

MitoFinder: efficient automated large-scale extraction of mitogenomic data in target enrichment phylogenomics

Rémi Allio¹, Alex Schomaker-Bastos^{2,†}, Jonathan Romiguier¹, Francisco Prosdocimi², Benoit Nabholz¹, and Frédéric Delsuc¹

¹*Institut des Sciences de l'Évolution de Montpellier (ISEM), CNRS, EPHE, IRD, Université de Montpellier, Montpellier, France.*

²*Laboratório Multidisciplinar para Análise de Dados (LAMPADA), Instituto de Bioquímica Médica Leopoldo de Meis, Universidade Federal do Rio de Janeiro, Rio de Janeiro, Brasil.*

[†] *In Memoriam (08/01/2015)*

Correspondence

Rémi Allio

Email: remi.allio@umontpellier.fr

Frédéric Delsuc

Email: frederic.delsuc@umontpellier.fr

Running head

Mitochondrial signal from UCE capture data

Abstract

Thanks to the development of high-throughput sequencing technologies, target enrichment sequencing of nuclear ultraconserved DNA elements (UCEs) now allows routinely inferring phylogenetic relationships from thousands of genomic markers. Recently, it has been shown that mitochondrial DNA (mtDNA) is frequently sequenced alongside the targeted loci in such capture experiments. Despite its broad evolutionary interest, mtDNA is rarely assembled and used in conjunction with nuclear markers in capture-based studies. Here, we developed MitoFinder, a user-friendly bioinformatic pipeline, to efficiently assemble and annotate mitogenomic data from hundreds of UCE libraries. As a case study, we used ants (Formicidae) for which 501 UCE libraries have been sequenced whereas only 29 mitogenomes are available. We compared the efficiency of four different assemblers (IDBA-UD, MEGAHIT, MetaSPAdes, and Trinity) for assembling both UCE and mtDNA loci. Using MitoFinder, we show that metagenomic assemblers, in particular MetaSPAdes, are well suited to assemble both UCEs and mtDNA. Mitogenomic signal was successfully extracted from all 501 UCE libraries allowing confirming species identification using COI barcoding. Moreover, our automated procedure retrieved 296 cases in which the mitochondrial genome was assembled in a single contig, thus increasing the number of available ant mitogenomes by an order of magnitude. By leveraging the power of metagenomic assemblers, MitoFinder provides an efficient tool to extract complementary mitogenomic data from UCE libraries, allowing testing for potential mito-nuclear discordance. Our approach is potentially applicable to other sequence capture methods, transcriptomic data, and whole genome shotgun sequencing in diverse taxa.

Keywords

Bioinformatics/Phyloinformatics, DNA Barcoding, Invertebrates, Metagenomics,

Systematics, Insects

Introduction

Next generation phylogenomics in which phylogenetic relationships are inferred from thousands of genomic markers gathered through high-throughput sequencing (HTS) is on the rise. More specifically, targeted enrichment or DNA sequence capture methods are becoming the gold standard in phylogenetic analyses because they allow subsampling the genome efficiently at reduced cost (Lemmon & Lemmon, 2013; McCormack, Hird, Zellmer, Carstens, & Brumfield 2013a). The field has witnessed the rapid parallel development of exon capture from transcriptome-derived baits (Bi *et al.* 2012), anchored hybrid enrichment techniques (Lemmon, Emme, & Lemmon 2012), and the capture of ultraconserved DNA elements (UCEs; McCormack, Hird, Zellmer, Carstens, & Brumfield 2013b). All hybridization capture methods target a particular portion of the genome corresponding to the defined probes plus flanking regions. Prior knowledge is required to generate sequence capture probes, but ethanol preserved tissues, old DNA extractions, and museum specimens can be successfully sequenced (Faircloth *et al.* 2012; Guschanski *et al.* 2013; Blaimer *et al.* 2015). The first UCEs were identified by Bejerano *et al.* (2004) in the human genome and have been shown to be conserved in mammals, birds, and even ray-finned fish. Thanks to their large-scale sequence conservation, UCEs are particularly well-suited for sequence capture experiments and have become popular for phylogenomic reconstruction of diverse animals groups (Guschanski *et al.* 2013; Blaimer *et al.* 2015; Esselstyn, Oliveros, Swanson, & Faircloth 2017). Initially restricted to a few vertebrate groups such as mammals (McCormack *et al.* 2012) and birds (McCormack *et al.* 2013a), new UCE probe sets have been designed to target thousands of loci in arthropods such as hymenopterans (Blaimer *et al.* 2015; Branstetter *et al.* 2017a; Faircloth, Branstetter, White, & Brady 2015), coleopterans (Baca, Alexander, Gustafson, & Short 2017), and arachnids (Starrett *et al.* 2017).

It has been shown that complete mitochondrial genomes could be retrieved as by-products of sequence capture/enrichment experiments such as whole exome capture in human (Picardi & Pesole, 2012). Indeed, mitogenomes can in most cases be assembled from off-target sequences of UCE capture libraries in amniotes (do Amaral *et al.* 2015). Despite its well-acknowledged limitations (Galtier, Nabholz, Glémin, & Hurst 2009), mitochondrial DNA (mtDNA) remains a marker of choice for phylogenetic inference (e.g. Hassanin *et al.* 2012), for species identification or delimitation through barcoding (e.g. Coissac *et al.* 2016), and to reveal potential cases of mito-nuclear discordance resulting from introgression and/or hybridization events (e.g. Zarza *et al.* 2016, 2018; Grummer, Morando, Avila, Sites Jr, & Leaché 2018). MtDNA could also be used to taxonomically validate the specimens sequenced for UCEs using COI barcoding (Ratnasingham & Hebert, 2007) and to control for potential cross-contaminations in HTS experiments (Ballenghien, Faivre, & Galtier 2017). In practice, the few studies that have extracted mtDNA signal from UCEs (e.g. Meiklejohn *et al.* 2014; Pie *et al.* 2017; Wang, Hosner, Liang, Braun, & Kimball 2017, Zarza *et al.* 2018) and anchored phylogenomics (Caparroz *et al.* 2018) have done so manually for only few taxa. Most studies assembling mitogenomes from UCE libraries have used contigs produced by the Trinity RNAseq assembler (Grabherr *et al.* 2011) as part of the PHYLUCE pipeline (Faircloth, 2016), which was specifically designed to extract UCE loci. Indeed, RNAseq assemblers such as Trinity allow dealing with the uneven coverage of target reads in sequence-capture libraries, but also multi-copy genes such as the ribosomal RNA cluster, and organelles (chloroplasts and mitochondria). However, this strategy is likely not scaling well with hundreds of taxa because of the high computational demand required by Trinity. Metagenomic assemblers could provide a powerful alternative because they have been designed for an efficient *de novo* assembly of complex read populations by explicitly dealing with uneven read coverage and are computationally and memory efficient. Comparisons

based on empirical bulk datasets of known composition (Vollmers, Wiegand, & Kaster 2017) have identified IDBA-UD (Peng, Leung, Yiu, & Chin 2012), MEGAHIT (Li *et al.* 2016), and MetaSPAdes (Nurk, Meleshko, Korobeynikov, & Pevzner 2017) as the most efficient current metagenomic assemblers.

As a case study, we focused on ants (Hymenoptera: Formicidae) for which a only 29 mitogenomes were available on GenBank compared to 501 UCE captured libraries as of March 29th, 2018 (**Appendix S1**). This contrasts sharply with the other most speciose group of social insects, termites (Isoptera), for which almost 500 reference mitogenomes have been produced (Bourguignon *et al.* 2017) and no UCE study has been conducted so far. Sequencing and assembling difficulties stemming from both the AT-rich composition (Foster, Jermin, & Hickey 1997) and a high rate of mitochondrial genome rearrangements in hymenopteran (Dowton, Castro, & Austin 2002) might explain the limited number of mitogenomes currently available for ants. It is only recently that a few ant mitogenomes have been assembled out from UCE data (Ströher *et al.* 2017; Meza-Lázaro, Poteaux, Bayona-Vásquez, Branstetter, & Zaldívar-Riverón 2018; Vieira & Prosdocimi, 2019). Here, we built a pipeline called MitoFinder designed to automatically assemble and extract mitogenomic data from raw UCE capture libraries. Using publicly available UCE libraries for 501 ants, we show that complementary mitochondrial phylogenetic signal can be efficiently extracted using metagenome assemblers along with targeted UCE loci.

Materials and methods

Data acquisition

We used UCE raw sequencing data for 501 ants produced in 10 phylogenomic studies (Blaimer *et al.* 2015; Faircloth *et al.* 2015; Blaimer *et al.* 2016; Branstetter *et al.* 2017a,b,c; Jesovnik *et al.* 2017; Pierce *et al.* 2017; Prebus *et al.* 2017; Ward & Branstetter 2017). This

dataset includes representatives of 15 of 16 subfamilies (Ward 2014) and 30 tribes. Raw sequence reads were downloaded from the NCBI Short Read Archive (SRA) on March 29th, 2018 (**Appendix S1**). For the 501 ant UCE libraries, raw reads were cleaned with Trimmomatic v0.36 (Bolger, Lohse, & Usadel 2014) using the following parameters: LEADING:3 TRAILING:3 SLIDINGWINDOW:4:15 MINLEN:50. A reference database with the 29 complete mitochondrial genomes available for ants on GenBank at the time was constructed.

