CrossMark

# Transcriptome and target DNA enrichment sequence data provide new insights into the phylogeny of vespid wasps (Hymenoptera: Aculeata: Vespidae)

Sarah Bank[a,1], Manuela Sann[a,b,1], Christoph Mayer[a], Karen Meusemann[a,b], Alexander Donath[a], Lars Podsiadlowski[c], Alexey Kozlov[d], Malte Petersen[a], Lars Krogmann[e], Rudolf Meier[f], Paolo Rosa[g], Thomas Schmitt[h], Mareike Wurdack[b,h], Shanlin Liu[i,j,k], Xin Zhou[l,m], Bernhard Misof[a], Ralph S. Peters[n,*], Oliver Niehuis[a,b,*]

[a] Center for Molecular Biodiversity Research, Zoological Research Museum Alexander Koenig, Adenauerallee 160, 53113 Bonn, Germany
[b] Department of Evolutionary Biology and Ecology, Institute for Biology I (Zoology), University of Freiburg, Hauptstraße 1, 79104 Freiburg, Germany
[c] Institute of Evolutionary Biology and Ecology, University of Bonn, An der Immenburg 1, 53121 Bonn, Germany
[d] Scientific Computing Group, Heidelberg Institute for Theoretical Studies, Schloss-Wolfsbrunnenweg 35, 69118 Heidelberg, Germany
[e] Department of Entomology, State Museum of Natural History, Rosenstein 1, 70191 Stuttgart, Germany
[f] National University of Singapore, 14 Science Dr 4, Singapore 117543, Singapore
[g] Via Belvedere 8/d, 20044 Bernareggio MI, Italy
[h] Department of Animal Ecology and Tropical Biology, Biocenter, University of Würzburg, Am Hubland, 97074 Würzburg, Germany
[i] China National GeneBank-Shenzhen, BGI-Shenzhen, Shenzhen, Guangdong Province 518083, People's Republic of China
[j] BGI-Shenzhen, Shenzhen, Guangdong Province 518083, People's Republic of China
[k] Centre for GeoGenetics, Natural History Museum of Denmark, University of Copenhagen, Øster Voldgade 5–7, 1350 Copenhagen, Denmark
[l] Beijing Advanced Innovation Center for Food Nutrition and Human Health, China Agricultural University, Beijing 100193, People's Republic of China
[m] Department of Entomology, China Agricultural University, Beijing 100193, People's Republic of China
[n] Center of Taxonomy and Evolutionary Research, Arthropoda Department, Zoological Research Museum Alexander Koenig, Adenauerallee 160, 53113 Bonn, Germany

## ARTICLE INFO

## ABSTRACT

The wasp family Vespidae comprises more than 5000 described species which represent life history strategies ranging from solitary and presocial to eusocial and socially parasitic. The phylogenetic relationships of the major vespid wasp lineages (i.e., subfamilies and tribes) have been investigated repeatedly by analyzing behavioral and morphological traits as well as nucleotide sequences of few selected genes with largely incongruent results. Here we reconstruct their phylogenetic relationships using a phylogenomic approach. We sequenced the transcriptomes of 24 vespid wasp and eight outgroup species and exploited the transcript sequences for design of probes for enriching 913 single-copy protein-coding genes to complement the transcriptome data with nucleotide sequence data from additional 25 ethanol-preserved vespid species. Results from phylogenetic analyses of the combined sequence data revealed the eusocial subfamily Stenogastrinae to be the sister group of all remaining Vespidae, while the subfamily Eumeninae turned out to be paraphyletic. Of the three currently recognized eumenine tribes, Odynerini is paraphyletic with respect to Eumenini, and Zethini is paraphyletic with respect to Polistinae and Vespinae. Our results are in conflict with the current tribal subdivision of Eumeninae and thus, we suggest granting subfamily rank to the two major clades of "Zethini": Raphiglossinae and Zethinae. Overall, our findings corroborate the hypothesis of two independent origins of eusociality in vespid wasps and suggest a single origin of using masticated and salivated plant material for building nests by Raphiglossinae, Zethinae, Polistinae, and Vespinae. The inferred phylogenetic relationships and the open access vespid wasp target DNA enrichment probes will provide a valuable tool for future comparative studies on species of the family Vespidae, including their genomes, life styles, evolution of sociality, and co-evolution with other organisms.

* Corresponding authors at: Center of Taxonomy and Evolutionary Research, Arthropoda Department, Zoological Research Museum Alexander Koenig, Adenauerallee 160, 53113 Bonn, Germany (R.S. Peters). Department of Evolutionary Biology and Ecology, Institute for Biology I (Zoology), University of Freiburg, Hauptstraße 1, 79104 Freiburg, Germany (O. Niehuis).
  E-mail addresses: r.peters@leibniz-zfmk.de (R.S. Peters), oliver.niehuis@biologie.uni-freiburg.de (O. Niehuis).
  [1] These two authors contributed equally to this work. Authors in alphabetic order.

# 1. Introduction

Vespids (Hymenoptera: Vespidae) represent a well-characterized group of more than 5000 described species of stinging wasps (Aculeata) (Carpenter, 1982; Brothers and Carpenter, 1993; Pickett and Carpenter, 2010). Most vespid wasp species are solitary and exhibit a predatory lifestyle providing their offspring with larvae of either moths (Lepidoptera), beetles (Coleoptera), or sawflies (Hymenoptera: Tenthredinidae) (Iwata, 1976; Krombein, 1979; Carpenter and Cumming, 1985; Budriene, 2003). Species of the subfamily Masarinae show a behavioral switch to collecting pollen and nectar as food source for their offspring (Gess, 1996). Besides solitary forms, vespids encompass obligatorily and facultatively eusocial species, presocial forms, and social parasites (Crespi and Yanega, 1995; Hunt, 2007; Archer, 2012). These extraordinary behavioral features have fueled many studies on the evolution of sociality within insects, but the basic question, how often eusociality evolved within Vespidae, has still remained controversial due to conflicting hypotheses regarding the phylogenetic relationships among major vespid wasp lineages (*e.g.*, Carpenter, 1982, 2003; Schmitz and Moritz, 1998; Hines et al., 2007; Pickett and Carpenter, 2010).

Phylogenetic inferences regarding vespid relationships have primarily been based on analyzing morphological and behavioral characters (*e.g.*, Carpenter, 1982, 1987, 1988a,b, 1991, 1993, 1996; Carpenter and Cumming, 1985; Carpenter and Rasnitsyn, 1990; Vernier, 1997; Gess, 1998; Krenn et al., 2002; Arévalo et al., 2004; Carpenter and Perera, 2006; Hermes et al., 2013; Perrard et al., 2017). The results from these studies led to the widely accepted recognition of six subfamilies, whose phylogenetic relationships are hypothesized to be as follows: Euparagiinae + (Masarinae + (Eumeninae + (Stenogastrinae + (Polistinae + Vespinae)))) (see Fig. 1, left diagram). According to this system, the three eusocial groups Stenogastrinae, Polistinae, and Vespinae constitute a monophylum, which implies that eusociality evolved only once in the family Vespidae.

The phylogenetic relationships of vespid wasps inferred from molecular sequence data are largely incongruent with those based on morphological and behavioral traits (Schmitz and Moritz, 1998; Hines et al., 2007; Peters et al., 2017), implying two origins of eusociality and challenging the monophyly of the Eumeninae (see Fig. 1, right diagram). Studying DNA sequence data of a mitochondrial and of a nuclear ribosomal gene, Schmitz and Moritz (1998) were the first to show that Polistinae and Vespinae are likely more closely related to Eumeninae than to Stenogastrinae. However, the authors' conclusions were rejected by Carpenter (2003) who argued that a combined analysis of the molecular sequence data with available morphological and behavioral trait information supports the traditional concept of vespid wasp relationships. Yet, Hines et al. (2007) inferred the same phylogenetic relationships as Schmitz and Moritz (1998) by studying a set of four nuclear encoded genes (including one analyzed also by Schmitz and Moritz (1998)) and a significantly improved taxon sample. The conclusions drawn by Hines et al. (2007) were later contradicted by Pickett and Carpenter (2010). In a recent phylogenomic study of all major

lineages of Hymenoptera (Peters et al., 2017), Stenogastrinae were proposed to be sister group of the remaining Vespidae. However, as this study included only few representatives of major vespid wasp lineages, it did not assess the phylogenetic position of the enigmatic eumenine tribe Zethini, which Hines et al. (2007) inferred as sister lineage of Polistinae and Vespinae (but see also Pickett and Carpenter, 2010).

Leaving the controversy about the phylogenetic position of Stenogastrinae aside, the phylogenetic relationships inferred by Hines et al. (2007) challenged the concept of monophyletic Eumeninae, the largest vespid wasp subfamily comprising more than 3500 species (Pickett and Carpenter, 2010). The Eumeninae *sensu* Carpenter (1982) (or Eumenidae, as the group was formerly given family status; Richards, 1962) unites the former subfamilies Eumeninae, Raphiglossinae, and Zethinae. Recently, Hermes et al. (2013) conducted a comprehensive phylogenetic study by analyzing morphological characters of species of the above three lineages. The results led the authors to subdivide the subfamily Eumeninae into three tribes: Eumenini (including part of the former Eumeninae), Odynerini (including the remaining part of the former Eumeninae), and Zethini (comprising the former Raphiglossinae and Zethinae). While the results of Hines et al. (2007) are compatible with two monophyletic tribes Eumenini and Odynerini within a subfamily Eumeninae, they argued that Zethini are more closely related to Polistinae and Vespinae than to the remaining Eumeninae (but see also Pickett and Carpenter, 2010). Therefore, Hines et al. (2007) suggested granting Zethini again subfamily status. However, the taxonomic sampling available to Hines et al. (2007) did not include samples of the species-poor former subfamily Raphiglossinae. It thus remained unclear whether these should be included in the subfamily Zethinae.

