

Bulk *de novo* mitogenome assembly from pooled total DNA elucidates the phylogeny of weevils (Coleoptera: Curculionoidea)

**Conrad P.D.T. Gillett^{1,2}, Alex Crampton-Platt^{1,3}, Martijn J. T. N. Timmermans^{1,4},
Bjarte Jordal⁵, Brent C. Emerson^{2,6} and Alfried P. Vogler^{1,4}**

¹Department of Life Sciences, The Natural History Museum, London, SW7 5BD,
United Kingdom

²School of Biological Sciences, Centre for Ecology, Evolution and Conservation,
University of East Anglia, Norwich, NR4 7TJ, United Kingdom

³Department of Genetics, Evolution and Environment, University College London,
Gower Street, London, WC1E 6BT, United Kingdom

⁴Division of Biology, Silwood Park Campus, Imperial College London, Ascot, SL5
7PY, United Kingdom

⁵The Natural History Museum, University Museum of Bergen, P.B. 7800, NO-5020
Bergen, Norway

⁶Island Ecology and Evolution Research Group, Instituto de Productos Naturales y
Agrobiología, C/Astrofísico Francisco Sánchez 3, La Laguna, Tenerife, Canary
Islands, 38206, Spain

ABSTRACT

Complete mitochondrial genomes have been shown to be reliable markers for phylogeny reconstruction among diverse animal groups. However, the relative difficulty and high cost associated with obtaining *de novo* full mitogenomes has frequently led to conspicuously low taxon sampling in ensuing studies. Here we report the successful use of an economical and accessible method for assembling complete or near-complete mitogenomes through shot-gun next generation sequencing of a single library made from pooled total DNA extracts of numerous target species. To avoid the use of separate indexed libraries for each specimen, and an associated increase in cost, we incorporate standard PCR-based 'bait' sequences to identify the assembled mitogenomes. The method was applied to study basal relationships in the weevils (Coleoptera: Curculionodea), producing 92 newly assembled mitogenomes obtained in a single Illumina MiSeq run, which were used to analyse the higher-level phylogenetic relationships of weevils. The analysis supported a separate origin of wood-boring behaviour by the subfamilies Scolytinae, Platypodinae and Cossoninae. This finding contradicts morphological hypotheses proposing a close relationship between the first two of these, but is congruent with previous molecular studies, reinforcing the utility of mitogenomes in phylogeny reconstruction. Our methodology provides a technically simple procedure for generating densely sampled trees from whole mitogenomes, and is widely applicable to groups of animals for which bait sequences are the only required prior genome knowledge.

INTRODUCTION

With the advent of high-throughput next generation sequencing (NGS) technologies and their ability to generate large amounts of data suitable for genomic assembly, systematists are increasingly adopting such methods to reconstruct complete mitochondrial genomes (mitogenomes) to infer phylogenies across a diverse range of taxa. Such research has provided compelling insights in studies ranging from the investigation of deep-level metazoan relationships (Osigus et al. 2013) to those within single phyla (e.g. Cnidaria; Kayal et al. 2013), orders (e.g. Primates; Finstermeier et al. 2013), families (e.g. Braconidae wasps; Wei et al. 2010) and genera (e.g. *Architeuthis* giant squid; Winkelmann et al. 2013). Mitogenomes have an intrinsic suitability for phylogenetic analysis due to their unambiguous orthology (Botero-Castro et al. 2013), varying nucleotide substitution rates that contribute to phylogenetic signal at diverse taxonomic ranks, and their uniparental inheritance consistent with bifurcating phylogenetic trees (Curole and Kocher 1999). In addition, mitochondrial DNA (mtDNA) is present in multiple copies per cell, facilitating its amplification and sequencing, which has undoubtedly contributed to the wide use of mitochondrial markers in phylogeny reconstruction. However, in spite of these advantages, complete mitogenome sequencing has been comparatively labour intensive and costly, resulting in often conspicuously few newly-generated mitogenomes per study (e.g. 17 bird mitogenomes in Pacheco et al. (2011), four complete Cnidarian mitogenomes in Kayal et al. (2013) and one cockroach and 13 termite mitogenomes in Cameron et al. (2012). Techniques have almost always included either shotgun sequencing of expensive multiple indexed-libraries (Botero-Castro et al. 2013) or a target-enrichment step such as primer walking using standard PCR amplification of overlapping fragments (Botero-Castro et al. 2013), long-range PCR followed by either sequencing-primer

walking (Roos et al. 2007) or NGS (Timmermans et al. 2010), and hybrid-capture using sheared long-range PCR products as ‘baits’ immobilised on magnetic beads (Winkelmann et al. 2013). While these techniques can generate full mitochondrial genomes, each of them has limitations that generally limit the number of taxa or samples that can be incorporated economically within a study.

The present study aims to address this sampling bottleneck by testing the possibility of parallel *de novo* mitogenome assembly from a single library of pooled genomic DNA from a bulk sample consisting of many species. This method has recently been applied to sequencing of environmental samples of arthropods from a rainforest canopy (Crampton-Platt et al., in review). Here, we apply this technique to investigate the higher-level phylogeny of an extremely diverse superfamily of insects, the superfamily of weevils (Curculionoidea), composed of no fewer than 62,000 described species distributed wherever terrestrial plants grow (Oberprieler et al. 2007). The current higher-level classification proposed by Bouchard et al. (2011) recognises 9 extant families, amongst which the Curculionidae *s.str.* is by far the largest, containing at least 51,000 species in 17 subfamilies and 292 tribes and subtribes. The phylogenetic classification of the weevils was recognised by the eminent beetle taxonomist Crowson (1955) as “...probably the largest and most important problem in the higher classification of Coleoptera...”. Since that time there have been considerable advances in our understanding of the phylogeny of this group, with significant morphological analyses by Kuschel (1995) and Marvaldi (1997). More recently, molecular data have contributed towards reconstructing weevil higher-level relationships, including studies by McKenna et al. (2009), Hundsdoerfer et al. (2009) and Jordal et al. (2011), which each incorporated between two and six gene markers. A recent analysis of 27 weevil mitogenomes using 12 protein-coding genes (Haran et al.,

2013) supported the paraphyly of Curculionidae *s.str.* as currently defined because the subfamily Platypodinae was recovered in a distant position, in a clade with members of the families Dryophthoridae and Brachyceridae, that together were sister to all other Curculionidae. Although undertaken with limited taxon sampling within the Curculionidae *s.str.* (18 tribes), this last study also supported the division of the family into two large clades; one comprising the ‘broad-nosed’ weevils (subfamilies Entiminae, Cyclominae and Hyperinae) and another containing the remaining subfamilies (except for Platypodinae). In the same study a tRNA^{Ala} to tRNA^{Arg} gene order rearrangement was identified in a cluster of six tRNA genes, located between *nad3* and *nad5*, which appears to be a synapomorphy for the ‘broad-nosed’ weevil subfamilies, further supporting their monophyly. This topology was consistent with that proposed by McKenna et al. (2009), who concluded that the initial diversification of weevils occurred on gymnosperm plants during the Early to early Middle Jurassic.

The Platypodinae is one of several weevil subfamilies that are specialist wood-borers, together with the bark-beetles (Scolytinae) and the subfamily Cossoninae, although other subfamilies also contain xylophagous members (e.g. Molytinae, Cryptorhynchinae and Conoderinae). The evolution of wood-boring behaviour was investigated in detail by Jordal et al. (2011), whose analyses incorporated morphological characters together with molecular data, concluding that both Scolytinae and Platypodinae are derived lineages within the Curculionidae *sensu* Oberprieler et al. (2007). However several important head characters that underpin this relationship are likely to be homoplasious and associated with tunnelling habit (Jordal et al. 2011). Thompson (1992) identified distinct characters of the platypodine eighth abdominal sternite and male genitalia, which indicated a distant relationship to Scolytinae and a possible justification for their inclusion in a separate curculionoid

family. Therefore, the question about the polyphyly of wood-boring lineages remains open, and the failure of previous mitogenome studies to recover the platypodine and scolytine lineages as monophyletic (Haran et al. 2013) may be due to limited taxon sampling. The issue therefore may only be resolved if Jordal et al.'s (2011) comprehensive taxon sampling of wood-boring lineages could be matched using mitochondrial genomes.

MATERIALS AND METHODS

Taxon sampling, DNA extraction and quantification

Throughout this study the most recent higher-level classification of Curculionoidea, proposed by Bouchard et al. (2011) is adhered to, whilst the assignment of genera to higher taxa follows the catalogue of Alonzo-Zarazaga and Lyal (1999). A total of 173 weevil specimens identified to species or a higher-level taxon and obtained through collecting or loans were selected for sequencing, including seven different families, 16 subfamilies and 104 tribes within the Curculionidae. DNA was extracted from each ethanol-preserved specimen individually using DNeasy blood and tissue extraction kits (Qiagen). As specimens were selected to represent a wide taxonomic coverage they were acquired from various sources and in different stages of preservation, leading to variable DNA quality, as is common in phylogenetic studies that involve lineages for which DNA-ready material is difficult to obtain. The DNA extracts included in the sequencing pool were not characterised in great detail, but based on bait PCR success are likely to differ in the degree of degradation and purity. Aliquots from 31 specimens had already been extracted for a previous study (Jordal et al. 2011). The concentration of double-stranded DNA (dsDNA) in most extractions (139 of 173) was assayed on a Qubit fluorometer using a dsDNA high-sensitivity kit (Invitrogen).

