# Supplementary Materials for

## The Genome of the Ctenophore *Mnemiopsis leidyi* and Its Implications for Cell Type Evolution

Joseph F. Ryan, Kevin Pang, Christine E. Schnitzler, Anh-Dao Nguyen, R. Travis Moreland, David K. Simmons, Bernard J. Koch, Warren R. Francis, Paul Havlak, NISC Comparative Sequencing Program, Stephen A. Smith, Nicholas H. Putnam, Steven H. D. Haddock, CaseyW. Dunn, Tyra G. Wolfsberg, James C.Mullikin, Mark Q. Martindale, Andreas D. Baxevanis

Corresponding author. E-mail: andy@mail.nih.gov

**The PDF file includes:**

Materials and Methods
Figs. S1 to S10
Tables S1 to S31
References

3
4
5
6   **Supplementary Materials:**
7   Materials and Methods
8   Tables S1 to S31
9   Figure S1 to S10
10  References
11
12
13  **Material and Methods**
14
15  <u>Source material and sequencing</u>
16
17  We used cteno-dippers (i.e., a beaker secured to a broom handle) to carefully collect adult
18  *Mnemiopsis leidyi* from the Vineyard Sound near Woods Hole, Massachusetts, USA. We
19  isolated genomic DNA from the embryos of two self-fertilized adults. DNA from one
20  embryo pool was used to construct a library for sequencing using a Roche 454 Genome
21  Sequencer FLX machine (Roche Applied Science, Indianapolis, IN). One picotiter plate
22  was run in the 100-cycle mode using FLX chemistry yielding an average read length of
23  236 bases. Eight more plates were run in the 200-cycle mode using Titanium chemistry
24  with an average read length of 334 bases. We generated 7,334,972 raw reads from these
25  nine runs, which yielded 2.5 Gb of sequence. The initial libraries were made with 'GS
26  FLX Titanium Rapid Library Preparation Kit' and subsequent libraries were made with
27  the same kit along with the 'GS FLX Titanium Library Paired End Adaptors Kit'.
28
29  <u>Genome assembly</u>
30
31  Using the Phusion assembler (*54*), we assembled this data into 24,884 contigs,
32  constituting 150,340,428 bases of sequence. The N50 of our contig assembly is 11,936
33  bases.
34
35  DNA from the other embryo pool was used to create two mate pair libraries for Illumina
36  GA-iiX sequencing, one with a 3-kb insert and the other with a 4-kb insert. Paired reads
37  were 51 bases each. After removing duplicate read-pairs, 4.2 million and 2.6 million pairs
38  remained for the 3 and 4-kb libraries, respectively. We mapped these mate-pair reads to
39  the assembly using Illumina's short read aligner ELAND and integrated these mappings
40  into Phusion's scaffolding process. The final assembly consists of 5,100 scaffolds
41  resulting in 160-fold physical coverage and an N50 of 187 kb. See Table S5 for scaffold
42  size frequencies and Table S6 for scaffold gap frequencies.
43
44  <u>Estimation of variation</u>
45

46  To estimate the level of genetic variation between animals, we sequenced a lane of
47  genomic material from the embryos of another self-fertilized animal using an Illumina
48  Mi-Seq and aligned these reads to our assembly. We identified 589,252 positions with
49  single nucleotide variations (SNVs) across 118,613,222 bases that were covered well
50  enough to call SNVs. Therefore, we conservatively estimate the occurrence of SNVs
51  between these two animals to occur once every 200 bases. This result suggests that there
52  is ample variation between animals for population studies.
53
54  Evaluation of completeness and correctness of genome assembly
55
56  We estimated the completeness of our genome assembly by aligning publicly available
57  (GenBank) *M. leidyi* ESTs (15,752) to the assembly using BLAT (*55*) (version 34x12)
58  with default parameters. We used baa.pl (*56*) (version 0.20) with default settings to
59  generate the following statistics: (1) 99.4% of the transcripts were mapped with BLAT;
60  (2) 98.2% of the positions in the mapped transcripts were aligned; and (3) 95.2% of the
61  transcripts mapped to a single scaffold.
62
63  We also generated 79 Mb of paired-end reads and 83 Mb of single-end reads of RNA-seq
64  data from mixed-stage *M. leidyi* embryos (15-30h post-fertilization) through Illumina
65  GA-II sequencing. We assembled these two runs of RNA-seq reads using Trinity (*57*)
66  (version r2012-10-05) into 32,630 and 27,315 transcripts respectively. As was done with
67  the public ESTs, we used BLAT to align these transcripts to our scaffolds. We used
68  baa.pl (*56*) with default settings to generate the following statistics for the RNA-seq
69  transcripts: (1) 99.2% of the transcripts were mapped with BLAT; (2) 98.1% of the
70  positions in the mapped transcripts were aligned; and (3) 92.7% of the transcripts mapped
71  to a single scaffold.
72
73  Large repeats and segmental duplications can be particularly problematic for assemblies
74  based on next-generation sequencing data (*58*). To be sure that we have not collapsed
75  large amounts of identical repeats, we generated a genome-size estimation based on the
76  occurrence of kmers in our reads. Using a kmer of size 17 and excluding low-count
77  kmers ($<= 4x$) we estimate the genome size to be 200,738,694 base pairs. This is 45MB
78  larger than our genome assembly. We suspect that most of this difference in size is due to
79  SNVs in our sequencing data due to the variation introduced by our sample material
80  (embryos from a self-fertilized hermaphrodite). To test this possibility, we simulated 454
81  data using the Art next generation sequencing read simulator (*59*), with one haplotype
82  coming from our assembly and another haplotype from a version of the assembly that
83  was permuted once every 200 bases and coverage being 12x. When we estimate the
84  genome size of the simulated data with two haplotypes, we see over-estimates of size in
85  the range 50MB, which is the difference that we see between our assembly and the
86  estimated genome size. Therefore, we suspect that our 156MB genome assembly closely
87  resembles the true genome size and that hidden genome complexity based on collapsed
88  repeats is minimal.
89
90  Detecting Novel repeats
91

