol. Syst. 36, 445–466.

Suzuki, Y., Glazko, G. V., Nei, M., 2002. Overcredibility of molecular phylogenies obtained by Bayesian phylogenetics. Proc. Natl. Acad. Sci. U. S. A. 99, 16138–16143.

Szucsich, N.U., Bartel, D., Blanke, A., Böhm, A., Donath, A., Fukui, M., Grove, S., Liu, S., Macek, O., Machida, R., Misof, B., Nakagaki, Y., Podsiadlowski, L., Sekiya, K., Tomizuka, S., Von Reumont, B.M., Waterhouse, R.M., Walzl, M., Meng, G., Zhou, X., Pass, G., Meusemann, K., 2020. Four myriapod relatives – but who are sisters? No end to debates on relationships among the four major myriapod subgroups. BMC Evol. Biol. 20, 144.

Thomas, G.W.C., Dohmen, E., Hughes, D.S.T., Murali, S.C., Poelchau, M., Glastad, K., Anstead, C.A., Ayoub, N.A., Batterham, P., Bellair, M., Binford, G.J., Chao, H., Chen, Y.H., Childers, C., Dinh, H., Doddapaneni, H.V., Duan, J.J., Dugan, S., Esposito, L.A., Friedrich, M., Garb, J., Gasser, R.B., Goodisman, M.A.D.,

Gundersen-Rindal, D.E., Han, Y., Handler, A.M., Hatakeyama, M., Hering, L., Hunter, W.B., Ioannidis, P., Jayaseelan, J.C., Kalra, D., Khila, A., Korhonen, P.K., Lee, C.E., Lee, S.L., Li, Y., Lindsey, A.R.I., Mayer, G., McGregor, A.P., McKenna, D.D., Misof, B., Munidasa, M., Munoz-Torres, M., Muzny, D.M., Niehuis, O., Osuji-Lacy, N., Palli, S.R., Panfilio, K.A., Pechmann, M., Perry, T., Peters, R.S., Poynton, H.C., Prpic, N.-M., Qu, J., Rotenberg, D., Schal, C., Schoville, S.D., Scully, E.D., Skinner, E., Sloan, D.B., Stouthamer, R., Strand, M.R., Szucsich, N.U., Wijeratne, A., Young, N.D., Zattara, E.E., Benoit, J.B., Zdobnov, E.M., Pfrender, M.E., Hackett, K.J., Werren, J.H., Worley, K.C., Gibbs, R.A., Chipman, A.D., Waterhouse, R.M., Bornberg-Bauer, E., Hahn, M.W., Richards, S., 2020. Gene content evolution in the arthropods. Genome Biol. 21, 15.

Thomsen, P.F., Elias, S., Gilbert, M.T.P., Haile, J., Munch, K., Kuzmina, S., Froese, D.G., Sher, A., Holdaway, R.N., Willerslev, E., 2009. Non-destructive sampling of ancient insect DNA. PLoS One 4, e5048.

Tong, K.J., Duchêne, S., Ho, S.Y.W., Lo, N., 2015. Comment on "Phylogenomics resolves the timing and pattern of insect evolution". Science. 349, 487.

Tyagi, K., Chakraborty, R., Cameron, S.L., Sweet, A.D., Chandra, K., Kumar, V., 2020. Rearrangement and evolution of mitochondrial genomes in Thysanoptera (Insecta). Sci. Rep. 10, 695.

Vasilikopoulos, A., Balke, M., Beutel, R.G., Donath, A., Podsiadlowski, L., Pflug, J.M., Waterhouse, R.M., Meusemann, K., Peters, R.S., Escalona, H.E., Mayer, C., Liu, S., Hendrich, L., Alarie, Y., Bilton, D.T., Jia, F., Zhou, X., Maddison, D.R., Niehuis, O., Misof, B., 2019. Phylogenomics of the superfamily Dytiscoidea

(Coleoptera: Adephaga) with an evaluation of phylogenetic conflict and systematic error. Mol. Phylogenet. Evol. 135, 270–285.

Vasilikopoulos, A., Gustafson, G.T., Balke, M., Niehuis, O., Beutel, R.G., Misof, B., 2020a. Resolving the phylogenetic position of Hygrobiidae (Coleoptera: Adephaga) requires objective statistical tests and exhaustive phylogenetic methodology: a response to Cai et al. (2020). Mol. Phylogenet. Evol. 106923.

