

Additional file

What is the phylogenetic signal limit from mitogenomes? The reconciliation between mitochondrial and nuclear data in the Insecta Class phylogeny

Gerard Talavera & Roger Vila

Table S1 - List of mitochondrial genomes used in this study

Organism	GenBank Code	Taxonomical group
Holometabola		
<i>Drosophila simulans</i>	NC_005781.1	Diptera
<i>Drosophila sechellia</i>	NC_005780.1	Diptera
<i>Drosophila mauritiana</i>	NC_005779.1	Diptera
<i>Drosophila melanogaster</i>	NC_001709.1	Diptera
<i>Drosophila yakuba</i>	NC_001322.1	Diptera
<i>Chrysomya putoria</i>	AF352790.1	Diptera
<i>Cochliomyia hominivorax</i>	NC_002660.1	Diptera
<i>Haematobia irritans</i>	NC_007102.1	Diptera
<i>Dermatobia hominis</i>	NC_006378.1	Diptera
<i>Bactrocera dorsalis</i>	NC_008748.1	Diptera
<i>Bactrocera oleae</i>	NC_005333.1	Diptera
<i>Ceratitis capitata</i>	NC_000857.1	Diptera
<i>Simosyrphus grandicornis</i>	NC_008754.1	Diptera
<i>Tricophtalma punctata</i>	NC_008755.1	Diptera
<i>Cydistomyia duplonata</i>	NC_008756.1	Diptera
<i>Anopheles gambiae</i>	NC_002084.1	Diptera
<i>Anopheles quadrimaculatus</i>	NC_000875.1	Diptera
<i>Aedes albopictus</i>	NC_006817.1	Diptera
<i>Bombyx mandarina</i>	NC_003395.1	Lepidoptera
<i>Bombyx mori</i>	NC_002355.1	Lepidoptera
<i>Antheraea pernyi</i>	AY242996.1	Lepidoptera
<i>Coreana raphaelis</i>	DQ102703.1	Lepidoptera
<i>Ostrinia furnacalis</i>	NC_003368.1	Lepidoptera
<i>Ostrinia nubilalis</i>	NC_003367.1	Lepidoptera
<i>Adoxophyes honmai</i>	DQ073916.1	Lepidoptera
<i>Bombus ignitus</i>	DQ870926.1	Hymenoptera
<i>Melipona bicolor</i>	NC_004529.1	Hymenoptera
<i>Apis mellifera</i>	NC_001566.1	Hymenoptera
<i>Vanhornia eucnemidarum</i>	NC_008323.1	Hymenoptera
<i>Primeuchroeus sp.</i>	DQ302102.1 DQ302101.1	Hymenoptera
<i>Perga condei</i>	AY787816.1	Hymenoptera
<i>Anoplophora glabripennis</i>	NC_008221	Coleoptera
<i>Crioceris duodecimpunctata</i>	NC_003372.1	Coleoptera
<i>Tribolium castaneum</i>	NC_003081.1	Coleoptera
<i>Rhagophthalmus lufengensis</i>	DQ888607.1	Coleoptera
<i>Rhagophthalmus ohbai</i>	AB267275.1	Coleoptera
<i>Pyrocoelia rufa</i>	NC_003970.1	Coleoptera

<i>Xenos vesparum</i>	DQ364229.1	Strepsiptera
Paraneoptera		
<i>Neomaskellia andropogonis</i>	NC 006159.1	Hemiptera
<i>Vasdavidius concursus</i>	AY648941.2	Hemiptera
<i>Aleurochiton aceris</i>	NC 006160.1	Hemiptera
<i>Bemisia tabaci</i>	NC 006279.1	Hemiptera
<i>Tetraleurodes acaciae</i>	NC 006292.1	Hemiptera
<i>Trialeurodes vaporariorum</i>	NC 006280.1	Hemiptera
<i>Aleurodicus dugesii</i>	NC 005939.1	Hemiptera
<i>Daktulosphaira vitifoliae</i>	DQ021446.1	Hemiptera
<i>Schizaphis graminum</i>	NC 006158.1	Hemiptera
<i>Pachypsylla venusta</i>	AY278317.1	Hemiptera
<i>Philaenus spumarius</i>	AY630340.1	Hemiptera
<i>Homalodisca coagulata</i>	AY875213.1	Hemiptera
<i>Triatoma dimidiata</i>	NC 002609.1	Hemiptera
<i>Heterodoxus macropus</i>	NC 002651.1	Phthiraptera
<i>Campanulotes bidentatus</i>	NC 007884.1	Phthiraptera
<i>Thrips imaginis</i>	NC 004371.1	Thysanoptera
<i>Lepidopsocid RS-2001</i>	NC 004816.1	Psocoptera

Table S2 - Number of characters in the final alignments for each phylogenetic reconstruction method tested. Resulting trees shown in figures 2-4 are indicated.

