

## Additional molecular support for the new chordate phylogeny.

Frédéric Delsuc, Georgia Tsagkogeorga, Nicolas Lartillot, Hervé Philippe

► **To cite this version:**

Frédéric Delsuc, Georgia Tsagkogeorga, Nicolas Lartillot, Hervé Philippe. Additional molecular support for the new chordate phylogeny.. *Genesis*, Wiley-Blackwell, 2008, 46 (11), pp.592-604. <10.1002/dvg.20450>. <halsde-00338411>

**HAL Id: halsde-00338411**

**<https://hal.archives-ouvertes.fr/halsde-00338411>**

Submitted on 13 Nov 2008

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

# **Additional Molecular Support for the New Chordate Phylogeny**

**Frédéric Delsuc<sup>1,2\*</sup>, Georgia Tsagkogeorga<sup>1,2</sup>, Nicolas Lartillot<sup>3</sup> and Hervé Philippe<sup>4</sup>**

<sup>1</sup>Université Montpellier II, CC064, Place Eugène Bataillon, 34095 Montpellier Cedex 5, France

<sup>2</sup>CNRS, Institut des Sciences de l'Evolution (UMR5554), CC064, Place Eugène Bataillon, 34095 Montpellier Cedex 5, France

<sup>3</sup>Laboratoire d'Informatique, de Robotique et de Microélectronique de Montpellier, CNRS-Université Montpellier II, 161 Rue Ada, 34392 Montpellier Cedex 5, France

<sup>4</sup>Département de Biochimie, Université de Montréal, Succursale Centre-Ville, Montréal, Québec, Canada H3C 3J7

\*Correspondence to: Frédéric Delsuc, CC064, Institut des Sciences de l'Evolution, UMR5554-CNRS, Université Montpellier II, Place Eugène Bataillon, 34095 Montpellier Cedex 5, France. Email: [Frederic.Delsuc@univ-montp2.fr](mailto:Frederic.Delsuc@univ-montp2.fr).

Contract grant sponsor: Research Networks Program in Bioinformatics from the High Council for Scientific and Technological Cooperation between France and Israel.

**SUMMARY**

**Recent phylogenomic analyses have suggested tunicates instead of cephalochordates as the closest living relatives of vertebrates. In direct contradiction with the long accepted view of Euchordates, this new phylogenetic hypothesis for chordate evolution has been the object of some scepticism. We assembled an expanded phylogenomic dataset focused on deuterostomes. Maximum-likelihood using standard models and Bayesian phylogenetic analyses using the CAT site-heterogeneous mixture model of amino-acid replacement both provided unequivocal support for the sister-group relationship between tunicates and vertebrates (Olfactores). Chordates were recovered as monophyletic with cephalochordates as the most basal lineage. These results were robust to both gene sampling and missing data. New analyses of ribosomal rRNA also recovered Olfactores when compositional bias was alleviated. Despite the inclusion of 25 taxa representing all major lineages, the monophyly of deuterostomes remained poorly supported. The implications of these phylogenetic results for interpreting chordate evolution are discussed in light of recent advances from evolutionary developmental biology and genomics.**

**Key words:** Phylogenomics – Deuterostomes – Chordates – Tunicates – Cephalochordates – Olfactores – Ribosomal RNA – Jackknife – Evolution.

## INTRODUCTION

Besides its fundamental role in systematics, phylogenetic reconstruction is a prerequisite for understanding the evolution of organisms. The essential contribution of phylogenetics for understanding morphological diversity has perhaps been best exemplified in the case of animal evolution (Telford and Budd, 2003). The Cambrian explosion has produced a bewildering diversity of body plans whose origins and evolution can only be apprehended by undertaking an integrative approach through evolutionary developmental biology (Evo-Devo) (Conway-Morris, 2003). The knowledge of phylogenetic relationships, by allowing the polarisation of character transformations, sheds light on the extent of morphological convergence and reversal. A phylogenetic framework is therefore required for distinguishing ancestral characters from those representing morphological innovations. Comparative genomics is now providing the opportunity to track these morphological innovations back to the molecular level by revealing the patterns of gene acquisition/loss, and giving clues to the molecular adaptations that underline the evolution of body plans (Cañestro *et al.*, 2007).

Animal taxonomy has deep roots. The study of morphological and embryological characters has allowed the definition of the major phyla, but left their interrelationships almost unresolved (Nielsen, 2001). The advent of molecular data during the 1990s has revolutionized the traditional classification through a series of phylogenetic analyses of the 18S ribosomal RNA (rRNA) gene for an ever increasing number of key taxa (Aguinaldo *et al.*, 1997; Halanych *et al.*, 1995). This period culminated with the proposition of a new view of animal phylogeny at odds with the traditional paradigm of a steady increase towards morphological complexity, and revealing instead the major role played by secondary simplification from complex ancestors (Adoutte *et al.*, 2000; Lwoff, 1944). Despite these undeniable achievements, the resolving power provided by 18S rRNA and other single genes is nevertheless limited and a number of open questions in animal phylogeny remained to be answered (Halanych, 2004).

The most recent advances in animal phylogeny have come from phylogenomics (Delsuc *et al.*, 2005) which considerably increases the resolving power by considering numerous concatenated genes from Expressed Sequence Tags (ESTs) and complete genome projects (Philippe and Telford, 2006). Despite some troubled beginnings due to the shortcomings of using only a restricted set of taxa (Philippe *et al.*, 2005a), phylogenomics has provided strong corroborating support for the new animal phylogeny, essentially confirming the monophyly of Protostomia, Ecdysozoa and Lophotrochozoa (Baurain *et al.*, 2007; Dunn *et al.*, 2008; Lartillot and Philippe, 2008; Philippe *et al.*, 2005b). Phylogenomics has also helped solving some longstanding mysteries such as the position of chaetognaths which finally appear to

belong to Protostomia (Marlétaz *et al.*, 2006b; Matus *et al.*, 2006) and also proposed unexpected phylogenetic affinities for enigmatic taxa such as *Buddenbrockia plumatellae* recently unmasked as a cnidarian worm (Jimenez-Guri *et al.*, 2007), or *Xenoturbella bocki*, representing a fourth deuterostome phylum on its own (Bourlat *et al.*, 2006).

Among the most groundbreaking results from recent phylogenomic studies was the identification of tunicates (or urochordates) as the closest living relatives of vertebrates, instead of cephalochordates as traditionally accepted (Delsuc *et al.*, 2006). Some hints of this unexpected result had been observed in previous large-scale phylogenetic studies including a single tunicate representative (Blair and Hedges, 2005; Philippe *et al.*, 2005b; Vienne and Pontarotti, 2006). However, a substantial increase in taxon sampling turned out to be required for recovering convincing support in favour of such an unorthodox relationship. In particular, the fact that the inclusion of the divergent appendicularian tunicate *Oikopleura dioica* did not disrupt the sister-group relationship between tunicate and vertebrates gave a good indication about the strength of the phylogenetic signal in its favour (Delsuc *et al.*, 2006). The grouping of tunicates and vertebrates had already been proposed on morphological grounds by Richard P.S. Jefferies who coined the name Olfactores after the presence a putatively homologous olfactory apparatus in fossils that were proposed to be precursors of tunicates and vertebrates (Jefferies, 1991). This phylogenetic result has nevertheless been the object of some scepticism. One reason for this maybe that it further invalidates the traditional textbook view of chordate evolution as a steady increase towards morphological complexity culminating with vertebrates, as betrayed by the use of the term Euchordates (literally “true chordates”) for denoting the grouping of cephalochordates and vertebrates (Gee, 2001). The lack of obvious morphological synapomorphies for Olfactores, apart from the presence of migratory neural crest-like cells (Jeffery, 2007; Jeffery *et al.*, 2004), and the apparent conflict with analyses of rRNA data which tend to favour Euchordates (Cameron *et al.*, 2000; Mallatt and Winchell, 2007; Winchell *et al.*, 2002) might also partly explain the caution with which this result has been considered at first.

Phylogenomics, despite being a powerful approach, is not immune to potential reconstruction artefacts however. The possible pitfalls associated with phylogenomic studies include systematic errors that can be traced back to some kind of model misspecifications (Philippe *et al.*, 2005a) and caused mainly by heterogeneity of evolutionary rates among taxa (Lartillot *et al.*, 2007; Philippe *et al.*, 2005b) and compositional bias (Blanquart and Lartillot, 2008; Jeffroy *et al.*, 2006; Lartillot and Philippe, 2008; Phillips *et al.*, 2004). Empirical protocols have been designed to detect and reduce the impact of systematic error in genome-scale studies (Rodríguez-Ezpeleta *et al.*, 2007) but the ultimate solution lies in the

development of improved models of sequence evolution (Felsenstein, 2004; Philippe *et al.*, 2005a; Steel, 2005). The reliance of current phylogenomic studies on a relatively limited number of highly expressed genes (Philippe and Telford, 2006) and the potential impact of missing data on phylogenomic inference (Hartmann and Vision, 2008; Philippe *et al.*, 2004) are also regularly cited as limitations of the phylogenomic approach.

The aim of this paper is to evaluate the current evidence for the new chordate phylogeny by: (1) reanalyzing previous phylogenomic data using improved models of amino-acid replacement, (2) assembling and analyzing an updated phylogenomic dataset with more genes and more taxa, (3) assessing the impact of missing data and gene sampling on phylogenomic results, and (4) performing new analyses of rRNA data taking compositional bias into account.

## MATERIALS AND METHODS

### Phylogenomic Dataset Assembly

We built upon previous phylogenomic datasets assembled in the Philippe lab (Delsuc *et al.*, 2006; Jimenez-Guri *et al.*, 2007; Lartillot and Philippe, 2008; Philippe *et al.*, 2005b; Philippe *et al.*, 2004) to select a set of 179 orthologous markers showing sufficient conservation across metazoans to be useful for inferring the phylogeny of metazoans. Alignments were built and updated with available sequences downloaded from the Trace Archive (<http://www.ncbi.nlm.nih.gov/Traces/>) and the EST Database (<http://www.ncbi.nlm.nih.gov/dbEST/>) of GenBank at the National Center for Biotechnology Information (<http://www.ncbi.nlm.nih.gov/>) using the program ED from the MUST package (Philippe, 1993). Unambiguously aligned regions were identified and excluded for each individual gene using the program GBLOCKS (Castresana, 2000) with a few manual refinements using NET from the MUST package. The complete list of genes with corresponding final numbers of amino-acid sites is available as Supplementary Material.

The concatenation of the 179 genes was constructed with the program SCAFOS (Roure *et al.*, 2007) by defining 51 metazoan operational taxonomic units (OTUs) including 25 deuterostomes representing all major lineages. When several sequences were available for a given OTU, the slowest evolving one was selected according to their degree of divergence using ML distances computed by TREE-PUZZLE (Schmidt *et al.*, 2002) under a WAG+F model (Whelan and Goldman, 2001) within SCAFOS. The percentage of missing data per taxon was reduced by creating some chimerical sequences for species belonging to the same OTU. The complete alignment consists of 179 genes and 51 taxa for 53,799 unambiguously aligned

amino-acid sites with 32% missing data. In order to study the potential impact of missing data on phylogenetic inference (Hartmann and Vision, 2008; Philippe *et al.*, 2004; Wiens, 2006), a concatenation of the 106 genes with sequences available for at least 41 of the 51 OTUs was also constructed with SCAFOS. This reduced alignment consists of 106 genes and 51 taxa for 25,321 amino-acid sites and contains only 20% of missing data. The list of defined OTUs, chimerical sequences and percentages of missing data are available as Supplementary Material. Individual gene alignments and their concatenations are available upon request.

### **Phylogenomic Analyses**

Bayesian phylogenetic analyses of the two phylogenomic datasets were conducted using the program PHYLOBAYES 2.3c (<http://www.atgc-montpellier.fr/phylobayes/>) under the CAT+ $\Gamma_4$  site-heterogeneous mixture model (Lartillot and Philippe, 2004). For each dataset, four independent Monte Carlo Markov Chains (MCMCs) starting from a random topology were run in parallel for 20,000 cycles (1,500,000 generations), saving a point every cycle, and discarding the first 2000 points as the burnin. Bayesian posterior probabilities (PP) were obtained from the 50% majority-rule consensus tree of the 18,000 MCMC sampled trees using the program READPB of PHYLOBAYES.