Vasilikopoulos, A., Misof, B., Meusemann, K., Lieberz, D., Flouri, T., Beutel, R.G., Niehuis, O., Wappler, T., Rust, J., Peters, R.S., Donath, A., Podsiadlowski, L., Mayer, C., Bartel, D., Böhm, A., Liu, S., Kapli, P., Greve, C., Jepson, J.E., Liu, X., Zhou, X., Aspöck, H., Aspöck, U., 2020b. An integrative phylogenomic approach to elucidate the evolutionary history and divergence times of Neuropterida (Insecta: Holometabola). BMC Evol. Biol. 20, 64.

Verlinden, H., Sterck, L., Li, J., Li, Z., Yssel, A., Gansemans, Y., Verdonck, R., Holtof, M., Song, H., Behmer, S.T., Sword, G.A., Matheson, T., Ott, S.R., Deforce, D., Van Nieuwerburgh, F., Van de Peer, Y., Vanden Broeck, J., 2020. First draft genome assembly of the desert locust, *Schistocerca gregaria* [version 1; peer review: 2 approved, 1 approved with reservations]. F1000Research 9, 775.

von Reumont, B.M., Jenner, R.A., Wills, M.A., Dell'Ampio, E., Pass, G., Ebersberger, I., Meyer, B., Koenemann, S., Iliffe, T.M., Stamatakis, A., Niehuis, O., Meusemann, K., Misof, B., 2012. Pancrustacean phylogeny in the light of new phylogenomic data: support for Remipedia as the possible sister group of Hexapoda. Mol. Biol. Evol. 29, 1031–1045.

Wahlberg, N., Braby, M.F., Brower, A.V.Z., De Jong, R., Lee, M.M., Nylin, S., Pierce, N.E., Sperling, F.A.H., Vila, R., Warren, A.D., Zakharov, E., 2005. Synergistic effects of combining morphological and molecular data in resolving the phylogeny of butterflies and skippers. Proc. R. Soc. B Biol. Sci. 272, 1577–1586.

Walker, J.F., Brown, J.W., Smith, S.A., 2018. Analyzing contentious relationships and outlier genes in phylogenomics. Syst. Biol. 67, 916–924.

Wang, H.-C., Minh, B.Q., Susko, E., Roger, A.J., 2018. Modeling site heterogeneity with posterior mean site frequency profiles accelerates accurate phylogenomic estimation. Syst. Biol. 67, 216–235.

Wang, H.-C., Susko, E., Roger, A.J., 2019. The relative importance of modeling site pattern heterogeneity versus partition-wise heterotachy in phylogenomic inference. Syst. Biol. 68, 1003–1019.

Wang, Y., Liu, X., Garzón-Orduña, I.J., Winterton, S.L., Yan, Y., Aspöck, U., Aspöck, H., Yang, D., 2017. Mitochondrial phylogenomics illuminates the evolutionary history of Neuropterida. Cladistics 33, 617–636.

Whelan, N. V., Kocot, K.M., Moroz, L.L., Halanych, K.M., 2015. Error, signal, and the placement of Ctenophora sister to all other animals. Proc. Natl. Acad. Sci. U. S. A. 112, 5773–5778.

Whitfield, J.B., Kjer, K.M., 2008. Ancient rapid radiations of insects: challenges for phylogenetic analysis. Annu. Rev. Entomol. 53, 449–472.

Whiting, M.F., Carpenter, J.C., Wheeler, Q.D., Wheeler, W.C., 1997. The Strepsiptera problem: phylogeny of the holometabolous insect orders inferred from 18S and 28S ribosomal DNA sequences and morphology. Syst. Biol. 46, 1–68.

Wickett, N.J., Mirarab, S., Nguyen, N., Warnow, T., Carpenter, E., Matasci, N., Ayyampalayam, S., Barker, M.S., Burleigh, J.G., Gitzendanner, M.A., Ruhfel, B.R., Wafula, E., Der, J.P., Graham, S.W., Mathews, S., Melkonian, M., Soltis, D.E., Soltis, P.S., Miles, N.W., Rothfels, C.J., Pokorny, L., Shaw, A.J., DeGironimo, L., Stevenson, D.W., Surek, B., Villarreal, J.C., Roure, B., Philippe, H., DePamphilis, C.W., Chen, T., Deyholos, M.K., Baucom, R.S., Kutchan, T.M., Augustin, M.M., Wang, J., Zhang, Y., Tian, Z., Yan, Z., Wu, X., Sun, X., Wong, G.K.-S., Leebens-Mack, J., 2014. Phylotranscriptomic analysis of the origin and early diversification of land plants. Proc. Natl. Acad. Sci. U. S. A. 111, E4859–E4868.