	Paraneoptera	Holometabola	Eumetabola
ML (Protein)	2731 (<i>Fig.3</i>)	3501 (<i>Fig.2</i>)	3288
BI (Protein)	2731	3501	3288
BI (DNA 1st and 2nd position)	7010	7368	7232
BI (DNA 1st and 2nd position + RNA)	9536 (<i>Fig.3</i>)	10202 (<i>Fig.2</i>)	9548
BI (DNA + RNA) - Site specific rate model	13068	13889	-
BI (DNA 1st and 2nd position) - CAT model	7010	7368	7232
BI (Protein) - CAT model	2731 (<i>Fig.3</i>)	3501 (<i>Fig.2</i>)	3288 (<i>Fig.4</i>)

Simulations methods

We compared the efficiency of the protein-based phylogenetic reconstruction strategies using simulations. We performed simulated protein alignments with 1000 amino acid positions conducted along a reference phylogenetic tree with eight tips, with a global divergence equivalent to the Holometabola BI-DNA tree. In order to imitate possible LBA effects, two unrelated branches were forced to be six times longer than the average length of the rest (in the Holometabola BI-DNA tree the longest branches were four times bigger). One hundred simulations were performed using Seq-Gen [1] with six categories of rate heterogeneity ($\alpha = 0.872$) and the MtRev evolutionary model. From these simulations, maximum likelihood trees with six categories of rate heterogeneity were

inferred with Phym1 2.4.4, Bayesian inference with Mr.Bayes 3.1.2 and MtRev model, and PhyloBayes 2.3 under the CAT model. After that, we calculated the scale-factor, a relative value for comparing branch lengths between two trees, and the Robinson-Foulds distance, which calculates the topological differences between two trees, using Ktreedist 1.0 software [2] from each resulting tree versus the reference tree used to conduct the simulations.

Simulations results and discussion

Simulations were performed to confirm the ability of the CAT model to suppress the LBA bias compared to ML-AA and BI-DNA. We created two unrelated long branches in a tree with eight terminals with a similar divergence to the Holometabola dataset, and also exaggerated this divergence 2 and 3 times. The general tendency of the simulation test results was the same than the one observed in real data. For the three divergence levels explored, BI-AA-CAT produced the lowest percentage of trees grouping the two long branches as sister taxa. In divergence x1, BI-CAT did not group the long branches in any of the simulations (0%), while we obtained a 9% for ML-AA and a 10% for BI-DNA. For divergence x2, long branches were grouped together in 3% of the cases for BI-AA-CAT, 12% for ML-AA and 26% for BI-DNA. Finally for divergence x3, the values increased to 10% for BI-AA-CAT, 34% for ML-AA and 38 % for BI-DNA. These percentages do not evaluate intermediate LBA effects, where both branches might be closer than they should, but not strictly sisters. To evaluate the topological differences between the simulated trees and the reference topology, presumably a product of LBA, we calculated Robinson-Foulds distances and calculated the average of the 100 simulations for each divergence type and method. This revealed again the better performance of BI-AA-CAT, which obtained the lowest values, followed by ML-AA and BI-DNA respectively (Figure S1). A better performance of the amino acid sequences versus DNA was also reflected in these results, and their use together with a site-heterogeneous mixture model under a Bayesian framework is the suggested combination to avoid LBA artefacts.

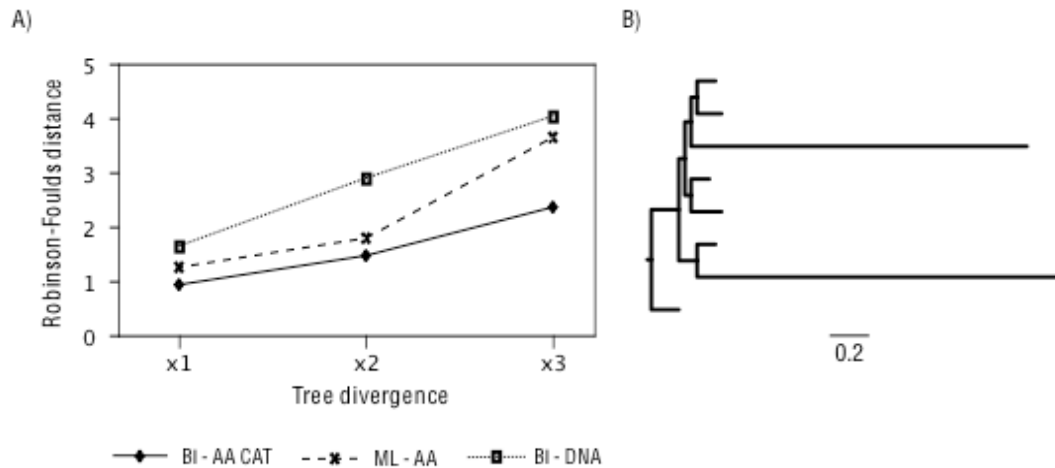


Figure S1 - Simulations

A) Average Robinson-Foulds distances relative to the reference tree calculated with ML - protein sequences (dotted line with squared symbols), BI - DNA excluding third codon positions (dashed line with cross symbols) and BI - protein with CAT model (solid line with diamonds) for three different tree divergences. B) Reference tree (divergence x1) used to conduct simulations. Scale bar represents 0.2 substitutions/site.

References:

1. Rambaut A, Grassly NC: **Seq-Gen: an application for the Monte Carlo simulation of DNA sequence evolution along phylogenetic trees.** *Comput Appl Biosci* 1997, **13**:235-238
2. Soria-Carrasco V, Talavera G, Igea J, Castresana J: **The K tree score: quantification of differences in the relative branch length and topology of phylogenetic trees.** *Bioinformatics* 2007, **23**:2954-2956