Maximum likelihood (ML) reconstruction on the new phylogenomic dataset was also performed using the program TREEFINDER version of March 2008 (Jobb *et al.*, 2004) under the empirical WAG+F+ $\Gamma_8$  model of amino-acid substitution. The  $\alpha$  shape parameter of the  $\Gamma$  distribution was estimated along with the topology and the branch lengths. Reliability of nodes was estimated by bootstrap resampling with 100 pseudo-replicate datasets generated by the program SEQBOOT of the PHYLIP package (Felsenstein, 2001). The 100 corresponding ML heuristic searches were run in parallel on a computing cluster and the majority-rule consensus of the 100 resulting trees was computed using TREEFINDER.

### **Jackknife Procedure**

The robustness of our phylogenomic inference with respect to gene sampling was assessed by applying a jackknife procedure. Fifty jackknife replicates of 50 genes drawn randomly from the full pool of 179 genes were generated. The only condition we imposed to this jackknife procedure was to require that each taxon is represented by at least one gene in each replicate. The 50 jackknife supermatrices ranging from 11,163 to 17,181 amino-acid sites with 27 to 35% missing data were then analyzed using PHYLOBAYES under the CAT+ $\Gamma_4$  model. To ensure correct convergence of the MCMC on each replicate, an automated stopping rule was used. Specifically, for each jackknife replicate, two independent parallel (and

synchronous) MCMC were run, until the posterior probability discrepancy between the two chains was less than 0.15 (maximum discrepancy over all bipartitions), and after removing the first 1000 sampled trees of each chain as the burnin. A global majority-rule consensus tree was obtained from the 50 replicates as follows: for each jackknife replicate ( $D_r$ ) taken in turn, we computed the frequency-table of all bipartitions (splits) observed in the sample collected from the posterior distribution  $p(T|D_r)$ . The frequencies associated to each bipartition were then averaged over the 50 replicates, and the resulting frequency table was used to build a consensus tree. The support values displayed by this Bayesian consensus tree are thus jackknife-resampled posterior probabilities ( $PP_{JK}$ ). High  $PP_{JK}$  values indicate nodes that have high probability support in most jackknife replicates.

### **Phylogenetic Analyses of Ribosomal RNA**

The 46-taxa dataset of combined 18S+28S rRNA genes assembled by Mallatt and Winchell (Mallatt and Winchell, 2007) for studying deuterostome phylogeny was reanalyzed. This alignment contains a total of 3,925 unambiguously aligned nucleotide sites. A principal component analysis (PCA) of nucleotide composition was realized using the R statistical package (Team, 2007). The best fitting model of nucleotide sequence evolution was evaluated using MODELTEST 3.7 (Posada and Crandall, 1998). The TIM+ $\Gamma_4$ +I transitional model (Posada and Crandall, 2001) was selected according to the Akaike information criterion. ML phylogenetic analysis of this nucleotide dataset was conducted with PAUP\* 4.0b10 (Swofford, 2002) using a heuristic search with Tree Bisection Reconnection (TBR) branch swapping starting from a Neighbor-Joining (NJ) tree.

The nucleotide dataset was RY-coded by pooling puRines ( $AG \Rightarrow R$ ) and pYrimidines ( $CT \Rightarrow Y$ ) in an attempt to alleviate both compositional heterogeneity and substitutional saturation of transition events. This RY-coded dataset was then analyzed by conducting a ML heuristic search with TBR branch swapping on a NJ starting tree using PAUP\* under the CF+ $\Gamma_8$  model for discrete characters (Cavender and Felsenstein, 1987). The  $\alpha$  shape parameter of the  $\Gamma$  distribution was previously estimated during a ML heuristic search on the nucleotide dataset conducted with TREEFINDER under the GTR2+ $\Gamma_8$  two-state model.

Reliability of nodes was estimated for each dataset by non-parametric bootstrap resampling using 100 pseudo-replicates generated by SEQBOOT. The 100 corresponding ML heuristic searches using PAUP\* with the previously estimated ML parameters, NJ starting trees, and TBR branch swapping were parallelized on a computing cluster. ML bootstrap percentages were obtained from the 50% majority-rule consensus tree of the 100 bootstrap ML trees using TREEFINDER.



## RESULTS AND DISCUSSION

### Effect of an Improved Model of Sequence Evolution

Our initial assessment of deuterostome phylogenetic relationships was based on a phylogenomic dataset encompassing 146 nuclear genes (33,800 amino-acids) from 38 taxa including 14 deuterostomes (Delsuc *et al.*, 2006). ML and Bayesian phylogenetic analyses conducted under the standard WAG+F+ $\Gamma_4$  model provided strong support for grouping tunicates with vertebrates (including cyclostomes), but also disrupted chordate monophyly because cephalochordates grouped with echinoderms, albeit with non-significant statistical support (Delsuc *et al.*, 2006). The limited taxon sampling available at the time for Ambulacraria (echinoderms and hemichordates), i.e. a single echinoderm, prompted us to be cautious about this result and to call for the inclusion of xenoturbellidans, hemichordates, and more echinoderms before drawing definitive conclusions. In fact, a subsequent phylogenomic study did exactly what we pleaded for by adding a representative species for each of these three groups (Bourlat *et al.*, 2006). The inclusion of these taxa allowed retrieving the monophyly of chordates in Bayesian analyses using standard amino-acid models, although the alternative hypothesis of chordate paraphyly was still not statistically rejected by ML non-parametric tests (Bourlat *et al.*, 2006). Importantly, the strong statistical support for the monophyly of Olfactores was unaffected by taxon addition (Bourlat *et al.*, 2006).

Models accounting for site-specific modulations of the amino-acid replacement process, such as the CAT mixture model (Lartillot and Philippe, 2004), seem to offer a significantly better fit to real data than empirical substitution matrices currently used in standard models of amino-acid sequence evolution. Accounting for site-specific amino-acid propensities has also been shown to induce a significant improvement of phylogenetic reconstruction in difficult cases such as long-branch attraction (Baurain *et al.*, 2007; Lartillot *et al.*, 2007; Lartillot and Philippe, 2008). This improvement essentially lays in the ability of the CAT model to detect multiple conservative substitutions more efficiently than standard amino-acid models (Lartillot *et al.*, 2007).

In order to test for an eventual effect of model misspecification on previous phylogenomic analyses, we reanalyzed our previous dataset (Delsuc *et al.*, 2006) under the CAT+ $\Gamma_4$  model. This analysis provides strong corroborating support for the grouping of tunicates and vertebrates (Fig. 1). However, in contrast with previous analyses using empirical amino-acid replacement matrices, which favoured a sister-group relationship between cephalochordates and echinoderms, the use of the CAT+ $\Gamma_4$  mixture model strongly supports the classical view of monophyletic chordates and deuterostomes (Fig. 1). The fact that chordate polyphyly is

disrupted both by a richer taxon sampling (Bourlat *et al.*, 2006), or upon the use of a more elaborate model, suggests that the previously observed grouping of cephalochordates and echinoderms (Delsuc *et al.*, 2006) was probably a phylogenetic reconstruction artefact. On the other hand, the fact that the grouping of tunicates and vertebrates is insensitive to the model used, adds further credence to the Olfactores hypothesis.

### An Updated Phylogenomic Dataset

The continuously growing genomic databases allowed us to build an updated phylogenomic dataset that includes both more genes and more taxa than previously considered to address the question of deuterostome phylogeny. This new dataset of 179 genes for 51 taxa includes 25 deuterostomes representing all major lineages: Xenoturbellida (1 taxon), Hemichordata (1), Echinodermata (5), Cephalochordata (1), Tunicata (6), Cyclostomata (2) and Vertebrata (9), plus 26 selected slow evolving metazoan taxa including Cnidarians and Poriferans as the most distant outgroups. Chordates are particularly well sampled with the inclusion, for the first time, of six tunicate species covering the four major clades evidenced by 18S rRNA studies (Swalla *et al.*, 2000). This diverse taxon sampling is essential to further test the new chordate phylogeny recently revealed by phylogenomics (Bourlat *et al.*, 2006; Delsuc *et al.*, 2006).

Bayesian (CAT+ $\Gamma_4$ ) and ML (WAG+F+ $\Gamma_8$ ) phylogenetic reconstructions conducted on this updated dataset (179 genes, 53,799 amino-acid sites, 51 taxa) resulted in a highly resolved tree (Fig. 2a). These analyses provided strong support for Ambulacraria (PP<sub>CAT+ $\Gamma_4$</sub>  = 1.0 / BP<sub>WAG+F+ $\Gamma_8$</sub>  = 97), chordates (1.0 / 69) and olfactores (1.0 / 100). Xenambulacraria (*Xenoturbella* + Ambulacraria) and a basal position for the chaetognath *Spadella* among Protostomia were also moderately supported by our analyses (Fig. 2a). These results are compatible with a recent phylogenomic analysis which also found strong support for Ambulacraria, chordates and Olfactores when using the CAT mixture model (Dunn *et al.*, 2008). However, the monophyly of Deuterostomes is unresolved in both Bayesian and ML phylogenetic reconstructions (Fig. 2a).

The complete dataset obtained by concatenating all 179 genes contains 32% missing data. Previous studies of the impact of missing data on the accuracy of phylogenetic inference have concluded that probabilistic methods are relatively tolerant to missing data (Hartmann and Vision, 2008; Philippe *et al.*, 2004; Wiens, 2003, , 2005), the most important factor being the absolute amount of available data for a given taxon. In phylogenomics, even incomplete taxa are usually represented by thousand of sites and the impact of missing data on accuracy is therefore relatively limited (Philippe *et al.*, 2004). Nonetheless, incomplete taxa might still be

difficult to place with confidence especially when they represent isolated lineages such as *Xenoturbella* (65% missing data) and *Spadella* (75%) in our dataset. In order to control for a potential effect of missing data on our phylogenomic results, we restricted our dataset to the 106 genes with sequences available for at least 41 of the 51 taxa. The concatenation of these 106 genes produces a matrix with 25,321 amino-acid sites that contains only 20% of missing data.

Bayesian and ML phylogenetic inference on this reduced dataset produced a tree fully congruent with the phylogenetic picture given by the complete dataset (Fig. 2b). In particular, the support for the monophyly of chordates was still maximal in terms of  $PP_{CAT+\Gamma_4}$ , but  $BP_{WAG+F+\Gamma_8}$  increased from 69 to 88%. The monophyly of Olfactores received maximal support in both cases and appeared not affected by missing data. Statistical support in terms of  $PP_{CAT+\Gamma_4}$  and  $BP_{WAG+F+\Gamma_8}$  was generally increased especially for locating incomplete taxa such as the *Xenoturbella* as the sister-group to Ambulacraria (from 0.98 / 69 to 1.0 / 88) and *Spadella* at the base of Protostomia (from 0.80 / 56 to 1.0 / 80). Altogether, reducing the amount of missing data, despite also reducing the total number of available sites, seems to result in a slight increase in bootstrap proportions. The only real noticeable difference between the two models concerns the monophyly of deuterostomes. The Bayesian inference under the CAT mixture model suggests deuterostome paraphyly by supporting a basal position of chordates within Bilateria (Fig. 2b) as previously reported (Lartillot and Philippe, 2008). However, ML retrieved the monophyly of deuterostomes, but with  $BP_{WAG+F+\Gamma_8}$  of only 50%, leaving the monophyly of deuterostomes unresolved by our data.

### Robustness of Phylogenomics to Gene Sampling

A legitimate question that can be directed to the phylogenomic approach is the degree to which the results are robust to the sample of genes used to infer phylogenetic trees. This potential concern was addressed by applying a jackknife statistical resampling protocol: fifty datasets were assembled by randomly drawing 50 genes from the total 179 genes, and subjected to Bayesian phylogenetic reconstruction using the  $CAT+\Gamma_4$  mixture model (see Methods). The resulting majority-rule consensus tree shows that the vast majority of inferred phylogenetic relationships are highly repeatable across the 50 jackknife replicates (Fig. 3). Olfactores, Chordata, and Ambulacraria all received  $PP_{JK}$  of more than 90% indicating that phylogenetic support is not dependent upon a particular gene combination. Xenambulacraria appears slightly more affected by gene sampling ( $PP_{JK} = 80\%$ ), but this relative instability

might be explained by the poor gene representation available for *Xenoturbella* with only 98 genes over 179 (55%). The same kind of reasoning could apply to the relatively unstable positions of the chaetognath *Spadella* within protostomes (PP<sub>JK</sub> = 62%) and of *Holothuria* within echinoderms (PP<sub>JK</sub> = 51%) (Fig. 3).