Mitogenomic data extraction with MitoFinder

To extract mitogenomic data from UCE libraries, we developed a dedicated bioinformatic pipeline called MitoFinder (**Fig. 1**). This pipeline was designed to assemble sequencing reads from target enrichment libraries, assemble, extract, and annotate mitochondrial contigs. To evaluate the impact of assembler choice, contigs were assembled with IDBA-UD v1.1.1, MEGAHIT v1.1.3, and MetaSPAdes v3.13.0 within MitoFinder, and with Trinity v2.1.1 within PHYLUCE using default parameters. Mitochondrial contigs were then identified by similarity search using blastn with e-value $\geq 1e-06$ against our ant reference mitogenomic database. Each detected mitochondrial contig was then annotated with tblastx for protein-coding genes (CDS) and blastn for both 16S and 12S taking advantage of the geneChecker module of mitoMaker (Schomaker-Bastos & Prosdocimi, 2018) that we incorporated in MitoFinder. Finally, we used ARWEN v1.2 (Laslett & Canbäck, 2007) to detect and annotate tRNA genes.

Considering possible rearrangements in ant mitogenomes, each mitochondrial CDS was first aligned with MAFFT v7.271 (Katoh & Standley, 2013) algorithm FFT-NS-2 with option *--adjustdirection*. Then, to take into account potential frameshifts and stop codons, mitochondrial CDSs were re-aligned with MACSE v2.03 (Ranwez *et al.* 2018) with option -

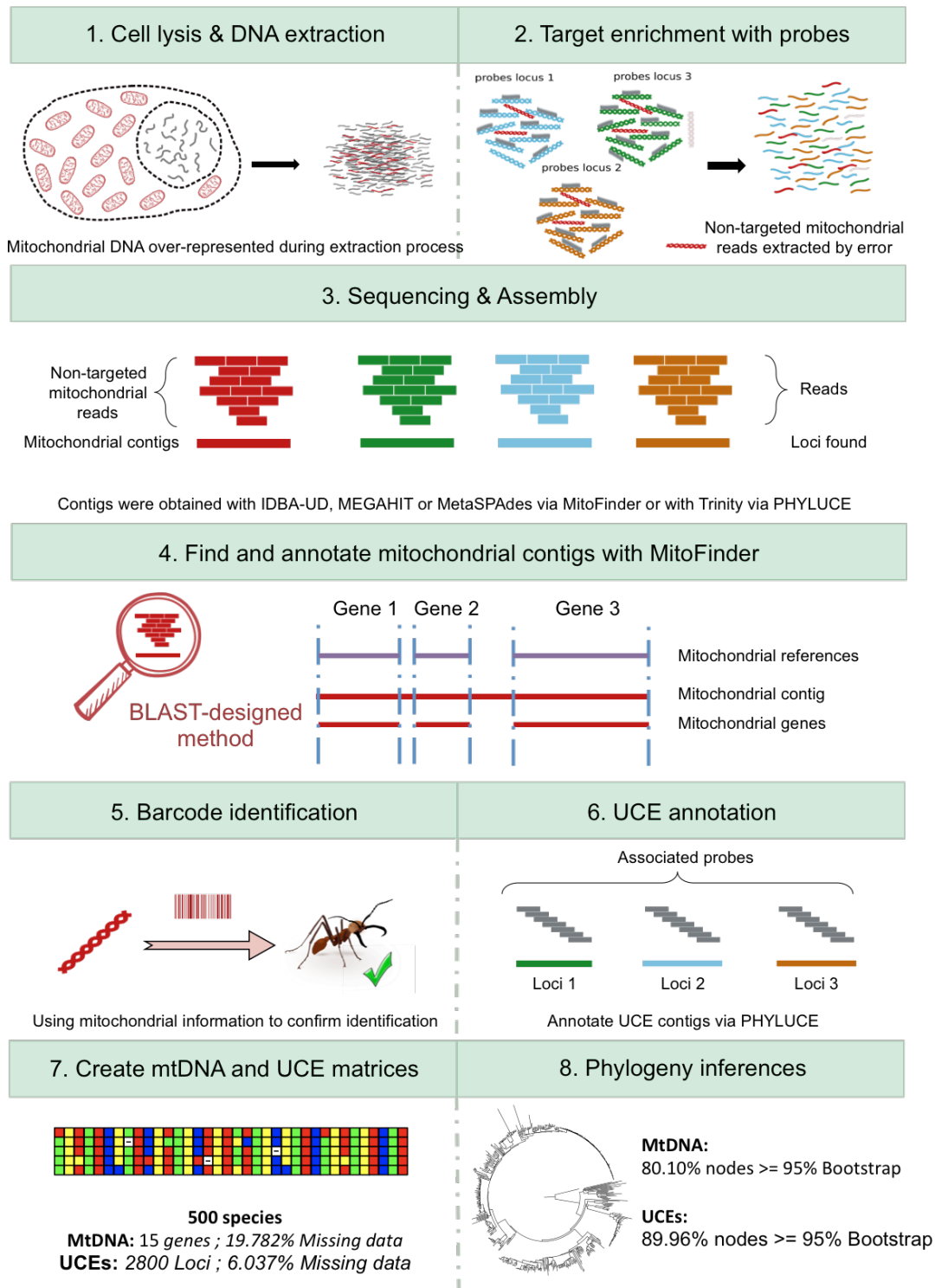


Figure 1. Conceptualization of the pipeline used to assemble and extract UCE and mitochondrial signal from ultraconserved element sequencing data.

prog alignSequences, which produces both nucleotide and amino acid alignments. To improve alignment accuracy and reduce calculation time, we used sequences from available mitogenomes as references for each CDS (option *-seq_lr*). Sequences with internal stop codons were excluded to remove incorrectly annotated fragments potentially corresponding to nuclear mitochondrial DNA segments (NUMTs) in each protein-coding gene alignment. Then, individual gene alignments were eye-checked to manually remove remaining aberrant sequences. Finally, a nucleotide supermatrix was created by concatenating protein-coding and ribosomal RNA genes. Considering mitochondrial signal saturation with high divergence, an amino acid supermatrix with the 13 mitochondrial CDSs was also assembled.

DNA barcoding

To verify species identification of the 501 ant UCE libraries, COI sequences extracted by MitoFinder using MetaSPAdes (mtDNA recovered for all species) were compared with Species Level Barcode Records (3,328,881 COI sequences including more than 100,000 ants) through the identification server of the Barcode Of Life Data System v4 (Ratnasingham & Hebert, 2007). The same COI sequences were also compared against the NCBI nucleotide database using Megablast with default parameters.

Assembly of UCEs

As recommended by Faircloth (2016), we first relied on Trinity to assemble UCE contigs using the *phyluce_assembly_assemblo_trinity* module of PHYLUCE. To assess the impact of assembler choice on UCE retrieval, we also used the assemblies obtained with IDBA-UD, MEGAHIT, and MetaSPAdes as implemented in MitoFinder. PHYLUCE scripts *phyluce_assembly_get_match_counts* and *phyluce_assembly_get_fastas_from_match_counts* were used to match contigs obtained for each sample to the bait set targeting 2590 UCE loci

for Hymenoptera (Branstetter *et al.* 2017b). The resulting alignments were then cleaned using Gblocks (Castresana 2000) with the *phyluce_align_get_gblocks_trimmed_alignments_from_untrimmed* script. Finally, loci found in at least 75% of species were selected to create the four corresponding UCE supermatrices using the *phyluce_align_get_only_loci_with_min_taxa* script.

Phylogenetic analyses

Phylogenetic relationships of ants were inferred from a total of 16 different supermatrices corresponding to the four supermatrices constructed from contigs obtained with each of the four assemblers (IDBA-UD, MEGAHIT, MetaSPAdes, and Trinity). The four supermatrices are as follows: (i) a UCE nucleotide supermatrix built from the concatenation of UCE loci retrieved for at least 75% of species, (ii) a mitochondrial nucleotide supermatrix consisting of the concatenation of the 13 protein-coding genes and the two rRNA genes, (iii) a mitochondrial amino-acid supermatrix of the 13 protein-coding genes, and (iv) a mixed supermatrix of UCE nucleotides and mitochondrial amino-acid protein-coding genes. For all supermatrices, phylogenetic inference was performed with Maximum Likelihood (ML) as implemented in IQ-TREE v1.6.8 (Nguyen, Schmidt, von Haeseler, & Minh 2015) using a GTR+ Γ_4 +I model for UCE and mitochondrial nucleotide supermatrices, a mtART+ Γ_4 +I model partitioned by gene for mitochondrial amino acids matrices, and a partitioned model mixing a GTR+ Γ_4 +I model for UCE nucleotides and a mtART+ Γ_4 +I model for mitochondrial amino acids for the mixed supermatrices. Statistical node support was estimated using ultrafast bootstrap (UFBS) with 1000 replicates (Hoang, Chernomor, von Haeseler, Minh, & Vinh 2018). Nodes with UFBS values higher than 95% were considered as strongly supported. For all supermatrices, the congruence between the different topologies obtained

with the four assemblers was evaluated by calculating quartet distances with Dquad (Ranwez, Criscuolo, & Douzery 2010).