92   Repeats were identified in the *M. leidyi* genome using a large-scale sequence substring-
93   matching program called VMatch (*60*). Repeats were selected according to their length
94   (length >= 50bp), percent identity (100%), and orientation (both direct and reverse-
95   complemented). In a post-processing step, we required five matches to qualify as a
96   repeat. 10.07% of the *M. leidyi* genome was identified as being comprised of repeats
97   (Table S3).
98
99   The VMatch algorithm was also used to evaluate the genomic repeat levels for several
100  other species of interest, including *A. queenslandica, M. brevicollis, N. vectensis, T.*
101  *adhaerens,* and *S. rosetta* (Table S3). *H. magnipapillata* repeats were not evaluated using
102  VMatch due to the excessive number of transposable elements in the *H. magnipapillata*
103  genome (~57%). Total genomic repeats range from 2.10% (for *T. adhaerens*) to 24.69%
104  (for *N. vectensis*) of each individual genome.
105
106  Detecting Known repeats
107
108  RepeatMasker (*61*), is a tool that identifies known repetitive DNA elements including
109  low-complexity sequences and interspersed repeats. We used RepeatMasker (with
110  crossmatch) to screen the *M. leidyi* genome for known repeats (Table S4).
111
112  We determined that 2.49% of the *M. leidyi* genome is comprised of previously classified
113  repeats (Table S4). We did not scan the *H. sapiens, D. melanogaster*, or *N. vectensis*
114  genomes since these genomes were used to generate the Repbase database (*62*) used by
115  RepeatMasker.
116
117  RNA sequencing and assembly
118
119  We generated 162 MB of RNA-seq data from mixed stage embryos. We mapped this
120  RNA-seq data to our genome assembly using the TopHat (*63*) read-mapping software
121  (length = 51; inner-distance = 295) (version_1.2.0). We assembled these tophat mappings
122  into 49,850 transcript fragments using Cufflinks (*64*).
123
124  Overview of gene prediction pipeline
125
126  Our gene prediction pipeline consisted of several rounds of automated and manual steps.
127  We loaded the Cufflink fragments, along with 15,775 publicly available EST sequences
128  and 161 publicly available cDNA sequences into PASA (*65*) (Program to Assemble
129  Spliced Alignments, version_08_22_2010). Prior to submitting the *M. leidyi* genome for
130  gene prediction analysis, we masked the genome with our repeat library, except for
131  regions where RNA-seq mappings overlapped. The RNA-seq overlap consisted of 53,244
132  regions and 7,387,140 base pairs of sequence. We generated protein-coding gene models
133  using various gene prediction programs, including FGENESH (*66*); AUGUSTUS (*67*)
134  (version_2.3.1); HMMgene (*68*) (version_1.1); and GenomeScan (*69*) (version_0.1). We
135  submitted all predicted gene models to EvidenceModeler (*70*) (EVM,
136  version_r03062010), which also considers EST and RNA-seq evidence (through PASA)