In fact, the only major clade whose monophyly appears to be influenced by gene sampling is deuterostomes for which PP<sub>JK</sub> was less than 50% (Fig. 3). In practise, this means that depending on the particular combination of 50 genes considered, deuterostomes might appear either monophyletic or paraphyletic, with the three possible topological alternatives retrieved in almost similar proportions: Deuterostomes (38%), basal chordates (28%), and basal Xenambulacraria (22%). Despite the inclusion of 25 taxa representing all major lineages in our dataset, these results confirm deuterostomes as one of the most difficult groups to resolve in the animal phylogeny despite its wide acceptance (see (Lartillot and Philippe, 2008)).

### **New Analyses of Ribosomal RNA Genes**

The sister-group relationship between tunicates and vertebrates (Olfactores) observed in phylogenomics is in conflict with most (if not all) analyses of rRNA which favour cephalochordates as the closest relatives of vertebrates (Euchordates) (Cameron *et al.*, 2000; Mallatt and Winchell, 2007; Swalla *et al.*, 2000; Wada and Satoh, 1994; Winchell *et al.*, 2002). However, the statistical support for Euchordates in rRNA-based phylogenetic studies is moderate. Indeed, the first 18S rRNA study, based on a limited taxon sampling of deuterostomes, reported a bootstrap value of only 45% for Euchordates (Wada and Satoh, 1994). A subsequent 18S rRNA study considering only slowly evolving sequences for 16 deuterostomes found only a moderate bootstrap support of 71% for grouping cephalochordates and vertebrates (Cameron *et al.*, 2000). A study focused on tunicates also obtained moderate support for Euchordates (58 to 85% depending on the dataset and reconstruction method) but failed to support chordate monophyly likely because tunicate 18S rRNA sequences are rapidly evolving (Swalla *et al.*, 2000).

The next studies used the combination of 18S and 28S rRNAs. An investigation using 28 taxa for the two rRNA subunits found strong bootstrap support (89 to 97% depending on the method) for Euchordates (Winchell *et al.*, 2002). However, this study again failed to support chordate monophyly. Detailed analyses confirmed that tunicate genes have evolved rapidly and showed that they are compositionally biased towards AT, rendering tunicates virtually impossible to locate convincingly in the tree on the basis of rRNA data (Winchell *et al.*, 2002). Finally, increasing the sampling to 46 taxa for this 18S+28S rRNA data did not help in further resolving the relationships among the major groups of deuterostomes and even

decreased the ML bootstrap support for Euchordates from 97% in the previous study to 50% (Mallatt and Winchell, 2007).

In order to gauge the extent to which the rRNA data conflicts with our phylogenomic results, we reanalyzed the 46-taxa dataset of Mallatt and Winchell (Mallatt and Winchell, 2007). The heterogeneity of base composition in this dataset is well illustrated by the PCA presented in Figure 4a. At one extreme, tunicates (especially *Oikopleura*) are particularly AT-rich, and at the other extreme, Myxinidae (*Myxine* and *Eptatretus*) and the pterobranch hemichordate *Cephalodiscus* are highly GC-rich. We therefore compared phylogenetic reconstructions conducted on nucleotides and on RY-coded data, a coding scheme allowing reducing both substitutional saturation and nucleotide compositional bias (Fig. 4b). The two inferred ML trees appear mostly congruent except for two major topological shifts.

The strongest topological change occurred within hemichordates (Fig. 2b). Whereas the use of a standard DNA model strongly supports the paraphyly of enteropneusts by grouping the pterobranch *Cephalodiscus* with *Saccoglossus* and *Harrimania* (BP = 95), RY-coding allows recovering the monophyly of enteropneusts with high bootstrap support (BP = 90). This helps in understanding the conflict between 18S rRNA that supports enteropneust paraphyly (Cameron *et al.*, 2000; Halanych, 1995) and 28S rRNA that rather favours their monophyly (Mallatt and Winchell, 2007; Winchell *et al.*, 2002). This result is of particular importance because it potentially invalidates the controversial hypothesis that pterobranchs evolved from an enteropneust (Cameron *et al.*, 2000; Halanych *et al.*, 1995) by suggesting that it is likely an artefact of 18S rRNA-based phylogenetic reconstructions due to shared nucleotide compositional bias between pterobranchs and Harrimaniidae (Fig. 4a).

Second, the support for the monophyly of Euchordates observed with nucleotides (BP = 84) disappeared in favour of the monophyly of Olfactores in the RY-coding dataset, albeit with no statistical support (BP = 44). This nevertheless strongly suggests that the high composition bias of tunicate sequences has blurred the phylogenetic signal for Olfactores in previous analyses. Thus, according to our interpretation, reducing compositional bias and substitutional saturation by RY-recoding allows recovering a limited signal for Olfactores in agreement with our phylogenomic analysis of amino-acid data. It is worth noting however, that rRNA does not statistically support chordate monophyly in both cases (Fig. 2b).

### **Molecular Phylogenetic Conclusions**

Our aim was to reanalyse the phylogenetic relationships among chordates. The revision of the position of tunicates proposed by recent phylogenomic studies (Bourlat *et al.*, 2006; Delsuc *et al.*, 2006; Dunn *et al.*, 2008) by concluding in favour of the monophyly of

Olfactores, has not yet been considered as totally convincing, essentially because it is at odds with both the traditional view based on embryological and morphological characters (Rowe, 2004; Schaeffer, 1987), and with earlier molecular phylogenetic analyses based on rRNA (Cameron *et al.*, 2000; Mallatt and Winchell, 2007; Swalla *et al.*, 2000; Wada and Satoh, 1994; Winchell *et al.*, 2002). The unexpected sister-group relationship between echinoderms and cephalochordates observed in one of these studies (Delsuc *et al.*, 2006) may also have suggested the possibility that the monophyly of Olfactores was due to an artefactual attraction of cephalochordates with echinoderms (Bourlat *et al.*, 2006).

In the present analysis, we have tried to address these points, essentially by reanalyzing both phylogenomic and rRNA data, under better taxonomic sampling and using more elaborate methods and probabilistic models. First, we demonstrate that, although the grouping of echinoderms and cephalochordates was indeed a probable artefact, disappearing upon the addition of several taxa or using an improved model of sequence evolution, the monophyly of Olfactores appears to be robust with respect to taxon sampling and model choice. Second, our reanalysis of rRNA data using RY-recoding also reveals a weak signal in favour of Olfactores, and suggests that the grouping of vertebrates and cephalochordates in former studies may have been an artefact driven by compositional biases. Altogether, our analyses allows a coherent interpretation of all empirical results observed thus far concerning chordate phylogeny, yielding further evidence in favour of the monophyly of both chordates and Olfactores.

At larger scale, however, we observe an overall lack of support for the monophyly of deuterostomes. Deuterostomes have nearly unanimously been considered as an unquestionable monophyletic group, a hypothesis backed up by traditional comparative analyses of embryological characters such as the fate of the blastopore (Nielsen, 2001), and morphological traits such as gill slits (Schaeffer, 1987). However, in our analyses, the status of deuterostomes seems to be sensitive to the model used, with CAT slightly favouring the paraphyletic configuration, and WAG the more traditional monophyly. In either case, the support measured by non-parametric resampling procedures (site-wise bootstrap or gene-wise jackknife) is weak.

Other phylogenomic studies (Dunn *et al.*, 2008; Lartillot and Philippe, 2008) also failed to obtain strong support for the relative phylogenetic positions of chordates and Ambulacraria. Moderate support for the monophyly of Deuterostomes was only obtained under empirical matrix models, the support disappearing when the CAT model was used instead (Dunn *et al.*, 2008; Lartillot and Philippe, 2008). Although profile mixture models such as CAT, whereas having a better fit than empirical matrices such as WAG, may have some inherent weaknesses

as to their phylogenetic accuracy, the WAG empirical matrix fails in many cases, particularly when confronted to a high level of saturation (Lartillot *et al.*, 2007). This casts doubts on results that seem to receive support exclusively under this model, as is the case for deuterostome monophyly. More observations are needed to better gauge the relative merits of either type of model. Overall, although deuterostome monophyly still remains a reasonable working hypothesis to date, more work is needed before the question can be settled.

### **Corroborative Evidence for Olfactores Monophyly**

The monophyly of Olfactores receives strong support from sequence-based phylogenomic inference. Rare genomic changes has also provided some evidence in its favour: the domain structure of cadherins (Oda *et al.*, 2002), a unique amino-acid insertion in fibrillar collagen (Wada *et al.*, 2006), and the distribution of micro RNAs (miRNAs) (Heimberg *et al.*, 2008). Finally, the *Branchiostoma floridae* genome helps confirming the sister-group relationship between tunicates and vertebrates in offering additional evidence from analyses of intron dynamics (Putnam *et al.*, 2008).

Cadherins are a superfamily of highly conserved adhesion molecules mediating cell communication and signalling that are pivotal for developmental processes of multicellular organisms. Their recent detection in the closest unicellular relatives of metazoans, the choanoflagellates, has highlighted their potential role in the origin of multicellularity (Abedin and King, 2008). Comparative studies on the classic cadherin subfamily has revealed that the structural element called Primitive Classic Cadherin Domain (PCCD) complex, otherwise termed non-chordate classic cadherin domain, is also present in cephalochordates, but has been lost in both tunicates and vertebrates (Oda *et al.*, 2002). The most parsimonious scenario is that this particular protein domain complex has been lost in the common ancestor of tunicates and vertebrates and constitutes a synapomorphy of Olfactores. However, cephalochordates possess two classic cadherin genes which originated by lineage-specific tandem duplication and that have a particular structure in lacking extracellular repeats found in all other investigated metazoans (Oda *et al.*, 2004). This derived state renders difficult to ascertain domain homology among chordate classic cadherin genes and casts doubt on its phylogenetic significance.

Further potential evidence for the clade Olfactores has been inferred from the evolution of fibrillar collagen genes within chordates. These genes represent important components of the notochord, the cartilage and mineralized bones in vertebrates. Phylogenetic analyses suggested that three ancestral fibrillar collagens gave rise to the gene diversity observed in living deuterostomes (Wada *et al.*, 2006). Comparative sequence analyses showed that

tunicates and vertebrates share a molecular signature in the form of a six to seven amino-acid insertion in the C-terminus noncollagenous domain of one type of fibrillar collagens, that is absent in cephalochordates and echinoderms (Wada *et al.*, 2006). This insertion was interpreted as supporting the idea that vertebrates are more closely related to tunicates than to cephalochordates (Wada *et al.*, 2006). The homology of the insertion appears nevertheless difficult to assert with certainty given the high degree of sequence divergence observed in this region of the molecule. More tunicate fibrillar collagen sequences might help in better understanding the dynamics of this peculiar amino-acid insertion and the phylogenetic signal it conveys.

The comparison of miRNA repertoires in metazoans has also recently unearthed some potential signatures for the sister-group relationship of tunicates and vertebrates (Heimberg *et al.*, 2008). miRNAs are small non-coding RNAs involved in regulation of gene expression in eukaryotes and playing an important role in the development of metazoans. Comparative genomic studies of miRNAs underlined that, during the evolution of metazoans, major body-plan innovations seemed to coincide with dramatic expansions of miRNA repertoires, suggesting a potential role in the increase of morphological complexity (Hertel *et al.*, 2006; Sempere *et al.*, 2006). The most recent study unveiled that three miRNA families (*mir-126*, *mir-135* and *mir-155*) were likely acquired in the common ancestor of tunicates and vertebrates (Heimberg *et al.*, 2008). Taking into consideration that miRNAs might be only rarely secondarily lost once they have been recruited, this finding provides corroborative evidence for the clade Olfactores. It should be noted however that, of these three miRNA families, only *mir-126* constitutes an exclusive synapomorphy for Olfactores without subsequent secondary lost in descendant taxa confirmed by Northern analysis. Moreover, the profound reorganization of miRNA repertoire undergone by tunicates requires being cautious when interpreting acquisition of miRNAs as potential signatures for reconstructing their phylogenetic relationships (Fu *et al.*, 2008).