‘Bait’ sequence PCR

Standard PCR reactions to amplify 4 different fragments of mitochondrial DNA (*cox1* 5’ ‘barcode region’, *cox1* 3’ region, *rrnL* and *cytb*) were undertaken for each of the 173 samples. Primers and reaction conditions are listed in Supplementary Table S1. PCR products were first cleaned with a size-exclusion filter (Merck Millipore) and then Sanger-sequenced; the resulting bait sequences were subsequently employed to identify mitogenomic assemblies in the manner detailed below.

Sample pooling and sequencing

To minimise the effects of DNA concentration on assembly success across all samples, approximately equimolar quantities of genomic DNA from each of the samples were pooled aiming for 10 ng of dsDNA per sample, resulting in a DNA pool of approximately 1.5 µg. This calculation did not consider 31 samples which were not quantified because of limited sample volume. For each of these, a fixed volume of either 5 or 8 µl was added to the pool. Based on the findings of Crampton-Platt et al. (in review), where longer insert size was found to result in longer mitochondrial contigs, a TruSeq library was prepared from the pool aiming for an insert size of 800 bp. Quantification of the final library indicated that the average insert size was 790 bp and this was sequenced on a single Illumina MiSeq run (500-cycle, 250 bp paired-end reads, version 2 reagent kit).

Mitogenomic assembly pipeline

The bioinformatics assembly pipeline used in this study was developed by Crampton-Platt et al. (in review) and is followed here with minor modifications. A list of the

software required (most freely available) is given in Table 1 and a schematic overview of the principal steps is presented in Figure 1. In brief, the raw data were trimmed of adapters using Trimmomatic (Lohse et al. 2012), and putative mitochondrial reads were identified in a BLAST search (Altschul et al. 1990) against a custom reference database of 258 Coleoptera mitogenomes ($E=1e-5$; no restriction in length overlap). The extracted mtDNA reads were subjected to whole-genome shotgun assembly using Celera Assembler (Myers et al. 2000) and IDBA-UD (Peng et al. 2011), and the resulting contigs were filtered again for mtDNA hits against the Coleoptera reference library for sequences of >1000 bp overlap at $E=1e-5$. Both assemblies were merged using Minimus2 (Sommer et al. 2007) to combine overlapping sequences from both assemblers into longer scaffolds.

To investigate the relationship between the number of generated sequencing reads and assembly success, all reads were mapped onto the obtained contigs using Geneious, allowing for 2% mismatches, a maximum gap size of 3 bp and requiring a minimum overlap of 100 bp. Annotations of each assembly were conducted by first mapping tRNA genes with COVE (Eddy and Durbin 1994), after which the intervening protein and rRNA coding genes were extracted with FeatureExtract 1.2 (Wernersson 2005). To identify these genes, the resulting sequences were mapped to the *Tribolium castaneum* mitogenome (GenBank accession NC_003081) using Geneious, and were afterwards exported, by gene, into separate FASTA files. Sequences of less than 1/3 of total gene length were discarded.

Identification of mitogenomic assemblies using ‘bait’ sequences

To identify the mitogenomic assemblies by association with their respective originating specimen, BLAST searches were conducted for each bait sequence

reference against all corresponding gene sequences extracted from the mitogenome assemblies (separately for *coxI* 5' and 3' regions, *cytB* and *rrnL*). Only hits with 100% pairwise identity and >100 bp overlap were considered a successful identification. Where multiple bait sequences from a single specimen were available, each bait was checked to have hit the same long assembly unequivocally to test for possible chimeras. If baits from a single specimen matched multiple, non-overlapping assemblies they presumably correspond to the same incompletely assembled mitogenome. These assemblies were combined and retained if they included eight or more genes in total. Once mitogenomic assemblies were identified, the tRNA gene order in the cluster of six tRNA genes located between *nad3* and *nad5* was visually recorded for all assemblies in order to test, with our greater taxon sampling, Haran et al.'s (2013) hypothesis that a ARNSEF to RANSEF tRNA gene rearrangement in this region is a synapomorphy for the Entiminae + Cyclominae + Hyperinae clade.

Sequence alignment and dataset concatenation

The sequences for the genes *nad5*, *nad4*, *nad4L* and *nad1*, which are transcribed on the reverse strand of the mitochondrial genome, were reverse complemented prior to alignment. Twenty-eight additional curculionoid mitogenome sequences were obtained from GenBank (primarily those generated by Haran et al. 2013; Supplementary Table S2) in order to maximise taxon sampling. Two members of Chrysomeloidea were included as outgroups, following Haran et al. (2013). The combined sequences from each of the separated 13 protein-coding and two ribosomal RNA genes were individually aligned using the MAFFT 7.0 online server, under the FFT-NS-I slow iterative refinement strategy (Kato et al. 2002). Alignments were thereafter checked manually in Geneious for quality and to ensure that protein-coding

genes were in the correct reading frame. Genes were concatenated together to make 6 different data matrices as follows: all genes (A), only protein-coding genes (B), all genes with 3rd codon positions removed from protein coding genes (C), protein-coding genes only with 3rd codon positions removed (D), all genes with 3rd codon positions removed from protein-coding genes and 1st codon positions R-Y coded (E) and only protein-coding genes with 3rd codon positions removed and 1st codon positions R-Y coded (F).

Phylogenetic analyses

Each of the six datasets were analysed under the maximum likelihood (ML) optimality criterion using RAxML 7.6.6 (Stamatakis 2006) run on the CIPRES web-based server (Miller et al. 2010). To assess nodal support, a rapid bootstrap analysis (BS) with 1000 iterations was run in parallel with tree-building. The datasets were each analysed both partitioned by gene and unpartitioned (i.e. a single partition). Additionally, three of the datasets (A, B and E) were first tested using PartitionFinder (Lanfear et al. 2012) in order to objectively select the best-fitting partitioning scheme and model of molecular evolution for each alignment. This was performed using the Bayesian Information Criterion from an initial partitioning of each of the three codon positions for each amino acid-coding sequence and each ribosomal RNA gene being separate partitions. The resulting ML trees were made ultrametric using the *chronos* function of the *ape* package in R (Paradis et al. 2004), which uses penalised likelihood to fit a chronogram to a phylogenetic tree (Paradis 2013). In order to obtain a measure of the suitability of the mitogenomic data to robustly support relationships across different nodal ages (putative taxonomic ranks) we investigated the distribution of nodal support across trees by calculating the branch length from the root for each node using a custom R

script and plotting this against its respective RAxML BS support. We also constructed a strict consensus tree from the 15 ML trees to visualise the distribution of consistent nodes across all our analyses. We performed additional RAxML analyses on datasets A and B partitioned by gene and separate codon positions for each protein-coding gene (41 and 39 partitions respectively) and various RAxML analyses on these two datasets with different combinations of partitioning schemes and topological constraints, as summarised in Table 3, in order to calculate the Akaike information criterion (AIC) as a means for preferred model selection (Posada & Buckley 2004).

RESULTS

Mitogenomic assembly

Following adapter trimming, approximately 5% of the Illumina reads resembled mitochondrial sequences after BLAST filtering (from a total of 18,341,901 paired-end reads). The Celera and IDBA-UD assemblies resulted in 348 and 336 assemblies of >1000 bp respectively, rising to 361 assemblies when combined using Minimus2. Of these, 105 were >10 kb in length and potentially represented (largely) complete mitogenomes. The cumulative distribution of the assemblies by sequence length is shown in Figure 2, whilst Figure 3 represents the frequency distribution of assembly lengths for each of the Celera, IDBA-UD and Minimus2 assemblies. The latter produced a shift towards longer contigs, especially for the critical contig length of >15kb that corresponds to the full-length of insect mitogenomes. All subsequent analyses were conducted on the Minimus2 assemblies. We were able to newly assemble and identify a total of 92 complete or near complete mitogenomes comprising at least eight genes, including 72 contigs containing the full complement of 15 genes and a further 16 with ≥ 12 genes (Supplementary Table S2). Three near-

complete mitogenomes contained sequences from two non-overlapping assemblies that each matched at least one bait from the same specimen. Those falling short of a full gene complement were mainly lacking the rRNA genes, in particular *rrnS*, which was the least common gene, present in only 56 of the assemblies, whilst *nad6* and *cytB* were present in all 92 assemblies.