137  as well as sequence similarity (through BLASTP of predicted proteomes). The total
138  number of gene predictions for each program is reported in Table S22.
139
140  Gene prediction with FGENESH
141
142  The FGENESH pipeline, consisting of three main steps, was implemented using partially
143  repeat-masked *M. leidyi* genomic sequence (5100 scaffolds). First, known mRNAs,
144  namely *M. leidyi* RNA-seq data (49,850 Cufflinks transcripts) and publicly available
145  mRNAs (161) and ESTs (15,752), were mapped to the *M. leidyi* genome. Next, genes
146  were predicted using ProtMap, based on homology to known proteins from the following
147  organisms: *H. sapiens, D. melanogaster, A. thaliana, C. intestinalis, B. floridae, C. teleta,*
148  *D. purpureum, L. gigantea, M. brevicollis, N. vectensis,* and *T. adhaerens*. Finally,
149  FGENESH performed an *ab initio* prediction of genes in the remaining regions having
150  neither mapped mRNAs nor genes predicted based on protein homology. FGENESH
151  predicted a total of 16,367 genes in the *M. leidyi* genome (Table S22).
152
153  Gene prediction with AUGUSTUS
154
155  AUGUSTUS uses species-specific parameters that can be trained on sets of annotated
156  genes (e.g., the Markov chain transition probabilities of coding and non-coding regions).
157  We trained AUGUSTUS with the *A. queenslandica* gene set. We then ran AUGUSTUS
158  gene modeling by incorporating known gene structure evidence ("hints") from extrinsic
159  sources, including *M. leidyi* RNA-seq reads (49,850 Cufflinks transcripts) and publicly
160  available mRNAs (161) and ESTs (15,752). Hints were used to search against the *M.*
161  *leidyi* genome using BLAT (*55*), and the results were combined into a single hints file.
162  AUGUSTUS was executed using this gene structure evidence and partially repeat-
163  masked *M. leidyi* genomic sequence (5100 scaffolds). AUGUSTUS generated a total
164  number of 29,359 *M. leidyi* gene predictions (Table S22).
165
166  Gene prediction with HMMgene
167
168  HMMgene, a program based on a hidden Markov model of gene structure, was also used
169  to generate *M. leidyi* gene models. HMMgene was run against the partially repeat-masked
170  (VMatch repeats not overlapping RNA-seq reads) *M. leidyi* genomic sequence (5100
171  scaffolds). HMMgene predicted a total of 13,948 *M. leidyi* gene models (Table S22).
172
173  Gene prediction with GenomeScan
174
175  We used the GenomeScan gene prediction program to generate gene models using the
176  partially repeat-masked (VMatch repeats not overlapping RNA-seq reads) *M. leidyi*
177  genomic sequence (5100 scaffolds). For the GenomeScan run we supplied sequence
178  similarity information based on BLASTX homology of known proteins from the
179  following organisms: *H. sapiens, D. melanogaster, A. thaliana, C. intestinalis, B.*
180  *floridae, C. teleta, D. purpureum, L. gigantea, M. brevicollis, N. vectensis,* and *T.*
181  *adhaerens*. BLASTX results (using the parameters -G 9, -E 2, and -e 0.05) were

182 converted to a "genoa" file in order to submit them to GenomeScan. GenomeScan
183 predicted a total of 6,443 *M. leidyi* gene models (Table S22).
184
185 <u>Gene prediction with PASA</u>
186
187 We implemented the PASA annotation pipeline. We first used Genomic Alignment and
188 Mapping Program (*71*) (GMAP) version_9_28_2007) to align the known ESTs, mRNAs
189 and RNA-seq reads to the partially-masked *M. leidyi* genome and transcript assemblies
190 were generated. We then trained PASA using the program options for longest ORFs and
191 extraction of terminal exons. The output of PASA was a GFF-formatted set of 153,004
192 "validated" transcripts, which we used as transcript evidence for consensus gene
193 modeling using EVM.
194
195 <u>Choosing between gene predictions with EVM</u>
196
197 We used the EvidenceModeler software (EVM) to combine *ab initio* gene predictions and
198 protein and transcript alignments, from all other gene modeling programs outlined above,
199 into weighted consensus gene structures. EVM requires an evidence weights file as input
200 where each piece of evidence is manually assigned a weight (1 to 10) based on data
201 confidence levels. To determine initial weights, we compared the gene predictions from
202 each program (through manual inspection) to 92 experimentally verified transcripts. We
203 determined these 92 transcripts using RACE-PCR (Clontech SMART™ RACE kit) from
204 cDNA of mixed developmental stages of several individuals. Our initial comparisons of
205 gene models from each program to our RACE sequences showed that FGENESH and
206 Augustus greatly outperformed the HMMgene and GenomeScan programs. This is
207 perhaps not surprising, since our runs of HMMgene and GenomeScan did not take into
208 consideration transcript information. Based on these initial evaluations, we set the initial
209 weighting of FGENESH and Augustus higher than HMMgene and GenomeScan.
210
211 We performed several rounds of weighting adjustments based on subsequent manual
212 inspections after each run (Table S23). We compared EVM predictions to the 92 full-
213 length RACE transcripts by visually evaluating their gene structures using their
214 designated transcript tracks in the JBrowse genome browser. Specifically, we inspected
215 the gene structure characteristics (e.g., exons, introns, splice junctions, and other genomic
216 features) with regard to quantity, length, and location and categorized any discrepancies
217 as joins, splits, partials, or missed predictions (Table S24).
218
219 In our initial EVM run, we generated 14,537 predicted genes, 85,446 exons, and 2,037
220 scaffolds. These were determined using four gene prediction programs (FGENESH,
221 AUGUSTUS, GenomeScan, and HMMgene), the PASA validated transcripts, and RNA-
222 seq (Cufflinks transcripts) data. Manual evaluation of EVM consensus gene models was
223 then conducted using JBrowse as described above. (Table S25). In the course of our
224 manual inspection of gene prediction results, we noticed that many of the mispredictions
225 corresponded with incorrect GenomeScan and HMMgene predictions, despite the output
226 from these two programs being heavily down-weighted. Removing GenomeScan and
227 HMMgene gene models from the next round of EVM runs led to improved results.