Results

Assembly of UCE and mitochondrial datasets

De novo assembly of 501 UCE capture sequencing libraries was performed with four different assemblers: IDBA-UD, MEGAHIT, and MetaSPAdes via MitoFinder and Trinity via PHYLUCe. All assemblers provided different numbers of contigs (**Table 1**) ranging from 30,544 (IDBA-UD) to 114,392 (MEGAHIT) on average. The average computational time per assembly was highly variable among assemblers with Trinity being by far the slowest (35 CPUs, median total-time: 1h:06m:22s) and IDBA-UD the fastest (5 CPUs, median total-time: 0h:11m:01s), MEGAHIT (5 CPUs, median total-time: 0h:12m:35s) being slightly slower, and MetaSPAdes (5 CPUs, median total-time: 0h:25m:44s) being about twice slower than the other two metagenomic assemblers (**Table 1 & Fig. 2A**).

The UCE supermatrices created by PHYLUCe for each of the four assemblers contained on average 2580/2590 UCE loci for Hymenoptera (**Table 1**). All matrices contained 501 species, but the size of the supermatrix and the percentage of missing data varied depending on the assembler (**Table 1**). Trinity, which is generally used as the default assembler in PHYLUCe, resulted in the shortest and most incomplete supermatrix with 2579 loci representing 127,803 sites (40.5% variable) and 17.8% missing data. Among metagenomic assemblers, MetaSPAdes provided the largest and most complete supermatrix with 2582 loci representing 156,456 sites (44.5% variable) and only 6.0% missing data. IDBA-UD retrieved 2581 loci representing 132,403 sites (43.9% variable) with only 6.7% missing data, and MEGAHIT resulted in a supermatrix with 2579 loci representing 147,589 sites (43.2% variable) but with 12.4% missing data. Note that less than 30 loci were retrieved

A) Summary statistics for UCE assemblies and supermatrices

| Assembler | Assembly time | UCEs | | | | |
|---------------------|---------------|-------------------|----------------|-------------|-----------------|---------------|
| | | Number of contigs | Number of loci | Matrix size | %Variable sites | %Missing data |
| IDBA-UD (5 CPUs) | 0h:11m:02s | 30,544 | 2581 | 132,403 | 43.9 | 6.7 |
| MEGAHIT (5 CPUs) | 0h:12m:35s | 114,392 | 2579 | 147,589 | 43.2 | 12.5 |
| MetaSPAdes (5 CPUs) | 0h:25m:42s | 113,303 | 2582 | 156,456 | 44.3 | 6.1 |
| Trinity (35 CPUs) | 1h:06m:22s | 43,481 | 2579 | 127,803 | 40.5 | 17.8 |

B) Summary statistic for mtDNA assemblies and supermatrices

| Assembler | Mitogenomes | | | | | | | | |
|---------------------|-------------------|-------------------|-----------------|----------------|----------------|------------------|----------------|----------------|------------------|
| | Number of contigs | Number of species | Number of genes | AA matrix size | % missing data | % Variable sites | NT matrix size | % missing data | % Variable sites |
| IDBA-UD (5 CPUs) | 4.2 | 499 | 13.04 | 3764 | 20.9 | 86.7 | 13635 | 26.1 | 85.8 |
| MEGAHIT (5 CPUs) | 3.9 | 499 | 13.61 | 3757 | 15.3 | 87.5 | 13718 | 20.6 | 86.4 |
| MetaSPAdes (5 CPUs) | 3.8 | 501 | 13.73 | 3766 | 14.6 | 88.9 | 13713 | 19.8 | 87.1 |
| Trinity (35 CPUs) | 4.2 | 500 | 13.37 | 3760 | 18.0 | 86.9 | 13648 | 26.7 | 86.1 |

Table 1. Summary statistics on assembly results according to the assembler used. The values are averages over the 501 assemblies, except for the assembly time, which is a median value. The two tables report specific statistics for A) ultraconserved elements data, and B) mitochondrial data. Note that 35 CPUs were used for Trinity whereas 5 CPUs were used for others assemblers.

for *Phalacromyrmex fugax* (between 4 and 27 loci depending on the assembler). This is congruent with the original publication in which this library was not included in phylogenetic analyses (Branstetter *et al.* 2017a). Accordingly, we removed the *Phalacromyrmex fugax* library (SRR5437956) from the dataset.

Depending on the assembler, mitochondrial signal was recovered in 499, 500, and 501 libraries out of a total of 501 (**Table 1, Fig. 2B**). Overall, mitochondrial signal thus was detected in all libraries but only MetaSPAdes retrieved it in all species (**Appendix S2**). On average, 3.8 contigs per species was identified (**Table 1, Fig. 2B**) and 13.7 genes were annotated with MitoFinder (**Fig. 2C**). In 296/501 cases, MitoFinder was able to assemble a contig of more than 15,000 bp containing at least 13 annotated genes that likely represents the complete mitochondrial genome. In 52 of these cases, all 15 genes were annotated. In the remaining, the putative mitogenome contigs were missing one or two genes, mostly the short and divergent ATP8 (131/296), the 12S rRNA (29/296) and the 16S rRNA (10/296), which were present but not directly annotated by our blast-based procedure.

After alignment and cleaning, mitochondrial genes were used to create nucleotide and amino acid supermatrices. To be consistent with UCE analyses, and despite the recovery of some mitochondrial signal, we ignored *Phalacromyrmex fugax* in further analyses. In the nucleotide supermatrices (13 protein-coding + 12S and 16S rRNAs), we obtained 13 genes on average per species, which resulted in supermatrices with 13,679 nucleotide sites (86.4% variable) and 23.3% missing data on average (**Table 1**). In the amino acid matrices (13 protein-coding genes), we obtained supermatrices with 3762 amino acid sites (87.4% variable) and 17.2% missing data on average (**Table 1**).

Barcoding analyses

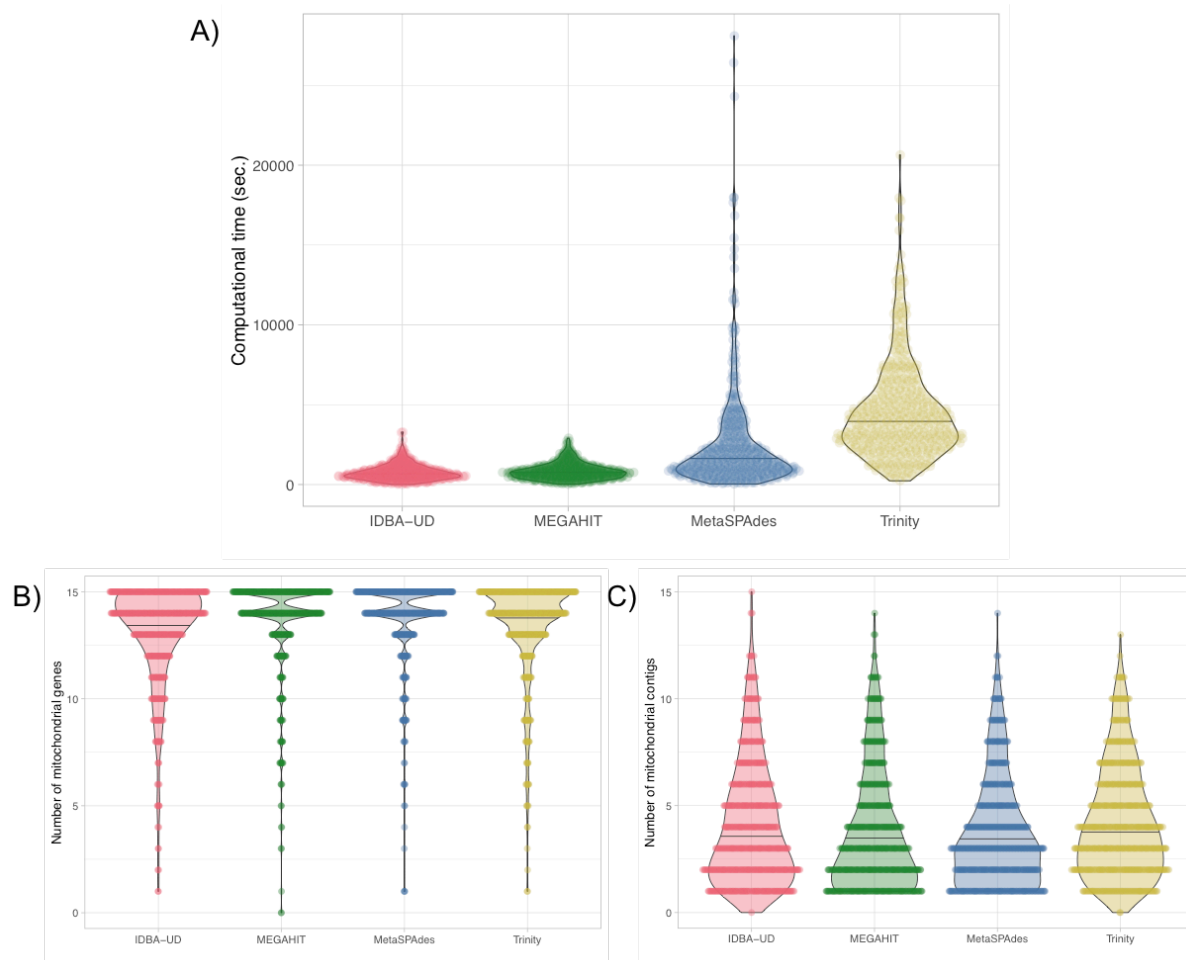


Figure 2. Comparison of the efficiency of the assemblers in terms of A) computational time, B) number of potentially mitochondrial contigs identified, and C) number of mitochondrial genes annotated. Violin plots reflect the data distribution with a horizontal line indicating the median. Note that for the three metagenomic assemblers, 5 CPUs were used compared to 35 CPUs for Trinity. Plots were obtained using PlotsOfData (Postma & Goedhart 2019).