Identification of mitogenomic assemblies using ‘bait’ sequences

From the set of 361 partial and complete contigs obtained with Minimus2, a total of 163 *coxI* (529-1560 bp), 154 *cytB* (218-1147 bp) and 162 *rrnL* (211-1340 bp) gene sequences were extracted. Sequences from each gene were grouped into libraries and used as queries in a BLAST search against each corresponding bait sequence reference library. The latter was composed of all successful PCR-based sequences from the 173 original DNA extractions and included 84 *coxI*-5’, 115 *coxI*-3’, 133 *cytB* and 107 *rrnL* sequences (Fig. 4). All samples used in the bulk sequencing were represented by at least one bait (36 samples), while 42, 60 and 36 samples were represented by two, three and four bait sequences, respectively. Matching these bait sequences to the 92 long mitogenomic assemblies, 14 assemblies showed a match to one bait, 32 assemblies matched two baits, 31 assemblies matched three baits and 15 assemblies matched all four baits. Out of the remaining 81 weevil samples, there were 33 instances where baits hit a short contig that was not included in the collection of near-complete or complete mitogenome assemblies, but in 44 instances the baits did not hit any of the assembled contigs. Additionally one divergent assembly was rejected because it was found to match Coleoptera other than weevils in the reference database, possibly present in the sample due to a contamination. Supplementary Table S3 summarises the bait-matching identification results, by bait, for each pooled sample,

with matching contigs given by their unique number. Total number of baits available per sample, the total number of bait hits per sample and the reasons for identification failures are also listed. Of the final set of mitogenomes, 2 belonged to the family Anthribidae, 5 to Attelabidae, 3 to Brachyceridae, 4 to Brentidae, 4 to Dryophthoridae, 1 to Nemonychidae and 101 belonged to 67 identified tribes within the Curculionidae, including 19 tribes of the wood-boring Scolytinae. Overall the different baits contributed fairly equally to the final identifications, with 56% of all *coxI*-3' baits leading to a successful identification, 53% of *cytB*, 53% of *rrnL* and 47% of *coxI*-5'. Proportions of total number of baits, bait hits and hits leading to assembly identifications by gene are illustrated in Figure 4.

The total number of reads making up each of the 92 mitogenomes (which were made up of 97 separate contigs) was used to calculate the sequencing depth (Fig. 5). The majority of sequences showed a 10-50x coverage that generally resulted in contigs of 15 – 20 kb. Coverage reached over 200x in a few cases but this does not appear to closely correlate with contig length. For example, two contigs of high coverage were <5kb in length and corresponded to two non-contiguous fragments from the same species (*Dryocoetes autographus*) linked by multiple baits obtained from a single specimen. In addition, read coverage was not closely correlated with the initial DNA concentration in the sequencing pool. Most samples were present at 10 ng, yet their coverage varied by more than an order of magnitude, while coverage for samples present at a concentration up to 4x lower varied over the same range (Fig. 5).

Phylogenetic analyses

The 92 new assemblies were combined with existing data, for an aligned data matrix of 122 samples and 13792 positions. The optimal partitioning scheme was established

using PartitionFinder, starting with a total of 39 partitions (41 partitions with the two rRNA genes included) that split all 13 genes (15 in datasets A, C and E) and three codon positions in each protein-coding gene. PartitionFinder selected five partitions for the ‘only protein-coding genes’ dataset and six partitions for the ‘all genes’ dataset, whereby the two rRNA genes were grouped with the first codon positions of *nad2*, *nad3* and *nad6* and the second codon position of *atp8* (Table 2). For both datasets the 1st and 3rd codon positions on forward and reverse strands were split into separate partitions, while all 2nd positions were collapsed into a single partition. Forward and reverse genes mainly differed in base frequencies, with a shift from A to T and G to C in the reverse strand partitions, and rates shifted accordingly (normalised to the time-reversible G-T changes: Supplementary Figure S4). The dataset containing ‘only protein-coding genes R-Y coded’ resulted in only 2 partitions, separating 1st and 2nd codon position for both strands combined (3rd positions are removed from this dataset). The findings are in accordance with previous observations on Curculionoidea that also showed a great improvement in likelihood values when partitioning by both codon position and strand (Haran et al. 2013), reflecting the great differences in codon usage in genes coded on either strand. However, this does not extend to produce differences in variation in amino acid changes, as forward and reverse strands were consistently grouped into a single partition for the dataset using 2nd position only and for the R-Y coded matrix (eliminating 1st codon synonymous changes).

The ML trees were greatly improved using six partitions over an unpartitioned analysis, but the benefit of using a model with 41 or 39 separate partitions was low, as seen from the small additional improvement in the AIC values (Table 3). Interestingly, the improvement in ML from using the partitioned models was very similar whether the trees were obtained directly under the partitioned model or obtained under the

unpartitioned model but with the likelihood calculated under partitioning (Table 3). Hence, despite the greatly improved likelihood scores after partitioning, the resulting trees differ only slightly in parameters of greatest impact on the likelihood. This suggests that the topologies are little changed between the unpartitioned model, six-partition model (five-partition model without rRNA genes) and the 41 (39) partition model, given the small increase in likelihood if the simpler model is imposed on the tree obtained with the more complex model.

ML trees obtained with the various coding schemes (including or excluding rRNA genes; R-Y coding; presence of 3rd codon position: Supplementary Table S5) also resulted in highly congruent topologies based upon strongly supported (>80% BS) nodes. Figure 6 depicts the best RAxML tree obtained with the ‘all genes’ dataset under six partitions. Indicated on this tree are nodes that are retained in the strict consensus of trees obtained from all different treatments of the data (Table X), and those nodes unresolved in the strict consensus, i.e. the nodes whose resolution is consistent with the strict consensus. Nodes with high nodal support (80-100% BS) occurred throughout the entire span of nodal ages and this pattern is found across all analyses (Figure 7).

Family-level relationships

All 15 analyses recovered the monophyletic ‘ambrosia beetles’, Platypodinae (100% BS) outside the other ‘true weevils’ (= Curculionidae *sensu* Bouchard et al. 2010), which would otherwise be monophyletic. In most analyses, except those including R-Y coded protein-coding genes, Platypodinae was placed in the sister clade to the rest of Curculionidae, together with the Dryophthoridae (palm weevils) and the brachycerid genus *Ocladius*, with moderate to strong support for this adelphic

relationship (62-95% BS). In all analyses the monophyletic Brentidae (100% BS) were recovered as the sister taxon to a Curculionidae + Dryophthoridae + Brachyceridae clade with very strong nodal support (100% BS). The sister relationship between the monophyletic (100% BS) Attelabidae (leaf-rolling weevils) and this latter clade plus Brentidae was similarly very strongly supported (100% BS) across all analyses. The Nemonychidae was consistently recovered as sister to the clade containing Attelabidae and all other weevil families mentioned so far. Support for this relationship was very high, ranging from 98-100% BS across analyses. The two taxa belonging to the Anthribidae were always recovered as monophyletic (100% BS). Within the Attelabidae, the subfamilies Apoderinae and Rhynchitinae were recovered as monophyletic with BS support of 100% and 83-97% respectively across analyses.

Relationships within Curculionidae *s.str.*

In most analyses the subfamily Bagoinae, represented only by a single *Bagous*, was recovered as the sister to all other Curculionidae (excepting Platypodinae as noted above), with BS support between 66 and 91%. Similarly, most analyses resulted in the recovery of both a monophyletic Entiminae + Cyclominae + Hyperinae clade (marked A in Figure 6; 100% BS) and a strongly supported sister relationship between this clade and a second clade (marked B in Figure 6) containing all other Curculionidae subfamilies (100% BS). Within the entimine clade, the Entiminae itself is not recovered as monophyletic because the tribe Sitonini is consistently recovered (100% BS) either as sister to the clade containing Hyperinae + Cyclominae + the rest of Entiminae, or in a sister clade also containing the Hyperinae (with generally weak nodal support for this relationship). Three entimine tribes are consistently recovered as monophyletic, with strong nodal support; the Otiorhynchini (100% BS), Brachyderini

(100% BS) and the Naupactini (100% BS). The tribe Tropiphorini is apparently paraphyletic because a well-supported clade (95% BS), containing two monophyletic Australian members (*Catasarcus* and *Leptopius*), is itself sister to the Naupactini with strong support (96% BS) and is only distantly related to the other Tropiphorini species in the dataset (*Tropiphorus*), which is sister to the Otorhynchini with strong nodal support (100% BS). All Entiminae (except *Sitona*) are marked by an ARNSEF to RANSEF rearrangement in the tRNA cluster, discovered in earlier studies (Haran et al., 2013; Song et al., 2010) and corroborated here (Fig. 6). One taxon, *Dichotrachelus manueli*, classified in Cyclominae by Alonso-Zarazaga and Lyal (1999), also possesses this same rearrangement, whilst the remaining Cyclominae taxa possess the common gene order, ARNSEF. *Sitona* and *Hypera* were characterised by unique RNSAEF and REANSF gene orders, respectively, observed initially by Haran et al. (2013) and hypothesized to constitute an initial step in the evolution of the derived gene order of the Entiminae. Here, *Hypera* + *Sitona* form a clade that is sister to all others in clade A, while the Cyclominae (minus *Dichotrachelus*), not represented in Haran et al. (2013), and exhibiting the ancestral gene order, occupy the next node as sister to the remaining Entiminae characterised by the derived gene order. This demonstrates that the gene order changes in *Hypera* and *Sitona* are independent of those in Entiminae.