228
229 The second EVM run with only FGENESH and AUGUSTUS predictions produced
230 14,835 predicted genes, 86,712 exons, and 1,998 scaffolds. When we compared these
231 results to our RACE sequences we found that these results were an improvement over the
232 previous run (Table S25).
233
234 Since FGENESH consistently outperformed AUGUSTUS we tried removing the
235 AUGUSTUS predictions from consideration and ran EVM with FGENESH only. These
236 new set of results resulted in a 12% increase in the total number of gene predictions
237 (16,845 predicted genes, 89,564 exons, and 1,915 scaffolds) and resulted in an overall
238 improvement in terms of accuracy towards predicting the RACE sequences (Table S25).
239 We were not able to improve upon this run.
240
241 To make sure our RACE sequences were accurately represented in our gene models, we
242 processed one final run that included our test set (the RACE sequences) in PASA and
243 reran EVM with FGENESH as the only set of predictions. This run produced 16,545 gene
244 predictions, 91,482 exons, and 1,748 scaffolds. In the few cases where RACE sequences
245 were wrongly predicted or missed completely, we manually replaced or added the correct
246 sequence. We later added a few additional manual additions and changes based on
247 manual refinements made during subsequent studies (*72-74*). Statistics for the final gene
248 set are in Table S9. Statistics for the genome, in terms of coding vs. non-coding
249 sequences, are in Table S10.
250
251 <u>Nested intronic genes</u>
252
253 Genes that occur in the introns of other genes are called nested intronic genes. Table S2
254 shows that *M. leidyi* has a high number of nested intronic genes compared to other
255 genomes where these genes have been characterized. We see some examples of
256 interesting functional relationships between nested and "host" genes. For example, a gene
257 likely involved in the cell cycle, ML41156a (BLAST hit to human cyclin B2) is
258 embedded within the first intron of a gene likely involved in DNA Repair, ML41157a
259 (BLAST hit to human DNA-repair protein RAD52). Similarly, we found a gene that
260 codes for a likely RNA-binding protein, ML431110a (BLAST hit to human RAD52)
261 embedded within the intron of a likely RNA-editing enzyme, uridylyltransferase (BLAST
262 hit to human ZCCHC11). A more in-depth analysis of these genes may lead to some
263 interesting regulatory interactions between nested genes and their host genes.
264
265 We suspect that some of the predicted nested intronic genes are incorrect predictions. We
266 manually examined many of these nested gene models in some depth and found examples
267 of (1) seemingly bona fide nested genes, (2) likely spurious predictions based on rare
268 isoforms of the "host" gene, and (3) many cases that are not clear.
269
270 An example of bona fide nested genes includes ML000128a and ML000127a, which are
271 situated within the fourth and seventh introns of ML000126a respectively (Fig. S5).
272 ML000126a has an E2F_TDP domain that spans the fourth and fifth intron, indicating
273 that it truly encompasses the ML00128a gene. ML000128a is a G-protein coupled

274  receptor with an Ldl_recept_a domain. According to Pfam, the Ldl_recept_a domain does
275  not co-occur with E2F_TDP domain, suggesting that the 2-exon ML00128a gene is not
276  part of an ML00126a isoform. ML00127a is translated on the opposite strand from
277  ML00126a, and it has two clearly defined Glyco_hydro domains, suggesting that it is a
278  bona fide gene.
279
280  ML000317a is a single-exon prediction within the 16th intron of the ML000316a gene
281  model. These two predictions are in the same orientation (Fig. S6). We find evidence,
282  based on paired end reads from independent transcriptomic sequencing (Dr. William
283  Browne, University of Miami, personal communication), suggesting that the single-exon
284  gene model ML000317a is actually an exon within an isoform of the "host" ML000316a.
285  This makes ML000317a likely a misprediction and not a bona fide nested gene. These
286  findings are noted on the gene wiki pages of both the ML000316a
287  (http://research.nhgri.nih.gov/mnemiopsis/wiki/index.php/ML000316a ) and ML000317a
288  (http://research.nhgri.nih.gov/mnemiopsis/wiki/index.php/ML000317a), and will be
289  corrected in the next version (ML2.3) of the gene models. We encourage users to make
290  annotations on the gene wiki pages, so that they can be incorporated into future versions
291  of the *M. leidyi* gene models.
292
293  As additional *M. leidyi* RNA-seq data and sequence data from other ctenophores become
294  available, we will be able to improve upon our existing gene models. We expect the
295  number of nested genes to decrease, but to remain high compared to other animals. A
296  more in depth look at these nested genes, will likely lead to interesting discoveries.
297
298  Testing for potential non-metazoan contaminants in ML2.2 gene models
299
300  To address concerns about possible contamination in the *M. leidyi* gene models, we
301  computed an alien index for each ML2.2 gene model, as described in Gladyshev et al.,
302  2008 (*75*). An alien index is computed as the log-transformed difference between the best
303  BLAST e-value to a metazoan hit in the NCBI non-redundant (NR) database, and the best
304  BLAST e-value to a non-metazoan hit in NR.  Of the 16,548 M. leidyi gene models, we
305  found 136 models (<1%) with alien indices above the threshold of 45 (i.e., possible non-
306  metazoan contaminants or horizontally transferred sequences). Ninety-six of these had
307  corresponding transcript data suggesting through independent sequencing evidence that
308  these were not the results of contamination. We included notes in the corresponding wiki
309  pages of the remaining 40 models at the *Mnemiopsis* Genome Portal warning that these
310  sequences could possibly be non-metazoan contaminants.
311
312  Ortholog assignments
313
314  We generated sets of genes with putative orthology using sequence similarity based on
315  BLAST (*76*) and relative position in a predetermined phylogenetic tree. We used as our
316  input tree Table S26. We assigned bit scores to hits between each pair of genes by
317  summing those for initial BLASTP high-scoring segments found on the same pair of
318  genes, in consistent order, and overlapping less than five percent (with bit scores
319  penalized proportional to the amount of overlap). We determined orthologous sets of