The lack of a robust phylogeny of vespids and in particular of the subfamily Eumeninae is not only a major obstacle for the stability of the classification of vespid wasps, but also for interpreting the group's evolutionary history. A poor understanding of the vespid wasp phylogenetic relationships makes it furthermore difficult to understand the evolution of those cleptoparasites and parasitoids (*e.g.*, cuckoo wasps; Hymenoptera: Chrysididae) that use vespid wasps as hosts (Kimsey and Bohart, 1991; Wurdack et al., 2015). In the present study, we address the most pressing and unresolved questions regarding phylogenetic relationships within the vespid family and establish a basis for future investigations that rely on a robust phylogeny of Vespidae and its subordinated groups. We seek to achieve this goal by two means: (1) simultaneous phylogenetic analyses of transcript and enriched target nucleotide sequence data of a total of 49 vespid wasp species covering all major lineages, except for Euparagiinae and Gayellini (Masarinae) that we were unable to sequence, to (1a) reassess the hypothesis of eusociality having evolved twice in the family Vespidae and (1b) evaluate the monophyly of the subfamily Eumeninae as well as of its tribes; and (2) design and publish a universal set of baits for enrichment of more than 900 single-copy protein-coding genes from Next Generation Sequencing (NGS) libraries of vespids to foster future in-depth phylogenomic analyses in subordinated vespid wasp lineages.
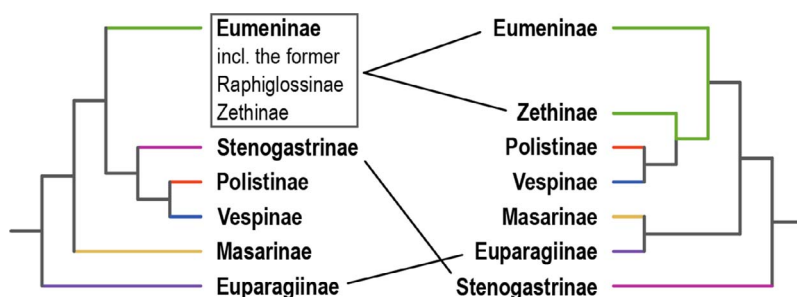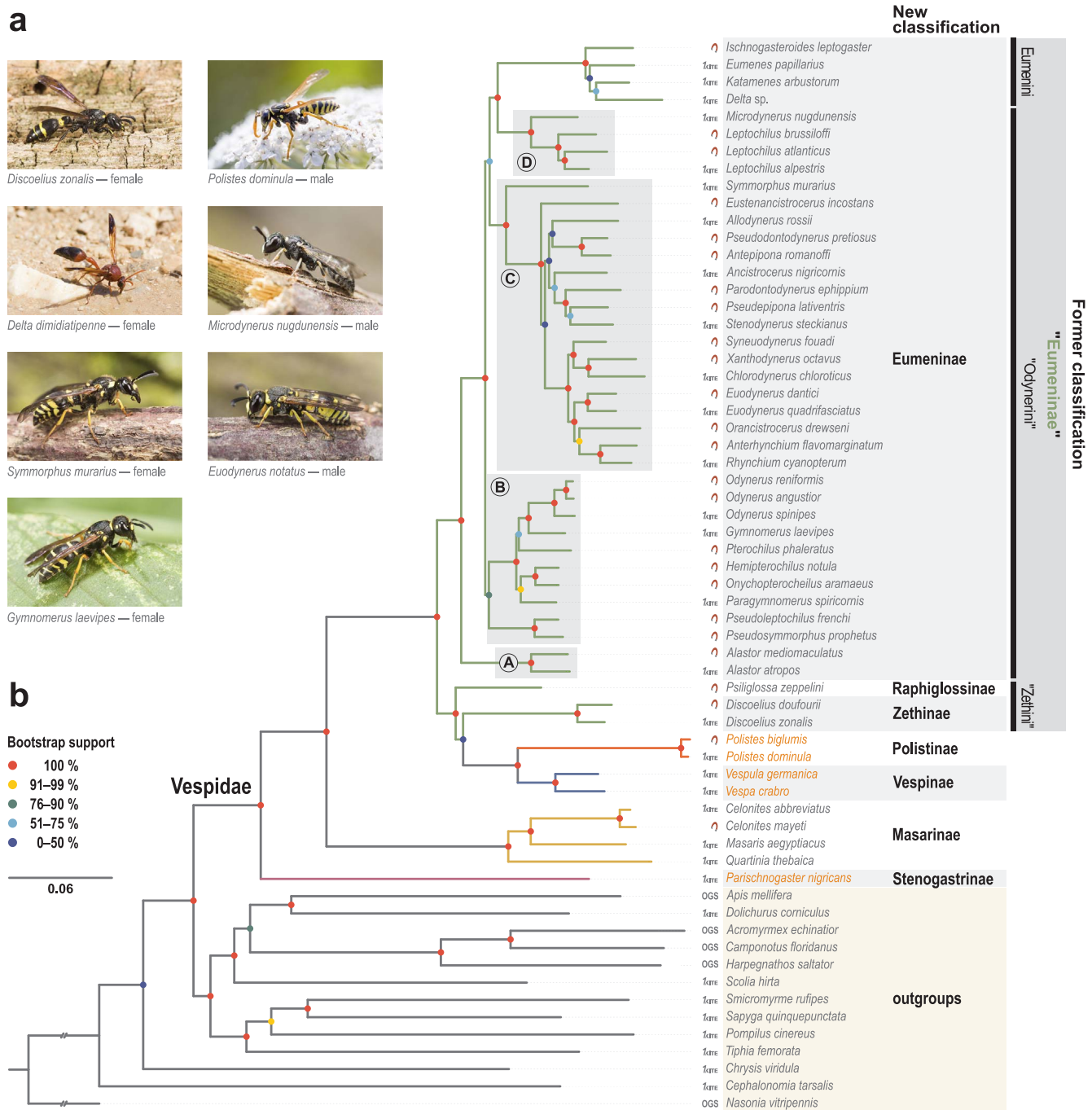


**Fig. 1.** Conflicting hypotheses on vespid subfamily relationships. The left cladogram was obtained by Carpenter (1982) from studying morphological data, the cladogram on the right was obtained by Hines et al. (2007) from studying molecular data. Carpenter (1982) inferred Stenogastrinae as the sister group to Polistinae + Vespinae and included the former subfamilies Raphiglossinae and Zethinae in the subfamily Eumeninae. Hines et al. (2007) inferred Stenogastrinae as sister group to all remaining Vespidae and found Zethinae to be the sister group of Polistinae and Vespinae. The position of the former subfamily Raphiglossinae remained unclear as they were not included in the study by Hines et al. (2007). Branch color-codes adopted from Hines et al. (2007) indicate subfamilies in the classificatory system of Vespidae proposed by Carpenter (1982): blue (Vespinae), green ("Eumeninae"), pink (Stenogastrinae), red (Polistinae), yellow (Masarinae).

(For interpretation of the references to color in this figure legend, the reader is referred to the web version of this article.)

**Fig. 2.** Vespid wasps and their phylogenetic relationships. (A) Representatives of vespid wasps analyzed in the present investigation. All photographs by O. Niehuis. (B) Phylogenetic relationships of major vespid wasp lineages and proposed changes of the taxonomic classification at the subfamily level. The tree was inferred with ExaML, analyzing transcript (1KITE) and enriched (horseshoe magnet) genomic nucleotide sequences plus corresponding nucleotide sequences from five genome projects (OGS) on the translational level (dataset A1/a; 1,004,596 amino acid sites, 511 partitions, see Section 2.9 and Table 1). Support values are inferred from 150 non-parametric bootstrap replicates. The phylogenetic tree was rooted with *Nasonia vitripennis*. Note that the branches connecting *N. vitripennis* with the rest of the topology have been truncated (//). Capitalized letters (A–D) specify clades referred to in the main text. Species names printed in orange letters indicate that the species is eusocial. Branch color-codes adopted from Hines et al. (2007) indicate subfamilies in the classificatory system of Vespidae proposed by Carpenter (1982): blue (Vespinae), green ("Eumeninae"), pink (Stenogastrinae), red (Polistinae), yellow (Masarinae). (For interpretation of the references to color in this figure legend, the reader is referred to the web version of this article.)

## 2. Material and methods

### 2.1. Taxon sampling and sample preservation

We studied a total of 49 species of vespid wasps representing the subfamilies Eumeninae (40 species, including three of the tribe Zethini and representing both the former Raphiglossinae and the former Zethinae), Masarinae (four species), Polistinae (two species),

Stenogastrinae (one species), and Vespinae (two species) (Fig. 2A and B; Supplementary Table 1). Our sampling did not include samples of the vespid wasp lineages Euparagiinae and Gayellini (Masarinae). We also included available transcriptomes of one species of each of the following aculeate wasp families for outgroup comparison (Peters et al., 2017): Ampulicidae, Bethylidae, Chrysididae, Mutillidae, Pompilidae, Sapygidae, Scoliidae, and Tiphiidae (Supplementary Table 1). Finally, we incorporated genomic sequences of the three ant species *Camponotus*

*floridanus*, *Harpegnathos saltator* (Bonasio et al., 2010), *Acromyrmex echinatior* (Nygaard et al., 2011), the honeybee *Apis mellifera* (Honeybee Genome Sequencing Consortium, 2006), and the jewel wasp *Nasonia vitripennis* (Werren et al., 2010). Note that a recently published study, which became available to us after having had completed our analyses, provided new evidence that the wasp family Rhopalosomatidae, which was not part of our taxonomic sampling, is likely the extant sister lineage of Vespidae (Branstetter et al., 2017; see also Pilgrim et al., 2008).

All wasps were hand-collected with an insect net. Samples collected for enriching the DNA of target genes were preserved and stored in 96% ethanol at − 20 °C. Samples collected for transcriptome sequencing were transferred into 2 ml Eppendorf vials containing 0.5 ml of RNAlater (Qiagen GmbH, Hilden, Germany) and were immediately ground with a disposable plastic pestle. Each Eppendorf vial was subsequently filled up to the lid with additional RNAlater and stored at 4 °C for subsequent procedures. Due to the destructive nature of the sample preservation in RNAlater, we preferentially sampled species that are easily identifiable in the field. However, in one instance (*Delta* sp.) the species of a collected sample remained unclear. We exclusively collected adult wasps and focused our sampling on representatives of Central European genera, as their transcriptomes were meant to facilitate future enrichment of target DNA of species occurring especially in this geographic region.

### 2.2. Transcriptome sequencing, assembly, and contamination check

RNA extraction, NGS library preparation, and sequencing of the prepared libraries on Illumina HiSeq sequencers followed the protocols given by Peters et al. (2017) and were conducted by BGI-Shenzhen (China). All cDNA libraries were paired-end (PE) sequenced on Illumina HiSeq2000 sequencing platforms (Illumina Inc., San Diego, CA, USA) with a read length of 150 base pairs (bp). Per species, we obtained about 2.5 Gbp of raw sequence data.

All raw reads were trimmed, assembled, and screened for possible contaminant sequences (which were then removed) as described by Peters et al. (2017). Both raw reads and the assembled transcriptomes are deposited at the Sequence Read Archive (SRA), respectively the Transcriptome Shotgun Assembly (TSA) of the National Center for Biotechnology Information (NCBI) under the Umbrella BioProject accession PRJNA183205 ("The 1KITE project: evolution of insects") (Supplementary Table 2).

### 2.3. Identification and alignment of single-copy genes in the sequenced transcriptomes

We identified contigs of putative single-copy genes in the transcriptome assemblies with Orthograph version 0.5.6 (https://github.com/mptrsen/Orthograph/; Petersen et al., 2017). The applied ortholog set comprised 3260 genes listed by OrthoDB version 7 (Waterhouse et al., 2013) to be single-copy in Holometabola. For the orthology identification in Orthograph, we used the official gene sets of six reference species with well-sequenced and annotated genomes (*A. echinatior*, Official Gene Set (OGS) version 3.8, Nygaard et al., 2011; *C. floridanus* and *H. saltator*, each OGS version 3.3, Bonasio et al., 2010; *A. mellifera*, OGS version 3.2, Honeybee Genome Sequencing Consortium, 2006; *N. vitripennis*, OGS version 2.0, Werren et al., 2010; *Tribolium castaneum*, OGS version 3.0, Tribolium Genome Sequencing Consortium, 2008). For details on the ortholog set and the applied Orthograph settings, see Peters et al. (2017). We included all five hymenopterans, whose amino acid and nucleotide sequences were part of the ortholog set in our analyses, while data of the flour beetle *T. castaneum* was only considered when identifying orthologous transcripts. The amino acid and nucleotide sequences of all 37 species whose transcriptomes (32 species) or official gene sets (five species) we exploited were further processed by removing terminal stop codons and

masking internal stop codons with 'X' and 'NNN' in the amino acid and nucleotide and sequences, respectively.

The orthologous amino acid sequences of each of the 3260 single-copy genes were aligned with MAFFT version 7.123 (Katoh and Standley, 2013) applying the L-INS-i alignment algorithm. The resulting alignments were checked for outlier amino acid sequences and underwent a refinement procedure described by Misof et al. (2014) except for one difference: when aligning outlier amino acid sequences to respective best matching amino acid sequences of a reference species, we called MAFFT L-INS-i with the "–addfragments" option, since this method is especially suited for aligning short amino acid sequences to an existing alignment. Refined alignments were rechecked for outlier amino acid sequences and remaining outliers were permanently removed from the amino acid alignments as well as from the corresponding nucleotide sequence datasets. We subsequently deleted all gap-only sites (columns) from the resulting amino acid alignments. Finally, we inferred nucleotide sequence alignments from the nucleotide sequence datasets with a modified version of Pal2Nal version 14.1 (Suyama et al., 2006; see Misof et al., 2014 for details on the modification), using the amino acid sequence alignments as blueprints.

### 2.4. Design of baits for enriching genomic DNA of target genes

In order to enlarge our taxonomic sampling, we not only used 24 vespid wasp transcriptomes, but also included additional 25 ethanol-preserved vespid species from which nucleotide sequence data was sampled by enriching and sequencing a set of 913 single-copy genes. For this purpose, we exploited the nucleotide sequence alignments of the orthologous single-copy protein-coding genes (Section 2.3) to design baits for target DNA enrichment (see Section 2.5). We analyzed the aligned transcript sequences of 23 out of the 24 sequenced vespid wasps with the BaitFisher software, version 1.2.7 (Mayer et al., 2016). Note that the transcriptome of *Paragymnomerus spiricornis* was not yet available when bait design was conducted. Using BaitFisher, aligned transcripts were split into individual coding sequence (CDS) sections using the honeybee gene models (OGS version 3.2) and the corresponding genome assembly (version 4.5) as a guide (Elsik et al., 2014). We specified a bait length of 120 bp to optimize the probes for the SureSelect^XT2 Target Enrichment System (Agilent Technologies) for enriching target DNA. Based on preliminary results from phylogenetic analyses of amino acid sequence data obtained from 24 transcriptome assemblies, we demanded that the nucleotide sequence of at least one representative of each of the following taxonomic groups (each group is enclosed by parentheses) was present in full length in all candidate bait regions with the length of the tiling design: (*Parischnogaster nigricans*), (*Quartinia thebaica*), (*Masaris aegyptiacus*), (*Celonites abbreviatus*), (*Discoelius zonalis*), (*Polistes dominula*), (*Vespa crabro*, *Vespula germanica*), (*Alastor atropos*), (*Allodynerus rossii*), (*Microdynerus nugdunensis*, *Leptochilus alpestris*), (*Delta* sp., *Eumenes papillarius*, *Katamenes arbustorum*), (*Gymnomerus laevipes*, *Odynerus spinipes*), (*Symmorphus murarius*), (*Ancistrocerus nigricornis*, *Stenodynerus steckianus*), (*Chlorodynerus chloroticus*, *Euodynerus quadrifasciatus*, *Rhynchium cyanopterum*). Depending on the length of the nucleotide sequence alignment suitable for bait design, we designed seven, five, three, or one bait(s) per CDS, with an offset between consecutive baits of 20 bp. Baits were inferred using the heuristic implementation of the unweighted Hamming 1-center DNA sequence search algorithm, specifying a maximum Hamming distance of 0.15 for clustering nucleotide sequences (Mayer et al., 2016).