Within the second main curculionid clade, the scolytine taxon *Coptonotus* (Coptonotini) is never recovered together with the bulk of the scolytines, which except for Scolytini (monophyletic with 100% BS), are consistently recovered in a clade with moderate to high support values of 66-100%. The scolytine tribes Corthylini and Ipinini are always recovered as monophyletic (100% BS support) within this. The following higher-level taxa from the second main Curculionidae clade are recovered as

monophyletic across all analyses (BS supports follow taxon name): Ceutorhynchinae (100%), Lixinae (100%), Conoderinae Lobotrachelini (100%) and Curculioninae Cionini (100%). The Cryptorhynchini appears to be paraphyletic owing to the presence of a sample (Cryptorhynchini sp. from Cameroon) falling outside the well supported clade (98% BS) comprising all four other genera analysed.

DISCUSSION

Contig formation from pooled total DNA sequencing

Our results provide a clear demonstration of economic, efficient and reliable sequencing, assembly and identification of large numbers of mitogenomes from a pool of total DNA of numerous samples, without any enrichment or PCR amplification.

Other recent papers attempting to generate full mitochondrial genomes from total DNA either generated a separate library for each taxon (Williams et al. 2014) or pooled a small number of distantly related taxa only (Rubinstein et al. 2013). We have been able to employ the resulting sequence data to reconstruct a higher-level phylogeny of the superfamily Curculionoidea that is highly congruent with recent molecular phylogenies and provides additional evidence for the convergent evolution of specialised wood-boring behaviour and morphology in weevils. The method has been explored previously for the analysis of bulk insect samples from a forest canopy (Crampton-Platt et al., in review), applied to nearly 500 individuals from >200 species. They found that the assembly of mitogenomes from bulk samples is hampered by substantial differences in DNA concentration for species in the pool, due to variation in both body size and number of specimens representing a species. In addition, intra-specific variation was found to cause difficulties with assembly due to polymorphisms, mirroring the well-known problem with genome assembly from heterozygotes (e.g.

Langley et al. 2011). The design of the current study was expected to avoid these problems by normalising the DNA concentration in the pool and by selecting a single individual per species. However, we find that there is no close correlation of sequencing depth and assembly success (Fig. 5), in accordance with Crampton-Platt et al. (in review). Our study excludes the presence of intra-specific variation, but indicates that there is a sequencing depth at which assemblers no longer operate optimally, possibly due to the larger numbers of individual sequencing errors contributed by overlapping reads.

A concern of pooled assemblies is the formation of chimeras by the mis-assembly of different mitogenomes. The potential for this is expected to increase if closely related samples that may not differ in conserved regions of the mitogenomes are included in the pool. The prevalence of chimeras was tested using 77 taxa for which multiple baits were available. In many cases these tests involved both the *cytb* or *rrnL* and the two fragments of the *cox1* gene that map to distant positions in the mitogenome. We did not observe a single case of chimera formation. In addition, the tree topology gave no reason to suggest chimeras, because of the monophyly of the smaller families of Curculionoidea, while chimera formation would also have produced great differences in the length of terminal branches that were not observed.

Phylogenetic analysis from densely sampled mitogenomes

Together with existing mitogenome sequences, a total of 120 terminals were included in the phylogenetic analysis. As mitogenome data sets increase with the numbers of taxa needed for dense sampling, this may produce problems with tree searches and model choice. Specifically, the most complex models, such as the amino acid based CAT model used by Timmermans et al. (2010) that was required for resolving the

deep-level relationships within the Coleoptera are not practical when the number of taxa becomes larger. This raises the question of what is the value of using complex models. Haran et al. (2013) have shown that likelihood trees of weevils can be substantially improved under model partitioning according to (i) codon position and (ii) forward vs. reverse strand, the latter presumably due to the well-established differences in codon usage on either strand. We conducted a formal analysis to test if this partitioning scheme by strand and codon captures the most important aspects of the nucleotide variation using the PartitionFinder software, starting from 41 potential partitions of each codon position within each gene. This could be reduced to the codon positions for all genes on either strands, similar to Haran et al. (2013), but maintaining a single partition for the 2nd codon position on either strand, while adding a separate partition for the rRNA genes not included in that study. The use of these six partitions over the full set of 41 partitions led only to a small reduction in likelihood, while the unpartitioned models were substantially worse (Table 3).

A general difficulty for comparing models is that comparisons are only possible for a single topology, but searches under different partitions favour different topologies. We therefore used the optimal trees obtained under no partitioning and the six and 41-partition schemes to assess likelihoods of the alternative partitioning schemes on those three topologies. The likelihoods on all trees for the three models were almost identical (Table 3), indicating that tree topology is not a major deciding factor for the best model. Taken at face value, the 41 partition wins out over the six partition scheme in all three analyses, but the likelihood gain is minor. As likelihood values become very large with the use of numerous whole mitogenomes, AIC values may not be an appropriate approach to avoid over-parameterisation, unless they are normalised for the total likelihood values (Castoe et al. 2005). We therefore believe

the six-partition scheme is fully adequate. In addition, the practicalities of tree searches on increasingly large datasets from full mitogenomes, as generated with the proposed methodology, also strongly argue for parameter reduction.

Implications for the systematics of weevils

The close relationship linking Platypodinae with Dryophthoridae, as sister to the Curculionidae *s.str.*, has been demonstrated multiple times (Marvaldi et al. 1997, McKenna et al. 2009 and Haran et al. 2013) and indicates that the family Curculionidae, as presently classified, is paraphyletic. The simplified classification system proposed by Oberprieler et al. (2007), recognising a broader Curculionidae also containing the presently defined Brachyceridae and Dryophthoridae as respective subfamilies (*sensu* Alonso-Zarazaga and Lyal 1999) would be consistent with our family-level results. Our results strongly support the relationships amongst the curculionoid families at the base of the tree, which are consistent with most previous molecular analyses, with the exception of the placement of Nemonychidae. This family has previously been suggested to be split off at the most basal node (e.g. McKenna et al. 2009), as opposed to Anthribidae in our results, but our sampling lacks two of the ‘primitive’ weevil families (Belidae and Caridae), prohibiting a definitive conclusion. Our results are also consistent with the previously suggested hypothesis that the Brentidae are the sister family to all the ‘true weevils’, Curculionidae, if we include Brachyceridae and Dryophthoridae in the latter.

A previously described deep split within the true weevils was confirmed by our substantially increased sampling. One strongly supported clade contains the Entiminae + Cyclominae + Hyperinae, and represents the monophyletic and diverse ‘broad-nosed’ weevils, so named because of their relatively short and blunt rostrums.

Rearrangements within the cluster of six tRNA genes are restricted to this clade, even with our increased taxon coverage, further supporting its distinctiveness. The cyclomine genus *Dichotrachelus*, containing the same RANSEF rearrangement as all other Entiminae (except *Sitona*) in our analysis, has been treated as belonging to the Entiminae by some authors (Meregalli and Osella 2007) on morphological grounds. Combined with the low nodal support for its inclusion in a monophyletic Cyclominae (< 50% BS), our tRNA rearrangement data are consistent with this opinion. The second clade containing all other curculionoid subfamilies, with the exception of Bagoinae, which is placed outside of the two main clades, is much less satisfactorily resolved, with only two of its constituent subfamilies (Lixinae and Ceutorhynchinae) being monophyletic. It contains a number of very large subfamilies including the Curculioninae, Molytinae, Baridinae, Cryptorhynchinae and Conoderinae, whose relationships remain obscure due to a lack of strong nodal support. Whilst the recovery of two tribes within this group being monophyletic (Lobotrachelini and Cionini) is encouraging, in order to further investigate the confusing topology of this clade, significantly more representative taxon sampling will be required. Indeed, limitations in taxon sampling are often cited as potentially limiting factors in higher-level phylogenetics (Franz and Engel 2010), and this is certainly an important consideration in such a large group as the Curculionoidea.

An interesting finding is that strong nodal support spans the full depth of the tree and differing taxonomic ranks (families, subfamilies and tribes; Fig. 7). This pattern was seen in analyses of all datasets and under all partitioning models. A potential criticism of mitochondrial sequence data is that due to accelerated evolutionary rates, saturation of sites may obscure or distort phylogenetic signal at deeper nodes (Talavera and Vila 2011). It is clear from our data that at least at the

intra-superfamily level in weevils, this is not necessarily the case, with phylogenetic signal being evenly distributed across the estimated 170 million year diversification history of the weevils (McKenna *et al.*, 2009).