320 genes at each tree node in two steps. First, if a set or gene from one child of the node was
321 in a mutual best hit-relation with a set or gene from the other child, they were combined
322 into a new set. Second, we considered in descending order all hits within this node's
323 subtree and between the subtree, and all outgroup genes. A hit to an outgroup gene
324 blocked any further merging of a gene or set (until another tree node was visited), while a
325 hit between two sets or genes within the subtree, neither previously blocked, resulted in
326 these being merged into a new set. This orthology computation was based on that
327 described in Putnam and coworkers, 2007 (*77*) with further refinement of the blocking
328 rules.
329
330 The use of a predetermined phylogenetic tree in the clustering algorithm makes these
331 clusters unsuitable for definitive phylogenetic analyses. As such, we did not use the
332 resulting clusters for gene selection in the phylogenetic analyses.
333
334 <u>Lineage-specific genes</u>
335
336 We used the clusters of orthologous genes and identified 7,171 clusters containing only
337 *M. leidyi* genes. These 7,171 clusters contained 9,288 genes. An RNA transcript is
338 associated with 7,798 (84%) of these 9,928 genes. We calculated the number of single-
339 species clusters for each species in our analysis and found that *M. leidyi* is eleventh out of
340 23 species in terms of numbers of lineage-specific clusters (Fig. S7).
341
342 <u>Gene duplications</u>
343
344 We looked in the clusters of orthologous genes generated in for lineage-specific
345 duplications. By definition, a cluster with more than one gene from a particular cluster
346 contains N-1 duplications from that lineage, where N is the number of genes in the cluster
347 from that species. We looked at clusters of genes that must have been in the metazoan
348 ancestor by only considering where at least one or both of the non-metazoan eukaryotes
349 *C. owczarzaki* and *M. brevicolis* are present and at least one gene from a metazoan is
350 present. *M. leidyi* has fewer duplicates than all 21 animals in the study other than
351 *Schistosoma mansoni* (Fig. S8)
352
353 <u>Conserved synteny</u>
354
355 Conservation of long-range gene linkage has been observed among diverse metazoan
356 genomes (*77-80*). We compared the gene composition of the 16 largest *M. leidyi*
357 scaffolds, where each contains genes representing from 20 to 41 ortholog groups to the
358 reconstructed linkage groups of the bilaterian ancestor (*80*) and found no scaffolds
359 showing significant conservation of synteny (relative to the null hypothesis of a random
360 distribution of genes across scaffolds). While a future genome assembly with longer
361 scaffolds would allow a more sensitive search for conservation, the current scaffolds are
362 long enough to detect conservation of the level observed in other studies. For
363 comparison, between chordates and the mollusk *Lottia gigantea*: 38% of scaffolds in the
364 same size range (20-41 ortholog groups) showed significant conservation with at least
365 one ancestral bilaterian linkage group (*81*). Our analysis takes into account the relatively

366 small number of gene families that *M. leidyi* shares with bilaterians, in that the metric we
367 used to compare scaffold sizes between genomes was the number of distinct ortholog
368 groups. If a large percentage of the *M. leidyi* genes that did not cluster with other animal
369 genes represent true orthologs that we were unable to detect, correct identification of
370 these difficult orthologs would increase our statistical power to detect a lower level of
371 synteny conservation. However, if there was extensive conserved synteny, we would
372 expect to detect this given the data and methods employed.
373
374 OrthoMCL clusters
375
376 Our ortholog assignments required a starting phylogenetic tree and therefore were not
377 suitable for phylogenetic inference. We therefore generated a set of clusters using
378 OrthoMCL (*82*) to use in phylogenetic analysis using default parameters. We used the
379 same set of species that were used to construct the ortholog clusters above, plus added 6-
380 frame translations of a de novo assembly of our *M. leidyi* RNA-seq (35,203 transcripts
381 assembled with Trinity (*57*)) to identify genes present in the transcriptome but missing
382 from the ML2.2 protein set. Of the 57,620 *M. leidyi* Trinity transcript translations
383 incorporated into the clusters, 57,161 were incorporated into clusters that also included
384 ML2.2 proteins or that consisted only of other *M. leidyi* Trinity transcript translations.
385 There were 459 transcript translations were incorporated into 132 clusters that did not
386 include an ML2.2 protein. These transcripts represent gene-structure improvements to
387 existing predictions or genes that were missed completely in ML2.2 and will ultimately
388 be incorporated into the next set of gene predictions (ML2.3).
389
390 Gene losses
391
392 Using the OrthoMCL clusters we calculated putative losses by identifying absences from
393 clusters where a gene from at least one of the non-metazoan eukaryotes *C. owczarzaki*
394 and *M. brevicolis* and one metazoan are present. We calculated this value for each animal
395 in our analysis. *M. leidyi* has 2,129 putative losses, comparable with *D. melanogaster*
396 with 2,123. The losses in *M. leidyi* are above average but fall within one half a standard
397 deviation of the median of 1,875  (Fig. S9).
398
399 Pfam domain content
400
401 We used the HMMER suite v3.0b (*83*) and Pfam v24.0 database (*84*) to identify protein
402 domains in the *M. leidyi* genome. We also searched the predicted proteomes of several
403 other animal and related non-animals. We have compiled a list of domains that appear to
404 have been lost in *M. leidyi* (Tables S27 and S28). Each proteome was scanned using the
405 HMMER program hmmscan on default settings. For each domain model in Pfam, if the
406 domain was detected in at least one peptide sequence with a HMMER independent
407 expect-value equal to or below .01, the domain was considered to be present in that
408 proteome. We also detected an additional 94 domains in the *M. leidyi* genome (with E-
409 Values better than 0.001) that did not occur in any of our gene predictions; we count
410 these as present in Table S27 and they are not included in Table S28). These regions may
411 be pseudogenes or coding genes that were not annotated.