Additional sequence-based phylogenomic reconstructions and analyses of rare genomic changes have been issued along with the recently published draft sequence of a cephalochordate (*Branchiostoma floridae*) genome (Putnam *et al.*, 2008). The phylogenetic analysis of a concatenation of 1,090 orthologs from 12 complete genomes retrieved maximal Bayesian support for Olfactores and chordates, whereas the corresponding bootstrap support was maximal for Olfactores but of only 78% for chordate monophyly (Putnam *et al.*, 2008). Moreover, the analysis of individual gene phylogenies revealed twice more cases where Olfactores was favoured over Euchordates than the reverse (Putnam *et al.*, 2008). Further evidence was obtained by analysing the phylogenetic signal deduced from the dynamics of



intron gain and loss among chordate genomes. Despite extensive intron losses along the tunicate lineage, a number of shared intron gain/loss events can be identified as a signature of tunicates and vertebrates common ancestry (Putnam *et al.*, 2008). Overall, the new evidence brought by the analysis of the *Branchiostoma floridae* genome essentially corroborates our present phylogenetic results.

### **Implications for Chordate Evo-Devo**

The additional evidence presented for the new chordate phylogeny provides a robust phylogenetic framework for (re)interpreting the evolution of morphological characters and developmental features. Inverting the phylogenetic position of tunicates and cephalochordates within monophyletic chordates highlights the prevalence of morphological simplification with characters that are likely ancestral for chordates, such as metameric segmentation, being lost secondarily in the tunicate lineage. On the other hand, the loss of pre-oral kidney and the presence of multiciliated epithelial cells might in fact constitute morphological synapomorphies for olfactores (Ruppert, 2005). The new chordate phylogeny further portrays tunicates as highly derived chordates with specialized lifestyles and developmental modes, whereas cephalochordates might have retained more ancestral chordate characteristics. We will use two examples to illustrate the importance of considering the new phylogenetic status of tunicates as the sister-group of vertebrates in the context of evolutionary developmental biology.

The first illustration concerns evolutionary origin of such fundamental structures as the neural crest and olfactory placodes. Migratory neural crest cells and sensory placodes have long been considered as vertebrate innovations. Implicated respectively in the development of major tissues and sensory organs, their origin is generally correlated with the increase in morphological complexity of vertebrates. However, recent molecular developmental studies have revealed the presence in tunicates of migratory neural crest-like cells (Jeffery, 2006; Jeffery *et al.*, 2004) and olfactory placodes (Bassham and Postlethwait, 2005; Mazet and Shimeld, 2005). When reinterpreted in light of the new chordate phylogeny, these results implied that both of these features did not evolve *de novo* in the vertebrate lineage, but rather evolved from specialized pre-existing structures in the common ancestor of vertebrates and tunicates.

The second example illustrates how the new phylogenetic context helps in understanding the genomic and developmental peculiarities of tunicates within chordates. The new phylogenetic picture implied that tunicate genomes have undergone significant genome reduction from the ancestral chordate genome (Holland, 2007). This genome compaction is

also associated with a high rate of genomic evolution at the levels of both primary sequences (Delsuc *et al.*, 2006; Edvardsen *et al.*, 2004) and genome organisation (Holland and Gibson-Brown, 2003). One of the most spectacular rearrangements of tunicate genomes is the lost of several Hox genes, the disintegration of the Hox cluster, and the lost of temporal colinearity in Hox gene expression during development (Ikuta *et al.*, 2004; Seo *et al.*, 2004). These observations raise the question of how tunicates, with their altered Hox clusters, are still able to develop a chordate body plan. In chordates, and deuterostomes more generally, temporal colinearity is regulated by the Retinoic-Acid (RA) signalling pathway which controls the antero-posterior patterning of the embryo (Cañestro *et al.*, 2006; Marlétaz *et al.*, 2006a). However, axial patterning in tunicates seems to have become independent of RA-signaling, with the genes of the RA machinery even being lost in *Oikopleura* (Cañestro and Postlethwait, 2007). Functional studies have shown that if “*Oikopleura* can be considered as a classical RA-signaling knock-down mutant naturally produced by evolution”, it is still capable of developing a typical chordate body plan (Cañestro and Postlethwait, 2007). With cephalochordates, which possess the RA genomic toolkit, being basal among chordates, RA-signalling must have been present in the tunicate ancestor and secondarily lost in *Oikopleura* suggesting that appendicularians use alternative mechanisms for the development of chordate features (Cañestro *et al.*, 2007; Holland, 2007).

The new chordate phylogeny strengthens the view that tunicates and cephalochordates represent complementary models for studying vertebrate Evo-Devo (Schubert *et al.*, 2006). Tunicates are phylogenetically closer to vertebrates but appear both morphologically and molecularly highly derived. The diversity of their developmental modes offers the opportunity to study the evolution of alternative adaptive solutions to the typical chordate development. In having retained more ancestral features, cephalochordates provide an ideal outgroup for polarizing evolutionary changes that occurred in tunicates and vertebrates. With the cephalochordate *Branchiostoma floridae* genome (Putnam *et al.*, 2008) and the upcoming genome sequence of the appendicularian *Oikopleura dioica*, the newly established phylogenetic framework makes chordate comparative genomics appearing full of promises for the Evo-Devo community as exemplified in a recent work on the origin and evolution of the Pax gene family (Bassham *et al.*, 2008).

## ACKNOWLEDGMENTS

We thank the associate editors Billie Swalla and José Xavier-Neto for the opportunity to write this paper. We also thank John Mallatt for kindly providing his 18S-28S rRNA alignment, Julien Claude for help in using R, and two anonymous reviewers for comments.

The extensive phylogenomic calculations benefited from the ISEM computing cluster. This is the contribution ISEM 2008-062 of the Institut des Sciences de l'Evolution.

**LITERATURE CITED**

- Abedin M, King N. 2008. The premetazoan ancestry of cadherins. *Science* 319:946-948.
- Adoutte A, Balavoine G, Lartillot N, Lespinet O, Prud'homme B, de Rosa R. 2000. The new animal phylogeny: reliability and implications. *Proc Natl Acad Sci USA* 97:4453-4456.
- Aguinaldo AM, Turbeville JM, Linford LS, Rivera MC, Garey JR, Raff RA, Lake JA. 1997. Evidence for a clade of nematodes, arthropods and other moulting animals. *Nature* 387:489-493.
- Bassham S, Cañestro C, Postlethwait JH. 2008. Evolution of developmental roles of Pax2/5/8 paralogs after independent duplication in urochordate and vertebrate lineages. *BMC Biol* 6:35.
- Bassham S, Postlethwait JH. 2005. The evolutionary history of placodes: a molecular genetic investigation of the larvacean urochordate *Oikopleura dioica*. *Development* 132:4259-4272.
- Baurain D, Brinkmann H, Philippe H. 2007. Lack of resolution in the animal phylogeny: closely spaced cladogeneses or undetected systematic errors? *Mol Biol Evol* 24:6-9.
- Blair JE, Hedges SB. 2005. Molecular phylogeny and divergence times of deuterostome animals. *Mol Biol Evol* 22:2275-2284.
- Blanquart S, Lartillot N. 2008. A site- and time-heterogeneous model of amino acid replacement. *Mol Biol Evol* 25:842-858.
- Bourlat SJ, Juliusdottir T, Lowe CJ, Freeman R, Aronowicz J, Kirschner M, Lander ES, Thorndyke M, Nakano H, Kohn AB, Heyland A, Moroz LL, Copley RR, Telford MJ. 2006. Deuterostome phylogeny reveals monophyletic chordates and the new phylum Xenoturbellida. *Nature* 444:85-88.

- Cameron CB, Garey JR, Swalla BJ. 2000. Evolution of the chordate body plan: new insights from phylogenetic analyses of deuterostome phyla. *Proc Natl Acad Sci USA* 97:4469-4474.
- Cañestro C, Postlethwait JH. 2007. Development of a chordate anterior-posterior axis without classical retinoic acid signaling. *Dev Biol* 305:522-538.
- Cañestro C, Postlethwait JH, Gonzalez-Duarte R, Albalat R. 2006. Is retinoic acid genetic machinery a chordate innovation? *Evol Dev* 8:394-406.
- Cañestro C, Yokoi H, Postlethwait JH. 2007. Evolutionary developmental biology and genomics. *Nat Rev Genet* 8:932-942.
- Castresana J. 2000. Selection of conserved blocks from multiple alignments for their use in phylogenetic analysis. *Mol Biol Evol* 17:540-552.
- Cavender JA, Felsenstein J. 1987. Invariants of phylogenies in a simple case with discrete states. *J Classif* 4:57-71.
- Conway-Morris S. 2003. The Cambrian "explosion" of metazoans and molecular biology: would Darwin be satisfied? *Int J Dev Biol* 47:505-515.
- Delsuc F, Brinkmann H, Chourrout D, Philippe H. 2006. Tunicates and not cephalochordates are the closest living relatives of vertebrates. *Nature* 439:965-968.
- Delsuc F, Brinkmann H, Philippe H. 2005. Phylogenomics and the reconstruction of the tree of life. *Nat Rev Genet* 6:361-375.
- Dunn CW, Hejnol A, Matus DQ, Pang K, Browne WE, Smith SA, Seaver E, Rouse GW, Obst M, Edgecombe GD, Sorensen MV, Haddock SH, Schmidt-Rhaesa A, Okusu A, Kristensen RM, Wheeler WC, Martindale MQ, Giribet G. 2008. Broad phylogenomic sampling improves resolution of the animal tree of life. *Nature* 452:745-749.
- Edwardsen RB, Lerat E, Maeland AD, Flat M, Tewari R, Jensen MF, Lehrach H, Reinhardt R, Seo HC, Chourrout D. 2004. Hypervariable and highly divergent intron-exon organizations in the chordate *Oikopleura dioica*. *J Mol Evol* 59:448-457.

- Felsenstein J. 2001. PHYLIP (Phylogenetic Inference Package) version 3.6. In: Distributed by the author, Department of Genetics, University of Washington, Seattle.
- Felsenstein J. 2004. *Inferring phylogenies*. Sunderland, MA, USA: Sinauer Associates, Inc. 645 p.
- Fu X, Adamski M, Thompson EM. 2008. Altered miRNA repertoire in the simplified chordate, *Oikopleura dioica*. *Mol Biol Evol* 25:1067-1080.
- Gee H. 2001. Deuterostome phylogeny: the context for the origin and evolution of the vertebrates. In: Ahlberg PE, editor. *Major events in early vertebrate evolution: palaeontology, phylogeny, genetics, and development*. London: Taylor and Francis. p 1-14.
- Halanych KM. 1995. The phylogenetic position of the pterobranch hemichordates based on 18S rDNA sequence data. *Mol Phylogenet Evol* 4:72-76.
- Halanych KM. 2004. The new view of animal phylogeny. *Annu Rev Ecol Evol Syst* 35:229-256.
- Halanych KM, Bacheller JD, Aguinaldo AM, Liva SM, Hillis DM, Lake JA. 1995. Evidence from 18S ribosomal DNA that the lophophorates are protostome animals. *Science* 267:1641-1643.
- Hartmann S, Vision TJ. 2008. Using ESTs for phylogenomics: can one accurately infer a phylogenetic tree from a gappy alignment? *BMC Evol Biol* 8:95.
- Heimberg AM, Sempere LF, Moy VN, Donoghue PC, Peterson KJ. 2008. MicroRNAs and the advent of vertebrate morphological complexity. *Proc Natl Acad Sci USA* 105:2946-2950.
- Hertel J, Lindemeyer M, Missal K, Fried C, Tanzer A, Flamm C, Hofacker IL, Stadler PF. 2006. The expansion of the metazoan microRNA repertoire. *BMC Genomics* 7:25.
- Holland LZ. 2007. Developmental biology: a chordate with a difference. *Nature* 447:153-155.