Evolution of wood-boring behaviour

The wood-boring weevil subfamilies are highly adapted to excavate galleries, either subcortically or in woody tissue, and feed on ligneous matter directly or cultivate symbiotic fungi in the tunnels as a food source, and for this reason many are widespread pests of forestry (Oberprieler *et al.* 2007). The taxon density of the current analysis nearly matched the extensive sampling of the wood-boring groups by Jordal *et al.* (2011), a study that is the basis for suggesting their close affinity. However, in contrast to Jordal *et al.* (2011) our results support the conclusions of Haran *et al.* (2013) and McKenna *et al.* (2009), indicating that wood-boring lineages are clearly not monophyletic, with Platypodinae consistently retrieved as closely related to the Dryophthoridae (and Brachyceridae) in a clade sister to all other Curculionidae *sensu* Bouchard *et al.* (2011). Although our analyses recovered neither the Scolytinae nor the Cossoninae as monophyletic, and they were never recovered as sister taxa or nested within the same clade, we cannot confidently conclude as to the relationship between them because only a series of weakly supported nodes separate the cossonine taxa and *Coptonotus* from the rest of the Scolytinae. The latter genus is interesting for consistently not being recovered in our analyses within the generally well-supported Scolytinae clade (excepting Scolytini). Based upon morphological characters, *Coptonotus* has been considered to be a transitional taxon between Platypodinae and other Curculionidae (Jordal *et al.* 2011) or alternatively as an intermediate form between Cossoninae and Scolytinae (Thompson 1992), whilst also containing

morphological characters linking it with Cossoninae. Thompson (1992) has suggested a close relationship between Coptonotini and the scolytine tribe Hylastini based on structures of the aedeagus. However our results argue against this because the Hylastini sample (*Hylastesopacus*) was retrieved with strong support as the sister of Tomicini, and this clade itself was strongly supported as sister to the Hylesini, within the main Scolytinae clade.

Conclusions

We have demonstrated the relative ease of efficiently and economically obtaining a large number of mitogenome DNA sequences from a pooled mixture of DNA extracts, without the need for enrichment or species specific tagging prior to genome pooling. Mitogenome sequences are confidently identified to specimen with a limited amount of prior mtDNA sequence data for each sample, and exhibit no error with regard to these bait sequences. Our mtDNA genome data yields phylogenetic relationships that are highly congruent with prior expectations, and provides phylogenetic signal with robustly supported nodes across a broad range of lineage divergence times and taxon diversity, from family-level to generic-level, which are consistent across different data partitioning schemes.

It is evident that the efficiency of our approach will be a function of the relative concentration of mitochondrial to nuclear DNA within a focal group. The average coleopteran genome size is estimated to be approximately 0.65 Gb +/- 0.05 (<http://www.genomesize.com>). Under the assumption that the copy number of mtDNA genomes does not differ substantially across organisms, our approach should be of broad utility within insect phylogenetics where mean nuclear genome size is estimated to be 1.22 Gb +/- 0.05. However, it may be less efficient for taxa with larger average

nuclear genome sizes (e.g. crustaceans: mean nuclear genome size = approximately 4.45 Gb +/- 0.45). A further consideration for the implementation of our approach is taxon sampling and the mitogenomic assembly pipeline. Our sampling for the higher-level taxonomic relationships within the Curculionoidea provides little challenge for the pipeline, as mtDNA genomes sampled from different genera exhibit high DNA sequence divergence. Genome divergence facilitates genome reassembly from a mixed pool of genome fragments, and the pipeline efficiency will eventually be compromised as mtDNA genome relatedness increases. Our data suggests this limit lies somewhere below an uncorrected divergence of 10% for *cox1* and *cytB* that characterises the two species of *Cionus* (*C. olens* and *C. griseus*) included in our sampling. To ascertain genome relatedness thresholds for the reassembly pipeline, simulation analyses can be employed. However, it is important to point out that as NGS technology and read lengths improve, relatedness thresholds will also become more favourable.

ACKNOWLEDGEMENTS

We are indebted to the following individuals who have lent or donated weevil specimens used in the present study: Max Barclay, Roberto Caldara, Christiana Faria, Michael Gillett, LeventGultekin, James Kitson, Christopher Lyal, Massimo Meregalli, Rolf Oberprieler, Charles O'Brien, Pedro Oromí, Li Ren and Clive Turner. This work was supported by a Natural Environment Research Council CASE PhD studentship to CPDTG, the University of East Anglia and the Natural History Museum, and the NHM Biodiversity Initiative. MJTNT is supported by NERC Postdoctoral Fellowship (NE/I021578/1).

REFERENCES

- Alonso-Zarazaga MA, Lyal CHC. 1999. A world catalogue of families and genera of Curculionoidea (Insecta: Coleoptera) (excepting Scolytidae and Platypodidae). Barcelona: Entomopraxis.
- Botero-Castro F, Tilak M-K, Justy F, Catzefflis F, Delsuc F, Douzery EJP. 2013. Next-generation sequencing and phylogenetic signal of complete mitochondrial genomes for resolving the evolutionary history of leaf-nosed bats (Phyllostomidae). *Molecular Phylogenetics and Evolution* 69:728-739.
- Bouchard P, Bousquet Y, Davies AE, Alonso-Zarazaga MA, Lawrence JF, Lyal CHC, Newton AF, Reid CAM, Schmitt M, Slipinski SA, Smith ABT. 2011. Family-group names in Coleoptera (Insecta). *Zookeys*: 1-895.
- Cameron SL, Lo N, Bourguignon T, Svenson GJ, Evans TA. 2012. A mitochondrial genome phylogeny of termites (Blattodea: Termitoidea): Robust support for interfamilial relationships and molecular synapomorphies define major clades. *Molecular Phylogenetics and Evolution* 65:163-173.
- Castoe TA, Sasa MM, Parkinson C. 2005. Modeling nucleotide evolution at the mesoscale: the phylogeny of the neotropical pitvipers of the *Porthidium* group (viperidae: crotalinae). *Molecular Phylogenetics and Evolution* 37:881-898.
- Crampton-Platt, AL, Timmermans MJTN, Gimmel ML, Narayanan Kutty S, Cockerill TD, Khen CV, Vogler AP. Pooled mitochondrial genome assembly for biodiversity discovery in a phylogenetic framework. *Systematic Biology*, in review.
- Crowson, RA. 1955. The natural classification of the families of Coleoptera. London: Nathaniel Lloyd & Co.
- Curole JP, Kocher TD. 1999. Mitogenomics: digging deeper with complete mitochondrial genomes. *Trends in Ecology & Evolution* 14:394-398.

- Finstermeier K, Zinner D, Brameier M, Meyer M, Kreuz E, Hofreiter M, Roos C. 2013. A mitogenomic phylogeny of living primates. *PLoS ONE* 8:e69504.
- Franz NM, Engel MS. 2010. Can higher-level phylogenies of weevils explain their evolutionary success? A critical review. *Systematic Entomology* 35:597-606.
- Haran J, Timmermans MJTN, Vogler AP. 2013. Mitogenome sequences stabilize the phylogenetics of weevils (Curculionoidea) and establish the monophyly of larval ectophagy. *Molecular Phylogenetics and Evolution* 67:156-166.
- Hundsdoerfer AK, Rheinheimer J, Wink M. 2009. Towards the phylogeny of the Curculionoidea (Coleoptera): Reconstructions from mitochondrial and nuclear ribosomal DNA sequences. *Zoologischer Anzeiger* 248:9-31.
- Hunt T, Bergsten J, Levkanicova Z, Papadopoulou A, John OS, Wild R, Hammond PM, Ahrens D, Balke M, Caterino MS, Gomez-Zurita J, Ribera I, Barraclough TG, Bocakova M, Bocak L, Vogler AP. 2007. A comprehensive phylogeny of beetles reveals the evolutionary origins of a superradiation. *Science* 318:1913-1916.
- Jordal BH, Sequeira AS, Cognato AI. 2011. The age and phylogeny of wood boring weevils and the origin of subsociality. *Molecular Phylogenetics and Evolution* 59:708-724.
- Katoh K, Misawa K, Kuma K, Miyata T. 2002. MAFFT: a novel method for rapid multiple sequence alignment based on fast Fourier transform. *Nucleic Acids Research* 30:3059-3066.
- Kayal E, Roure B, Philippe H, Collins AG, Lavrov DV. 2013. Cnidarian phylogenetic relationships as revealed by mitogenomics. *BMC Evolutionary Biology* 13.
- Kuschel G. 1995. A phylogenetic classification of Curculionoidea to families and subfamilies. *Memoirs of the Entomological Society of Washington* 14:5-33.