BaitFisher designs baits at every potential start position of a bait region. The resulting redundancy (*i.e.*, having more bait start positions than required for realizing a specific tiling design at a given locus) was useful, since we used BaitFilter version 1.0.5 (part of the BaitFisher package) to search for and exclude suggested baits that possibly enrich no-target loci. BaitFilter was run with the following options: "-m blast-l –blast-min-hit-coverage-of-baits-in-tiling-stack 0.84 –blast-first-hit-evalue 0.000001". With these options, BaitFilter searched with the aid

of Blast+ software suite version 2.2.29 (Camacho et al., 2009) all potential baits against a reference genome, in this study an early draft genome assembly of the spiny mason wasp, *O. spinipes* (unpublished data). The blast result was used first to remove bait sets at start positions at which not at least one bait of a given bait stack showed a hit coverage of at least 84% with the target sequence in the reference genome. Second, we removed bait sets at start positions if one (or more) bait(s) in the bait set exhibited a significant sequence similarity with more than one position in the reference genome (*i.e.*, the best and second best hit had e-values smaller than 0.000001; see the BaitFisher and BaitFilter manual for more details). Finally, we used BaitFilter in a separate run to choose the optimal bait region among all remaining bait regions. Specifically, we chose the start position within a given CDS region at which the highest number of transcript sequences was available for designing baits. With BaitFisher, we assessed different tiling designs for 120-bp-long baits, since not all CDS regions contain sufficiently long and suitable alignment segments which can host full tiling designs and contain all required taxa: (1) seven baits tiled across 240 bp with an offset of 20 bp between baits, (2) five baits tiled across 200 bp with an offset of 20 bp between baits, (3) three baits tiled across 160 bp with an offset of 20 bp between baits, (4) two baits tiled across 140 bp with an offset of 20 bp between baits, or (5) a single bait. If tiling designs of different lengths fit into a given CDS, we chose the tiling design with the highest number of tiled baits.

## 2.5. Target DNA enrichment, sequencing, assembly, and contamination check

Genomic DNA (gDNA) was extracted from muscle tissue of 25 vespid wasp species (Supplementary Table 1) using the DNeasy Blood & Tissue Kit (Qiagen, Hilden, Germany) and eluted in 100 μl nuclease-free water. Quality and quantity of the extracted gDNA was assessed with a Fragment Analyzer (Advanced Analytical Technologies GmbH, Heidelberg, Germany) and a Quantus Fluorometer (Promega, Fitchburg, Wisconsin, USA) (Supplementary Table 3).

During library preparation, we followed the SureSelect^XT2 Target Enrichment System Protocol for Illumina Paired-End Multiplexed Sequencing Version E1 published in June 2015 by Agilent Technologies Inc., with some minor modifications. First, gDNA was cut into fragments of 150–400 bp using the Next dsDNase Fragmentase Kit (New England Biolabs Inc., Ipswich, USA) by incubating 100 ng gDNA with 2 μl NEB Next dsDNA Fragmentase and 2 μl 10x Fragmentase Reaction Buffer v2 for 20–25 min. The fragmented gDNA was purified with AMPure XP beads (Beckman Coulter GmbH, Krefeld, Germany) in a ratio of 1:1. Purified fragmented gDNA was subsequently eluted in 30 μl nuclease-free water. The quality and quantity of the fragmented gDNA was assessed by using again a Fragment Analyzer and a Quantus Fluorometer. In the library preparation steps "End Repair", "A-tailing", "Ligation of indexed adapter", and "Pre-amplification of indexed libraries", we reduced the reaction volume specified in Agilent's protocol (pages 43–54 for 100 ng DNA samples) by 50%. For the pre-amplification reaction, we applied the following PCR program: initial denaturation temperature of 98 °C for 2 min, followed by 12 cycles of 30 s at 98 °C, 30 s at 60 °C, and 60 s at 72 °C, followed by a 10 min final extension at 72 °C.

For enriching the target gDNA in the indexed libraries, we continued following the procedure outlined in Agilent's SureSelect^XT2 Target Enrichment System Protocol for Illumina Paired-End Multiplexed Sequencing Version E1 published in June 2015 (pages 55–74) with minor modifications. Briefly, we used a SureSelect^XT2 Custom 59.1 Mbp capture library comprising 49,226 different baits (Supplementary Table 4) and pooled the indexed libraries before the hybridization reaction as follows: one pool (A) comprised the libraries of 14 samples [plus two additional ones not included in this study] (~1.5 μg in total), with each library contributing 93 ng. Two additional pools comprised the indexed libraries of six [plus two not included in

this study] (B1) and of seven [plus one not included in this study] (B2) samples (each with ~750 ng in total), again with each library contributing 93 ng. For more information, see Supplementary Table 5. The two specimens of the species *Odynerus reniformis* and *Eustenancistrocerus inconstans* were enriched twice, once in a 16-samples pool and once in an 8-samples pool for quality control reasons. After pooling the libraries, the total volume of the pools was reduced to 7.0 μl (pool of 14 [total 16] samples) and 3.5 μl (pools of seven [total eight] and six [total eight] samples) with a SpeedVac R SPD 111V (ThermoFisher Scientific, Waltham, MA; USA). Hybridization with the baits was allowed for 48 h at 65 °C in a GeneAmp PCR System 2720. We then initiated the physical separation of the target DNA fragments from the remaining DNA fragments by adding 50 μl Dynabeads MyOne Streptavidin T1 beads and incubating the mixture for 30 min at room temperature. After washing of the beads, the captured DNA was re-suspended in 30 μl nuclease-free water and post-amplified in an on-bead PCR reaction. For the post-amplification, we followed Agilent's protocol by applying the recommended PCR cycling program for a capture library size of > 1.5 Mb with a slightly increased cycle number: initial denaturation temperature of 98 °C for 2 min, followed by 12 cycles of 30 s at 98 °C, 30 s at 60 °C, and 60 s at 72 °C, followed by a 10 min final extension at 72 °C. We purified the amplicons with AMPure XP beads in a ratio of 1:0.7 to remove oligonucleotide primer dimers and to further select for fragments with a size between 200 and 500 bp. Each of the three processed library pools was eluted in 30 μl nuclease-free water and checked for quality and quantity with a Fragment Analyzer and a Quantus Fluorometer.

The three pools of enriched gDNA libraries were sequenced on an Illumina NextSeq 500 Serious sequencer (Illumina Inc., San Diego, CA, USA) with 150 bp PE generating about 0.7 Gbp of raw data per sample (the total amount of raw data of the twice-sequenced samples, *O. reniformis* and *E. inconstans*, was 1.27 Gbp and 1.54 Gbp, respectively). All obtained raw reads were trimmed with Trimmomatic version 0.35 (Bolger et al., 2014) and *de novo*-assembled with IDBA-UD version 1.1.1 (Peng et al., 2012) as described by Mayer et al. (2016). Finally, we searched all contigs sequenced on the same lane against each other using the program blastn of the Blast+ software suite version 2.2.31 (Camacho et al., 2009) in order to identify possible contaminant contigs. Contigs identified as contaminants were removed following the procedure outlined in Mayer et al. (2016), except that we selected a 10-fold expression difference between contigs rather than a 2-fold difference for distinguishing between contaminants and non-contaminants (see Mayer et al., 2016 for details).

## 2.6. Post-processing of assembled gDNA sequences

We used Orthograph version 0.5.6 (Petersen et al., 2017) to search the assembled gDNA data for contigs containing sections of enriched target genes. Orthograph concatenates by default contigs referring to different CDS regions of the same gene and provides the predicted amino acid and corresponding coding nucleotide sequences. However, since Orthograph is optimized to process cDNA rather than gDNA sequences, it translates into intronic sequence sections if possible by chance. This can severely bias downstream analyses, because sequences obtained from applying target DNA enrichment could share a small fraction of erroneously predicted amino acid residues (in contrast to sequences obtained from transcriptome sequencing; see Section 2.2). To remove such erroneously predicted sequence sections, we first mapped the predicted amino acid sequences of the target genes onto the aligned amino acid sequences of transcript origin using MAFFT version 7.273 (Katoh and Standley, 2016) applying the L-INS-i alignment algorithm. This was done by selecting the following alignment options: (1) "–*add*" for adding sequence fragments to an existing multiple sequence alignment, (2) "–*keeplength*" to not allow adding gaps to an existing multiple sequence alignment by removing any extra amino acids from the added sequences, and (3) "–*mapout*" to record information about where amino

acids were removed from the added sequences in order to keep the alignment length fixed. We subsequently used the recorded information of how the extra amino acid sequences were mapped onto the amino acid alignments to edit the nucleotide sequences of the enriched exons and remove corresponding codons. Finally, we aligned the corrected nucleotide sequences of the 25 added vespid wasps to the transcript nucleotide sequences with a modified version of Pal2Nal version 14.1 (Suyama et al., 2006; see Misof et al., 2014 for details on the modification), using the amino acid sequence alignments from the preceding step as blueprints. Next, we identified with custom Perl scripts the individual CDS sections in the amino acid sequence alignments, using the honeybee gene models of OGS version 3.2 in the draft genome assembly version 4.5 as a guide (Elsik et al., 2014). We then removed all amino acid residues that were aligned to non-target CDS sections from sequences obtained via target DNA enrichment. We additionally and conservatively removed with custom Perl scripts all amino acid sequence sections covering less than 95% of the honeybee target exon sequence in each multiple sequence alignment from sequences obtained via target DNA enrichment to ensure that no erroneously translated intronic sequence sections, which could bias the phylogenetic analyses, remained in the dataset. The nucleotide sequence alignments were subsequently processed accordingly, using the amino acid sequence alignments as blueprints and custom Perl scripts.

### 2.7. Enrichment statistics

We calculated the base coverage depth of all full-length or near full-length target exons as well as of the bait-binding sites on each enriched exon by mapping the raw reads onto the respective contig with the software segemehl version 0.2.0 (Hoffmann et al., 2009, 2014). The mapped data was subsequently exploited with SAMtools version 1.2 (Li et al., 2009) to infer base-coverage depth estimates of specific sequence sections (target coding exons and bait-binding sites). We further assessed the extent to which target DNA was enriched, applying the approach suggested by Mayer et al. (2016) for analyzing species with known genome size. While the genome size of none of the enriched species is currently known, we hypothesized that the genome sizes of *P. dominula* (246.3 Mbp; Standage et al., 2016) and *O. spinipes* (197.1 Mbp, inferred by analyzing the k-mer coverage distribution in paired-end sequenced libraries of this species; Niehuis, pers. comm.) are reasonable estimates to those of *Polistes biglumis* and of *Odynerus angustior* and *O. reniformis* for which we estimated the enrichment success. We acknowledge that the genome size even of closely related species can differ. However, significant genome size discrepancies should (in most instances) result in vastly disparate enrichment coefficient estimates when assessing different species, while similar enrichment coefficient estimates would be consistent with the idea of similar genome sizes of these species. Following Mayer et al. (2016), we compared the average base-coverage depth ($C_t$) of the bait-binding sites on sequenced target exons to the average base-coverage depth ($C_g$) expected for the sequenced and assembled fragments of a given genome in the absence of enrichment. $C_g$ was calculated by dividing the total number of nucleotides considered for assembling the library of a respective species by the estimated size of the species' genome. Since we applied different tiling designs for enriching target loci, we also investigated whether or not the tiling design had an impact on the base-coverage depth of the bait-binding sites of enriched target exons, using the base-coverage estimates inferred with SAMtools. However, to reduce edge effects (*i.e.*, the base-coverage depth of one exon influencing the base-coverage depth of a flanking target exon), we restricted our analyses to genes for which we enriched only a single coding exon.