- Holland LZ, Gibson-Brown JJ. 2003. The *Ciona intestinalis* genome: when the constraints are off. *Bioessays* 25:529-532.
- Ikuta T, Yoshida N, Satoh N, Saiga H. 2004. *Ciona intestinalis* Hox gene cluster: Its dispersed structure and residual colinear expression in development. *Proc Natl Acad Sci USA* 101:15118-15123.
- Jefferies RPS. 1991. Two types of bilateral symmetry in the Metazoa: chordate and bilaterian. In: Bock GR, Marsh J, editors. *Biological Asymmetry and Handedness*. Chichester: Wiley. p 94-127.
- Jeffery WR. 2006. Ascidian neural crest-like cells: phylogenetic distribution, relationship to larval complexity, and pigment cell fate. *J Exp Zool B Mol Dev Evol* 306:470-480.
- Jeffery WR. 2007. Chordate ancestry of the neural crest: new insights from ascidians. *Semin Cell Dev Biol* 18:481-491.
- Jeffery WR, Strickler AG, Yamamoto Y. 2004. Migratory neural crest-like cells form body pigmentation in a urochordate embryo. *Nature* 431:696-699.
- Jeffroy O, Brinkmann H, Delsuc F, Philippe H. 2006. Phylogenomics: the beginning of incongruence? *Trends Genet* 22:225-231.
- Jimenez-Guri E, Philippe H, Okamura B, Holland PWH. 2007. *Buddenbrockia* is a Cnidarian worm. *Science* 317:116-118.
- Jobb G, von Haeseler A, Strimmer K. 2004. TREEFINDER: a powerful graphical analysis environment for molecular phylogenetics. *BMC Evol Biol* 4:18.
- Lartillot N, Brinkmann H, Philippe H. 2007. Suppression of long-branch attraction artefacts in the animal phylogeny using a site-heterogeneous model. *BMC Evol Biol* 7 Suppl 1:S4.
- Lartillot N, Philippe H. 2004. A Bayesian mixture model for across-site heterogeneities in the amino-acid replacement process. *Mol Biol Evol* 21:1095-1109.
- Lartillot N, Philippe H. 2008. Improvement of molecular phylogenetic inference and the phylogeny of Bilateria. *Philos Trans R Soc Lond B Biol Sci* 363:1463-1472.

- Lwoff A. 1944. L'évolution physiologique : étude des pertes de fonctions chez les microorganismes. Paris: Hermann. 308 p.
- Mallatt J, Winchell CJ. 2007. Ribosomal RNA genes and deuterostome phylogeny revisited: more cyclostomes, elasmobranchs, reptiles, and a brittle star. *Mol Phylogenet Evol* 43:1005-1022.
- Marlétaz F, Holland LZ, Laudet V, Schubert M. 2006a. Retinoic acid signaling and the evolution of chordates. *Int J Biol Sci* 2:38-47.
- Marlétaz F, Martin E, Perez Y, Papillon D, Caubit X, Lowe CJ, Freeman B, Fasano L, Dossat C, Wincker P, Weissenbach J, Le Parco Y. 2006b. Chaetognath phylogenomics: a protostome with deuterostome-like development. *Curr Biol* 16:R577-578.
- Matus DQ, Copley RR, Dunn CW, Hejnal A, Eccleston H, Halanych KM, Martindale MQ, Telford MJ. 2006. Broad taxon and gene sampling indicate that chaetognaths are protostomes. *Curr Biol* 16:R575-576.
- Mazet F, Shimeld SM. 2005. Molecular evidence from ascidians for the evolutionary origin of vertebrate cranial sensory placodes. *J Exp Zool B Mol Dev Evol* 304:340-346.
- Nielsen C. 2001. Animal evolution, interrelationships of the living phyla. Oxford, UK: Oxford University Press.
- Oda H, Akiyama-Oda Y, Zhang S. 2004. Two classic cadherin-related molecules with no cadherin extracellular repeats in the cephalochordate amphioxus: distinct adhesive specificities and possible involvement in the development of multicell-layered structures. *J Cell Sci* 117:2757-2767.
- Oda H, Wada H, Tagawa K, Akiyama-Oda Y, Satoh N, Humphreys T, Zhang S, Tsukita S. 2002. A novel amphioxus cadherin that localizes to epithelial adherens junctions has an unusual domain organization with implications for chordate phylogeny. *Evol Dev* 4:426-434.



- Philippe H. 1993. MUST, a computer package of Management Utilities for Sequences and Trees. *Nucleic Acids Res* 21:5264-5272.
- Philippe H, Delsuc F, Brinkmann H, Lartillot N. 2005a. Phylogenomics. *Annu Rev Ecol Evol Syst* 36:541-562.
- Philippe H, Lartillot N, Brinkmann H. 2005b. Multigene analyses of bilaterian animals corroborate the monophyly of ecdysozoa, lophotrochozoa, and protostomia. *Mol Biol Evol* 22:1246-1253.
- Philippe H, Snell EA, Bapteste E, Lopez P, Holland PW, Casane D. 2004. Phylogenomics of eukaryotes: impact of missing data on large alignments. *Mol Biol Evol* 21:1740-1752.
- Philippe H, Telford MJ. 2006. Large-scale sequencing and the new animal phylogeny. *Trends Ecol Evol* 21:614-620.
- Phillips MJ, Delsuc F, Penny D. 2004. Genome-scale phylogeny and the detection of systematic biases. *Mol Biol Evol* 21:1455-1458.
- Posada D, Crandall KA. 1998. MODELTEST: testing the model of DNA substitution. *Bioinformatics* 14:817-818.
- Posada D, Crandall KA. 2001. Selecting the best-fit model of nucleotide substitution. *Syst Biol* 50:580-601.
- Putnam NH, Butts T, Ferrier DE, Furlong RF, Hellsten U, Kawashima T, Robinson-Rechavi M, Shoguchi E, Terry A, Yu JK, Benito-Gutierrez EL, Dubchak I, Garcia-Fernandez J, Gibson-Brown JJ, Grigoriev IV, Horton AC, de Jong PJ, Jurka J, Kapitonov VV, Kohara Y, Kuroki Y, Lindquist E, Lucas S, Osoegawa K, Pennacchio LA, Salamov AA, Satou Y, Sauka-Spengler T, Schmutz J, Shin IT, Toyoda A, Bronner-Fraser M, Fujiyama A, Holland LZ, Holland PW, Satoh N, Rokhsar DS. 2008. The amphioxus genome and the evolution of the chordate karyotype. *Nature* 453:1064-1071.
- R Development Core Team. 2007. R: A language and environment for statistical computing. Vienna, Austria: R Foundation for Statistical Computing.

- Rodríguez-Ezpeleta N, Brinkmann H, Roure B, Lartillot N, Lang BF, Philippe H. 2007. Detecting and overcoming systematic errors in genome-scale phylogenies. *Syst Biol* 56:389-399.
- Roure B, Rodriguez-Ezpeleta N, Philippe H. 2007. SCaFoS: a tool for selection, concatenation and fusion of sequences for phylogenomics. *BMC Evol Biol* 7 Suppl 1:S2.
- Rowe T. 2004. Chordate phylogeny and development. In: Cracraft J, Donoghue MJ, editors. *Assembling the Tree of Life*. Oxford: Oxford University Press. p 384-409.
- Ruppert EE. 2005. Key characters uniting hemichordates and chordates: homologies or homoplasies? *Can. J. Zool.* 83:8-23.
- Schaeffer B. 1987. Deuterostome monophyly and phylogeny. *Evol Biol* 21:179-235.
- Schmidt HA, Strimmer K, Vingron M, von Haeseler A. 2002. TREE-PUZZLE: maximum likelihood phylogenetic analysis using quartets and parallel computing. *Bioinformatics* 18:502-504.
- Schubert M, Escriva H, Xavier-Neto J, Laudet V. 2006. Amphioxus and tunicates as evolutionary model systems. *Trends Ecol Evol* 21:269-277.
- Sempere LF, Cole CN, McPeck MA, Peterson KJ. 2006. The phylogenetic distribution of metazoan microRNAs: insights into evolutionary complexity and constraint. *J Exp Zool B Mol Dev Evol* 306:575-588.
- Seo HC, Edvardsen RB, Maeland AD, Bjordal M, Jensen MF, Hansen A, Flaatt M, Weissenbach J, Lehrach H, Wincker P, Reinhardt R, Chourrout D. 2004. Hox cluster disintegration with persistent anteroposterior order of expression in *Oikopleura dioica*. *Nature* 431:67-71.
- Steel M. 2005. Should phylogenetic models be trying to "fit an elephant"? *Trends Genet* 21:307-309.

- Swalla BJ, Cameron CB, Corley LS, Garey JR. 2000. Urochordates are monophyletic within the deuterostomes. *Syst Biol* 49:52-64.
- Swofford DL. 2002. PAUP\*: Phylogenetic Analysis Using Parsimony and other methods version 4.0b10. In: Sinauer, Sunderland, MA.
- Telford MJ, Budd GE. 2003. The place of phylogeny and cladistics in Evo-Devo research. *Int J Dev Biol* 47:479-490.
- Vienne A, Pontarotti P. 2006. Metaphylogeny of 82 gene families sheds a new light on chordate evolution. *Int J Biol Sci* 2:32-37.
- Wada H, Okuyama M, Satoh N, Zhang S. 2006. Molecular evolution of fibrillar collagen in chordates, with implications for the evolution of vertebrate skeletons and chordate phylogeny. *Evol Dev* 8:370-377.
- Wada H, Satoh N. 1994. Details of the evolutionary history from invertebrates to vertebrates, as deduced from the sequences of 18S rDNA. *Proc Natl Acad Sci USA* 91:1801-1804.
- Whelan S, Goldman N. 2001. A general empirical model of protein evolution derived from multiple protein families using a maximum-likelihood approach. *Mol Biol Evol* 18:691-699.
- Wiens JJ. 2003. Missing data, incomplete taxa, and phylogenetic accuracy. *Syst Biol* 52:528-538.
- Wiens JJ. 2005. Can incomplete taxa rescue phylogenetic analyses from long-branch attraction? *Syst Biol* 54:731-742.
- Wiens JJ. 2006. Missing data and the design of phylogenetic analyses. *J Biomed Inform* 39:34-42.
- Winchell CJ, Sullivan J, Cameron CB, Swalla BJ, Mallatt J. 2002. Evaluating hypotheses of deuterostome phylogeny and chordate evolution with new LSU and SSU ribosomal DNA data. *Mol Biol Evol* 19:762-776.



**FIGURE LEGENDS**

**FIG. 1:** Reanalysis of previous phylogenomic data using an improved model of sequence evolution. The Delsuc *et al.* (2006) phylogenomic dataset of 146 genes (38 taxa and 33,800 sites) was analyzed under the CAT+ $\Gamma_4$  site-heterogeneous mixture model of amino-acid replacement. Values at nodes represent Bayesian posterior probabilities (PP). Circles indicate nodes with maximal support PP = 1.0. The scale bar represents the estimated number of substitutions per site.

**FIG. 2:** Phylogenetic analyses of an updated phylogenomic dataset. **(a)** Bayesian consensus tree obtained using the CAT+ $\Gamma_4$  mixture model on the complete dataset based on the concatenation of 179 genes (51 taxa and 53,799 amino-acid sites) containing 32% missing data. **(b)** Bayesian inference using the CAT+ $\Gamma_4$  mixture model on the dataset reduced to the concatenation of the 106 genes for which sequences were available for at least 41 of the 51 taxa (25,321 amino-acid sites) containing only 20% of missing data. Values at nodes indicate Bayesian posterior probabilities (PP) / Maximum-likelihood bootstrap percentages (BP; 100 replicates) obtained under the WAG+ $\Gamma_8$ . Circles indicate strongly supported nodes with PP  $\geq$  0.95 and BP  $\geq$  95. The scale bar represents the estimated number of substitutions per site.

**FIG. 3:** Assessing the robustness of phylogenetic results to gene sampling using a jackknife procedure. The Bayesian phylogenetic inference was conducted under the CAT+ $\Gamma_4$  mixture model on 50 jackknife replicates of 50 genes over a total of 179. The tree presented is the weighted majority-rule consensus of all trees sampled every 10 cycles across the 50 replicates after removing the first 1000 trees in each MCMC as the burnin. Values at nodes represent corresponding jackknife-resampled posterior probabilities indices (PP<sub>JK</sub>). Circles indicate highly repeatable nodes with PP<sub>JK</sub>  $\geq$  95%. The scale bar represents the number of substitutions per site.

**FIG. 4:** New phylogenetic analyses of ribosomal RNA genes. **(a)** Principal component analysis of nucleotide composition of the combined 18S+28S rRNA dataset. The graph represents the projection of individuals on the first two axes, which explain more than 98% of the total variance. **(b)** Maximum-likelihood analyses of the 18S+28S dataset using the best-fitting standard DNA model (TIM+ $\Gamma_4$ +I) on nucleotides (left) and a two-state model (CF+ $\Gamma_8$ ) after RY-coding of nucleotides (right). ML bootstrap percentages are given at nodes when greater than 70 except within vertebrates. Circles indicate strongly supported nodes with BP  $\geq$

95. Squares points to shifting nodes of interest between the two ML trees. Scale bars represent the number of substitutions per site.