- Lanfear R, Calcott B, Ho SYW, Guindon S. 2012. PartitionFinder: Combined selection of partitioning schemes and substitution models for phylogenetic analyses. *Molecular Biology and Evolution* 29:1695-1701.
- Langley CH, Crepeau M, Cerdeno C, Corbett-Detig R, Stevens K. 2011. Circumventing heterozygosity: sequencing the amplified genome of a single haploid *Drosophila melanogaster* embryo. *Genetics* 188:239-246.
- Marvaldi AE. 1997. Higher level phylogeny of curculionidae (Coleoptera : Curculionoidea) based mainly on larval characters, with special reference to broad-nosed weevils. *Cladistics* 13:285-312.
- McKenna DD, Sequeira AS, Marvaldi AE, Farrell BD. 2009. Temporal lags and overlap in the diversification of weevils and flowering plants. *Proceedings of the National Academy of Sciences of the United States of America* 106:7083-7088.
- Meregalli M, Osella G. 2007. *Dichotrachelus kahleni* sp. n., a new weevil species from the Carnian Alps, north-eastern Italy (Coleoptera, Curculionidae, Entiminae). *Deutsche Entomologische Zeitschrift* 54:169-177.
- Miller MA, Pfeiffer W, Schwartz T. 2010. Creating the CIPRES Science Gateway for inference of large phylogenetic trees. In: *Proceedings of the Gateway Computing Environments Workshop (GCE)*, 14 Nov. 2010, New Orleans, LA. p.1-8.
- Oberprieler RG, Marvaldi AE, Anderson RS. 2007. Weevils, weevils, weevils everywhere. *Zootaxa* 1668:491-520.
- Osigus H-J, Eitel M, Bernt M, Donath A, Schierwater B. 2013. Mitogenomics at the base of Metazoa. *Molecular Phylogenetics and Evolution* 69: 339-351.
- Pacheco MA, Battistuzzi FU, Lentino M, Aguilar RF, Kumar S, Escalante AA. 2011. Evolution of modern birds revealed by mitogenomics: timing the radiation and origin of major orders. *Molecular Biology and Evolution* 28:1927-1942.

- Paradis E. 2013. Molecular dating of phylogenies by likelihood methods: A comparison of models and a new information criterion. *Molecular Phylogenetics and Evolution* 67:436-444.
- Paradis E, Claude J, Strimmer K. 2004. APE: Analyses of phylogenetics and evolution in R language. *Bioinformatics* 20:289-290.
- Posada D, Buckley T. 2004. Model selection and model averaging in phylogenetics: advantages of Akaike information criterion and Bayesian approaches over likelihood ratio tests. *Systematic Biology* 53:793-808.
- Roos J, Aggarwal RK, Janke A. 2007. Extended mitogenomic phylogenetic analyses yield new insight into crocodylian evolution and their survival of the Cretaceous-Tertiary boundary. *Molecular Phylogenetics and Evolution* 45:663-673.
- Rubinstein ND, Feldstein T, Shenkar N, Botero-Castro F, Griggio F, Mastrototaro F, Delsuc F, Douzery EJP, Gissi C, Huchon D. 2013. Deep Sequencing of Mixed Total DNA without Barcodes Allows Efficient Assembly of Highly Plastic Ascidian Mitochondrial Genomes. *Genome Biology and Evolution* 5:1185-1199.
- Sommer DD, Delcher AL, Salzberg SL, Pop M. 2007. Minimus: a fast, lightweight genome assembler. *BMC Bioinformatics* 8:64.
- Song HJ, Sheffield NC, Cameron SL, Miller KB, Whiting MF. 2010. When phylogenetic assumptions are violated: base compositional heterogeneity and among-site rate variation in beetle mitochondrial phylogenomics. *Systematic Entomology* 35:429-448.
- Stamatakis A. 2006. RAxML-VI-HPC: Maximum likelihood-based phylogenetic analyses with thousands of taxa and mixed models. *Bioinformatics* 22:2688-2690.

- Talavera G, Vila R. 2011. What is the phylogenetic signal limit from mitogenomes? The reconciliation between mitochondrial and nuclear data in the Insecta class phylogeny. *BMC Evolutionary Biology* 11:315.
- Thompson RT. 1992. Observations on the morphology and classification of weevils (Coleoptera, Curculionoidea) with a key to major groups. *Journal of Natural History* 26:835-891.
- Timmermans MJTN, Dodsworth S, Culverwell CL, Bocak L, Ahrens D, Littlewood DTJ, Pons J, Vogler AP. 2010. Why barcode? High-throughput multiplex sequencing of mitochondrial genomes for molecular systematics. *Nucleic Acids Research* 38:1-14.
- Wei S-j, Shi M, Sharkey MJ, van Achterberg C, Chen X-X. 2010. Comparative mitogenomics of Braconidae (Insecta: Hymenoptera) and the phylogenetic utility of mitochondrial genomes with special reference to Holometabolous insects. *BMC Genomics* 11:371.
- Williams S, Foster PG, Littlewood, DTJ. 2014. The complete mitochondrial genome of a turbinid vetigastropod from MiSeq Illumina sequencing of genomic DNA and steps towards a resolved gastropod phylogeny. *Gene* 533:38-47.
- Winkelmann I, Campos PF, Strugnell J, Cherel Y, Smith PJ, Kubodera T, Allcock L, Kampmann M-L, Schroeder H, Guerra A, Norman M, Finn J, Ingrao D, Clarke M, Gilbert MTP. 2013. Mitochondrial genome diversity and population structure of the giant squid *Architeuthis*: genetics sheds new light on one of the most enigmatic marine species. *Proceedings of the Royal Society B-Biological Sciences* 280:1759

FIGURE LEGENDS

Figure 1. Schematic flowchart of the principal steps for the bulk *de novo* assembly of mitogenomes and identification with PCR-amplified ‘bait’ fragments.

Figure 2. Cumulative distribution of assembly lengths from the Celera, IDBA-UD and the combined Minimus2-generated assemblies.

Figure 3. Frequency distribution of assembly lengths from the Celera, IDBA-UD and the combined Minimus2-generated assemblies.

Figure 4. Relative proportions, by gene, of total ‘bait’ sequences available, ‘bait’ sequences with matching ‘hits’ to the assembled genes and matching hits that contributed to a successful mitogenome identification following a BLAST search.

Figure 5. Mean sequencing coverage versus A) assembly (contig) length (bp) and B) approximate mass of genomic DNA in the sample pool, for identified mitogenomic assemblies.

Figure 6. Maximum likelihood tree resulting from the analysis of the ‘all genes’ dataset partitioned according to PartitionFinder (see Table 2). Within Curculionidae *s.str.* (*sensu* Bouchard *et al.* 2010) branches are coloured according to subfamily. Other curculionoid families have their name labels coloured by family. Numbers adjacent to nodes are RAxML rapid bootstrap scores, with values >80% highlighted in red. The three principal wood-boring subfamilies are represented by dashed branches and the

nodes labelled A and B indicate the two large divisions within Curculionidae referred to in the text. Nodes indicated in green correspond to nodes present in the strict consensus tree and nodes indicated in blue are consistent with it. The positions of the three tRNA rearrangements are indicated. Scale bar represents substitution rate.

Family and subfamily codes precede taxa names as follows: Anthribidae (ANTH), Attelabidae (ATTE), Brachyceridae (BRAC), Brentidae (BREN), Dryophthoridae (DRYO), Nemonychidae (NEMO), Bagoinae (BAGO), Baridinae (BARI), Ceutorhynchinae (CEUT), Conoderinae (CONO), Cossoninae (COSS), Cryptorhynchinae (CRYP), Curculioninae (CURC), Lixinae (LIXI), Mesoptillinae (MESO), Molytinae (MOLY), Platypodinae (PLAT) and Scolytinae (SCOL).

Figure 7. Graph of RAxML nodal bootstrap support against branch length of nodes from the root for the analysis of all 15 concatenated genes under the six partition scheme (dataset A).

TABLES

Table 1. List of software used for the *de novo* assembly of mitogenomes, with their main function and source URL.

| Program | Function | URL |
|--------------------|--|---|
| FastQC | NGS quality assesment | http://www.bioinformatics.babraham.ac.uk/projects/fastqc/ |
| Trimmomatic | Adapter trimming | http://www.usadellab.org/cms/index.php?page=trimmomatic |
| Celera | Genome assembly | http://sourceforge.net/apps/mediawiki/wgs-assembler/index.php?title=Main_Page |
| IDBA-UD | Genome assembly | http://i.cs.hku.hk/~alse/hkubrg/projects/idba_ud/ |
| Minimus2 | Merging sequence sets | http://sourceforge.net/apps/mediawiki/amos/index.php?title=Minimus2 |
| Prinseq | Sequence quality control | http://edwards.sdsu.edu/cgi-bin/prinseq/prinseq.cgi |
| COVE | tRNA annotation | http://selab.janelia.org/software.html |
| FeatureExtract | Gene extraction | http://www.cbs.dtu.dk/services/FeatureExtract/ |
| Geneious | Gene annotation / sequence editing | http://www.geneious.com/ |
| MAFFT | Sequence alignment | http://mafft.cbrc.jp/alignment/software/ |
| BLAST | Local alignment search | http://blast.ncbi.nlm.nih.gov/Blast.cgi |
| PartitionFinder | Partitioning scheme selection | http://www.robertlanfear.com/partitionfinder/ |
| CIPRES | Phylogenetic analysis server | http://www.phylo.org/ |
| RAxML | Maximum Likelihood phylogenetic analysis | http://sco.h-its.org/exelixis/software.html |
| 'ape' package in R | Phylogenetic analysis | http://ape-package.ird.fr/ |