412
413 <u>Construction of the Genome Set amino acid matrix</u>
414
415 We started with the 242-gene data matrix (104,840 columns, 12 animals) from the *A.*
416 *queenslandica* genome paper (*80*). In this study, the authors generated mutual best hits
417 from 17 proteomes to genes in the *M. brevicolis* genome. Hits with E-values ≥ 0.001
418 were discarded. To avoid paralogs, hits were only kept if the score of the second best hit
419 in both directions was less than half of the score of the best hit. These hits are considered
420 filtered mutual best hits. We used this same criteria to produce filtered mutual best hits
421 from the predicted proteomes of the following species: *Mnemiopsis leidyi*, *Salpingoeca*
422 *rosetta*, *Capsaspora owczarzaki*, *Sphaeroforma arctica*, *Spizellomyces punctatus*, and
423 *Saccharomyces cerevisiae*. For each of the 242 genes in the original data matrix, a file
424 with the full-length amino acid sequence from each species with a filtered mutual best hit
425 was created. These sequences were aligned using MAFFT (*85*) with the following
426 command: `mafft --ep 0 --genafpair --maxiterate 1000`
427 `FILE.fasta > FILE.mafft; einsi FILE.mafft > FILE.einsi`. Next
428 we trimmed each alignment with Gblocks (*86*) using the following command (where
429 X=0.65 * number of seqs in FILE.einsi): `Gblocks FILE.einsi -b2=X -b3=10`
430 `-b4=5 -b5=a`. These 242 individual alignments were then concatenated to create the
431 Genome Set.
432
433 In subsequent analyses we varied the non-animal outgroups as follows. The complete
434 matrix (Opisthokonta) consists of 13 animals plus *M. brevicolis, S. rosetta*, *C.*
435 *owczarzaki*, *S. punctatus*, *S. arctica, and S. cerevisiae*. The Holozoa matrix consists of all
436 species in the Opisthokonta matrix except for *S. punctatus*, *S arctica, and S. cerevisiae*.
437 The Choanimalia matrix consists of all species in the Holozoa matrix except for *C.*
438 *owczarzaki*. The animalia matrix includes only the 13 animals.
439
440 <u>Construction of the EST Set amino acid matrix</u>
441
442 To identify homologs between many taxa, we BLASTed each sequence from each taxa-
443 specific dataset (transcriptome or predicted proteome) against OrthoDB (*87*). We
444 categorized hits by the ortholog hit in OrthoDB and created fasta files for each ortholog.
445 We generated alignments using MAFFT and trimmed these alignments with Gblocks.
446 Individual phylogenies were created for each of these alignments and orthologs were
447 identified using the method described in Hejnol et al. (*88*). We analyzed final orthologs
448 for completeness and created a concatenated matrix with approximately 50% occupancy.
449 The final dataset includes 88,384 sites and 406 gene regions from the following 70 taxa:
450 *Oscarella carmela, Asterina pectinifera, Strongylocentrotus purpuratus, Capitella telata,*
451 *Branchiostoma floridae, Helobdella robusta, Lottia gigantea, Saccoglossus kowalevskii,*
452 *Drosophila melanogaster, Anoplodactylus eroticus, Daphnia pulex, Petromyzon marinus,*
453 *Ciona intestinalis, Acropora palmata, Acropora millepora, Nematostella vectensis,*
454 *Anemonia viridis, Capsaspora owczarzaki ATCC 30864, Saccharomyces cerevisiae,*
455 *Cryptococcus neoformans, Monosiga brevicollis, Phycomyces blakesleeanus,*
456 *Batrachochytrium dendrobatidis, Spizellomyces punctatus, Salpingoeca rosetta,*
457 *Amphimedon queenslandica, Lubomirskia baicalensis, Mnemiopsis, Hydra*
458 *magnipapillata, Trichoplax adhaerens, Cyanea capillata, Crassostrea virginica,*