A total of 534 COI sequences retrieved from the 501 MetaSPAdes assemblies were used to verify species identification of the UCE libraries (**Appendix S3**). In 42 cases, two or three COI sequence fragments were retrieved from the same UCE library. In seven of these cases, the slightly-overlapping COI fragments most likely resulted from bad assembly or erroneous annotation. However, in the 35 remaining cases, a genuine complete COI sequence overlapped with shorter fragments suggesting either cross-contaminations, nuclear mitochondrial DNA segments (NUMTs), endoparasites, or bacterial symbionts. For instance, in *Temnothorax* sp. mmp11 (SRR5809551), a 391 bp fragment annotated as COI by MitoFinder was found to be 98.2% identical to both the *Wolbachia pipientis* wAlbB and *Wolbachia* Pel strain wPip genomes, which are bacterial endosymbionts of the mosquitoes *Aedes albopictus* and *Culex quinquefasciatus*, respectively. Also, in *Sericomyrmex bondari* (SRR5044901) and *Sericomyrmex mayri* (SRR5044856) short COI fragments best matched with nematodes. However, in the 312 cases for which COI barcoding allowed to confirm the species identity of the UCE libraries, we did not detect any obvious cases of cross-contaminations where the COI extracted from a given library would have been identical to the one of another library (**Appendix S3**).

Phylogenetic results

The ML topologies inferred from the different UCE supermatrices were very similar with an average quartet distance of 0.005 among assemblers (**Appendix S4**). However, the percentage of supported nodes (UFBS > 95) differed depending on the assembler: IDBA-UD (91.37%), MetaSPAdes (89.96%), MEGAHIT (89.56%), and Trinity (85.85%). In the following, we only discuss the phylogenetic results obtained with MetaSPAdes that provides the most comprehensive assemblies for both UCE and mitochondrial data (**Table 1**). The following 12 well-established subfamilies were retrieved with maximal UFBS support

(100%): Aneuretinae, Amblyoponinae, Dolichoderinae, Dorylinae, Ectatomminae, Formicinae, Heteroponerinae, Myrmeciinae, Myrmicinae, Pseudomyrmecinae, and Ponerinae (**Fig. 3A**). The two supergroups Formicoid and Poneroid were also retrieved with maximal UFBS support, as well as consensual phylogenetic relationships among Formicoid subfamilies (Ward 2014).

For mitochondrial matrices, the percentage of supported nodes (UFBS > 95) with nucleotides also differed depending on the assembler and was higher than with the amino acids: MetaSPAdes (84.5% vs. 80.1%), MEGAHIT (84.0% vs. 79.4%), Trinity (83.3% vs. 80.4%), and IDBA-UD (80.2% vs. 78.0%). However, ML mitogenomic trees inferred from amino acids were more congruent with UCE topologies than the ones inferred from the mitochondrial nucleotides (average quartet distance = 0.035 v.s. 0.063; **Appendix S4**). Among assemblers, the ML topologies inferred with amino acid matrices were highly congruent with an average quartet distance of 0.007 (**Appendix S4**). In the ML tree obtained with the MetaSPAdes supermatrix (**Fig. 3B**), all ant subfamilies were retrieved with maximal UFBS support values except for Myrmicinae (93%), Ponerinae (97%), and Proceratinae (99%) (**Fig. 3A**). However, relationships among subfamilies were not congruent with UCE phylogenomic inferences except for Heteroponerinae + Ectatomminae (UFBS = 100) and Dolichoderinae + Aneuretinae (UFBS = 96) (**Fig. 3A**).

Finally, phylogenetic inference carried on mixed supermatrices composed of UCes and mitochondrial amino acids resulted in ML topologies that were also highly similar among assemblers with an average quartet distance of 0.006 (**Appendix S4**). The percentage of supported nodes (UFBS > 95) were: IDBA-UD (91.2%), MEGAHIT (92.8%), MetaSPAdes (92.2%), and Trinity (90.4%). As with UCE matrices, the 12 well-established subfamilies, the two supergroups Formicoid and Poneroid and consensual Formicoid inter-subfamilies relationships (Ward 2014) were all retrieved with maximal UFBS support.

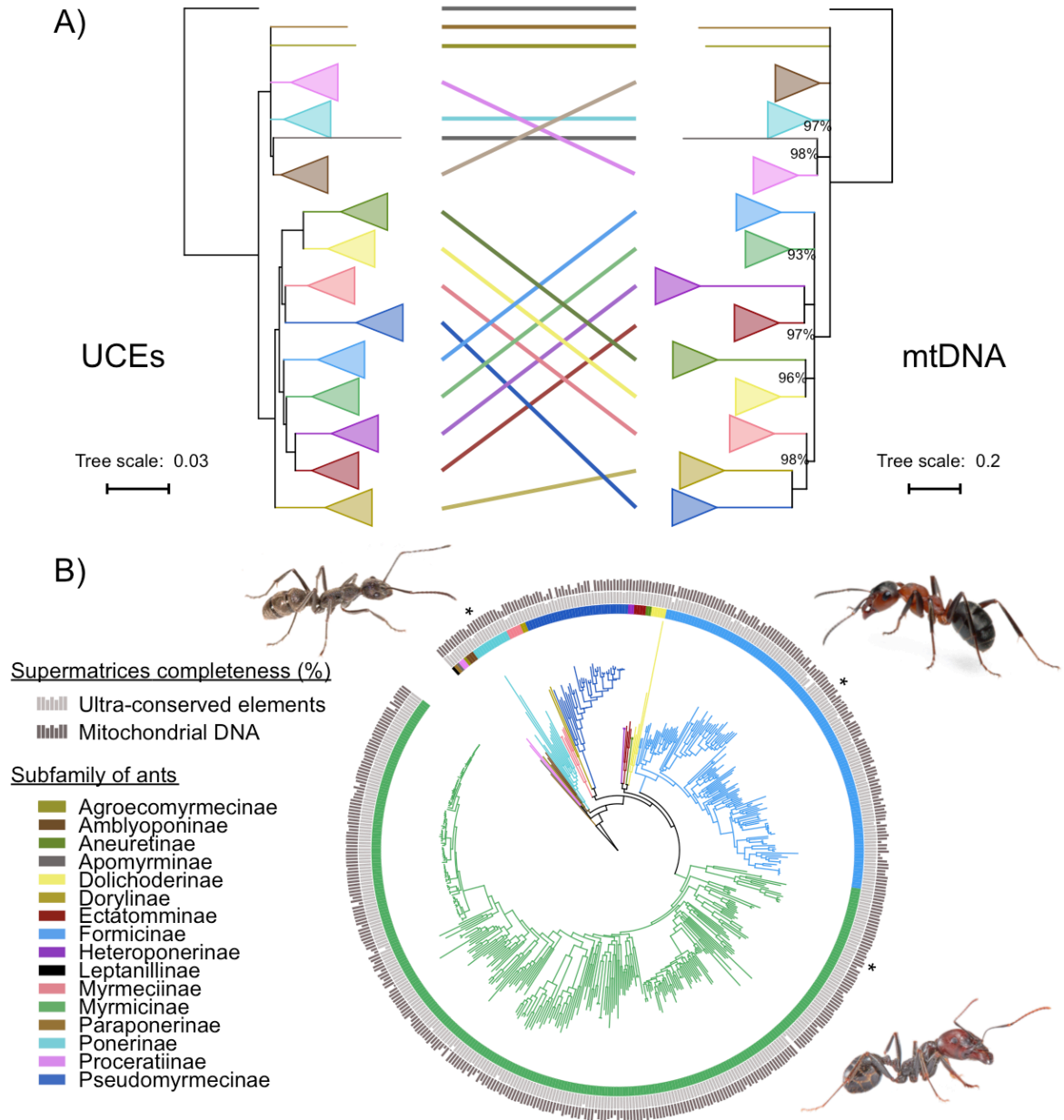


Figure 3. Phylogenomic relationships of Formicidae. A) The topology reflects the results of phylogenetic analyses based on UCEs amino acid mitochondrial supermatrix. Histograms reflect the percent of UCE (dark grey) and mitochondrial genes (light grey) recovered for each species. Illustrative pictures (*): *Diacamma* sp. (Ponerinae; top left), *Formica* sp. (Formicinae; top right), and *Messor barbarus* (Myrmicinae; bottom right). B) Mito-nuclear phylogenetic differences for ancient relationships. Clade corresponding to subfamilies were collapsed. Inter-subfamily relationships with UFBS < 95 were collapsed. Non-maximal node support values are reported.