Wiegmann, B.M., Regier, J.C., Mitter, C., 2002. Combined molecular and morphological evidence on the phylogeny of the earliest lepidopteran lineages. Zool. Scr. 31, 67–81.

Wiens, J.J., 2005. Can incomplete taxa rescue phylogenetic analyses from long-branch attraction? Syst. Biol. 54, 731–742.

Wiens, J.J., Kuczynski, C.A., Townsend, T., Reeder, T.W., Mulcahy, D.G., Sites  Jr, J.W., 2010. Combining phylogenomics and fossils in higher-level squamate reptile phylogeny: molecular data change the placement of fossil taxa. Syst. Biol. 59, 674–688.

Wiens, J.J., Morrill, M.C., 2011. Missing data in phylogenetic analysis: reconciling results from simulations and empirical data. Syst. Biol. 60, 719–731.

Wilbrandt, J., Misof, B., Panfilio, K.A., Niehuis, O., 2019. Repertoire-wide gene structure analyses: a case study comparing automatically predicted and manually annotated gene models. BMC Genomics 20, 753.

Winterton, S.L., 2003. Molecular phylogeny of Neuropterida with Emphasis on the lacewings (Neuroptera). Entomol. Abhandlungen 61, 158–160.

Winterton, S.L., Hardy, N.B., Wiegmann, B.M., 2010. On wings of lace: phylogeny and Bayesian divergence time estimates of Neuropterida (Insecta) based on morphological and molecular data. Syst. Entomol. 35, 349–378.

Winterton, S.L., Lemmon, A.R., Gillung, J.P., Garzon, I.J., Badano, D., Bakkes, D.K., Breitkreuz, L.C.V., Engel, M.S., Lemmon, E.M., Liu, X., Machado, R.J.P., Skevington, J.H., Oswald, J.D., 2018. Evolution of lacewings and allied orders using anchored phylogenomics (Neuroptera, Megaloptera, Raphidioptera). Syst. Entomol. 43, 330–354.

Winterton, S.L., Yang, L., Wiegmann, B.M., Yeates, D.K., 2001. Phylogenetic revision of Agapophytinae subf.n. (Diptera: Therevidae) based on molecular and morphological evidence. Syst. Entomol. 26, 173–211.

Wortley, A.H., Scotland, R.W., 2006. The effect of combining molecular and morphological data in published phylogenetic analyses. Syst. Biol. 55, 677–685.

Yang, Q., Makarkin, V.N., Winterton, S.L., Khramov, A. V, Ren, D., 2012. A remarkable new family of Jurassic insects (Neuroptera) with primitive wing venation and its phylogenetic position in Neuropterida. PLoS One 7, e44762.

Yang, Z., 2014. Molecular evolution: a statistical approach, 1st ed. Oxford University Press, Oxford.

Yang, Z., 1998. On the best evolutionary rate for phylogenetic analysis. Syst. Biol. 47, 125–133.

Young, A.D., Lemmon, A.R., Skevington, J.H., Mengual, X., Ståhls, G., Reemer, M., Jordaens, K., Kelso, S., Lemmon, E.M., Hauser, M., De Meyer, M., Misof, B., Wiegmann, B.M., 2016. Anchored enrichment dataset for true flies (order Diptera) reveals insights into the phylogeny of flower flies (family Syrphidae). BMC Evol. Biol. 16, 143.

Zhang, C., Scornavacca, C., Molloy, E.K., Mirarab, S., 2020. ASTRAL-Pro: quartet-based species-tree inference despite paralogy. Mol. Biol. Evol. 37, 3292–3307.

Zhang, S.-Q., Che, L.-H., Li, Y., Dan Liang, Pang, H., Ślipiński, A., Zhang, P., 2018. Evolutionary history of Coleoptera revealed by extensive sampling of genes and species. Nat. Commun. 9, 205.

Zhao, J., Li, H., Winterton, S.L., Liu, Z., 2013. Ancestral gene organization in the mitochondrial genome of *Thyridosmylus langii* (McLachlan, 1870) (Neuroptera: Osmylidae) and implications for lacewing evolution. PLoS One 8, e62943.

Zhou, X., Lutteropp, S., Czech, L., Stamatakis, A., von Looz, M., Rokas, A., 2020. Quartet-based computations of internode certainty provide robust measures of phylogenetic incongruence. Syst. Biol. 69, 308–324.

Zwickl, D.J., Hillis, D.M., 2002. Increased taxon sampling greatly reduces phylogenetic error. Syst. Biol. 51, 588–598.