459 *Schmidtea mediterranea, Cerebratulus lacteus, Ephydatia muelleri, Pleurobrachia*
460 *pileus, Suberites domuncula, Hydractinia echinata, Clytia hemisphaerica, Metridium*
461 *senile, Porites astreoides, Montastraea faveolata, Halocynthia roretzi, Euprymna*
462 *scolopes, Boophilus microplus, Isodiametra pulchra, Symsagittifera roscoffensis,*
463 *Monosiga ovata, Oscarella lobularis, Gallus gallus, Meara stichopi, Terebratalia*
464 *transversa, Euperipatoides kanangrensis, Xiphinema index, Oopsacas minuta,*
465 *Convolutriloba longifissura, Mertensiid, Paraplanocera oligoglena, Nemertoderma*
466 *westbladi, Echinoderes horni, Ptychodera flava, Sphaeroforma arctica, Rhizopus orizae,*
467 *Podocoryna carnea, Amoebidium parasiticum, Sycon raphanus, Xenoturbella bocki,*
468 *Carteriospongia foliascens, Aiptasia pallida,* and *Leucetta chagosensis.*
469
470 In subsequent analyses we varied the non-animal outgroups as follows. The complete
471 matrix (Opisthokonta) consists of 58 animals plus *M. brevicolis, S. rosetta, M. ovata, C.*
472 *owczarzaki, S. arctica, A. parasiticum, S. cerevisiae, S. punctatus, R. orizae, C.*
473 *neoformans, B. dendrobatidis, and C. neoformans.* The Holozoa matrix consists of all
474 species in the Opisthokonta matrix, except for *S. cerevisiae, S. punctatus, R. orizae, C.*
475 *neoformans, B. dendrobatidis,* and *C. neoformans.* The Choanimalia matrix consists of all
476 species in the Holozoa matrix except for *C. owczarzaki, S. arctica,* and *A. parasiticum.*
477 The Animalia matrix includes only the 58 animals.
478
479 Note: *E. muelleri* was inadvertently included twice in all matrices. The two *E. muelleri*
480 entries went together 100% of all tree runs with a near-zero length node. We pruned one
481 of the *E. muelleri* branches from all trees. Numbers above were adjusted accordingly.
482
483 <u>Testing for potential non-metazoan contaminants in the Genome or EST matrices</u>
484
485 The only gene model in our Genome or EST matrices with an alien index above the
486 suggested threshold of 45 was ML02315a (alien index=51.35), which was present in both
487 sets. This gene has both RNA-seq data and public EST data that aligned to the model and
488 is therefore likely not a contaminant.
489
490 <u>Maximum-likelihood phylogenetic analyses of amino acid matrices</u>
491
492 We ran a maximum-likelihood analysis on both the genome matrices and the EST
493 matrices using RaxML (version 7.2.8). For each matrix, we generated a best tree with the
494 following command: `raxmlHPC -m PROTGAMMAGTR -s MATRIX.phy -n`
495 `NAME -q MODELFILE`. Next, 100 bootstraps were generated with the following
496 command: `raxmlHPC -b RANDOMSEED -N 100 -m PROTGAMMAGTR -s`
497 `MATRIX.phy -n NAME.bs -q MODELFILE`. All trees are available in Figure S1.
498
499 <u>Bayesian analyses of amino acid matrices</u>
500
501 We ran two instances of Phylobayes (version 3.2e) on each matrix with the following
502 commands: `pb -d MATRIX.nex NAME.01; pb -d MATRIX.nex NAME.02`
503 For the Genome Set, we ran: `bpcomp -x BURN-IN NAME.01 NAME.02`

504  For the EST Set, we ran: `readpb -x BURN-IN NAME` for each individual run. We
505  set the burn-in to be approximately 1/3 of the total length of the shortest chain. Table S29
506  gives statistics including the burn-in used for each run. Bayesian trees resulting from the
507  genome matrices are available as Figure S1. Despite an average runtime of 205 days per
508  run, none of the runs on the EST Set converged. Consequently, we report trees from all
509  eight runs. An unexpected result in Bayesian analyses of the Genome Set is that *T.*
510  *adhaerens* forms a clade with *M. leidyi* and *A. queenslandica*, which is positioned sister
511  to the rest of animals. We are not aware of this result ever being reported in the literature.
512
513  Calculating support values for each scenario
514
515  For each phylogenetic scenario in Figure 2a, we constructed constraint trees and used
516  Paup* (*89*) to determine the percentage of bootstrap trees (in the case of maximum-
517  likelihood analysis) or post-burn-in trees (in the case of Bayesian analysis) that fit a
518  particular scenario (i.e., constraint tree). These values are reported in Table 1 of the main
519  paper.
520
521  Phylogenetic analysis of individual genes
522
523  It is well documented that phylogenies performed on a number individual genes will
524  produce phylogenetic trees that are incongruent with a tree produced from the
525  concatenation of these same set of genes (*90*). Nevertheless, it is interesting to see how
526  well individual genes support competing hypotheses. We have therefore, performed
527  analyses on each individual gene from both the Genome Set and the EST Set. In our
528  analyses of the 406 genes from the EST set (with all outgroups included), we found that,
529  of the 361 which included at least one sponge, one ctenophore, and one non-metazoan
530  outgroup, 19 of these produced a topology congruent with ctenophores as the sister group
531  to the rest of animals (Fig. 2D) and nine were congruent with sponges as the sister group
532  to the rest of animals (Fig. 2C). In our analyses of the 242 genes from the Genome set
533  (with all outgroups included), we found that of the 196 genes, which included at least one
534  sponge, one ctenophore and one non-metazoan outgroup, that 18 of these produced a
535  topology congruent with ctenophores as the sister group to the rest of animals (Fig. 2D)
536  and 15 were congruent with sponges as the sister group to the rest of animals (Fig. 2C).
537
538  For this analysis we generated a parsimony tree for each gene matrix using RAxML
539  (version 7.7.8) and then used this parsimony tree as a starting tree under RAxML-Light
540  (version 1.0.9) with the PROTGAMMAAUTO specified as the model. Under this model
541  setting RAxML conducts an ML estimate of all available pre-defined AA models
542  (excluding GTR) every time the model parameters are optimized during the tree search.
543
544  Branch length comparison
545
546  For every tree in our phylogenetic analyses of the Genome and EST Sets, except those
547  without a non-metazoan outgroup, we calculated the distance from the root in terms of
548  branch length for each taxon. These values are included in Tables S30 and S31.
549