Discussion

Metagenomic assemblers are efficient tools for assembling UCEs

Currently, genomic and transcriptomic *de novo* assemblers are commonly used to assemble UCE loci from DNA capture sequencing data (Faircloth 2016). Since metagenomic assemblers such as IDBA-UD, MEGAHIT, and MetaSPAdes have been designed to account for variance in sequencing coverage, they seem to be well adapted for targeted enrichment or DNA sequence capture data. Our results show that metagenomic assemblers are indeed more efficient at assembling UCE loci than the classically used, but computationally intensive, Trinity transcriptomic assembler. As a consequence, they could lead to datasets containing more variable sites, less missing data, and increased phylogenetic signal (**Table 1**). Indeed, the topologies obtained with the metagenomic assemblers are very similar to the topology obtained with the Trinity-based supermatrix, contain a higher number of supported nodes (UFBS \geq 95%), and are consistent to previous studies (Ward 2014). Furthermore, assemblies obtained with the three metagenomic assemblers provide variable numbers of contigs (ranging from 30,544 to 114,392) resulting in differences in the completeness of the matrices (6.0% to 17.8% of missing data for UCE matrices and 29.9% to 41.3% for mitochondrial matrices) and in numbers of variable sites (for UCE, 40.5% to 44.3%; for mtDNA, 77.2% to 79.0%). Interestingly, for both UCE matrices and mtDNA matrices, MetaSPAdes consistently provides more loci, more variable sites, and less missing data. In addition, mitochondrial signal was extracted from all libraries only using MetaSPAdes within Mitofinder. Despite a computation time on average twice that of the other two metagenomic assemblers, MetaSPAdes is the more efficient assembler for ant UCEs. This software therefore provides a much needed alternative to Trinity for efficiently assembling hundreds of UCE libraries.

Mitochondrial signal can systematically be extracted from UCE capture data

Ultraconserved elements are key loci exploited as target capture sequences in an increasing number of phylogenomic studies. DNA sequence capture methods are used to efficiently enrich targeted DNA regions in library preparation prior to sequencing, but non-targeted regions are always sequenced in the process resulting in so called “off-target reads”. Interestingly, off-target reads could represent up to 40% of the sequenced reads in exome capture experiments (Chilamakuri *et al.* 2014) and many contigs not belonging to targeted UCE loci are typically assembled from UCE capture data (e.g. Smith, Harvey, Faircloth, Glenn, & Brumfield, 2014; Faircloth *et al.* 2015). Given this high proportion of off-target reads, we can expect that mitochondrial DNA could be found as off-target sequences in many target enrichment data. Accordingly, several studies have succeeded in extracting mtDNA from UCE libraries (e.g. Smith *et al.* 2014; do Amaro *et al.* 2015). The development of MitoFinder allowed the automatic extraction of mitochondrial signal from all 501 ant UCE libraries. This maximum success rate indicates that this approach is highly efficient at least in Formicidae. However, the success in retrieving mitochondrial sequences, ultimately depends on the number of mitochondria contained in the tissue used for DNA extraction and library preparation. As expected, mitochondrial off-target reads are much more common in muscle and heart than in lung tissues in human (D’Erchia *et al.* 2015). Similarly, mitochondrial sequences are probably rare or absent in library constructed from vertebrate blood, even in birds in which nucleated red blood cells contain mitochondria, but in very low numbers (Reverter *et al.* 2016). In invertebrates, our case study with 100% success rate in ant UCES demonstrates that mitochondrial sequences could probably be easily retrieved for many arthropod taxa as a by product of target enrichment sequencing experiments.

The value of complementary mitochondrial signal

Mitochondrial sequences could provide interesting and important complementary information compared to nuclear sequences. First, mtDNA can be used to confirm the identity of the species sequenced for conserved UCE loci. Here, we were able to confirm the identification of 312 ant species out of the 501 UCE libraries using COI barcoding without revealing a single case of obvious species misidentification. Given that ant UCE libraries have been constructed from museum specimens, the 501 COI sequences we annotated could be used as reference barcoding sequences in future studies. Then, even though we did not detect such cases, the high mutation rate and the absence of heterozygous sites in mtDNA also make it well adapted for cross-contamination detection analyses (Ballenghien *et al.* 2017).

Nevertheless, mitochondrial markers also have some well identified limitations (Galtier *et al.* 2009). First, mtDNA could be inserted in the nuclear genome in the form of NUMTs (Bensasson, Zhang, Hartl, & Hewitt 2001). NUMTs could potentially be assembled as off-target contigs in DNA capture libraries and we might have indeed extracted some fragments corresponding to NUMTs for the COI gene using MitoFinder (**Appendix S2**). Theoretically, NUMTs could be picked up by analysing the coverage of putative mitochondrial contigs as they are expected to have a coverage comparable to other off-targets nuclear contigs, whereas genuine mitochondrial contigs should have a higher coverage. A second limitation of mtDNA exists in arthropods where maternally inherited intra-cellular bacteria are frequent. Among those bacteria, *Wolbachia* is particularly widespread and could distort the mitochondrial genealogy when a particular strain spreads within the host species hitchhiking its linked mitochondrial haplotype (Cariou, Duret, & Charlat 2017). *Wolbachia* is frequent among ants and could therefore be responsible of some mito-nuclear discordance (Wenseleers *et al.* 1998). We indeed discovered such an instance with a *Wolbachia* COI

sequence identified in *Temnothorax* sp. mmp11 (SRR5809551), which was confirmed by several assembled contigs matching to *Wolbachia* strain genomes in this sample.

Beyond the methodological aspects of species identification and potential cross-contamination detection, mitochondrial sequences could also be useful to tackle fundamental evolutionary questions. UCEs have also proved to be useful genetic markers for phylogeography and for resolving shallow phylogenetic relationships (Musher & Cracraft 2018; Smith *et al.* 2014). In this context, mtDNA could also bring complementary information. In most animals, mtDNA has a maternal inheritance without recombination, which means that all mitochondrial genes behave as a single locus. This simplifies the interpretation of the phylogenetic pattern between closely related species or within subdivided populations of a species. Mito-nuclear phylogenetic discordance could also reveal interesting phenomena involving hybridization, sex-biased dispersal, and introgression (Toews & Brelsford, 2012). In practise, hybridization events are often identified using mito-nuclear discordance (Li *et al.* 2016) and in some cases, the mitochondrial introgression events have proven to be adaptive (Seixas, Boursot & Melo-Ferreira 2018). Nevertheless, in our ant case study, a detailed comparison of mitochondrial and UCE phylogenies did not allow revealing convincing occurrences of such discordances.

Ant phylogenetic relationships from 500 UCE and mitochondrial data

Both nuclear and mitochondrial data retrieved the most consensual phylogenetic relationships in the ant phylogeny (Ward 2014; Branstetter *et al.* 2017b; Borowiec *et al.* 2019). Twelve Formicidae subfamilies were recovered as monophyletic in all analyses, both with the nuclear and mitochondrial datasets, confirming their robustness. However, the well-defined inter-subfamily relationships within Formicoids (Ward 2014; Branstetter *et al.* 2017; Borowiec *et al.* 2019) were only supported by the UCE dataset, but not by the mitochondrial amino acid

dataset. For example, the army ant subfamily (Dorylinae) was not retrieved as the sister-group of all other Formicoids, but was the closest relative of Pseudomyrmicinae (UFBS = 100). Similarly, contradicting the classical and well-defined relationship of Heteroponerinae+Ectatomminae as the sister-group of Myrmicinae (Ward 2014; Branstetter *et al.* 2017; Borrowiec *et al.* 2019), the mitochondrial dataset supported an alternative relationship with Dolichoderinae+Aneuretinae (UFBS = 96). These differences suggest that mitochondrial data might be not well-suited to resolve ancient phylogenetic relationships at the ant inter-subfamily level, even if they look suitable for more recent nodes such as intra-subfamily relationships.

Interestingly, these topological incongruences between UCEs and mitochondrial genes also featured different topologies regarding the existence of the Poneroid taxa, a controversial clade not always retrieved depending on the studies (Ward 2014), but that tends to be retrieved in the most recent studies (Branstetter *et al.* 2017; Borrowiec *et al.* 2019; UCE dataset in this study) and is not recovered by our mitochondrial amino acid dataset (**Fig. 3B**). The same applies to the phylogenetic placement of Apomyrminae, a subfamily either grouped with Leptanillinae or Amblyoponinae in past studies (Ward 2014), but that was grouped with Proceratinae in our mitochondrial dataset (UFBS = 98; **Fig. 3B**). For such controversial nodes, our study demonstrates that the nature of the phylogenetic markers can provide different results. Such differences between nuclear and mitochondrial data might be due to the substitutional saturation of mitochondrial data even at the amino acid level. This problem may actually be exacerbated in hymenopteran mitochondria that possess high AT content translating into strongly biased codon usage potentially leading to phylogenetic reconstruction artefacts (Foster, Jermin & Hickey 1997; Foster & Hickey 1999). Interestingly, such differences between mitochondrial and nuclear inference for ancient phylogenetic relationships, is not observed with insects with less AT-rich mitochondrial

genomes such as, for instance, swallowtail butterflies (Condamine, Nabholz, Clamens, Dupuis & Sperling, 2018; Allio *et al.* 2019) or tiger beetles (Vogler & Pearson 1996). This calls for additional studies on both controversial and consensual ant inter-subfamily relationships with more comprehensive genome-wide datasets.