Table 3. Maximum likelihood of trees under different partitioning schemes. Trees were obtained under no partitioning, under the 6-partition scheme selected by PartitionFinder, and by the maximum number of partitions tested (partitioning by gene and codon position). Each of the resulting trees were then assessed for their likelihood under the alternative models. Note the comparatively small difference in likelihood (Δ AIC) under each partitioning scheme regardless of the model used in the tree search.

| Data set | Partitioning Scheme | Topological constraint | No. partitions | Substitution model | No. Parameters | LnL | AIC | Δ AIC |
|---------------------------------|-------------------------------------|------------------------|----------------|--------------------|----------------|---------|---------|--------------|
| All genes (A) | Unpartitioned (1 partition) | None | 1 | GTR | 8 | -787773 | 1575562 | 62885 |
| | PartitionFinder (6 partitions) | on 1 partition tree | 6 | GTR | 48 | -758061 | 1516219 | |
| | Gene/codon-position (41 partitions) | on 1 partition tree | 41 | GTR | 328 | -756379 | 1513414 | 737 |
| | Gene/codon-position (41 partitions) | on 6 partition tree | 41 | GTR | 328 | -756272 | 1513199 | 522 |
| | PartitionFinder (6 partitions) | on 41 partition tree | 6 | GTR | 48 | -758010 | 1516097 | 3417 |
| | Gene/codon-position (41 partitions) | None | 41 | GTR | 328 | -756010 | 1512677 | n/a |
| Protein-coding genes (B) | Unpartitioned (1 partition) | None | 1 | GTR | 8 | -684161 | 1368339 | 34473 |
| | PartitionFinder (5 partitions) | on 1 partition tree | 5 | GTR | 40 | -668567 | 1337213 | |
| | PartitionFinder (5 partitions) | None | 5 | GTR | 40 | -668480 | 1337039 | 3173 |
| | Gene/codon-position | on 5 partition tree | 39 | GTR | 312 | -666678 | 1333981 | 115 |
| | PartitionFinder (5 partitions) | on 39 partition tree | 5 | GTR | 40 | -668523 | 1337043 | 3177 |
| | Gene/codon-position (39 partitions) | None | 39 | GTR | 312 | -666621 | 1333866 | n/a |

FIGURES

Figure 1.

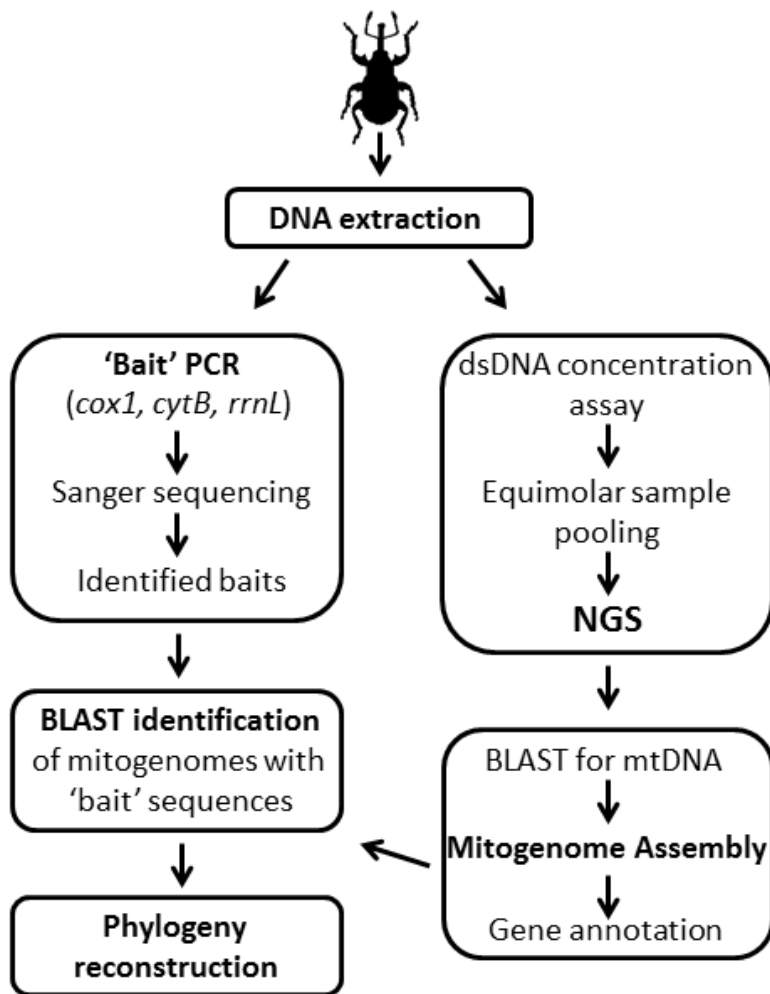


Figure 2.

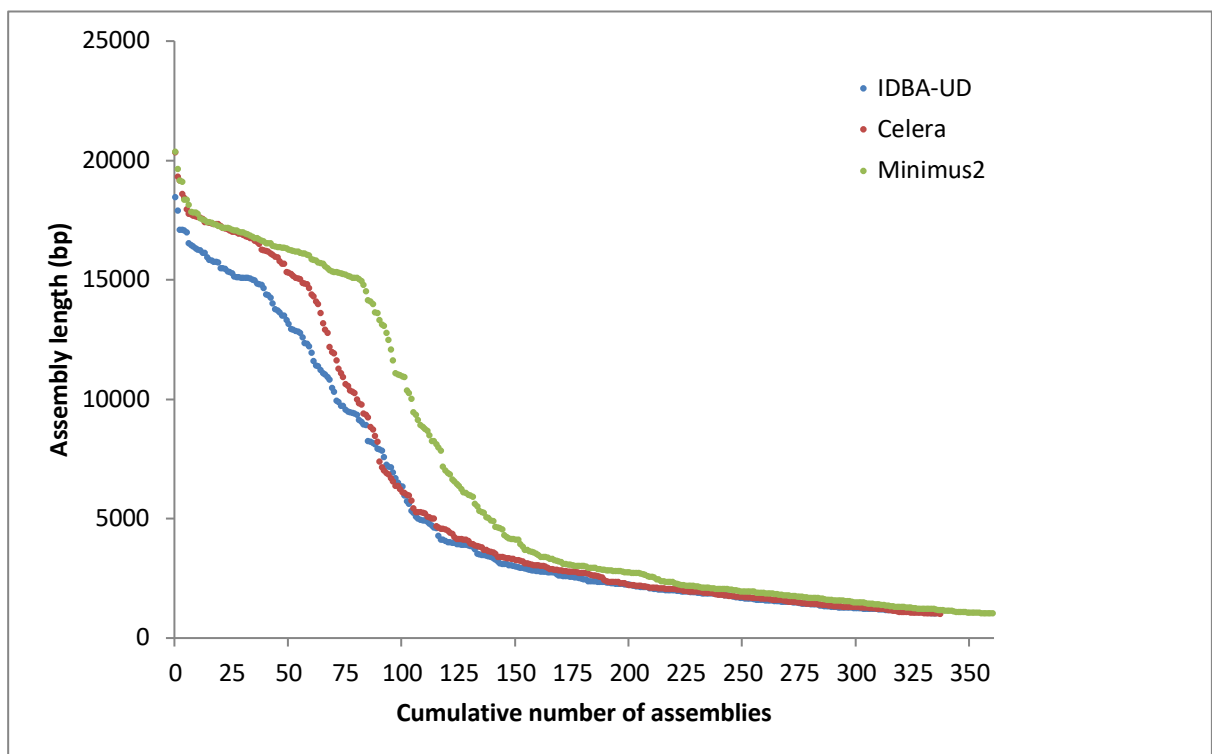


Figure 3.