### 2.8. Phylogenetic analyses of transcript sequences and genes from official gene sets

All amino acid alignments were searched for sequence sections

showing random similarity or ambiguously aligned residues with Aliscore version 1.2 (Misof and Misof, 2009; Kück et al., 2010). Aliscore was run with default parameters exept for using the '-e' option to cope with transcript sequence alignments containing many gaps (see Meusemann et al., 2010) and the '–r' option set to $10^{27}$ to compare all sequence pairs in each sliding window.

We decided to apply a protein domain-based partitioning scheme to improve the fit of substitution models for the amino acid and nucleotide sequence data, as suggested by Misof et al. (2014) when studying comparable transcriptome sequence data. We identified protein domains, families, and clans in each predicted (unmasked) transcript alignment on the amino acid level, exploiting information from the protein family databases Pfam-A (release 28; Finn et al., 2014) and Pfam-B (release 27; Finn et al., 2014). Domains were searched for with the aid of PfamScan software version 1.5 (released 2013-10-15, Finn et al., 2014) and HMMER version 3.1b2 (Eddy, 2011) as outlined in Misof et al. (2014) and Peters et al. (2017). The two Pfam databases were separately used to search for domains in the multiple sequence alignment (MSA) of each gene, and the domain with the highest number of hits across all species' sequences in the MSA was selected as the dominant domain. To merge the results of both databases and to avoid overlapping domains, we gave Pfam-A annotations priority over Pfam-B annotations. Please note that we did not consider any of the enriched target gene sequences when searching for protein domains (see Section 2.9). We then merged the coordinates received from the protein domain identification with the information on sites suggested to be removed by Aliscore. We deleted respective sections and concatenated the data blocks into a supermatrix on the amino acid level and generated a corresponding supermatrix on the nucleotide level. During this process, terminal gap symbols ('-') were masked for each data block with 'X' and 'N' in the amino acid and the nucleotide alignments, respectively. All sequence sections were concatenated according to the domain identification as follows: (i) sequence segments identified as Pfam-A domains belonging to the same clan were concatenated to clan-specific data blocks, (ii) sequence segments identified as the same Pfam-A domain (not associated with any clan) were concatenated to Pfam-A domain-specific data blocks, (iii) sequences segments identified as Pfam-B domains were concatenated to Pfam-B domain-specific data blocks, and (iv) sequence segments without any domain annotation were concatenated to the gene-specific data blocks.

The information content within the amino acid supermatrix was evaluated for each data block with the software MARE version 0.1.2-rc (Misof et al., 2013). All data blocks with zero information content were removed. In order to minimize non-random distribution of missing data, we only kept those data blocks that included sequences of each of the 37 species (*i.e.*, 32 species whose transcriptomes we analyzed plus five reference species). We kept the corresponding data blocks from the nucleotide supermatrix.

Having protein domain-based data blocks at hand, we next conducted a two-step heuristic approach to search for both an optimal partitioning scheme and best-fitting substitution models to the inferred partitions. To reduce the complexity of this task, we restricted the search for the best partitioning scheme to a subset of substitution models. Thus, we searched with PartitionFinder version 2.0.0 pre-release 10 (http://www.robertlanfear.com/partitionfinder/; Lanfear et al., 2014, 2016) in combination with RAxML 8.2.4 (Stamatakis, 2014) with the settings '–raxml –weights 1,1,0,1 –rcluster-max 10000 –rcluster-percent 100 –all-states –min-subset-size 50' for the best partitioning of the amino acid supermatrix, allowing only two different substitution models to be used (*i.e.*, LG + G and LG + G + F). Once the best partitioning scheme was found, we assessed in a second step the best fitting model for each partition. This was done with the help of the corrected Akaike information criterion (AICc; Hurvich and Tsai, 1989) and by comparing the fit of the following substitution models: WAG + G, WAG + G + F, BLOSUM62 + G, BLOSUM62 + G + F, DCMUT + G, DCMUT + G + F, JTT + G, JTT + G + F, LG + G, LG + G + F, LG4X. We

used exactly the same partitioning scheme, except that we additionally treated the three codon positions within each partition of the above inferred partitioning scheme as separate partitions, when analyzing the supermatrix at the nucleotide level and applied the GTR + G model to all partitions.

Phylogenetic relationships were inferred by applying the Maximum Likelihood (ML) optimality criterion as implemented in the software ExaML version 3.0.15 (Kozlov et al., 2015). We conducted 50 tree searches: 25 using randomized stepwise addition parsimony starting trees and 25 using completely random starting trees. All starting trees were inferred with RAxML version 8.2.7 (Stamatakis, 2014). The tree with the best log-likelihood score among the 50 evaluated ones was considered to be the best to reflect the phylogenetic hypotheses supported by the analyzed dataset.

We assessed support values for phylogenetic relationships by a partitioned non-parametric bootstrap analysis using a total of 50 (amino acid data set) and 100 (nucleotide data set) bootstrap replicates with ExaML version 3.0.15 (Kozlov et al., 2015). We determined whether or not the number of bootstrap replicates was sufficient for assessing support values for different hypothesis by applying the *a posteriori* bootstopping criterion (Pattengale et al., 2010) implemented in RAxML version 8.2.7 (Weighted Robinson Foulds distance building an extended majority-rule (MRE) consensus tree (autoMRE, threshold [0.03], with 1000 permutations; Stamatakis, 2014)). Bootstrap support values were mapped onto the two inferred best ML phylogenetic trees (one of which is based on the analysis of amino acids [analysis scheme A0], the other one is based on the analysis of nucleotides [analysis scheme N0]), which were subsequently drawn with FigTree version 1.4.3 (Rambaut, 2016) and rooted with the parasitoid wasp *N. vitripennis*. Exported vector graphics were edited with Inkscape version 0.91. Given that FigTree (and other tree visualization software) suffers from a significant software bug resulting in bootstrap support values being assigned to wrong nodes after re-rooting of a tree (Czech et al., 2017), we manually checked all bootstrap support values in the inferred illustrations.

We also searched for rogue taxa in the topologies inferred from analyzing the partitioned amino acid and the partitioned nucleotide sequence data, using the software RogueNaRok version 1.0 (Aberer et al., 2013) with the same wide array of settings as applied and specified by Peters et al. (2017). However, none of the species showed rogue behavior in the phylogenetic analyses.

To assess whether or not the dataset contained conflicting signal that is not obvious from the two inferred phylogenetic trees and to evaluate whether or not confounding signal due to compositional heterogeneity across taxa and/or non-random distribution of missing data (see Dell'Ampio et al., 2014) had an impact on the support of specific phylogenetic hypotheses, we applied the Four-Cluster Likelihood Mapping method (FcLM) on the original amino acid supermatrix as well as on permuted versions of it, following the strategy suggested by Misof et al. (2014). For more information on the approach, please consult Strimmer and von Haeseler (1997), Misof et al. (2014), as well as the legend to Supplementary Fig. 4. Note that we used the LG substitution matrix (Le and Gascuel, 2008) for permuting the supermatrices. FcLM was used to evaluate whether *D. zonalis* (the only representative of the tribe Zethini and of the former subfamily Zethinae whose transcriptome we sequenced), is closer related to Polistinae + Vespinae or to Eumeninae (excl. *D. zonalis*) (Supplementary Table 6). FcLM was done with ExaML version 3.0.17 (Kozlov et al., 2015) on the original amino acid supermatrix, using parsimony start trees, and applying the partitioning scheme and substitution models inferred when analyzing the complete supermatrix at the amino acid level. For the permutation approach, we used the same software and partition scheme, but replaced the original supermatrix with random data inferred with the aid of the LG substitution matrix (we consequently applied LG substitution model across all partitions when analyzing the permuted matrices via FcLM). Results were visualized in simplex graphs, using a custom Perl script.

## 2.9. Phylogenetic analysis of transcript and enriched target coding sequences

We analyzed nucleotide sequences obtained via target DNA enrichment in conjunction with the orthologous transcript sequences using the transcript sequence alignments (Section 2.3) onto which the gDNA sequences (Sections 2.4 and 2.5) had been mapped (Section 2.6). We generated a supermatrix mirroring the one inferred in Section 2.8 when analyzing the transcriptomic sequences alone. Thus, the supermatrix exhibited exactly the same number of sites and partitions. This was achieved by applying all previously acquired information about what sites in the transcript alignments to remove, mask, and combine (Section 2.8) onto the corresponding alignments containing the additional gDNA sequences. However, since the enriched sequences encompassed only a subset of the single-copy genes that were present (and analyzed) in the transcript sequences, we applied various filtering and modeling schemes to assess the impact of missing data and to improve the model fitting: (1) for a partition to be considered in the phylogenetic analysis, each partition previously inferred from analyzing only the transcript sequences (Section 2.8) had to contain the sequences of all 32 species, whose transcriptome we sequenced plus the sequences of the five reference species (same conditions we demanded when analyzing these species alone; Section 2.8), while no minimum number of sequences was specified for the DNA enrichment dataset (datasets based on amino acids [A1] and nucleotides [N1]). (2) Same conditions as in (1), except that sequence data obtained via DNA enrichment from at least one additional species had to be present in a given partition (datasets A2 and N2). (3) For a partition to be considered in the phylogenetic analysis, it had to contain the sequences of all 62 species, *i.e.*, those whose transcriptome we sequenced plus those of the reference species plus those which we sequenced via target DNA enrichment (datasets A3 and N3).

When analyzing the above three datasets at the amino acid level, we applied two substitution modeling schemes: (a) we applied the same substitution models as we did when analyzing the transcript sequences alone (Section 2.8) (analysis schemes A1/a, A2/a, and A3/a), and (b) we inferred the best fitting substitution model to each partition using PartitionFinder version 2.0.0 prerelease 10 (Lanfear et al., 2016) and testing the same substitution models as listed in Section 2.8 (analysis schemes A1/b, A2/b, and A3/b). We thus conducted a total of six phylogenetic inferences on the amino acid level using the maximum likelihood tree inference method implemented in ExaML (*i.e.*, A1/a, A1/b, A2/a, A2/b, A3/a, and A3/b; Supplementary Fig. 1). Since we consistently applied the GTR + G model when studying the dataset on the nucleotide level, the total number of phylogenetic inferences on the nucleotide level was three (*i.e.*, N1, N2, N3). A summary of the data processing workflow is given in Supplementary Fig. 1. Phylogenetic trees were inferred, branch support was assessed, and rogue taxa were identified in each of the analysis schemes as outlined in Section 2.8. As only species within subordinated lineages, whose relationships to each other we intended to infer, exhibited rogue behavior (*i.e.*, *E. papillarius* and *Ischnogasteroides leptogaster* within the tribe Eumenini; *A. rossii, E. incostans*, and *Orancistrocerus drewseni* within clade C of the tribe "Odynerini"; Supplementary Table 7), we refrained from excluding these species in any of our inferences. FcLM, as outlined in Section 2.8, was used to check for conflicting and/or confounding signal when exploring whether *Discoelius* spp. (Zethini and representatives of the former Zethinae) or *P. zeppelini* (Zethini and representative of the former Raphiglossinae) are the closest relatives of Polistinae and Vespinae, analyzing the amino acid datasets A1/a, A1/b, A2/a, A2/b, A3/a, and A3/b (Supplementary Table 8).

To assess the possible impact of the ML tree inference method on the inferred tree topology, we additionally conducted phylogenetic inferences in a Bayesian framework, using the software ExaBayes (Aberer

et al., 2014). We applied this approach exclusively to the datasets A1/a and N1, which contained all compiled sequence information (note that all other datasets represented subsets of the datasets A1/a and N1, and their analysis resulted in virtually identical tree topologies). ExaBayes was run as outlined by Peters et al. (2017), except that we generated Markov chain Monte Carlo chains (four coupled chains in three independent runs) for 1,000,000 generations each when analyzing the dataset A1/a, and for 3,000,000 generations each when analyzing dataset N1. Since one of the three runs got trapped in a local optimum when analyzing dataset N1 which prevented the three runs from converging (average standard deviation of split frequencies [ASDSF] = 12.25%), we additionally sampled trees from a fourth run, which converged with the two previously converging runs (ASDSF = 4.41%). We analyzed only the trees from the three runs that converged. The three runs from analyzing dataset A1/a also converged (ASDSF = 3.97%). While we applied the same data partitioning scheme that we used in ExaML, we enabled automatic substitution model detection when analyzing the amino acid dataset, since ExaBayes does not support the LG4X amino acid substitution model that PartitionFinder suggested to apply on several of the inferred data partitions. Trees were sampled every 500 generations and the first 25% of the sampled trees were discarded (burn-in phase). This resulted in a total of 4500 (dataset A1/a) and 13,500 (dataset N1) sampled trees based from which we calculated posterior probability values.