Conclusions

In this study, we developed the MitoFinder tool to automatically extract and annotate mitogenomic data from raw sequencing data in an efficient way. For the assembly step of our pipeline, we tested four different assemblers and showed that MetaSPAdes is the most efficient and accurate assembler for both UCE and mitochondrial data. Applying MitoFinder to ants, we were able to extract mitochondrial signal from 501 UCE libraries. This demonstrates that mitochondrial DNA can be found as off-target sequences in UCEs sequencing data. Interestingly, mitochondrial DNA extracted from UCE libraries can also be used to: (i) confirm species identification with barcoding methods, (ii) highlight potential sample cross-contamination, and (iii) reveal potential cases of mito-nuclear discordance caused by hybridization events leading to mitochondrial introgression. Finally, MitoFinder was developed with UCE libraries but our approach should also work with data obtained from other capture methods in which numerous off-targets reads are sequenced, as well as with transcriptomic and whole genome sequencing data, in which mitochondrial reads are overrepresented.

Acknowledgements

This paper is dedicated to the memory of graduate student Alex Schomaker-Bastos (1992-2015) who was assassinated by the time he was writing the mitoMaker program on which we built upon for the annotation module of MitoFinder. We also thank Fabien Condamine for

providing helpful comments on a previous version of the manuscript. This work has been supported by grants from Investissements d’Avenir of the Agence Nationale de la Recherche (CEBA: ANR-10-LABX-25-01; CEMEB: ANR-10-LABX-0004), and the European Research Council (ERC-2015-CoG-683257 project ConvergeAnt). This is contribution ISEM 2019-XXX of the Institut des Sciences de l’Evolution de Montpellier.

References

- Allio, R., Scornavacca, C., Nabholz, B., Clamens, A. L., Sperling, F. A., & Condamine, F. (2019). Whole genome shotgun phylogenomics resolves the pattern and timing of swallowtail butterfly evolution. *Systematic Biology*, syz030. doi:10.1093/sysbio/syz030
- Baca, S. M., Alexander, A., Gustafson, G. T., & Short, A. E. Z. (2017). Ultraconserved elements show utility in phylogenetic inference of Adephaga (Coleoptera) and suggest paraphyly of ‘Hydradephaga’. *Systematic Entomology*, 42(4), 786–795. doi:10.1111/syen.12244
- Ballenghien, M., Faivre, N., & Galtier, N. (2017). Patterns of cross-contamination in a multispecies population genomic project: detection, quantification, impact, and solutions. *BMC Biology*, 15(1), 25. doi:10.1186/s12915-017-0366-6
- Bejerano, G., Pheasant, M., Makunin, I., Stephen, S., Kent, W. J., Mattick, J. S., & Haussler, D. (2004). Ultraconserved elements in the human genome. *Science*, 304(5675), 1321–5. doi:10.1126/science.1098119
- Bensasson, D., Zhang, D.-X., Hartl, D. L., & Hewitt, G. M. (2001). Mitochondrial pseudogenes: evolution’s misplaced witnesses. *Trends in Ecology & Evolution*, 16(6), 314–321. doi:10.1016/S0169-5347(01)02151-6
- Bi, K., Vanderpool, D., Singhal, S., Linderoth, T., Moritz, C., & Good, J. M. (2012). Transcriptome-based exon capture enables highly cost-effective comparative genomic

data collection at moderate evolutionary scales. *BMC Genomics*, *13*(1), 403.

doi:10.1186/1471-2164-13-403

Blaimer, B. B., Brady, S. G., Schultz, T. R., Lloyd, M. W., Fisher, B. L., & Ward, P. S.

(2015). Phylogenomic methods outperform traditional multi-locus approaches in resolving deep evolutionary history: a case study of formicine ants. *BMC Evolutionary Biology*, *15*(1), 271. doi:10.1186/s12862-015-0552-5

Bolger, A. M., Lohse, M., & Usadel, B. (2014). Trimmomatic: a flexible trimmer for Illumina sequence data. *Bioinformatics*, *30*(15), 2114–2120. doi:10.1093/bioinformatics/btu170

Borowiec, M. L., Rabeling, C., Brady, S. G., Fisher, B. L., Schultz, T. R., & Ward, P. S.

(2019). Compositional heterogeneity and outgroup choice influence the internal phylogeny of the ants. *Molecular Phylogenetics and Evolution*, *134*, 111–121.

doi:10.1016/J.YMPEV.2019.01.024

Bourguignon, T., Lo, N., Šobotník, J., Ho, S. Y. W., Iqbal, N., Coissac, E., ... Evans, T. A.

(2016). Mitochondrial phylogenomics resolves the global spread of higher termites, ecosystem engineers of the tropics. *Molecular Biology and Evolution*, *34*(3), 589–597.

doi:10.1093/molbev/msw253

Branstetter, M. G., Danforth, B. N., Pitts, J. P., Faircloth, B. C., Ward, P. S., Buffington, M.

L., ... Brady, S. G. (2017a). Phylogenomic Insights into the Evolution of Stinging Wasps and the Origins of Ants and Bees. *Current Biology*, *27*(7), 1019–1025.

doi:10.1016/J.CUB.2017.03.027

Branstetter, M. G., Ješovnik, A., Sosa-Calvo, J., Lloyd, M. W., Faircloth, B. C., Brady, S. G.,

& Schultz, T. R. (2017b). Dry habitats were crucibles of domestication in the evolution of agriculture in ants. *Proceedings of the Royal Society B: Biological Sciences*,

284(1852), 20170095. doi:10.1098/rspb.2017.0095

- Branstetter, M. G., Longino, J. T., Ward, P. S., & Faircloth, B. C. (2017c). Enriching the ant tree of life: enhanced UCE bait set for genome-scale phylogenetics of ants and other Hymenoptera. *Methods in Ecology and Evolution*, *8*(6), 768–776. doi:10.1111/2041-210X.12742
- Breinolt, J. W., Earl, C., Lemmon, A. R., Lemmon, E. M., Xiao, L., & Kawahara, A. Y. (2018). Resolving relationships among the megadiverse butterflies and moths with a novel pipeline for anchored phylogenomics. *Systematic Biology*, *67*(1), 78–93. doi:10.1093/sysbio/syx048
- Caparroz, R., Rocha, A. V., Cabanne, G. S., Tubaro, P., Aleixo, A., Lemmon, E. M., & Lemmon, A. R. (2018). Mitogenomes of two neotropical bird species and the multiple independent origin of mitochondrial gene orders in Passeriformes. *Molecular Biology Reports*, *45*(3), 279–285. doi:10.1007/s11033-018-4160-5
- Capella-Gutierrez, S., Silla-Martinez, J. M., & Gabaldon, T. (2009). trimAl: a tool for automated alignment trimming in large-scale phylogenetic analyses. *Bioinformatics*, *25*(15), 1972–1973. doi:10.1093/bioinformatics/btp348
- Cariou, M., Duret, L., & Charlat, S. (2017). The global impact of *Wolbachia* on mitochondrial diversity and evolution. *Journal of Evolutionary Biology*, *30*(12), 2204–2210. doi:10.1111/jeb.13186
- Castresana, J. (2000). Selection of conserved blocks from multiple alignments for their use in phylogenetic analysis. *Molecular biology and evolution*, *17*(4), 540-552. doi:10.1093/oxfordjournals.molbev.a026334
- Chilamakuri, C. S., Lorenz, S., Madoui, M.-A., Vodák, D., Sun, J., Hovig, E., ... Meza-Zepeda, L. A. (2014). Performance comparison of four exome capture systems for deep sequencing. *BMC Genomics*, *15*(1), 449. doi:10.1186/1471-2164-15-449

- Coissac, E., Hollingsworth, P. M., Lavergne, S., & Taberlet, P. (2016). From barcodes to genomes: extending the concept of DNA barcoding. *Molecular Ecology*, *25*(7), 1423–1428. doi:10.1111/mec.13549
- Condamine, F. L., Nabholz, B., Clamens, A.-L., Dupuis, J. R., & Sperling, F. A. (2018). Mitochondrial phylogenomics, the origin of swallowtail butterflies, and the impact of the number of clocks in Bayesian molecular dating. *Systematic entomology*, *43*(3), 460–480. doi:10.1111/syen.12284
- D’Erchia, A. M., Atlante, A., Gadaleta, G., Pavesi, G., Chiara, M., De Virgilio, C., ... Pesole, G. (2015). Tissue-specific mtDNA abundance from exome data and its correlation with mitochondrial transcription, mass and respiratory activity. *Mitochondrion*, *20*, 13–21. doi:10.1016/J.MITO.2014.10.005
- Di Franco, A., Poujol, R., Baurain, D., & Philippe, H. (2019). Evaluating the usefulness of alignment filtering methods to reduce the impact of errors on evolutionary inferences. *BMC Evolutionary Biology*, *19*(1), 21. doi:10.1186/s12862-019-1350-2
- do Amaral, F. R., Neves, L. G., Resende Jr, M. F. R., Mobili, F., Miyaki, C. Y., Pellegrino, K. C. M., & Biondo, C. (2015). Ultraconserved Elements Sequencing as a Low-Cost Source of Complete Mitochondrial Genomes and Microsatellite Markers in Non-Model Amniotes. *PLoS One*, *10*(9), e0138446. doi:10.1371/journal.pone.0138446
- Dowton, M., & Austin, A. D. (1999). Evolutionary dynamics of a mitochondrial rearrangement ‘hot spot’ in the Hymenoptera. *Molecular Biology and Evolution*, *16*(2), 298–309. doi:10.1093/oxfordjournals.molbev.a026111
- Dowton, M., Castro, L. R., & Austin, A. D. (2002). Mitochondrial gene rearrangements as phylogenetic characters in the invertebrates: the examination of genome ‘morphology’. *Invertebrate Systematics*, *16*(3), 345. doi:10.1071/IS02003