Figure 1

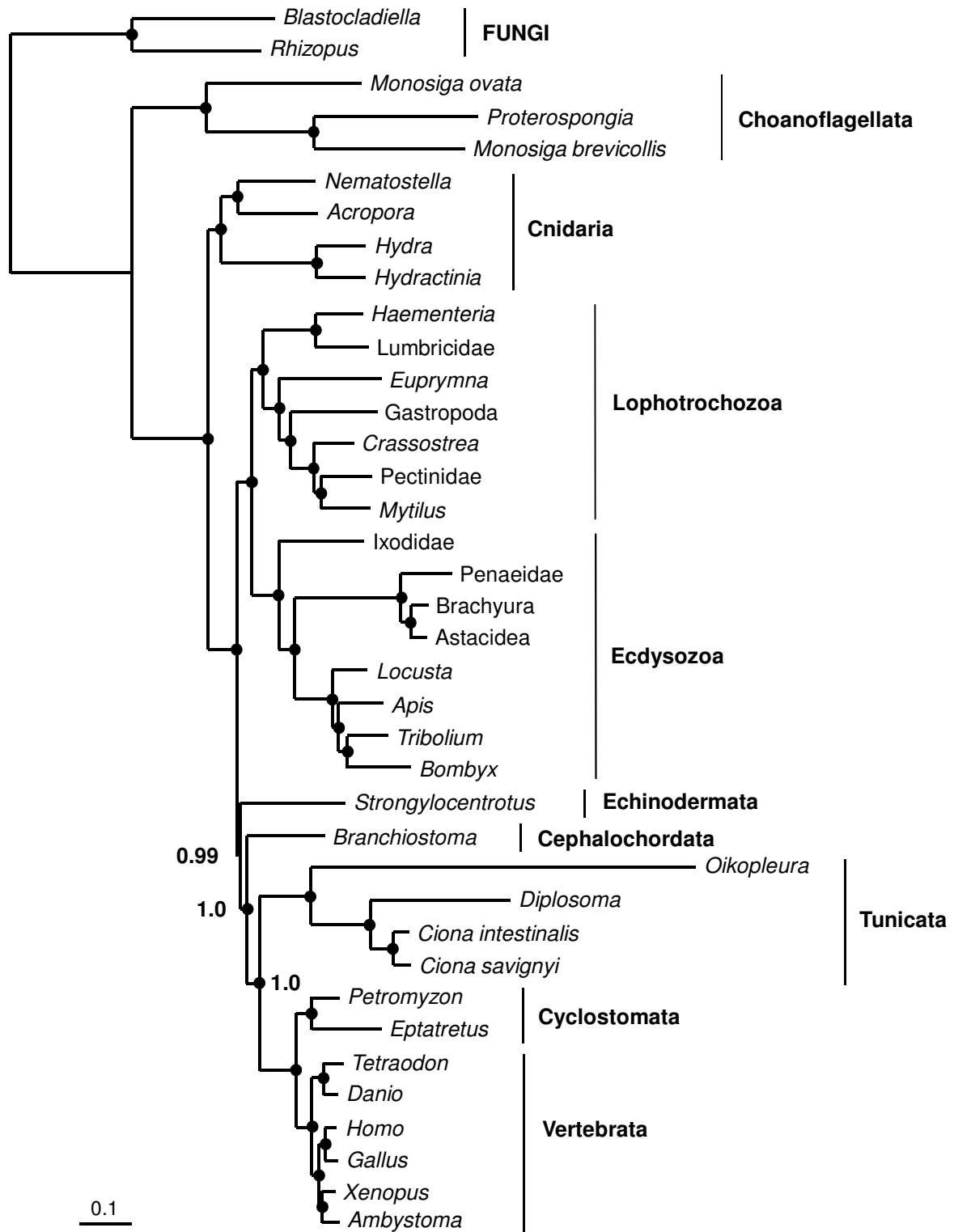
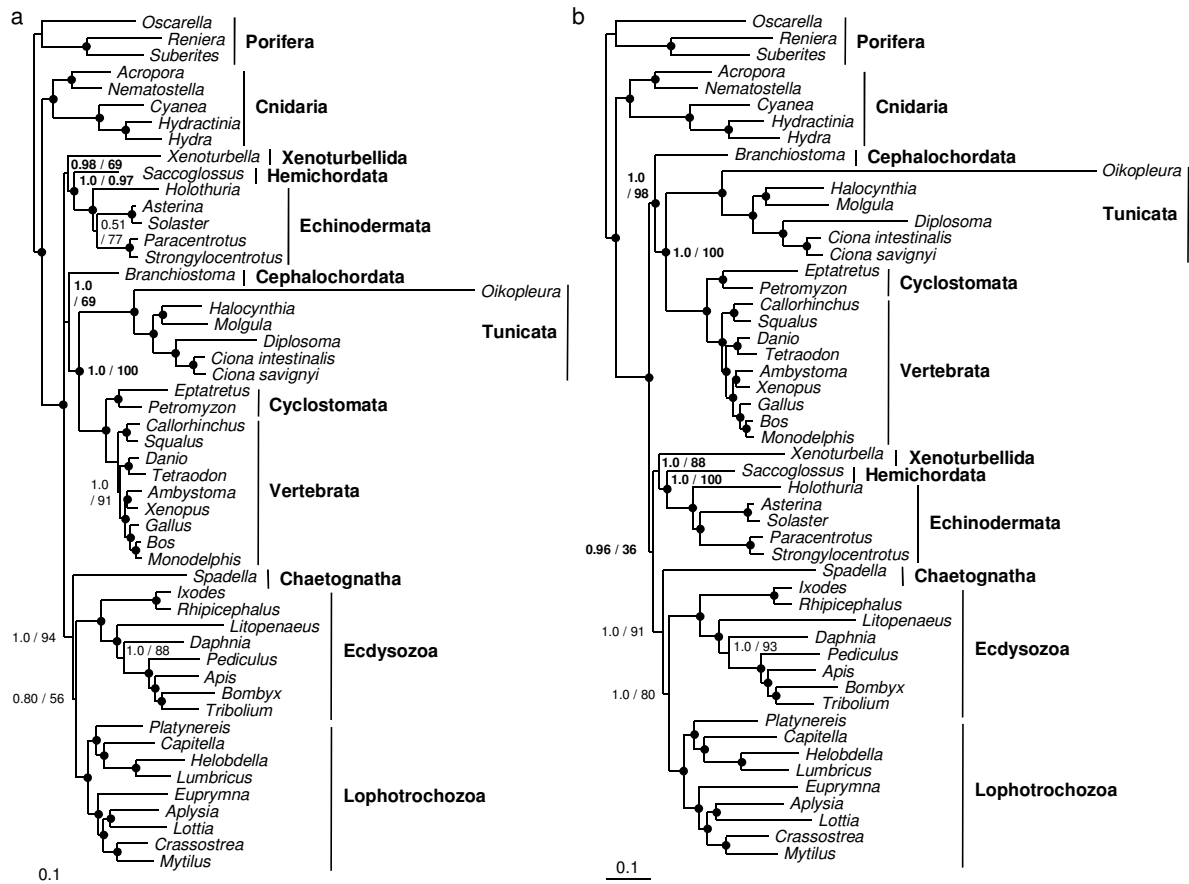


Figure 2





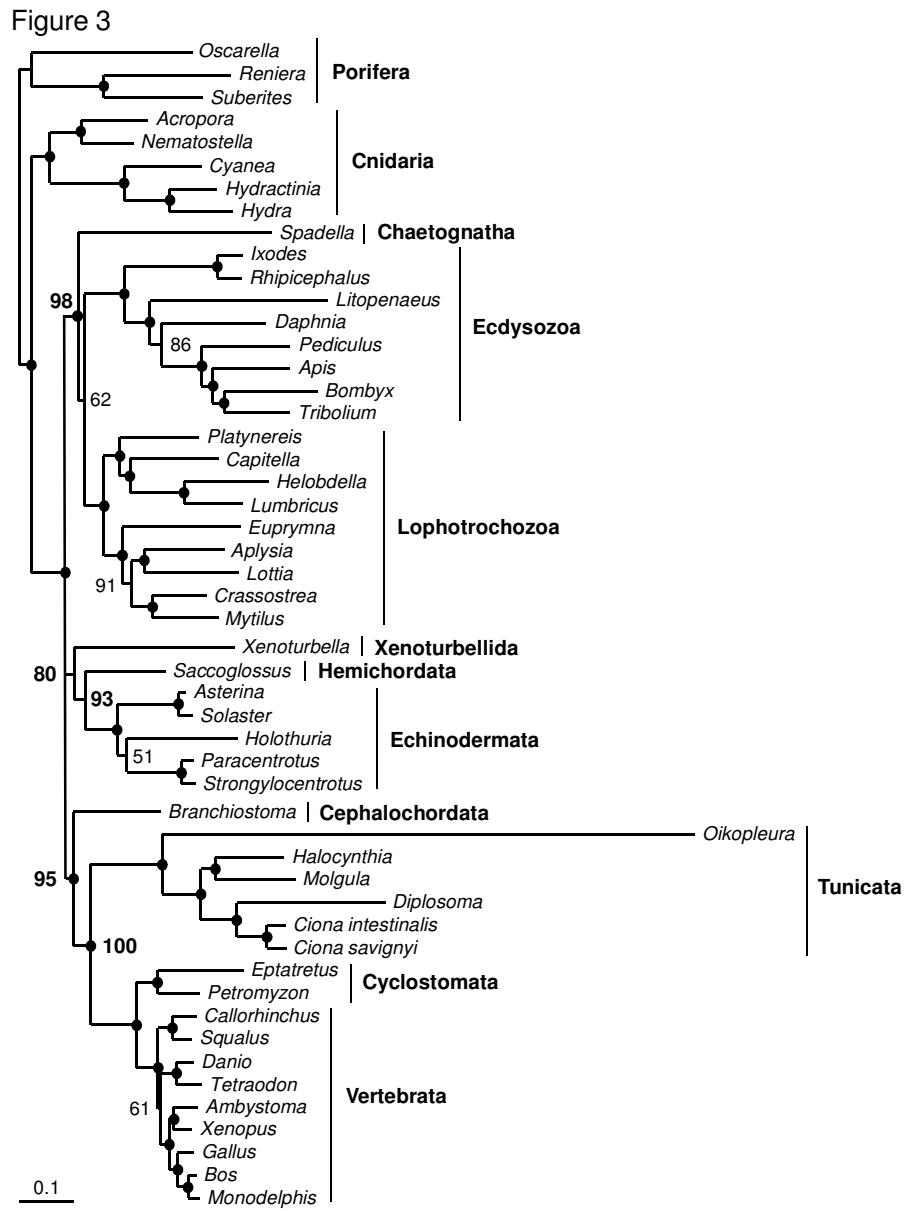
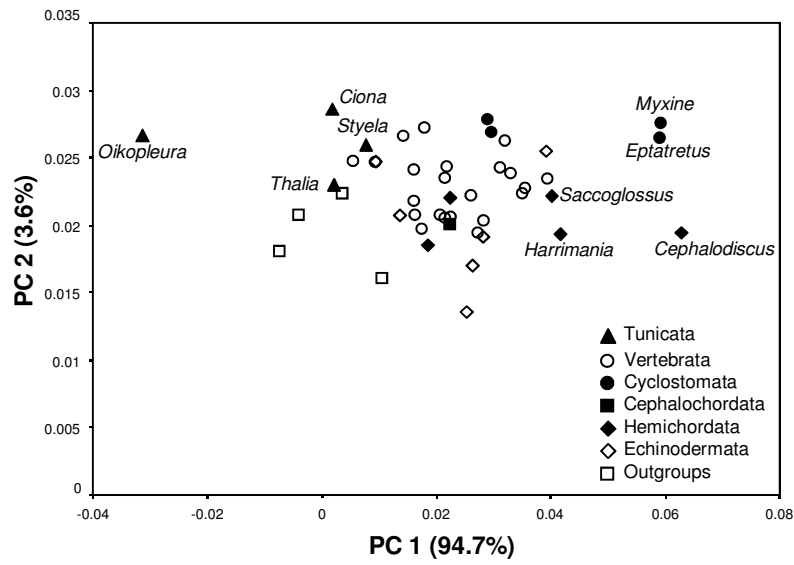
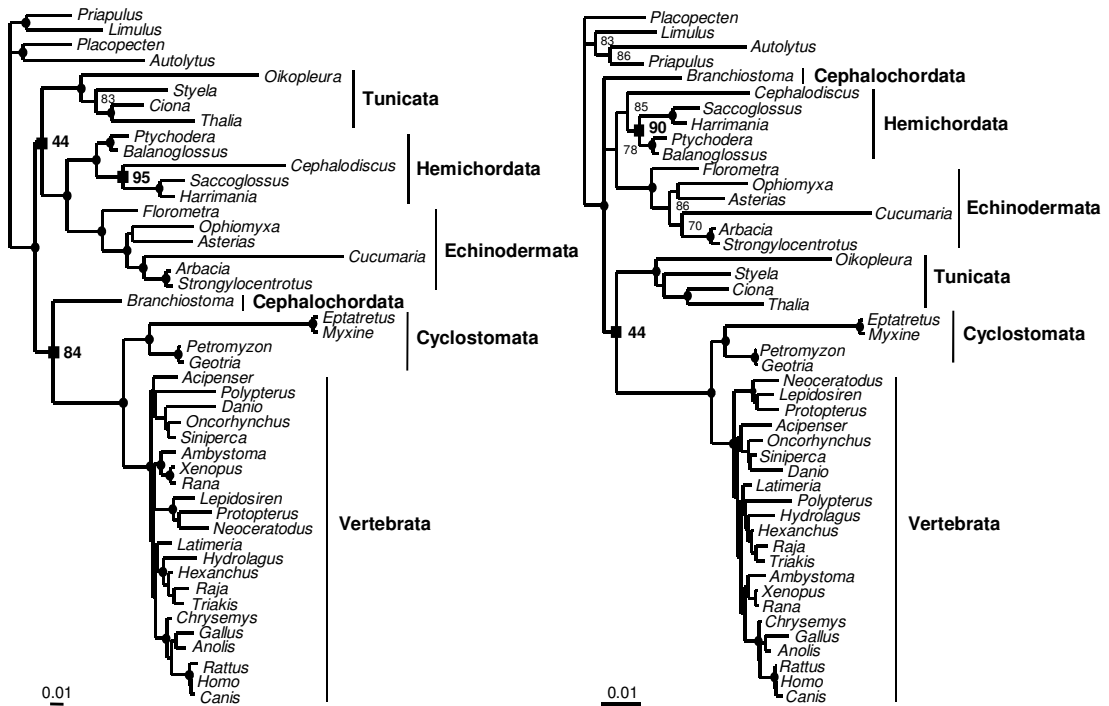