To assess the possible impact of species in our datasets, whose sequence evolution violated the assumption of global stationary, reversibility, and homogeneity (SRH conditions), on the tree topology (Jermiin et al., 2004; Ababneh et al., 2006), we conducted pairwise sequence comparisons using Bowker's matched-pairs tests of symmetry (Bowker, 1948) and generated heat maps based on the inferred p-values as implemented in SymTest version 2.0.47 (https://github.com/ottmi/symtest). We applied Bowker's test exclusively to the datasets A1/a and N1 (for the same reasons as given above in context of the Bayesian tree inference) and compared the results obtained from analyzing the two datasets with each other.

Branch support values were mapped onto the best corresponding phylogenetic tree. All phylogenetic trees were rooted with *N. vitripennis* as outgroup using FigTree version 1.4.3 (Rambaut, 2016). All bootstrap values in the rooted trees were visually checked (see above and Czech et al., 2017) before further editing the resulting vector graphics with Inkscape version 0.91 for publication.

## 3. Results

### 3.1. Transcriptome sequencing, assembly, contamination screening, and identification of single-copy protein-coding genes

We sequenced transcriptomes of 32 aculeate wasp species in context of the international 1KITE project (some of which had previously been released by Peters et al., 2017) comprising 24 representatives of the family Vespidae and eight outgroup species. All sequences have been submitted to the NCBI Transcriptome Shotgun Assembly (TSA) database (accession numbers are listed in Supplementary Table 2). Per species, we analyzed 5.9–19.0 M (median: 9.9 M) raw reads, which assembled after adapter clipping and quality trimming into 19,607–43,567 (median: 27,522) contigs. We removed between 56 and 2532 contigs identified as possible contaminants per assembly. The number of contigs in the cleaned assemblies consequently dropped to 19,309–43,415 (median: 27,150). We identified transcripts of 2766–3099 (median: 2990) of the 3260 protein-coding single-copy genes in the 32 transcriptomes. The number of different protein-coding single-copy genes identified in at least one of the 32 transcriptomes was 3251 which constitutes the number of gene alignments we obtained. These and additional assembly statistics are summarized in Supplementary Table 9.

### 3.2. Phylogenetic analysis of the transcript sequences

After identification of outlier sequences in the 3251 multiple sequence alignments at the amino acid level and subsequent alignment refinement, we removed 577 sequences referring to 217 single-copy genes. Search for protein domains in the refined amino acid sequence alignments assigned 30% of the alignment sites to Pfam-A domains and 6.1% of the alignment sites to Pfam-B domains. A total of 63.9% of the alignment sites consequently remained unannotated (voids). Based on the domain identification results, we split the 3251 multiple sequence alignments and rearranged their sites into 6066 different data blocks, with each block encoding a given protein domain or protein domain clan (comprising domains with a common evolutionary origin; Finn, 2006) or voids. Overall, 1669 data blocks referred to different Pfam-A domains (or domain clans), 1146 referred to different Pfam-B domains, and 3251 referred to voids of the 3251 analyzed genes. After removing ambiguously aligned sites identified by Aliscore (resulting in 5935 data blocks), removing data blocks that contained no phylogenetic information (resulting in 4939 data blocks), eliminating data blocks that did not encompass sequences of all 37 species (resulting in 2531 data blocks), and concatenating the supermatrix resulted in 1,004,596 amino acid and 3,013,788 nucleotide sites, respectively. Both supermatrices covered the sequences of 2531 data blocks and comprised 850 Pfam-A data blocks (incl. clan data blocks), 197 Pfam-B data blocks, and 1484 unannotated gene data blocks (voids). Finally, PartitionFinder suggested a best partitioning scheme integrating these data blocks into 511 partitions.

Phylogenetic analyses of the transcript sequences (1KITE) in combination with the corresponding sequences from five genome projects (OGS) on the amino acid (dataset A0) and on the nucleotide level (dataset N0) with ExaML resulted in two trees whose ingroup relationships were largely congruent (Supplementary Figs. 2 and 3). The results from analyzing the amino acid sequence data are identical to those illustrated in Fig. 2b. The only notable difference between the two obtained tree topologies was that eumenine clade C was inferred as sister to eumenine clade B when analyzing the nucleotide sequence data (Supplementary Fig. 3). Both analyses consistently inferred Stenogastrinae as sister lineage to all remaining Vespidae and confirmed Masarinae being the closest relatives of a clade ("Eumeninae" + (Polistinae + Vespinae)). Both analyses corroborate that the genus *Discoelius* (Eumeninae: Zethini; representative of the former Zethinae) is more closely related to Polistinae + Vespinae than to the remaining eumenine tribes. Finally, both analyses revealed that the tribe Odynerini (Eumeninae) is paraphyletic, with clade D of the tribe Odynerini being more closely related to Eumenini than to any of the remaining clades (A–C) of the tribe Odynerini. Our analyses also consistently inferred the genus *Alastor* (clade A) as the sister lineage to all remaining Eumeninae excluding Discoelius (Eumeninae: Zethini; representative of the former Zethinae).

When assessing the signal for the relationships of Eumeninae (16 species; excl. *Discoelius*), *Discoelius* (one species; Eumeninae: Zethini; representative of the former Zethinae), Polistinae and Vespinae (three species), and Masarinae and Stenogastrinae plus outgroup taxa (17 species in total) to each other via FcLM, we found the highest support (100% of the quartets in the analysis) for *Discoelius* and Polistinae + Vespinae being closest relatives (Supplementary Fig. 4). Permutation tests did not indicate that the FcLM results obtained when analyzing the original amino acid supermatrix were biased by confounding signal (*e.g.*, violation of SRH conditions, non-random distribution of [missing] data). We therefore consider the strong support for a possible sister group relationship of *Discoelius* to Polistinae + Vespinae in both the ML tree inference and in the FcLM results when analyzing the original supermatrix on the translational level (dataset A0) as reliable.

**Table 1**

Dataset characteristics and analysis schemes in eleven phylogenetic inferences outlined in Sections 2.8 and 2.9 (see also Supplementary Fig. 1). 1KITE: transcripts of 32 species whose transcript libraries were sequenced in the 1KITE project; OGS: genes from the official gene sets of five reference species with sequenced genome; enrichment: genomic sequences of target genes enriched in 25 species.

| Dataset/analysis scheme | Data origin | Number of species | Character type | Size of dataset | Number of partitions | Minimum number of species in each partition | Partition-specific substitution models | Number of bootstrap replicates until convergence |
| --- | --- | --- | --- | --- | --- | --- | --- | --- |
| A0 | 1KITE + OGS | 37 | Amino acids | 1,004,596 | 511 | 37 | Dataset-specific | 50 |
| N0 | 1KITE + OGS | 37 | Nucleotides | 3,013,788 | 1533[a] | 37 | GTR+G | 100 |
| A1/a | 1KITE + OGS + enrichment | 62 | Amino acids | 1,004,596 | 511 | All of A0 | As in A0 | 150 |
| A1/b | 1KITE + OGS + enrichment | 62 | Amino acids | 1,004,596 | 511 | All of A0 | Dataset-specific | 200 |
| A2/a | 1KITE + OGS + enrichment | 62 | Amino acids | 519,093 | 344 | All of A0 + at least 1 and up to 25 | As in A0 | 150 |
| A2/b | 1KITE + OGS + enrichment | 62 | Amino acids | 519,093 | 344 | All of A0 + at least 1 and up to 25 | Dataset-specific | 200 |
| A3/a | 1KITE + OGS + enrichment | 62 | Amino acids | 335,029 | 199 | 62 | As in A0 | 150 |
| A3/b | 1KITE + OGS + enrichment | 62 | Amino acids | 335,029 | 199 | 62 | Dataset-specific | 150 |
| N1 | 1KITE + OGS + enrichment | 62 | Nucleotides | 3,013,788 | 1533[a] | All of N0 | GTR+G | 100 |
| N2 | 1KITE + OGS + enrichment | 62 | Nucleotides | 1,557,279 | 1032[a] | All of N0 + at least 1 and up to 25 | GTR+G | 300 |
| N3 | 1KITE + OGS + enrichment | 62 | Nucleotides | 1,005,087 | 597[a] | 62 | GTR+G | 150 |

[a] Each of the three codon positions of a given partition in the inferred partitioning scheme was treated as separate partition.

### 3.3. Bait design

To extend the taxonomic sampling of vespid wasps by analyzing also ethanol-preserved samples, we inferred baits from the aligned transcripts of 23 vespid species for which transcript libraries were available. BaitFisher, using the parameters and specifications outlined in Section 2.4, suggested a set of 49,226 baits for enriching 2158 coding exons of a total of 913 genes. The 2158 coding exons were targeted using different tiling strategies: (i) seven baits tiled across 240 bp, with a new bait every 20 bp (663 exons referring to 506 genes); (ii) five baits tiled across 200 bp, with a new bait every 20 bp (390 exons referring to 366 genes); (iii) three baits tiled across 160 bp, with a new bait every 20 bp (468 exons referring to 458 genes); (iv) a single bait (637 exons referring to 320 genes).

### 3.4. Capture of target coding sequences

We applied the 49,226 designed baits to enrich coding exons of 913 single-copy genes in 25 vespid species. Per sample, we collected 2.6–7.7 M raw reads (median: 4.4 M). These assembled after adapter clipping and quality trimming into 7224–69,492 contigs (median: 24,884). After removing possible contaminated contigs (10–704 per assembly, median 160), the assemblies comprised 7186–69,361 contigs (median: 24,790). All sequences of the cleaned assemblies are available at Mendeley Data: http://dx.doi.org/10.17632/npht7b2426.2. The assembled transcripts of the species contained 895–911 (median: 904) target genes. After further data removal as outlined in Section 2.6, the number of target genes per species decreased to 671–733 (median: 709). Supplementary Table 10 provides an overview of the assembly statistics and target gene recovery rates.

The base-coverage depth of the enriched coding exons ranged between 355x and 1023x (median: 618x) in the 23 single-sequenced libraries and was 1577x and 1204x in the two double-sequenced libraries of *E. incostans* and *O. reniformis*, respectively. The base-coverage depth of the bait-binding sites ranged between 407x and 1245x (median: 730x) in the 23 single-sequenced libraries and was 1926x and 1468x in the two double-sequenced libraries (Supplementary Table 11).

Assuming congeneric species exhibiting a similar genome size and using the genome sizes of *O. spinipes* (197.1 Mbp) and *P. dominula* (246.3 Mbp) as references, we estimated enrichment coefficients ($C_t/C_g$) of 231, 260, and 298 when considering the average base-coverage depth of the bait-binding sites of *O. angustior* ($C_t$ = 1049; 896.4 Mbp sequenced), *O. reniformis* ($C_t$ = 1468; 1114.8 Mbp sequenced), and *P. biglumis* ($C_t$ = 885; 730.7 Mbp sequenced) (all sequence volumes after adapter clipping and quality trimming).

Comparing the base-coverage depth of bait-binding sites between genes, for which we enriched a single exon each, we found a median increase across the 25 species of 82%, 14%, and 25% when rising the number of tiled baits from one to three, from three to five, and from five to seven, respectively (Supplementary Table 11).

### 3.5. Phylogenetic analysis of transcript and captured target coding sequences

After removal of data blocks without phylogenetic signal from the combined transcriptomic and gDNA alignments, our dataset covered 4939 data blocks. Using this dataset as basis, we applied various filtering and modeling schemes (Section 2.9) that reduced the number of considered data blocks. After (1) eliminating data blocks that did not encompass sequences of all 37 species, the supermatrix consisted of 1,004,596 amino acid and 3,013,788 nucleotide sites, respectively (A1; N1). Both supermatrices covered the sequences of 2531 data blocks and comprised 850 Pfam-A data blocks (incl. clans), 197 Pfam-B data blocks, and 1484 unannotated gene data blocks (merged void regions); (2) eliminating data blocks that did not encompass sequences of all 37 species and have at least one additional sequence obtained via DNA

enrichment present in a given data partition, the supermatrix consisted of 519,093 amino acid and 1,557,279 nucleotide sites, respectively (A2; N2). Both supermatrices covered the sequences of 983 data blocks and comprised 452 Pfam-A data blocks (incl. clans), 79 Pfam-B data blocks, and 452 unannotated gene data blocks (voids); and (3) eliminating data blocks that did not encompass sequences of all analyzed 62 species, the supermatrix consisted of 335,029 amino acid and 1,005,087 nucleotide sites, respectively (A3; N3). Both supermatrices covered the sequences of 376 data blocks and comprised 220 Pfam-A data blocks (incl. clans), 16 Pfam-B data blocks, and 140 unannotated gene data blocks (voids).