- Esselstyn, J. A., Oliveros, C. H., Swanson, M. T., & Faircloth, B. C. (2017). Investigating Difficult Nodes in the Placental Mammal Tree with Expanded Taxon Sampling and Thousands of Ultraconserved Elements. *Genome Biology and Evolution*, *9*(9), 2308–2321. doi:10.1093/gbe/evx168
- Faircloth, B. C., McCormack, J. E., Crawford, N. G., Harvey, M. G., Brumfield, R. T., & Glenn, T. C. (2012). Ultraconserved Elements Anchor Thousands of Genetic Markers Spanning Multiple Evolutionary Timescales. *Systematic Biology*, *61*(5), 717–726. doi:10.1093/sysbio/sys004
- Faircloth, B. C., Branstetter, M. G., White, N. D., & Brady, S. G. (2015). Target enrichment of ultraconserved elements from arthropods provides a genomic perspective on relationships among Hymenoptera. *Molecular Ecology Resources*, *15*(3), 489–501. doi:10.1111/1755-0998.12328
- Faircloth, B. C. (2016). PHYLUCE is a software package for the analysis of conserved genomic loci. *Bioinformatics*, *32*(5), 786–788. doi:10.1093/bioinformatics/btv646
- Foster, P. G., Jermin, L. S., & Hickey, D. A. (1997). Nucleotide composition bias affects amino acid content in proteins coded by animal mitochondria. *Journal of Molecular Evolution*, *44*(3), 282–288. doi:10.1007/PL00006145
- Foster, P. G., & Hickey, D. A. (1999). Compositional bias may affect both DNA-based and protein-based phylogenetic reconstructions. *Journal of Molecular Evolution*, *48*(3), 284–290. doi:10.1007/PL000006
- Galtier, N., Nabholz, B., Glémin, S., & Hurst, G. D. D. (2009). Mitochondrial DNA as a marker of molecular diversity: a reappraisal. *Molecular Ecology*, *18*(22), 4541–4550. doi:10.1111/j.1365-294X.2009.04380.x

- Grabherr, M. G., Haas, B. J., Yassour, M., Levin, J. Z., Thompson, D. A., Amit, I., ... Regev, A. (2011). Full-length transcriptome assembly from RNA-Seq data without a reference genome. *Nature Biotechnology*, *29*(7), 644–652. doi:10.1038/nbt.1883
- Grummer, J. A., Morando, M. M., Avila, L. J., Sites, J. W., & Leaché, A. D. (2018). Phylogenomic evidence for a recent and rapid radiation of lizards in the Patagonian *Liolaemus fitzingerii* species group. *Molecular Phylogenetics and Evolution*, *125*, 243–254. doi:10.1016/J.YMPEV.2018.03.023
- Guschanski, K., Krause, J., Sawyer, S., Valente, L. M., Bailey, S., Finstermeier, K., ... Savolainen, V. (2013). Next-Generation Museomics Disentangles One of the Largest Primate Radiations. *Systematic Biology*, *62*(4), 539–554. doi:10.1093/sysbio/syt018
- Hassanin, A., Delsuc, F., Ropiquet, A., Hammer, C., Jansen van Vuuren, B., Matthee, C., ... Couloux, A. (2012). Pattern and timing of diversification of Cetartiodactyla (Mammalia, Laurasiatheria), as revealed by a comprehensive analysis of mitochondrial genomes. *Comptes Rendus Biologies*, *335*(1), 32–50. doi:10.1016/J.CRVI.2011.11.002
- Hoang, D. T., Chernomor, O., von Haeseler, A., Minh, B. Q., & Vinh, L. S. (2018). UFBoot2: Improving the Ultrafast Bootstrap Approximation. *Molecular Biology and Evolution*, *35*(2), 518–522. doi:10.1093/molbev/msx281
- Ješovnik, A., Sosa-Calvo, J., Lloyd, M. W., Branstetter, M. G., Fernández, F., & Schultz, T. R. (2017). Phylogenomic species delimitation and host-symbiont coevolution in the fungus-farming ant genus *Sericomyrmex* Mayr (Hymenoptera: Formicidae): ultraconserved elements (UCEs) resolve a recent radiation. *Systematic Entomology*, *42*(3), 523–542. doi:10.1111/syen.12228
- Katoh, K., & Standley, D. M. (2013). MAFFT Multiple Sequence Alignment Software Version 7: Improvements in Performance and Usability. *Molecular Biology and Evolution*, *30*(4), 772–780. doi:10.1093/molbev/mst010

- Laslett, D., & Canback, B. (2007). ARWEN: a program to detect tRNA genes in metazoan mitochondrial nucleotide sequences. *Bioinformatics*, *24*(2), 172–175.
doi:10.1093/bioinformatics/btm573
- Le, S. Q., & Gascuel, O. (2008). An Improved General Amino Acid Replacement Matrix. *Molecular Biology and Evolution*, *25*(7), 1307–1320. doi:10.1093/molbev/msn067
- Le, S. Q., Dang, C. C., & Gascuel, O. (2012). Modeling Protein Evolution with Several Amino Acid Replacement Matrices Depending on Site Rates. *Molecular Biology and Evolution*, *29*(10), 2921–2936. doi:10.1093/molbev/mss112
- Lemmon, A. R., Emme, S. A., & Lemmon, E. M. (2012). Anchored Hybrid Enrichment for Massively High-Throughput Phylogenomics. *Systematic Biology*, *61*(5), 727–744.
doi:10.1093/sysbio/sys049
- Lemmon, E. M., & Lemmon, A. R. (2013). High-Throughput Genomic Data in Systematics and Phylogenetics. *Annual Review of Ecology, Evolution, and Systematics*, *44*(1), 99–121. doi:10.1146/annurev-ecolsys-110512-135822
- Li, D., Luo, R., Liu, C.-M., Leung, C.-M., Ting, H.-F., Sadakane, K., ... Lam, T.-W. (2016). MEGAHIT v1.0: A fast and scalable metagenome assembler driven by advanced methodologies and community practices. *Methods*, *102*, 3–11.
doi:10.1016/J.YMETH.2016.02.020
- McCormack, J. E., Faircloth, B. C., Crawford, N. G., Gowaty, P. A., Brumfield, R. T., & Glenn, T. C. (2012). Ultraconserved elements are novel phylogenomic markers that resolve placental mammal phylogeny when combined with species-tree analysis. *Genome Research*, *22*(4), 746–54. doi:10.1101/gr.125864.111
- McCormack, J. E., Harvey, M. G., Faircloth, B. C., Crawford, N. G., Glenn, T. C., & Brumfield, R. T. (2013a). A Phylogeny of Birds Based on Over 1,500 Loci Collected

by Target Enrichment and High-Throughput Sequencing. *PLoS ONE*, 8(1), e54848.

doi:10.1371/journal.pone.0054848

McCormack, J. E., Hird, S. M., Zellmer, A. J., Carstens, B. C., & Brumfield, R. T. (2013b).

Applications of next-generation sequencing to phylogeography and phylogenetics.

Molecular Phylogenetics and Evolution, 66(2), 526–538.

doi:10.1016/J.YMPEV.2011.12.007

Meiklejohn, K. A., Danielson, M. J., Faircloth, B. C., Glenn, T. C., Braun, E. L., & Kimball,

R. T. (2014). Incongruence among different mitochondrial regions: A case study using complete mitogenomes. *Molecular Phylogenetics and Evolution*, 78, 314–323.

doi:10.1016/j.ympcv.2014.06.003

Meza-Lázaro, R. N., Poteaux, C., Bayona-Vásquez, N. J., Branstetter, M. G., & Zaldívar-

Riverón, A. (2018). Extensive mitochondrial heteroplasmy in the neotropical ants of the

Ectatomma ruidum complex (Formicidae: Ectatomminae). *Mitochondrial DNA Part A*,

29(8), 1203–1214. doi:10.1080/24701394.2018.1431228

Musher, L. J., & Cracraft, J. (2018). Phylogenomics and species delimitation of a complex

radiation of Neotropical suboscine birds (Pachyramphus). *Molecular Phylogenetics and*

Evolution, 118, 204–221. doi:10.1016/j.ympcv.2017.09.013

Nguyen, L.-T., Schmidt, H. A., von Haeseler, A., & Minh, B. Q. (2015). IQ-TREE: A Fast

and Effective Stochastic Algorithm for Estimating Maximum-Likelihood Phylogenies.