Figure 4

a



b



# Additional molecular evidence for the new chordate phylogeny

Frédéric Delsuc, Georgia Tsagkogeorga, Nicolas Lartillot & Hervé Philippe

## Supplementary Material

**Table S1: List of the 179 genes with names and numbers of amino-acid positions conserved for each gene alignment.**

<b>Gene Abbreviation</b>	<b>Gene name</b>	<b>Amino-acid positions</b>
<b>ar21</b>	Actin-related protein 2/3 complex subunit 3	155
<b>arc20</b>	Actin-related protein 2/3 complex subunit 4	166
<b>arp23</b>	Actin-related protein 2/3 complex subunit 1b	200
<b>atpsynthalph-a-mt</b>	F0F1 ATP synthase subunit alpha	464
<b>cct-A</b>	T complex protein 1 alpha subunit	515
<b>cct-B</b>	T complex protein 1 beta subunit	513
<b>cct-D</b>	T complex protein 1 delta subunit	489
<b>cct-E</b>	T complex protein 1 epsilon subunit	526
<b>cct-G</b>	T complex protein 1 gamma subunit	512
<b>cct-N</b>	T complex protein 1 eta subunit	528
<b>cct-T</b>	T complex protein 1 theta subunit	484
<b>cct-Z</b>	T complex protein 1 ? subunit	499
<b>cpn60-mt</b>	Heat shock protein HSP 60kDa mitochondrial	489
<b>crfg</b>	Nucleolar GTP binding protein 1	404
<b>ef1-RF3</b>	Release factor RF3	417
<b>ef1-RFS</b>	EF1alpha-like	399
<b>ef2-EF2</b>	Elongation factor EF2	806
<b>ef2-EF4</b>	EF2-like	567
<b>ef2-U5</b>	Elongation factor Tu family U5 snRNP specific protein	914
<b>fibri</b>	Fibrillarin	228
<b>fpps</b>	Farnesyl pyrophosphate synthase	250
<b>glcn</b>	N-acetyl glucosamine phosphotransferase	331
<b>grc5</b>	60S ribosomal protein L10 QM protein	214
<b>hsp70-E</b>	Heat shock 70kDa protein form ER	594
<b>hsp70-mt</b>	Heat shock 70kDa protein, mitochondrial form	600
<b>hsp90-C</b>	Heat shock 90kDa protein cytosolic form	641
<b>hsp90-E</b>	Heat shock 90kDa protein form ER	673
<b>hsp90-Z</b>	TNF receptor-associated protein 1	474
<b>if1a</b>	Eukaryotic translation initiation factor 1a	118
<b>if2b</b>	Eukaryotic translation initiation factor 2b	194
<b>if2g</b>	Eukaryotic translation initiation factor 2g	451
<b>if2p</b>	Eukaryotic translation initiation factor 2p	590
<b>if4a-a</b>	Translation initiation factor 4A	390
<b>if4a-b</b>	Translation initiation factor 4A	369
<b>if6</b>	Eukaryotic translation initiation factor 6	239
<b>ino1</b>	Myo-inositol-1-phosphate synthase	442
<b>I12e-A</b>	40S ribosomal Protein S12	123
<b>I12e-B</b>	High mobility group like nuclear protein 2 NHP2	127
	High mobility group like nuclear protein 2 NHP2-like protein	
<b>I12e-C</b>	1	124
<b>I12e-D</b>	60S ribosomal Protein L7a	246
<b>limif</b>	Small CTD phosphatase	186

<b>mcm-A</b>	minichromosome family maintenance protein 5	617
<b>mcm-B</b>	minichromosome family maintenance protein 2	708
<b>mcm-C</b>	minichromosome family maintenance protein 3	525
<b>mcm-D</b>	minichromosome family maintenance protein 7	641
<b>mcm-E</b>	minichromosome family maintenance protein 4	627
<b>mcm-F</b>	minichromosome family maintenance protein 6	686
<b>mcm-G</b>	minichromosome family maintenance protein 9	398
<b>mcm-H</b>	minichromosome family maintenance protein 8	461
<b>mra1</b>	Ribosome biogenesis protein NEP1 C2F protein	187
<b>nsf1-C</b>	Vacuolar protein sorting factor 4b	400
<b>nsf1-G</b>	26S proteasome AAA-ATPase regulatory subunit 8	385
<b>nsf1-I</b>	putative 26S proteasome ATPase regulatory subunit 7	422
<b>nsf1-J</b>	26S proteasome AAA-ATPase regulatory subunit 6	385
<b>nsf1-K</b>	26S proteasome AAA-ATPase regulatory subunit 6a	403
<b>nsf1-L</b>	26S proteasome AAA-ATPase regulatory subunit 6b	380
<b>nsf1-M</b>	26S proteasome AAA-ATPase regulatory subunit 4	441
<b>nsf1-N</b>	Katanin p60 subunit A	265
<b>nsf2-A</b>	Transitional endoplasmic reticulum ATPase TER ATPase	762
<b>nsf2-B</b>	Nuclear VCP-like	451
<b>nsf2-F</b>	Vesicular fusion protein nsf2	497
<b>orf2</b>	putative 28 kDa protein	177
<b>ornamtrans-a</b>	Ornithine aminotransferase	371
<b>pace2-A</b>	XPA binding protein 1	209
<b>pace2-B</b>	Conserved hypothetical ATP binding protein	239
<b>pace2-C</b>	Conserved hypothetical ATP binding protein	229
<b>pace4</b>	protein chromosome 2 ORF 4	230
<b>pace5-A</b>	Shwachman-Badian-Diamond syndrome protein	237
<b>pace6</b>	programmed cell death protein 5	106
<b>psma-A</b>	20S proteasome beta subunit macropain zeta chain	225
<b>psma-B</b>	20S proteasome alpha 1a chain	242
<b>psma-C</b>	20S proteasome alpha 1b chain	252
<b>psma-D</b>	20S proteasome alpha 2 chain	228
<b>psma-E</b>	20S proteasome alpha 1c chain	216
<b>psma-F</b>	20S proteasome alpha 3 chain	241
<b>psma-G</b>	20S proteasome alpha 6 chain	241
<b>psmb-H</b>	20S proteasome alpha 1d chain	185
<b>psmb-I</b>	20S proteasome alpha 1e chain	205
<b>psmb-J</b>	20S proteasome alpha 1f chain	212
<b>psmb-K</b>	20S proteasome beta 7 chain	234
<b>psmb-L</b>	20S proteasome beta 6 chain	196
<b>psmb-M</b>	20S proteasome beta 5 chain	195
<b>psmb-N</b>	20S proteasome beta 4 chain	194
<b>pyrdehydroe1b-B</b>	Branched chain ketoacid dehydrogenase E1 beta	321
<b>pyrdehydroe1b-mt</b>	Pyruvate dehydrogenase E1 beta subunit	322
<b>rad23</b>	UV excision repair protein RAD23	216
<b>rad51-A</b>	DNA repair protein RAD51	315
<b>rad51-B</b>	DNA repair protein DMC1	320
<b>rf1</b>	Eukaryotic peptide chain release factor subunit 1	402
<b>rla2-A</b>	60S acidic ribosomal protein P2	92
<b>rla2-B</b>	60S acidic ribosomal protein P1	76
<b>rpl1</b>	60S ribosomal Protein 1	213
<b>rpl11b</b>	60S ribosomal Protein 11b	169
<b>rpl12b</b>	60S ribosomal Protein 12b	163
<b>rpl13</b>	60S ribosomal Protein 13	179
<b>rpl14a</b>	60S ribosomal Protein 14a	123

<b>rpl15a</b>	60S ribosomal Protein 15a	204
<b>rpl16b</b>	60S ribosomal Protein 16b	173
<b>rpl17</b>	60S ribosomal Protein 17	175
<b>rpl18</b>	60S ribosomal Protein 18	183
<b>rpl19a</b>	60S ribosomal Protein 19a	190
<b>rpl2</b>	60S ribosomal Protein 2	248
<b>rpl20</b>	60S ribosomal Protein 20	160
<b>rpl21</b>	60S ribosomal Protein 21	154
<b>rpl22</b>	60S ribosomal Protein 22	96
<b>rpl23a</b>	60S ribosomal Protein 23a	139
<b>rpl24-A</b>	60S ribosomal Protein 24a	123
<b>rpl24-B</b>	60S ribosomal Protein 24b	137
<b>rpl25</b>	60S ribosomal Protein 25	126
<b>rpl26</b>	60S ribosomal Protein 26	135
<b>rpl27</b>	60S ribosomal Protein 27	137
<b>rpl3</b>	60S ribosomal Protein 3	393
<b>rpl30</b>	60S ribosomal Protein 30	105
<b>rpl31</b>	60S ribosomal Protein 31	108
<b>rpl32</b>	60S ribosomal Protein 32	129
<b>rpl33a</b>	60S ribosomal Protein 33a	106
<b>rpl34</b>	60S ribosomal Protein 34	108
<b>rpl35</b>	60S ribosomal Protein 35	123
<b>rpl36</b>	60S ribosomal Protein 36	85
<b>rpl37a</b>	60S ribosomal Protein 37a	81
<b>rpl38</b>	60S ribosomal Protein 38	65
<b>rpl39</b>	60S ribosomal Protein 39	51
<b>rpl42</b>	60S ribosomal Protein 4	102
<b>rpl43b</b>	60S ribosomal Protein 43b	88
<b>rpl4B</b>	60S ribosomal Protein 4b	326
<b>rpl5</b>	60S ribosomal Protein 5	275
<b>rpl6</b>	60S ribosomal Protein 6	186
<b>rpl7-A</b>	60S ribosomal Protein 7a	207
<b>rpl9</b>	60S ribosomal Protein 9	183
<b>rpo-A</b>	RNA polymerase alpha subunit	684
<b>rpo-B</b>	RNA polymerase beta subunit	1419
<b>rpo-C</b>	RNA polymerase gamma subunit	1091
<b>rpp0</b>	60S acidic ribosomal protein P0 L10E	289
<b>rps1</b>	40S ribosomal Protein 1	244
<b>rps10</b>	40S ribosomal Protein 10	119
<b>rps11</b>	40S ribosomal Protein 11	140
<b>rps13a</b>	40S ribosomal Protein 13a	151
<b>rps14</b>	40S ribosomal Protein 14	149
<b>rps15</b>	40S ribosomal Protein 15	139
<b>rps16</b>	40S ribosomal Protein 16	138
<b>rps17</b>	40S ribosomal Protein 17	119
<b>rps18</b>	40S ribosomal Protein 18	152
<b>rps19</b>	40S ribosomal Protein 19	132
<b>rps20</b>	40S ribosomal Protein 20	106
<b>rps22a</b>	40S ribosomal Protein 22a	130
<b>rps23</b>	40S ribosomal Protein 23	143
<b>rps24</b>	40S ribosomal Protein 24	122
<b>rps25</b>	40S ribosomal Protein 25	92
<b>rps26</b>	40S ribosomal Protein 26	101
<b>rps27</b>	40S ribosomal Protein 27	84
<b>rps27a</b>	40S ribosomal Protein 27a	155