Combined analysis of the transcript and enriched genomic nucleotide sequences and corresponding nucleotide sequences from five genome projects (OGS) on the translational (amino acid) and nucleotide level (see Table 1 for additional information on the datasets and analysis schemes) revealed, irrespective of the applied tree inference method, largely congruent topologies (Supplementary Figs. 5–14; see Section 2.9 for details on the various analyses schemes). Bowker's matched-pairs tests of symmetry revealed that the nucleotide dataset N1 strongly violates the SRH conditions (Supplementary Fig. 15). The amino acid dataset A1/a, by contrast, suffers much less from such violations (Supplementary Fig. 16). Most sequence comparisons violating the SRH conditions in the amino acid dataset A1/a include at least one outgroup taxon (*e.g.*, *Acromyrmex echinatior*, *Apis mellifera*, *Camponotus floridanus*, *Chrysis viridula*, *Dolichurus corniculus*, *Harpegnathos saltator*, *Pompilus cincereus*, *Sapyga quinquepunctata*, *Tiphia femorata*) and in particular the two ingroup taxa *Vespa crabro* and *Vespula germanica*.

Differences between the nine inferred topologies concern (i) the phylogenetic relationships of the genera *Discoelius* and *Psiliglossa* (both tribe Zethini and representing the former subfamilies Zethinae and Raphiglossinae) relative to Polistinae + Vespinae, (ii) the phylogenetic relationships of species within clade B and within clade C of the tribe Odynerini, (iii) the phylogenetic relationships of species within the tribe Eumenini, and (iv) the phylogenetic position of Scoliidae relative to Formicidae (both outgroup taxa). In context of the present study, only phylogenetic relationships of the genera *Discoelius* and *Psiliglossa* relative to Polistinae + Vespinae are of special interest (see below).

We inferred similar phylogenetic relationships of the major vespid wasp lineages to those obtained when analyzing the transcript sequence (plus the nucleotide sequences from five genome projects) alone (see Section 3.2; Fig. 2b): Stenogastrinae + (Masarinae + ("Eumeninae" + (Polistinae + Vespinae))). Within "Eumeninae", the tribe Odynerini is paraphyletic and comprises four major clades (A–D) of which clade D is sister to the Eumenini. The genus *Alastor* (clade A) was again inferred as sister lineage to all remaining "Eumeninae" (excl. Zethini). Finally, the obtained topologies strongly corroborate the hypothesis of Zethini being more closely related to Polistinae + Vespinae than to the remaining Eumeninae. However, in none of our analyses did the genera *Discoelius* and *Psiliglossa* cluster in a monophyletic clade Zethini. Instead, five of the inferred topologies suggest *Discoelius* (Zethini and representative of the former Zethinae) being more closely related to Polistinae + Vespinae than to *Psiliglossa* (Zethini and representative of the former Raphiglossinae), although with low bootstrap support (28–35%; Fig. 2; Supplementary Figs. 7–10). The remaining four topologies (including all three analyses on the nucleotide level) suggest *Psiliglossa* being more closely related to Polistinae + Vespinae than to *Discoelius*, but with weak (59%; Supplementary Fig. 6) to moderate (81–86 %; Supplementary Figs. 11, 13, 14) bootstrap support. Phylogenetic analysis in a Bayesian framework (datasets A1/a and N1) suggested *Discoelius* being more closely related to Polistinae + Vespinae with 100% posterior probability when analyzing dataset A1/a (amino acids; Supplementary Fig. 5) and suggest *Psiliglossa* more closely related to Polistinae + Vespinae with 100% posterior probability when analyzing dataset N1 (nucleotides; Supplementary Fig. 12).

We assessed the signal in the datasets A1/a, A1/b, A2/a, A2/b, A3/a, and A3b (Table 1; see also Section 2.9 for further details on the

datasets) for the possible phylogenetic relationships of *Discoelius* spp. (two species; Eumeninae: Zethini; representative of the former Zethinae), *Psiliglossa* (one species; Eumeninae: Zethini; representative of the former Raphiglossinae), Polistinae + Vespinae (four species), and all remaining species (55 species) via FcLM. We found few quartets supporting *Discoelius* and *Psiliglossa* being closely related (11–18% of the quartets in each of the six analyses) or *Discoelius* and Polistinae + Vespinae being closely related (10–22%). The majority of quartets support a closer relationship between *Psiliglossa* and Polistinae + Vespinae (60–72%; see Supplementary Figs. 17 and 18). The results from the FcLM permutation approaches I and II suggest that the support of a closer relationship between *Psiliglossa* and Polistinae + Vespinae when analyzing the original amino acid supermatrix cannot be explained by violation of SRH conditions or non-random distribution of (missing) data (or by a combination of both). Note that in the permutation approach I, which assessed violation of SHR conditions and non-random distribution of (missing) data, a sister group relationship *Discoelius* to Polistinae + Vespinae was supported by 30% of the quartets, indicating that the support for this relationship when analyzing the original supermatrix could be due to confounding signal in dataset A1/b (Supplementary Fig. 17-1b). Unexpected was the support of a close relationship of *Psiliglossa* to Polistinae + Vespinae by 39% of the quartets when applying permutation scheme III (Supplementary Fig. 17-1d). This result could be due to the low number of drawn quartets, which caused a random bias in the completely randomized dataset.

## 4. Discussion

We aimed to infer the phylogenetic relationships of the major vespid wasp lineages (*i.e.*, Eumeninae, Masarinae, Polistinae, Stenogastrinae, Vespinae; excl. Euparagiinae and Gayellini, which were not available to us). Specifically, we were interested in reassessing the hypothesis of eusociality having evolved twice in the family Vespidae and evaluating the monophyly of the subfamily Eumeninae as well as of its tribes (*i.e.*, Eumenini, Odynerini, Zethini). Our results are in line with previous molecular phylogenetic investigations which indicated that Stenogastrinae are likely to be the sister group of all remaining Vespidae (Schmitz and Moritz, 1998; Hines et al., 2007; Peters et al., 2017; see Figs. 1 and 2B), while earlier analyses of morphological and behavioral characters suggested a sister group relationship of Stenogastrinae to Polistinae + Vespidae (Carpenter, 1982, 2003; Pickett and Carpenter, 2010; Hermes et al., 2013). Each molecular phylogenetic study that included Stenogastrinae utilized largely different sets of molecular markers (the studies by Schmitz and Moritz, 1998 and Hines et al., 2007 had one gene in common), which contrast in substitution patterns and evolutionary constrains from each other. Yet, these studies obtained the same result in respect of the phylogenetic position of Stenogastrinae. While our current study builds on the same set of molecular markers as the one by Peters et al. (2017) (*i.e.*, single-copy protein-coding genes), our taxon sample is significantly denser (44 species vs. four species) in the lineage to which Stenogastrinae were previously thought to belong (*i.e.*, "Eumeninae" *sensu lato*, Polistinae, Vespinae; Carpenter, 2003; Pickett and Carpenter, 2010). Nevertheless, our analysis still lacks representatives of the subfamily Euparagiinae and of the tribe Gayellini of the subfamily Masarinae, which would be needed for an even more rigorous test of the relationships between the major vespid lineages. Given that Euparagiinae and Gayellini comprise exclusively solitary nesting species and assuming that the vespid wasp relationships inferred in our study have not been misled by long-branch attraction (Felsenstein, 1978), the specific phylogenetic positions of these two lineages have no impact on our conclusions on how often eusociality evolved within vespid wasps (see below).

The recent confirmation that Rhopalosomatidae likely represent the extant sister lineage of Vespidae (Branstetter et al., 2017; see also Pilgrim et al., 2008) opens up the possibility to even more accurately infer ancestral character states of the family Vespidae by including representatives of Rhopalosomatidae in phylogenetic studies. Having said that, we have currently no reason to assume that the rooting of Vespidae has been compromised by the omission of this outgroup taxon in our study: the number of substitutions that have to be hypothesized along the lineage leading to Vespidae has not been particularly high at the amino acid level. Furthermore, we obtained virtually the same tree topology when analyzing the nucleotide and amino acid datasets irrespective of the tree inference method. Finally, and despite of the fact that deviation from the assumptions of SRH conditions differs significantly between our most comprehensive dataset on the amino acid and on the nucleotide level, we inferred the same topology. This makes us presume that non-stationary processes across the analyzed taxa, which have been reported to have impacted phylogenetic inferences in other lineages of Hymenoptera (Romiguier et al., 2016; Bossert et al., 2017), likely had no major impact on our results. This assumption receives further support from the results of FcLM permutation tests which did not indicate that support for specific phylogenetic hypotheses was driven by compositional heterogeneity and/or non-random distribution of data.

We found *Discoelius* (representative of the former Zethinae) and *Psiliglossa* (representative of the former Raphiglossinae), currently united in the tribe Zethini within the subfamily Eumeninae (Hermes et al., 2013), to be more closely related to Polistinae + Vespinae than to the remaining Eumeninae. Hines et al. (2007) already suggested granting Zethini subfamily status, but our investigation indicates that *Discoelius* and *Psiliglossa* do not necessarily constitute a natural group, since we obtained such a relationship in none of our phylogenetic inferences. However, despite analyzing a significant amount of data, our results are unfortunately not fully conclusive in respect of whether *Discoelius* or *Psiliglossa* is closer related to Polistinae + Vespinae. In our ML tree inferences that suggested a sister group relationship of *Discoelius* to Polistinae + Vespinae, the bootstrap support for this relationship was negligible (28–35%). In those phylogenetic analyses obtained with ExaML that suggested a sister group relationship of *Psiliglossa* to Polistinae + Vespinae, the bootstrap support was 59–86%. The Bayesian phylogenetic inferences provided strong support (100% posterior probability) but contradictory results on whether *Discoelius* or *Psiliglossa* is more closely related to Polistinae + Vespinae. Future studies should improve the taxonomic sampling in this part of the phylogenetic tree (*e.g.*, via target DNA enrichment and exploitation of museum specimens; Mayer et al., 2016) in order to address the phylogenetic relationships between the representatives of the former Raphiglossinae, the representatives of the former Zethinae, and Polistinae + Vespinae. Given the distinct morphology of the former two lineages and their unclear phylogenetic relationship to each other, we propose granting both of them again subfamily status: Raphiglossinae and Zethinae.

The inferred close phylogenetic relationship between Raphiglossinae, Zethinae, and Polistinae + Vespinae substantiates the idea of two independent origins of eusociality within the family Vespidae (Hines et al., 2007): one in the Stenogastrinae and a second in the most recent common ancestor of Polistinae + Vespinae. As outlined by Hines et al. (2007), there are also morphological and behavioral differences between Stenogastrinae and Polistinae + Vespinae that would be consistent with two independent origins of eusociality (*e.g.*, differences in wing morphology, in the provisioning of the larvae, and in the eusocial behavior itself; Hunt, 1991, 2007; Strassmann et al., 1994; Turillazzi, 1991; Yoshikawa et al., 1969). The close phylogenetic relationship between Raphiglossinae, Zethinae, and Polistinae + Vespinae has also implications for the interpretation of the evolution of other traits, such as nest-building: Polistinae and Vespinae are well known for building nests from paper-like material (Evans and West-Eberhard, 1970). Intriguingly, Raphiglossinae and Zethinae apparently also exploit masticated and salivated plant material for constructing their nests (Ferton, 1920; Bischoff, 1927; Blüthgen, 1961 Bohart and

Stange, 1965; Krombein, 1991). Assuming that the use of moistened soil as nest building substrate represents the ancestral character state in Stenogastrinae, Euparagiinae, Masarinae, and Eumeninae *sensu stricto* (Hansell, 1985; Mauss, 2007), the use of plant material for nest-building could represent a synapomorphy of Zethinae, Raphiglossinae, Polistinae, and Vespinae, a hypothesis already discussed by Evans and West-Eberhard (1970). Utilizing plant material enables the eusocial Polistinae and Vespinae to overcome nest size constrains enforced by the limited availability of naturally occurring structures with individual chambers suitable for raising colonies. In this respect, the evolutionary success of Polistinae and Vespinae likely only became possible after their solitary ancestors evolved the ability to exploit masticated and salivated plant material for constructing nests. A similar reasoning has been put forth by Litman et al. (2011) for explaining the evolutionary success of bee lineages that include foreign material in their nest construction. Knowledge of the sister lineage of Polistinae + Vespinae furthermore provides the basis for testing hypotheses on the evolution of eusociality *per se* in vespid wasps (*e.g.*, Hunt and Amdam, 2005).

The phylogenetic relationships within the subfamily Eumeninae (excl. Raphiglossinae and Zethinae) do not support the idea of a monophyletic tribe Odynerini. Given that a tribe Zethini within the subfamily Eumeninae can no longer be justified (see above), the only monophyletic tribe within the subfamily Eumeninae is the Eumenini. We therefore suggest relinquishing a tribal subdivision of the subfamily Eumeninae until the phylogenetic relationships of all major lineages of Eumeninae have been satisfactorily inferred. The present study provides a strong basis for such efforts by delivering both a robust basic phylogenetic framework that can help guide future taxon sampling and designed target DNA enrichment baits.