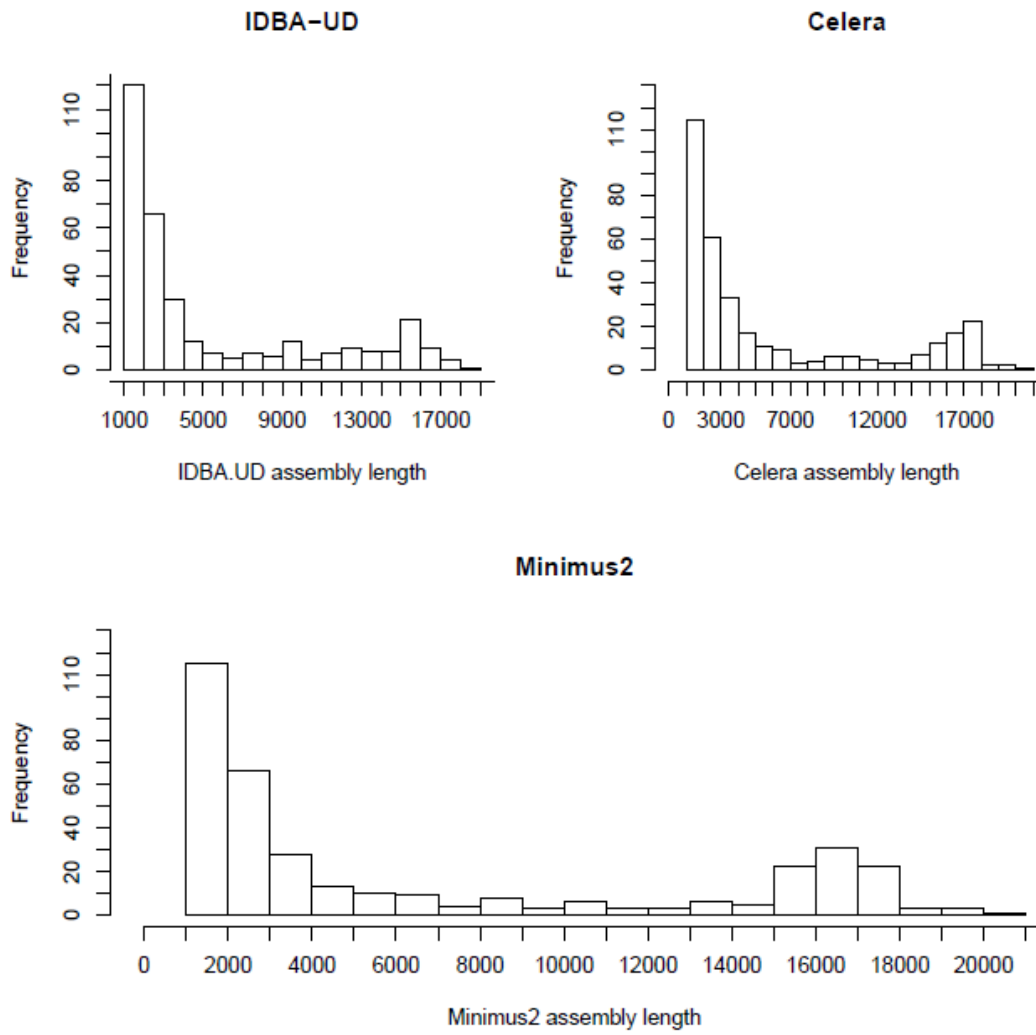


Figure 4.

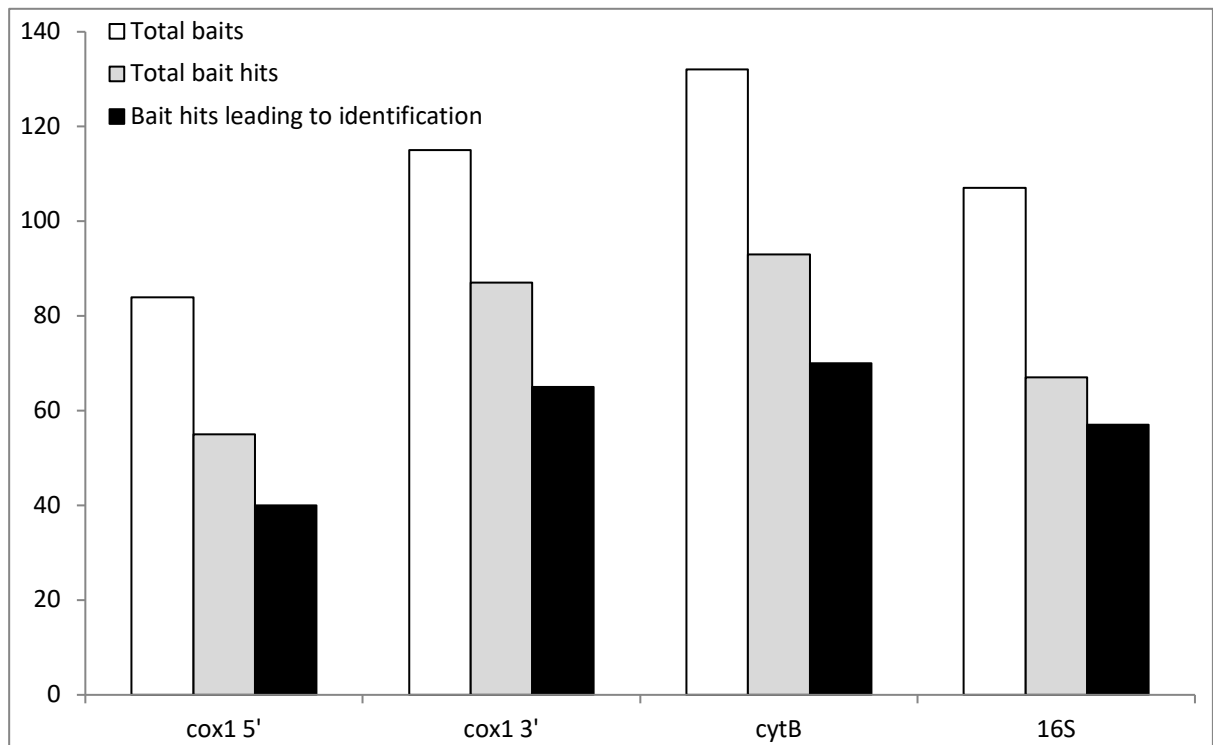


Figure 6.

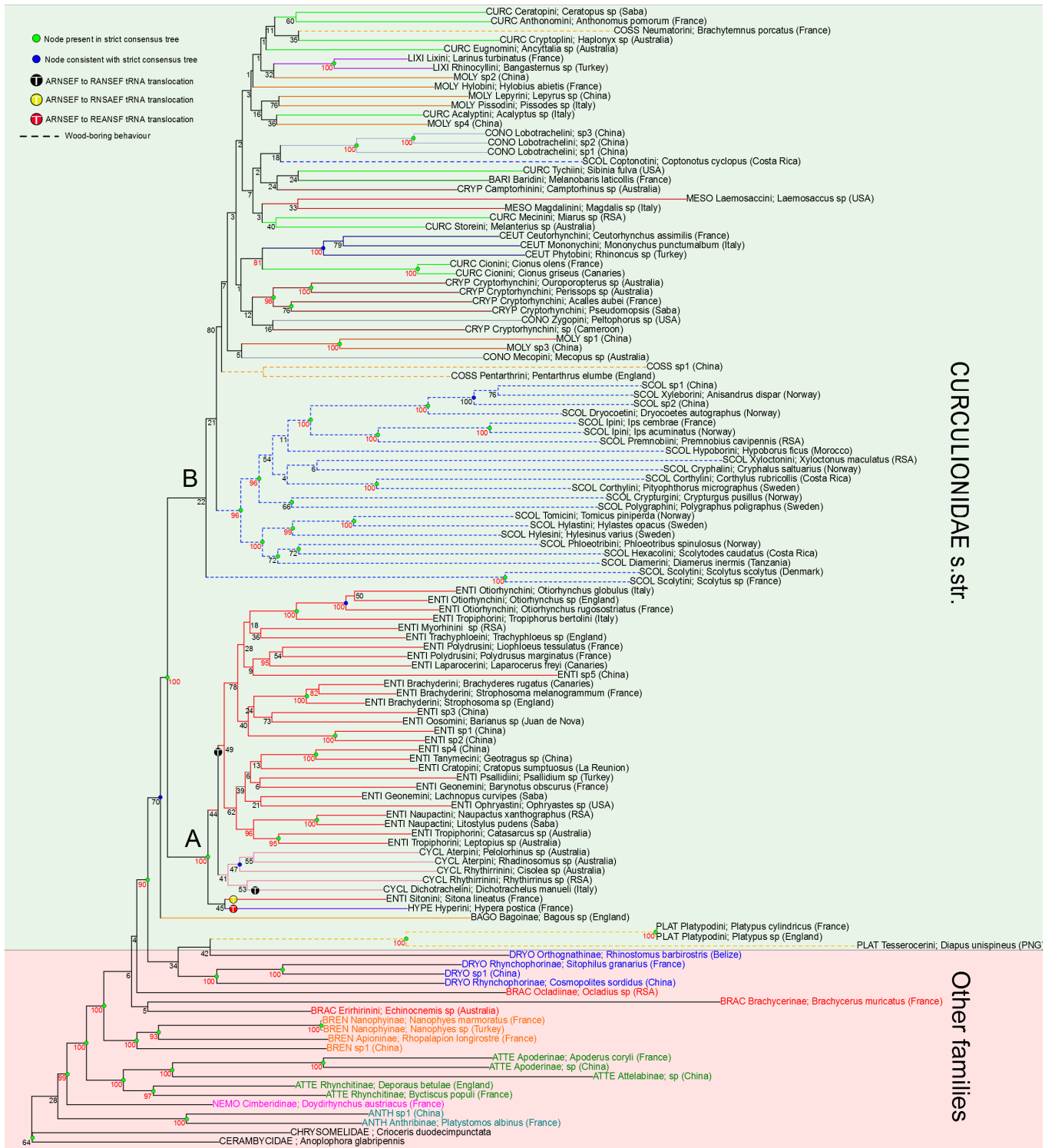


Figure 7.