Molecular Biology and Evolution, 32(1), 268–274. doi:10.1093/molbev/msu300

Nurk, S., Meleshko, D., Korobeynikov, A., & Pevzner, P. A. (2017). metaSPAdes: a new

versatile metagenomic assembler. *Genome Research*, 27(5), 824–834.

doi:10.1101/gr.213959.116

- Peng, Y., Leung, H. C. M., Yiu, S. M., & Chin, F. Y. L. (2012). IDBA-UD: a de novo assembler for single-cell and metagenomic sequencing data with highly uneven depth. *Bioinformatics*, 28(11), 1420–1428. doi:10.1093/bioinformatics/bts174
- Picardi, E., & Pesole, G. (2012). Mitochondrial genomes gleaned from human whole-exome sequencing. *Nature Methods*, 9(6), 523–524. doi:10.1038/nmeth.2029
- Pie, M. R., Ströher, P. R., Belmonte-Lopes, R., Bornschein, M. R., Ribeiro, L. F., Faircloth, B. C., & McCormack, J. E. (2017). Phylogenetic relationships of diurnal, phytotelm-breeding *Melanophryniscus* (Anura: Bufonidae) based on mitogenomic data. *Gene*, 628, 194–199. doi:10.1016/J.GENE.2017.07.048
- Pierce, M. P., Branstetter, M. G., & Longino, J. T. (2017). Integrative taxonomy reveals multiple cryptic species within Central American *Hylomyrma* FOREL, 1912 (Hymenoptera: Formicidae). *Myrmecological News*, 25, 131–143. doi:10.25849/myrmecol.news_025:131
- Prebus, M. (2017). Insights into the evolution, biogeography and natural history of the acorn ants, genus *Temnothorax* Mayr (hymenoptera: Formicidae). *BMC Evolutionary Biology*, 17(1), 250. doi:10.1186/s12862-017-1095-8
- Postma, M., & Goedhart, J. (2019). PlotsOfData—A web app for visualizing data together with their summaries. *PLoS biology*, 17(3), e3000202.
- Ranwez, V., Criscuolo, A., & Douzery, E. J. P. (2010). SuperTriplets: a triplet-based supertree approach to phylogenomics. *Bioinformatics*, 26(12), i115-i123. doi:10.1093/bioinformatics/btq196
- Ranwez, V., Douzery, E. J. P., Cambon, C., Chantret, N., & Delsuc, F. (2018). MACSE v2: toolkit for the alignment of coding sequences accounting for frameshifts and stop codons. *Molecular Biology and Evolution*, 35(10), 2582–2584. doi:10.1093/molbev/msy159

- Ratnasingham, S., & Hebert, P. D. N. (2007). BOLD: The Barcode of Life Data System (<http://www.barcodinglife.org>). *Molecular Ecology Notes*, 7(3), 355–364.
doi:10.1111/j.1471-8286.2007.01678.x
- Reverter, A., Okimoto, R., Sapp, R., Bottje, W. G., Hawken, R., & Hudson, N. J. (2017). Chicken muscle mitochondrial content appears co-ordinately regulated and is associated with performance phenotypes. *Biology Open*, 6(1), 50–58.
doi:10.1242/bio.022772
- Schomaker-Bastos, A., & Prosdocimi, F. (2018). mitoMaker: a pipeline for automatic assembly and annotation of animal mitochondria using raw NGS data.
doi:10.20944/preprints201808.0423.v1
- Seixas, F. A., Boursot, P., & Melo-Ferreira, J. (2018). The genomic impact of historical hybridization with massive mitochondrial DNA introgression. *Genome Biology*, 19(1), 91. doi:10.1186/s13059-018-1471-8
- Smith, B. T., Harvey, M. G., Faircloth, B. C., Glenn, T. C., & Brumfield, R. T. (2014). Target capture and massively parallel sequencing of ultraconserved elements for comparative studies at shallow evolutionary time scales. *Systematic Biology*, 63(1), 83–95.
doi:10.1093/sysbio/syt061
- Starrett, J., Derkarabetian, S., Hedin, M., Bryson, R. W., McCormack, J. E., & Faircloth, B. C. (2017). High phylogenetic utility of an ultraconserved element probe set designed for Arachnida. *Molecular Ecology Resources*, 17(4), 812–823. doi:10.1111/1755-0998.12621
- Ströher, P. R., Zarza, E., Tsai, W. L. E., McCormack, J. E., Feitosa, R. M., & Pie, M. R. (2017). The mitochondrial genome of *Octostruma stenognatha* and its phylogenetic implications. *Insectes Sociaux*, 64(1), 149–154. doi:10.1007/s00040-016-0525-8

- Tavaré, S. (1986). Some probabilistic and statistical problems in the analysis of DNA sequences. *Lectures on Mathematics in the Life Sciences*, 17(2), 57–86.
- Vieira, G. A., & Prosdocimi, F. (2019). Accessible molecular phylogenomics at no cost: obtaining 14 new mitogenomes for the ant subfamily Pseudomyrmecinae from public data. *PeerJ*, 7, e6271. doi:10.7717/peerj.6271
- Vogler, A. P., & Pearson, D. L. (1996). A molecular phylogeny of the tiger beetles (Cicindelidae): congruence of mitochondrial and nuclear rDNA data sets. *Molecular Phylogenetics and Evolution*, 6(3), 321-338.
- Vollmers, J., Wiegand, S., & Kaster, A.-K. (2017). Comparing and Evaluating Metagenome Assembly Tools from a Microbiologist's Perspective - Not Only Size Matters! *PloS One*, 12(1), e0169662. doi:10.1371/journal.pone.0169662
- Wang, N., Hosner, P. A., Liang, B., Braun, E. L., & Kimball, R. T. (2017). Historical relationships of three enigmatic phasianid genera (Aves: Galliformes) inferred using phylogenomic and mitogenomic data. *Molecular Phylogenetics and Evolution*, 109, 217–225. doi:10.1016/J.YMPEV.2017.01.006
- Ward, P. S. (2014). The Phylogeny and Evolution of Ants. *Annual Review of Ecology, Evolution, and Systematics*, 45(1), 23–43. doi:10.1146/annurev-ecolsys-120213-091824
- Ward, P. S., & Branstetter, M. G. (2017). The acacia ants revisited: convergent evolution and biogeographic context in an iconic ant/plant mutualism. *Proceedings of the Royal Society B: Biological Sciences*, 284(1850), 20162569. doi:10.1098/rspb.2016.2569
- Wenseleers, T., Ito, F., Van Borm, S., Huybrechts, R., Volckaert, F., & Billen, J. (1998). Widespread occurrence of the microorganism Wolbachia in ants. *Proceedings of the Royal Society B: Biological Sciences*, 265(1404), 1447–1452. doi:10.1098/rspb.1998.0456

Young, A. D., Lemmon, A. R., Skevington, J. H., Mengual, X., Ståhls, G., Reemer, M., ...

Wiegmann, B. M. (2016). Anchored enrichment dataset for true flies (order Diptera) reveals insights into the phylogeny of flower flies (family Syrphidae). *BMC Evolutionary Biology*, 16(1), 143. doi:10.1186/s12862-016-0714-0

Zarza, E., Faircloth, B. C., Tsai, W. L. E., Bryson, R. W., Klicka, J., & McCormack, J. E. (2016). Hidden histories of gene flow in highland birds revealed with genomic markers. *Molecular Ecology*, 25(20), 5144–5157. doi:10.1111/mec.13813

Zarza, E., Connors, E. M., Maley, J. M., Tsai, W. L., Heimes, P., Kaplan, M., & McCormack, J. E. (2018). Combining ultraconserved elements and mtDNA data to uncover lineage diversity in a Mexican highland frog (Sarcohyala; Hylidae). *PeerJ*, 6, e6045.

Data accessibility

The MitoFinder software is available from XXXXXX. Annotated mitogenomes and partial mitogenomic contigs have been deposited in GenBank (Accession Numbers XXXXX-XXXXX). The full analytical pipeline, phylogenetic datasets and corresponding trees can be retrieved from zenodo.org (DOI:10.5281/zenodo.3231390).

Authors' contributions

RA and FD conceived the ideas and designed methodology, analysed the data, and led the writing of the manuscript; RA implemented the MitoFinder software in part using code previously written by AS-B; JR, FP, and BN contributed to the writing of the manuscript. All authors contributed critically to the drafts and gave final approval for publication.

Supporting information

Appendix S1. List of the 501 UCE libraries (SRA accessions) and associated metadata.

Appendix S2. Summary statistics on mitochondrial signal recovered per species and depending on the assembler used. The table provides the number of contigs and genes recovered with MitoFinder and the size of each annotated gene.

Appendix S3. Summary statistics of barcoding analyses. Detailed results for both BOLDsystem and Megablast analyses are provided for each COX1 recovered with MitoFinder using MetaSPAdes.

Appendix S4. Detailed results of tree distance analyses realized with Dquad (Ranwez, Criscuolo, & Douzery 2010). Trees obtained with each assembler with mitochondrial amino acid supermatrix, mitochondrial nucleotide supermatrix, and UCE nucleotide supermatrix were compared with each others.