One aim of our study was to develop and test a dedicated set of target DNA enrichment baits for studying single-copy protein-coding genes in vespid wasps. Target DNA enrichment requires prior knowledge of the target nucleotide sequence in order to design baits for enriching target sites. Given that ingroup nucleotide sequence information for the design of enrichment baits is still often limited, one popular strategy has been to enrich ultra-conserved elements (UCEs) whose nucleotide sequences do not differ even among distantly related reference species for which (typically) sequenced genomes are available (Faircloth et al., 2014; Faircloth, 2017). A second strategy has been termed anchored hybrid-enrichment. It also targets conserved regions of the genome for enrichment, but it copes with known target locus nucleotide sequence variation by using a more diverse set of (reference species-specific) baits per locus (Lemmon et al., 2012). What both strategies have in common is that they primarily exploit the phylogenetic signal of the flanking regions of target loci. The main drawbacks of the two strategies are consequently (a) that it is unavoidable that phylogenetically uninformative sequence sections are enriched and sequenced (due to the fact that these sections serve as anchors for enrichment), (b) that it remains *a priori* uncertain whether or not the obtained flanking sequence sections are orthologous and phylogenetically informative among the analyzed species, and (c) that there is a low probability (primarily when enriching UCEs) that the flanking sequence sections can be analyzed on both the nucleotide and the amino acid level. Being able to study DNA sequences on the amino acid level typically allows to more reliably align the corresponding nucleotide sequences, to phylogenetically analyze more strongly diverged lineages, and to potentially circumvent problems associated with compositional heterogeneity on the nucleotide level (*e.g.*, Misof et al., 2014; present study). For these reasons, we followed a different approach proposed by Mayer et al. (2016). Thus, we first sequenced transcriptomes of representative ingroup species to obtain reliable nucleotide sequence information on potential protein-coding target loci as well as their variation among species. We then designed a set of baits to capture these protein-coding loci in additional species by exploiting all available nucleotide sequence information and optimizing bait design using the software BaitFisher (Mayer et al., 2016). Since the enriched and sequenced protein-coding loci represent a subset of the loci in the sequenced transcriptomes, the enriched protein-coding nucleotide sequences can seamlessly be aligned to the transcriptome sequence data. We consider this a major advantage of the applied approach. The main disadvantage of our approach is the necessity to first have to invest in obtaining ingroup sequence information.

Our sets of baits proved to be highly efficient (~ 231x to 298x), with a DNA sequence recovery of 98–99.8% of the 913 target genes being captured. Note that we conservatively discarded parts of the enriched sequences in downstream analyses due to the fact that the applied software for identifying and concatenating coding target DNA sequences (Orthograph; Petersen et al., 2017) is optimized for analyzing transcript sequences (cDNA) rather than genomic DNA (gDNA) (outlined in Section 2.6). Since we relied on gene models of the honeybee (Elsik et al., 2014) to identify and remove any possibly erroneously annotated coding sequence section that is not necessarily identical to those of the investigated vespid wasps, we focused the phylogenetic analyses on those exons that largely corresponded in length between the honeybee and vespid wasp. The recently published gene models of the European paper wasp, *Polistes dominula* (Standage et al., 2016), had unfortunately not been available for our study, but will allow future studies to use gene models for an ingroup lineage and will likely reduce the amount of discarded data.

The comprehensive set of baits for enriching single-copy protein-coding genes in vespid wasps will facilitate extending the taxonomic sampling considerably, because it allows for exploiting genomic information from ethanol-preserved samples and possibly also from older museum specimens (Mayer et al., 2016). Enrichment of hundreds of exons in closely related species could also enable coping with phylogenetic uncertainties that result from incomplete lineage sorting by applying shortcut coalescence approaches (Liu et al., 2009a,b; Liu et al., 2010; but see also Springer and Gatesy, 2016). At the same time, genome and transcriptome sequencing data will continue to accumulate (*e.g.*, Lopez-Osorio et al., 2017) and rapidly increase our knowledge of the evolutionary history of the family Vespidae.

## Author contributions

B.M., O.N., R.S.P. conceived the study. L.K., M.W., O.N., P.R., R.M., T.S. collected samples. A.D., K.M., L.P., O.N., R.S.P., S.L., X.Z. sequenced, assembled, and processed the transcriptomes. C.M., M.S., O.N. conducted the target DNA enrichment and sequencing experiments. A.K., B.M., C.M., K.M., M.P., M.S., O.N., R.S.P., S.B. phylogenetically analyzed the sequence data. All authors contributed to the writing of the manuscript, with M.S., O.N., R.S.P., S.B. taking the lead.

## Appendix A. Supplementary material

Supplementary data associated with this article can be found, in the online version, at http://dx.doi.org/10.1016/j.ympev.2017.08.020.

## 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., Kobert, K., Stamatakis, A., 2014. ExaBayes: massively parallel Bayesian tree inference for the whole-genome era. Mol. Biol. Evol. 31, 2553–2556.

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

Archer, M.E., 2012. Vespine Wasps of the World. Siri Scientific Press, Manchester, UK.

Arévalo, E., Zhu, Y., Carpenter, J.M., Strassmann, J.E., 2004. The phylogeny of the social wasp subfamily Polistinae: evidence from microsatellite flanking sequences, mitochondrial COI sequence, and morphological characters. BMC Evol. Biol. 4, 8.

Bischoff, H., 1927. Biologie der Hymenopteren. Springer, Berlin.

Blüthgen, P., 1961. Die Faltenwespen Mitteleuropas (Hymenoptera, Diploptera). Abh. Dt. Akad. Wiss. Berlin 2, 1–249.

Bohart, R.M., Stange, L.A., 1965. A revision of the genus *Zethus* in the Western Hemisphere (Hymenoptera, Eumenidae). Univ. Calif. Publ. Entomol. 40, 1–208.

Bolger, A.M., Lohse, M., Usadel, B., 2014. Trimmomatic: a flexible trimmer for Illumina sequence data. Bioinformatics 30, 2114–2120.

Bonasio, R., Zhang, G., Ye, C., Mutti, N.S., Fang, X., Qin, N., Donahue, G., Yang, P., Li, Q., Li, C., Zhang, P., Huang, Z., Berger, S.L., Reinberg, D., Wang, J., Liebig, J., 2010. Genomic comparison of the ants *Camponotus floridanus* and *Harpegnathos saltator*. Science 329, 1068–1071.

Bossert, S., Murraya, E.A., Blaimer, B.B., Danforth, B.N., 2017. The impact of GC bias on phylogenetic accuracy using targeted enrichment phylogenomic data. Mol. Phylogenet. Evol. 111, 149–157.

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

Branstetter, M.G., Danforth, B.N., Pitts, J.P., Faircloth, B.C., Ward, P.S., Buffington, M.L., Gates, M.W., Kula, R.R., Brady, S.G., 2017. Phylogenomic insights into the evolution of stinging wasps and the origins of ants and bees. Curr. Biol. 27, 1019–1025.

Brothers, D.J., Carpenter, J.M., 1993. Phylogeny of Aculeata: Chrysidoidea and Vespoidea. J. Hym. Res. 2, 227–302.

Budriene, A., 2003. Prey of *Symmorphus* wasps (Hymenoptera: Eumeninae) in Lithuania. Acta Zool. Lituanica 13, 306–310.

Camacho, C., Coulouris, G., Avagyan, V., Ma, N., Papadopoulos, J., Bealer, K., Madden, T.L., 2009. BLAST+: architecture and applications. BMC Bioinform. 10, 421.

Carpenter, J.M., 1982. The phylogenetic relationships and natural classification of the Vespoidea (Hymenoptera). Syst. Entomol. 7, 11–38.

Carpenter, J.M., 1987. Phylogenetic relationships and classification of the Vespinae (Hymenoptera: Vespidae). Syst. Entomol. 12, 413–431.

Carpenter, J.M., 1988a. The phylogenetic systems of the Gayellini (Hymenoptera: Vespidae, Masarinae). Psyche 95, 211–241.

Carpenter, J.M., 1988b. The phylogenetic system of the Stenogastrinae (Hymenoptera, Vespidae). J. New York Ent. Soc. 96, 140–175.

Carpenter, J.M., 1991. Phylogenetic relationships and the origin of social behaviour in the Vespidae. In: Ross, K.G., Matthews, R.W. (Eds.), The Social Biology of Wasps. Cornell University Press, Ithaca, New York, USA, pp. 7–32.

Carpenter, J.M., 1993. Biogeographic Patterns in the Vespidae (Hymenoptera): Two Views of Africa and South America. In: Goldblatt, P. (Ed.), Biological Relationships Between Africa and South America Proceedings of the 37th Annual Systematics Symposium, Held at Missouri Botanical Gardens, 4–6 October 1990. Yale Univ. Press, New Haven, London, pp. 139–155.

Carpenter, J.M., 1996. Generic classification of the Australian pollen wasps (Hymenoptera: Vespidae; Masarinae). J. Kans. Entomol. Soc. 69, 384–400.

Carpenter, J.M., 2003. On "Molecular Phylogeny of Vespidae (Hymenoptera) and the Evolution of Sociality in Wasps". Am. Mus. Novit. 3389, 1–20.

Carpenter, J.M., Cumming, J.M., 1985. A character analysis of the North American potter wasps (Hymenoptera: Vespidae; Eumeninae). J. Nat. Hist. 19, 877–916.

Carpenter, J.M., Perera, E.P., 2006. Phylogenetic relationships among yellowjackets and the evolution of social parasitism (Hymenoptera: Vespidae, Vespinae). Am. Mus. Novit. 3507, 1–19.

Carpenter, J.M., Rasnitsyn, A.P., 1990. Mesozoic Vespidae. Psyche 97, 1–20.

Crespi, B.J., Yanega, D., 1995. The definition of eusociality. Behav. Ecol. 6, 109–115.

Czech, L., Huerta-Cepas, J., Stamatakis, A., 2017. A critical review on the use of support values in tree viewers and bioinformatics toolkits. Mol. Biol. Evol. 34, 1535–1542.

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.

Eddy, S.R., 2011. Accelerated profile HMM searches. PLoS Comput. Biol. 7, e1002195.

Elsik, C.G., Worley, K.C., Bennett, A.K., Beye, M., Camara, F., Childers, C.P., de Graaf, D.C., Debyser, G., Deng, J., Devreese, B., Elhaik, E., Evans, J.D., Foster, L.J., Graur, D., Guigo, R., HGSC production teams, Hoff, K.J., Holder, M.E., Hudson, M.E., Hunt, G.J., Jiang, H., Joshi, V., Khetani, R.S., Kosarev, P., Kovar, C.L., Ma, J., Maleszka, R., Moritz, R.F., Muñoz-Torres, M.C., Murphy, T.D., Muzny, D.M., Newsham, I.F., Reese, J.T., Robertson, H.M., Robinson, G.E., Rueppell, O., Solovyev, V., Stanke, M., Stolle, E., Tsuruda, J.M., Vaerenbergh, M.V., Waterhouse, R.M., Weaver, D.B., Whitfield, C.W., Wu, Y., Zdobnov, E.M., Zhang, L., Zhu, D., Gibbs, R.A., 2014. Finding the missing honey bee genes: lessons learned from a genome upgrade. BMC Genom. 15, 86.

Evans, H.E., West-Eberhard, M.J., 1970. The Wasps. University of Michigan Press, Ann Arbor, Michigan.

Faircloth, B.C., 2017. Identifying conserved genomic elements and designing universal bait sets to enrich them. Methods Ecol. Evol. (early access). http://dx.doi.org/10.1111/2041-210X.12754.

Faircloth, B.C., Branstetter, M.G., White, N.D., Brady, S.G., 2014. Target enrichment of ultraconserved elements from arthropods provides a genomic perspective on relationships among Hymenoptera. Mol. Ecol. Res. 15, 489–501.

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

Ferton, C., 1920. Notes détachées sur l'instinct des Hyménoptères mellifères et ravisseurs avec la description deux espèces nouvelles. (9e Série). Ann. Soc. Ent. Fr. 89, 329–375.

Finn, R.D., 2006. Pfam: clans, web tools and services. Nucl. Acids Res. 34, D247–D251.

Finn, R.D., Bateman, A., Clements, J., Coggill, P., Eberhardt, R.Y., Eddy, S.R., Heger, A., Hetherington, K., Holm, L., Mistry, J., Sonnhammer, E.L.L., Tate, J., Punta, M., 2014. Pfam: the protein families database. Nucl. Acids Res. 42, D222–D230.