# List of Abbreviations

**aa**: amino acid

**ACSR**: ancestral character state reconstructions

**AICc**: corrected Akaike information criterion

**ASRV**: among-site rate variation

**AU test**: approximately unbiased test

**BI**: Bayesian inference

**BIC**: Bayesian information criterion

**bp**: basepairs

**BS**: bootstrap support

**BSHETU**: Bayesian site-heterogeneous unpartitioned model CAT+GTR+G4

**COG**: cluster of orthologous and single-copy genes

**DNA**: deoxyribonucleic acid

**ER**: equal-rates model

**FcLM**: Four-cluster likelihood mapping

**gDNA**: genomic DNA

**GTD**: Gene-tree discordance

**HGT**: horizontal gene transfer

**IC**: phylogenetic information content

**ILS**: incomplete lineage sorting

**LB score**: Long-branch score

**LM**: likelihood mapping (percentage of fully resolved quartets)

**MCMC**: Markov chain Monte Carlo

**ML**: maximum likelihood

**MSA**: multiple sequence alignment

**MSC**: multi-species coalescent

**Mya**: million years ago

**ng**: nanogram

**NGS**: next-generation sequencing

**OGS**: official gene set

**PI**: Parsimony informative (number of informative sites)

**PMSF**: posterior mean site frequency model

**RAD-seq**: restriction-site-associated DNA sequencing

**RCFV**: relative composition frequency variation

**RF**: Robinson-Foulds distances

**SCA**: summary coalescent analyses

**SCM**: stochastic character mapping

**SH-aLRT**: SH-like approximate likelihood ratio test

**SHETU**: site-heterogeneous unpartitioned model (ML-based, C60)

**SHOMP**: site-homogeneous partitioned model

**SHOMU**: site-homogeneous unpartitioned model

**SH**: SH-aLRT-based average branch support

**SRH**: stationary, homogeneous and (time)-reversible conditions

**SV**: Saturation values

**TBE**: bootstrap support by transfer

**UCE**: Ultraconserved element

**UFBoot1**: ultrafast bootstrap support 1

**UFBoot2**: ultrafast bootstrap support 2

**UFB**: ultrafast bootstrap

**WGS**: whole-genome shotgun sequencing

**µl**: microliter

# List of Supplementary Materials

## Chapter 2

Supplementary materials (i.e., Additional files 1–4) to chapter 2 are freely available online at:

- https://doi.org/10.1186/s12862-020-01631-6 and

- https://doi.org/10.1186/s12862-020-01695-4

The datasets of chapter 2 (and of the corresponding published article) are available online at the Dryad digital repository: https://doi.org/10.5061/dryad.1jwstqjrs.

## Chapter 3

Supplementary materials to chapter 3 (S1, S2, S3) are provided in the accompanying CD-ROM/electronic supplement:

**S1**: Supplementary Text 1 - Supplementary experimental procedures

**S2**: Supplementary Tables 1–4

1. S2A: Suppl. Table 1 - Statistics and source of the official gene sets used in the ortholog set.

2. S2B: Suppl. Table 2 - FcLM results at the amino-acid sequence level

3. S2C: Suppl. Table 3 - FcLM results at the nucleotide sequence level

4. S2D: Suppl. Table 4: Results of outlier sequence removal

## Chapter 3 (continued from previous page)

**Chapter 3** (continued from previous page)

**Chapter 3** (continued from previous page)

**Chapter 3** (continued from previous page)

The datasets of chapter 3 (and of the corresponding published article) are available online at the MENDELEY DATA repository:https://dx.doi.org/10.17632/j8xwxdtbyb.1.

## Chapter 4

Supplementary materials to chapter 4 (S4, S5, S6) are provided in the accompanying CD-ROM/electronic supplement:

**S4:** Supplementary Tables S1–S11

1. Table S1: Transcriptomes used for bait design and downstream phylogenetic analyses

2. Table S2: Species used for hybrid enrichment. Summarized statistics of the hybrid enrichment and sequencing as well as NCBI accession numbers are provided

3. Table S3: Summarized results of the different tiling design experiments

4. Table S4: Collection information for the specimens used for hybrid enrichment

5. Table S5: Results of Mann-Whitney tests for pairs of families concerning the calculated hybrid enrichment statistics

6. Table S6: Summarized results of cross-contamination checks for hybrid-enrichment data

7. Table S7: Convergence statistics of the Bayesian phylogenetic analyses

8. Table S8: Comparison of AICc and BIC scores for different schemes of evolutionary modeling

9. Table S9: Results of model selection in IQ-TREE for all amino-acid supermatrices

10. Table S10: Statistics and information on nucleotide sequence supermatrices

11. Table S11: LB score statistics per species in supermatrix G.


**S5:** Supplementary Text 1 - Supplementary materials and methods

## **Chapter 4** (continued from previous page)

**S6:** Figures S1–S105

1. Fig. S1: Flowchart of the steps used to generate the amino-acid supermatrices

2. Fig. S2: Box-plots of the Ct / Ca ratios separately for each family of Adephaga

3. Fig. S3: Box plots of the number of assigned orthologs for each family of

   Adephaga (hybrid-capture data)