<b>rps28a</b>	40S ribosomal Protein 28a	61
<b>rps29</b>	40S ribosomal Protein 29	56
<b>rrp46-A</b>	Exosome component 5	165
<b>rrp46-B</b>	Exosome complex exonuclease Rrp41	210
<b>sadhchydrolase-A</b>	Adenosylhomocysteinase 89E	455
<b>sadhchydrolase-E1</b>	S-adenosylhomocysteine hydrolase	424
<b>sap40</b>	40S ribosomal protein SA 40kDa laminin receptor 1	212
<b>Sra</b>	Signal recognition particle receptor alpha subunit SR alpha	422
<b>srp54</b>	Signal recognition particle 54 kDa protein	492
<b>Srs</b>	Seryl tRNA synthetase	454
<b>stbproptase2a-b</b>	Protein phosphatase-2A catalytic subunit beta	302
<b>stcproptase2a-c</b>	Protein phosphatase 6	294
<b>Suca</b>	Succinyl-CoA ligase alpha chain mitochondrial precursor?	291
<b>Tfiid</b>	TATA box binding protein related factor 2	175
<b>tif2a</b>	Translation initiation factor 2 alpha	304
<b>topo1</b>	DNA topoisomerase I, mitochondrial precursor	520
<b>u2snrnp</b>	Splicing factor U2AF 35 kDa subunit	179
<b>vacaatpasep121-a</b>	V-type ATP synthase subunit K	147
<b>Vata</b>	Vacuolar ATP synthase catalytic subunit A	610
<b>Vatb</b>	Vacuolar ATP synthase catalytic subunit B	479
<b>Vatc</b>	Vacuolar ATP synthase catalytic subunit C	336
<b>Vate</b>	Vacuolar ATP synthase catalytic subunit E	200
<b>Vatpased</b>	ATPase H <sup>+</sup> transporting V1 subunit D	210
<b>vdac2</b>	Voltage-dependent anion channel 2	276
<b>w09c</b>	TGF beta inducible nuclear protein	260
<b>Wrs</b>	tryptophanyl-tRNA synthetase	379
<b>Xpb</b>	Helicase XPB subunit 2	631
<b>yif1p</b>	homolog of Yeast Golgi membrane protein	188

**Table S2: List of Operational Taxonomic Units (OTUs) and Chimerical Sequences**

In the manuscript, OTUs were indicated by the genus name of the most frequent species indicated (underlined for chimerical OTUs).

## # Echinodermata

Strongylocentrotus: *Strongylocentrotus purpuratus*  
 Paracentrotus: *Paracentrotus lividus*  
 Asterina: *Asterina pectinifera*  
 Holothuria: *Holothuria glaberrima*, *Apostichopus japonicus*  
 Solaster: *Solaster stimpsonii*

## # Hemichordata

Saccoglossus: *Saccoglossus kowalevskii*  
 Xenoturbella: *Xenoturbella bocki*

## # Cephalochordata

Branchiostoma: *Branchiostoma floridae*, *Branchiostoma belcheri*, *Branchiostoma lanceolatum*

## # Tunicata

Ciona savignyi: *Ciona savignyi*  
 Ciona intestinalis: *Ciona intestinalis*  
 Oikopleura: *Oikopleura dioica*  
 Molgula: *Molgula tectiformis*  
 Halocynthia: *Halocynthia roretzi*  
 Diplosoma: *Diplosoma listerianum*

## # Myxini

Eptatretus: *Eptatretus burgeri*, *Myxine glutinosa*

## # Petromyzontiformes

Petromyzon: *Petromyzon marinus*

## # Chondrichthyes

Squalus: *Squalus acanthias*  
 Callorhynchus: *Callorhynchus milii*

## # Actinopterygii

Danio: *Danio rerio*  
 Tetraodon: *Tetraodon nigroviridis*

## # Amphibia

Xenopus: *Xenopus tropicalis*, *Xenopus laevis*  
 Ambystoma: *Ambystoma mexicanum*, *Ambystoma tigrinum*

## # Aves

Gallus: *Gallus gallus*

## # Metatheria

Monodelphis: *Monodelphis domestica*

## # Placentalia

Bos: *Bos taurus*, *Canis familiaris*, *Homo sapiens*, *Macaca mulatta*, *Pan troglodytes*, *Pongo pygmaeus*, *Macaca fascicularis*, *Mus musculus*, *Rattus norvegicus*

## # Annelida

Capitella: *Capitella species-2004*  
 Helobdella: *Helobdella robusta*, *Haementeria depressa*  
 Lumbricus: *Lumbricus rubellus*, *Eisenia fetida*, *Eisenia andrei*  
 Platynereis: *Platynereis dumerilii*

## # Mollusca

Aplysia: *Aplysia californica*  
 Lottia: *Lottia gigantea*

Euprymna: *Euprymna scolopes*, *Idiosepius paradoxus*  
Crassostrea: *Crassostrea virginica*, *Crassostrea gigas*  
Mytilus: *Mytilus galloprovincialis*, *Mytilus californianus*, *Mytilus edulis*

# Chaetognatha

Spadella: *Spadella cephaloptera*, *Flaccisagitta enflata*

# Arthropoda

Ixodes: *Ixodes scapularis*, *Ixodes pacificus*

Rhipicephalus: *Rhipicephalus microplus*, *Rhipicephalus appendiculatus*

Litopenaeus: *Litopenaeus vannamei*, *Litopenaeus setiferus*, *Penaeus monodon*, *Fenneropenaeus chinensis*,  
*Marsupenaeus japonicus*

Daphnia: *Daphnia pulex*, *Daphnia magna*

Tribolium: *Tribolium castaneum*

Apis: *Apis mellifera*

Bombyx: *Bombyx mori*

Pediculus: *Pediculus humanus*

# Cnidaria

Nematostella: *Nematostella vectensis*

Hydra: *Hydra magnipapillata*, *Hydra vulgaris*

Acropora: *Acropora millepora*, *Acropora palmata*, *Montastraea faveolata*

Hydractinia: *Hydractinia echinata*, *Podocoryne carnea*

Cyanea: *Cyanea capillata*

# Porifera

Reniera: *Reniera sp.*

Oscarella: *Oscarella carmela*, *Oscarella lobularis*, *Oscarella sp.*

Suberites: *Suberites domuncula*, *Suberites fuscus*















		vat				% missing	% missing	%
		vat	vat	pas	vda	positions	genes	chimeras
		c	e	ed	c2			
<i>Acropora</i>	<i>millepora</i>	28	5	10	0	25	7	2
<i>Ambystoma</i>	<i>mexicanum</i>	45	10	3	0	10	3	3
<i>Apis</i>	<i>mellifera</i>	0	0	0	0	1	1	0
<i>Aplysia</i>	<i>californica</i>	9	8	2	0	4	0	0
<i>Asterina</i>	<i>pectinifera</i>	100	100	100	0	43	19	0
<i>Bombyx</i>	<i>mori</i>	0	0	0	7	2	0	0
<i>Bos</i>	<i>taurus</i>	0	0	0	0	2	1	0
<i>Branchiostoma</i>	<i>floridae</i>	0	0	1	0	0	0	0
<i>Callorhinchus</i>	<i>mili</i>	52	100	63	47	55	30	0
<i>Capitella</i>	<i>sp.-2004</i>	18	5	24	0	7	0	0
<i>Ciona</i>	<i>intestinalis</i>	9	4	0	0	5	2	0
<i>Ciona</i>	<i>savignyi</i>	0	0	0	0	1	0	0
<i>Crassostrea</i>	<i>virginica</i>	51	100	13	1	26	8	8
<i>Cyanea</i>	<i>capillata</i>	100	28	100	82	72	46	0
<i>Danio</i>	<i>rerio</i>	0	0	0	0	1	2	0
<i>Daphnia</i>	<i>pulex</i>	0	0	0	0	0	0	0
<i>Diplosoma</i>	<i>listerianum</i>	100	30	100	0	66	44	0
<i>Eptatretus</i>	<i>burgeri</i>	68	2	100	1	18	6	0
<i>Euprymna</i>	<i>scolopes</i>	100	1	100	0	27	10	3
<i>Gallus</i>	<i>gallus</i>	0	0	0	0	8	8	0
<i>Halocynthia</i>	<i>roretzi</i>	1	13	4	0	36	42	0
<i>Helobdella</i>	<i>robusta</i>	60	100	27	55	18	3	8
<i>Holothuria</i>	<i>glaberrima</i>	100	100	100	100	60	36	0
<i>Hydra</i>	<i>magnipapillata</i>	0	0	0	0	0	0	0
<i>Hydractinia</i>	<i>echinata</i>	100	100	100	1	25	9	8
<i>Ixodes</i>	<i>scapularis</i>	0	6	0	0	3	0	0
<i>Litopenaeus</i>	<i>vannamei</i>	29	8	61	42	18	3	13
<i>Lottia</i>	<i>gigantea</i>	81	0	29	38	32	7	1
<i>Lumbricus</i>	<i>rubellus</i>	64	0	0	1	26	8	6
<i>Molgula</i>	<i>tectiformis</i>	1	0	0	0	4	3	0
<i>Monodelphis</i>	<i>domestica</i>	0	0	0	0	7	11	0
<i>Mytilus</i>	<i>galloprovincialis</i>	11	0	60	5	25	6	3
<i>Nematostella</i>	<i>vectensis</i>	1	0	1	0	0	0	0
<i>Oikopleura</i>	<i>dioica</i>	12	14	1	100	8	1	0
<i>Oscarella</i>	<i>carmela</i>	0	100	1	0	39	40	1
<i>Paracentrotus</i>	<i>lividus</i>	32	0	0	0	29	34	0
<i>Pediculus</i>	<i>humanus</i>	0	0	7	0	3	1	0
<i>Petromyzon</i>	<i>marinus</i>	29	0	4	0	10	2	0
<i>Platynereis</i>	<i>dumerilii</i>	34	100	100	0	48	48	0
<i>Reniera</i>	<i>sp.</i>	28	0	0	0	2	0	0
<i>Rhipicephalus</i>	<i>microplus</i>	0	20	0	0	4	5	1
<i>Saccoglossus</i>	<i>kowalevskii</i>	0	0	45	0	3	0	0
<i>Solaster</i>	<i>stimpsonii</i>	100	0	100	100	66	58	0
<i>Spadella</i>	<i>cephaloptera</i>	100	38	100	19	56	29	1
<i>Squalus</i>	<i>acanthias</i>	46	16	77	18	40	23	0
<i>Strongylocentrotus</i>	<i>purpuratus</i>	0	0	0	100	2	3	0
<i>Suberites</i>	<i>domuncula</i>	100	100	0	100	56	37	0
<i>Tetraodon</i>	<i>nigroviridis</i>	10	0	0	1	14	13	0
<i>Tribolium</i>	<i>castaneum</i>	0	0	0	0	0	0	0
<i>Xenopus</i>	<i>tropicalis</i>	0	0	0	0	0	0	0
<i>Xenoturbella</i>	<i>bocki</i>	80	100	0	0	36	21	0
% missing positions		33	24	28	16			
% missing OTUs		18	20	20	10			
% chimeras		0	0	0	2			
# amino-acid sites		336	200	210	276			