Gess, F.W., 1998. *Priscomasaris namibiensis* Gess, a new genus and species of Masarinae (Hymenoptera: Vespidae) from Namibia, southern Africa, with a discussion of its position within the subfamily. J. Hym. Res. 7, 296–304.

Gess, S.K., 1996. The Pollen Wasps – Ecology and Natural History of the Masarinae. Harvard University Press, Cambridge, Massachusetts, pp. 1–340.

Hansell, M.H., 1985. The nest material of Stenogastrinae (Hymenoptera, Vespidae) and its effect on the evolution of social behaviour and nest design. Actes Coll. Insectes Soc. 2, 57–63.

Hermes, M.G., Melo, G.A.R., Carpenter, J.M., 2013. The higher-level phylogenetic relationships of the Eumeninae (Insecta, Hymenoptera, Vespidae), with emphasis on *Eumenes* sensu lato. Cladistics 30, 1–32.

Hines, H.M., Hunt, J.H., O'Connor, T.K., Gillespie, J.J., Cameron, S.A., 2007. Multigene phylogeny reveals eusociality evolved twice in vespid wasps. Proc. Natl. Acad. Sci. USA 104, 3295–3299.

Hoffmann, S., Otto, C., Doose, G., Tanzer, A., Langenberger, D., Christ, S., Kunz, M., Holdt, L., Teupser, D., Hackermueller, J., Stadler, P.F., 2014. A multi-split mapping algorithm for circular RNA, splicing, trans-splicing, and fusion detection. Genome Biol. 15, R34.

Hoffmann, S., Otto, C., Kurtz, S., Sharma, C.M., Khaitovich, P., Vogel, J., Stadler, P.F., Hackermueller, J., 2009. Fast mapping of short sequences with mismatches, insertions and deletions using index structures. PLoS Comput. Biol. 5, e1000502.

Honeybee Genome Sequencing Consortium, 2006. Insights into social insects from the genome of the honeybee *Apis mellifera*. Nature 443, 931–949.

Hunt, J.H., 1991. Nourishment and the evolution of the social vespidae. In: Ross, K.G., Matthews, R.W. (Eds.), The Social Biology of Wasps. Cornell University Press, Ithaca, pp. 426–450.

Hunt, J.H., 2007. The Evolution of Social Wasps. Oxford University Press, New York, USA.

Hunt, J.H., Amdam, G.V., 2005. Bivoltinism as an antecedent to eusociality in the paper wasp genus *Polistes*. Science 308, 264–267.

Hurvich, C.M., Tsai, C.L., 1989. Regression and time series model selection in small samples. Biometrika 76, 297–307.

Iwata, K., 1976. Evolution of Instinct – Comparative Ethology of Hymenoptera. Amerind Publishing Co., New Dehli.

Jermiin, L., 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.

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

Katoh, K., Standley, D.M., 2016. A simple method to control over-alignment in the MAFFT multiple sequence alignment program. Bioinformatics 32, 1933–1942.

Kimsey, L.S., Bohart, R.M., 1991. [1990]: The Chrysidid Wasps of the World. Oxford University Press, Oxford, New York, Toronto.

Kozlov, A., Aberer, A.J., Stamatakis, A., 2015. ExaML version 3: a tool for phylogenomic analyses on supercomputer. Bioinformatics 31, 2577–2579.

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.

Krenn, H.W., Mauss, V., Plant, J., 2002. Evolution of the suctorial proboscis in pollen wasps (Masarinae, Vespidae). Arthropod Struct. Develop. 31, 103–120.

Krombein, K.V., 1979. Vespoidea. In: Krombein, K.V., Hurd, P.D., Smith, D.R., Burks, B.D. (Eds.), Catalog of Hymenoptera in America North of Mexico 2. Smithsonian Institution Press, Washington, pp. 1469–1522.

Krombein, K.V., 1991. Biosystematic studies of Ceylonese wasps XIX: Natural history notes in several families (Hymenoptera: Eumenidae, Vespidae, Pompilidae and Crabronidae). Smithson. Contrib. Zool. 515, 1–41.

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., 2016. PartitionFinder 2: new methods for selecting partitioned models of evolution for molecular and morphological phylogenetic analyses. Mol. Biol. Evol. 34, 772–773.

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

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

Li, H., Handsaker, B., Wysoker, A., Fennell, T., Ruan, J., Homer, N., Marth, G., Abecasis, G., Durbin, R., 1000 Genome Project Data Processing Subgroup, 2009. The sequence alignment/map format and SAMtools. Bioinformatics 25, 2078–2079.

Litman, J.R., Danforth, B.N., Eardley, C.D., Praz, C.J., 2011. Why do leafcutter bees cut leaves? New insights into the early evolution of bees. Proc. R. Soc. B 278, 3593–3600.

Liu, L., Yu, L., Kubatko, L.S., Pearl, D.K., Edwards, S.V., 2009a. Coalescent methods for estimating phylogenetic trees. Mol. Phylogenet. Evol. 53, 320–328.

Liu, L., Yu, L., Pearl, D.K., Edwards, S.V., 2009b. Estimating species phylogenies using coalescence times among sequences. Syst. Biol. 58, 468–477.

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-Osorio, F., Pickett, K.M., Carpenter, J.M., Ballif, B.A., Agnarsson, I., 2017. Phylogenomic analysis of yellowjackets and hornets (Hymenoptera: Vespidae, Vespinae). Mol. Phylogenet. Evol. 107, 10–15.

Mauss, V., 2007. Evolution verschiedener Lebensformtypen innerhalb basaler Teilgruppen der Faltenwespen (Hymenoptera, Vespidae). Denisia 20, 701–722.

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 multispecies target DNA enrichment probe design. Mol. Biol. Evol. 33, 1875–1886.

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., Böhm, 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., Schütte, 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., 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., Li, Y., Xu, X., Zhang, Y., Yang, H., Wang, J., 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 Bioinform. 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.

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., 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.

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.

Peng, Y., Leung, H.C., Yiu, S.M., Chin, F.Y., 2012. IDBA-UD: a de novo assembler for single-cell and metagenomic sequencing data with highly uneven depth. Bioinformatics 28, 1420–1428.

Perrard, A., Grimaldi, D., Carpenter, J.M., 2017. Early lineages of Vespidae (Hymenoptera) in Cretaceous amber. Syst. Entomol. 42, 379–386.

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, 1–6.

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 Bioinform. 18, 111.

Pickett, K.M., Carpenter, J.M., 2010. Simultaneous analysis and the origin of eusociality in the Vespidae (Insecta: Hymenoptera). Arthropod Syst. Phylo. 68, 3–33.

Pilgrim, E.M., von Dohlen, C.D., Pitts, J.P., 2008. Molecular phylogenetics of Vespoidea indicate paraphyly of the superfamily and novel relationships of its component families and subfamilies. Zool. Scr. 37, 539–560.

Rambaut, A., 2016. FigTree version 1.4.3 for Mac OS X. http://tree.bio.ed.ac.uk/software/figtree/.

Richards, O.W., 1962. A Revisional Study of the Masarid Wasps. British Museum (Natural History), London, UK.

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.

Schmitz, J., Moritz, R.F.A., 1998. Molecular phylogeny of Vespidae (Hymenoptera) and the evolution of sociality in wasps. Mol. Phylogenet. Evol. 9, 183–191.

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.

Standage, D.S., Berens, A.J., Glastad, K.M., Severin, A.J., Brendel, V.P., Toth, A.L., 2016. Genome, transcriptome and methylome sequencing of a primitively eusocial wasp reveal a greatly reduced DNA methylation system in a social insect. Mol. Ecol. 25, 1769–1784.

Strassmann, J.E., Hughes, C.R., Turillazzi, S., Solís, C.R., Queller, D.C., 1994. Genetic relatedness and incipient eusociality in stenogastrine wasps. Anim. Behav. 48, 813–821.

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.

Suyama, M., Torrents, D., Bork, P., 2006. PAL2NAL: robust conversion of protein sequence alignments into the corresponding codon alignments. Nucl. Acids Res. 34, W609–W612.

Tribolium Genome Sequencing Consortium, 2008. The genome of the model beetle and pest *Tribolium castaneum*. Nature 452, 949–955.

Turillazzi, S., 1991. The Stenogastrinae. In: Ross, K.G., Matthews, R.W. (Eds.), The Social Biology of Wasps. Cornell University Press, Ithaca, pp. 74–98.

Vernier, R., 1997. Essai d'analyse cladistique des genres d'Eumeninae (Vespidae, Hymenoptera) représentés en Europe septentrionale, occidentale et centrale. B. Soc. Neuchâteloise Sci. Nat. 120, 87–98.

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

Werren, J.H., Richards, S., Desjardins, C.A., Niehuis, O., Gadau, J., Colbourne, J.K., Beukeboom, L.W., Desplan, C., Elsik, C.G., Grimmelikhuijzen, C.J., Kitts, P., Lynch, J.A., Murphy, T., Oliveira, D.C., Smith, C.D., van de Zande, L., Worley, K.C., Zdobnov, E.M., Aerts, M., Albert, S., Anaya, V.H., Anzola, J.M., Barchuk, A.R., Behura, S.K., Bera, A.N., Berenbaum, M.R., Bertossa, R.C., Bitondi, M.M., Bordenstein, S.R., Bork, P., Bornberg-Bauer, E., Brunain, M., Cazzamali, G., Chaboub, L., Chacko, J., Chavez, D., Childers, C.P., Choi, J.H., Clark, M.E., Claudianos, C., Clinton, R.A., Cree, A.G., Cristino, A.S., Dang, P.M., Darby, A.C., de Graaf, D.C., Devreese, B., Dinh, H.H., Edwards, R., Elango, N., Elhaik, E., Ermolaeva, O., Evans, J.D., Foret, S., Fowler, G.R., Gerlach, D., Gibson, J.D., Gilbert, D.G., Graur, D., Gründer, S., Hagen, D.E., Han, Y., Hauser, F., Hultmark, D., Hunter 4th, H.C., Hurst, G.D., Jhangian, S.N., Jiang, H., Johnson, R.M., Jones, A.K., Junier, T., Kadowaki, T., Kamping, A., Kapustin, Y., Kechavarzi, B., Kim, J., Kim, J., Kiryutin, B., Koevoets, T., Kovar, C.L., Kriventseva, E.V., Kucharski, R., Lee, H., Lee, S.L., Lees, K., Lewis, L.R., Loehlin, D.W., Logsdon Jr, J.M., Lopez, J.A., Lozado, R.J., Maglott, D., Maleszka, R., Mayampurath, A., Mazur, D.J., McClure, M.A., Moore, A.D., Morgan, M.B., Muller, J., Munoz-Torres, M.C., Muzny, D.M., Nazareth, L.V., Neupert, S., Nguyen, N.B., Nunes, F.M., Oakeshott, J.G., Okwuonu, G.O., Pannebakker, B.A., Pejaver, V.R., Peng, Z., Pratt, S.C., Predel, R., Pu, L.L., Ranson, H., Raychoudhury, R., Rechtsteiner, A., Reese, J.T., Reid, J.G., Riddle, M., Robertson, H.M., Romero-Severson, J., Rosenberg, M., Sackton, T.B., Sattele, D.B., Schlüns, H., Schmitt, T., Schneider, M., Schüler, A., Schurko, A.M., Shuker, D.M., Simões, Z.L., Sinha, S., Smith, Z., Solovyev, V., Souvorov, A., Springauf, A., Stafflinger, E., Stage, D.E., Stanke, M., Tanaka, Y., Telschow, A., Trent, C., Vattathil, S., Verhulst, E.C., Viljakainen, L., Wanner, K.W., Waterhouse, R.M., Whitfield, J.B., Wilkes, T.E., Williamson, M., Willis, J.H., Wolschin, F., Wyder, S., Yamada, T., Yi, S.V., Zecher, C.N., Zhang, L., Gibbs, R.A., 2010. Functional and evolutionary insights from the genomes of three parasitoid *Nasonia* species. Science 327, 343–348.

Wurdack, M., Herbertz, S., Dowling, D., Kroiss, J., Strohm, E., Baur, H., Niehuis, O., Schmitt, T., 2015. Striking cuticular hydrocarbon dimorphism in the mason wasp *Odynerus spinipes* and its possible evolutionary cause (Hymenoptera: Chrysididae, Vespidae). Proc. Royal Soc. B 282, 20151777.

Yoshikawa, K., Ohgushi, R., Sakagami, S.F., 1969. Preliminary report on entomology of the Osaka City University 5th Scientific Expedition to Southeast Asia 1966 – with descriptions of two new genera of stenogastrine wasps by J. van der Vecht. Nat. Life Southeast Asia 6, 153–200.