4. Fig. S4: Phylogenetic tree from analysis of supermatrix B (SHETU model)

5. Fig. S5: Phylogenetic tree from analysis of supermatrix C (SHETU model)

6. Fig. S6: Phylogenetic tree from analysis of supermatrix E (SHETU model)

7. Fig. S7: Phylogenetic tree from analysis of supermatrix F (SHETU model)

8. Fig. S8: Phylogenetic tree from analysis of supermatrix G (SHETU model)

9. Fig. S9: Phylogenetic tree from analysis of supermatrix H (SHETU model)

10. Fig. S10: Phylogenetic tree from analysis of supermatrix I (SHETU model)

11. Fig. S11: Phylogenetic tree from analysis of supermatrix J (SHETU model)

12. Fig. S12: Phylogenetic tree from analysis of supermatrix D (BSHETU model)

13. Fig. S13: Phylogenetic tree from analysis of supermatrix D-recoded (BSHETU

    model)

14. Fig. S14: Phylogenetic tree from analysis of supermatrix F (BSHETU model)

15. Fig. S15: Phylogenetic tree from analysis of supermatrix G (BSHETU model)

16. Fig. S15: Phylogenetic tree from analysis of supermatrix H (BSHETU model)

17. Fig. S17: Phylogenetic tree from analysis of supermatrix I (BSHETU model)

18. Fig. S18: Phylogenetic tree from analysis of supermatrix J (BSHETU model)

19. Fig. S19: Phylogenetic tree from analysis of supermatrix nt_A (partitioned)

20. Fig. S20: Phylogenetic tree from analysis of supermatrix nt_B (partitioned)

**Chapter 4** (continued from previous page)

**Chapter 4** (continued from previous page)

**Chapter 4** (continued from previous page)

**Chapter 4** (continued from previous page)

**Chapter 4** (continued from previous page)

# Acknowledgments

First and foremost, I am grateful to Bernhard Misof for his guidance and support throughout my Ph.D. studies. I learned so may things by working close to him and feel I am now a much better scientist. I would also like to thank him for the refreshing scientific discussions and for offering me the opportunity to pursue a Ph.D. degree. Next, I would like to thank Oliver Niehuis because he has always been helpful whenever there were questions and for his critical approach to commenting my manuscripts that provided me with all necessary skills for becoming a critically thinking scientist. Lastly, I want to thank Albert Haas and Diana Imhof for their willingness to be part of my dissertation committee and for their time reading my thesis.

Special thanks I owe to my colleagues and friends Jan Philip Oeyen and Panagiotis Provataris for their support, help, discussions and motivating attitude throughout my Ph.D. studies. I think that things would have been very different for all of us in pursuing our Ph.D. if we had not received the right amount of mental support and advice from each other. I also really enjoyed the time that we spent together outside the office.

I also owe many thanks to Sandra Kukowka for her help and training in the molecular laboratory methods and to Claudia Etzbauer for allowing the molecular lab work in my project to run smoothly.

Thanks to the following people in no particular order for the interesting scientific discussions, wonderful collaborations and valuable help: Karen Meusemann, Michael Balke, Rolf Beutel, Doria Lieberz, Lars Podsiadlowski, Alexander Donath, Malte Petersen, Jeanne Wilbrandt, Robert Waterhouse, Tanja Ziesmann, Sebastian Martin, James Pflug, Jes Rust, Torsten Wappler, David Maddison, Ulrike Aspöck, Horst Aspöck, Xingyue Liu, Paschalia Kapli, Tomáš Flouri, Patrick Kück, Christoph Mayer, Niklas Noll, Lars Dietz, France Gimnich, Ralph Peters and anyone else I might have forgotten.

I would also like to thank the German Research Foundation for funding my Ph.D. project and for allowing me and my collaborators to get further insights into the phylogeny and evolution of Adephaga. Additionally, many thanks to the Zoological Research Museum Alexander Koenig for funding parts of my studies and for offering me the space and infrastructure to conduct my experiments.

Lastly, I would like to thank all members of my family for their mental and emotional support throughout my Ph.D. studies. Many times things did not go as I had planned but you have always been there for me and for this I am grateful.