550    To determine branch lengths, we opened each tree in FigTree v.1.3.1 (*91*) and rooted on
551    branch separating the metazoan clade from the non-metazoan clade. We exported this
552    rooted tree and opened the Nexus formatted tree in TreeStat v. 1.2 (*92*). We used TreeStat
553    to calculate the root-tip-lengths for each taxa. In all trees, the lengths of the *M. leidyi*
554    branches are similar to those of *Drosophila melanogaster*.
555
556    Maximum likelihood analysis of gene content
557
558    We assembled a presence/absence matrix of the OrthoMCL clusters and analyzed these
559    data with RAxML version 7.2.6 (SSE3 version) under the GTR gamma model of rate
560    heterogeneity as was done in previous studies (*93, 94*). The substitution probability
561    matrix depends on the branch lengths as well as the instantaneous rate matrix which, in
562    turn, depends on both the equilibrium state frequencies and parameters for rate changes
563    between states. Though only a single rate parameter is used in the RAxML binary model,
564    the equilibrium frequencies can, to a certain extent, account for differences in the rates of
565    gain and loss. Though the rate parameter is constant for gains and losses, the rate matrix
566    is therefore not necessarily symmetric. The equilibrium frequencies were estimated from
567    the data and do deviate from the observed frequencies as expected for asymmetric rates
568
569    Using input data and the NCBI taxonomy tree to initially infer a binary and taxonomy-
570    constrained tree topology as a starting tree we ran: (raxmlHPC-PTHREADS-SSE3 -T 4 -s
571    min2.I1.5.phy -m BINCAT -g known_relationships.tre -p 12345 -n S1). This returns a list
572    of files ended with '.S1', in which the file 'RAxML_bestTree.S1' is the binary constraint
573    starting tree. We used this resulting tree (via '-t' option) to infer the integer weights for
574    features (via '-f u' option), up-weighting the congruent ones while down-weighting those
575    that were incongruent: (raxmlHPC-PTHREADS-SSE3 -T 4 -f u -m BINGAMMA -s
576    min2.I1.5.phy -t RAxML_bestTree.S1 -# 100 -p 12345 -n S2). This returns a list of files
577    ending with '.S2', in which the file 'RAxML_weights.S2' is the weight vector.  Using the
578    weight vector to infer the binary and taxonomy-constrained best tree as a final tree
579    topology we ran: (raxmlHPC-PTHREADS-SSE3 -T 4 -m BINGAMMA -s min2.I1.5.phy
580    -t RAxML_bestTree.S1 -a RAxML_weights.S2 -p 12345 -n S3). This will return a list of
581    files ending with '.S3', in which the file 'RAxML_bestTree.S3' is the final tree.  To see if
582    the result is the same without the starting tree we ran:  (raxmlHPC-PTHREADS-SSE3 -T
583    4 -m BINGAMMA -s min2.I1.5.phy -a RAxML_weights.S2 -p 12345 -n S4) We also ran
584    a rapid bootstrap analysis: (raxmlHPC-PTHREADS-SSE3 -T 4 -m BINGAMMA -s
585    min2.I1.5.phy -a RAxML_weights.S2 -f a -# 100 -x 234544 -p 12345 -n S5).  To see the
586    effect of weighting we ran: (raxmlHPC-PTHREADS-SSE3 -T 4 -m BINGAMMA -s
587    min2.I1.5.phy -f a -# 100 -x 234544 -p 12345 -n S6). All analyses produced the topology
588    seen in Figure 4.
589
590    Maximum parsimony analysis of gene content
591
592    We ran a maximum parsimony analyses using Paup*, using the same character weighting
593    as the ML analyses, and recovered 100% bootstrap support for Ctenophora as the sister
594    group to all animals. This was the case whether characters were coded as ordered or
595    Dollo.

596    Bayesian analysis of gene content
597
598    We also performed Bayesian analyses of the unweighted gene presence/absence matrix
599    using Mr. Bayes. The data type was set to restriction (i.e., binary) and rates model was set
600    to invgamma. Two runs, each with one million generations, were completed. The average
601    standard deviation of split frequencies was 0.000000 by 80000 generations, indicating
602    quick convergence. These analyses recover Parahoxozoa and place *M. leidyi* as sister to
603    all other sampled animals, both with posterior probabilities of 100%. Some relationships
604    within Parahoxozoa were inconsistent with well-supported relationships (e.g., Cnidaria
605    was not recovered as monophyletic).
606
607    Hypothesis testing of gene content data using CONSEL
608
609    To see if this result is significantly better than the *a priori* alternative hypotheses
610    presented in Figure 2, we built constraint trees corresponding to each hypothesis and
611    generated trees with branch lengths for each starting tree (`raxmlHPC-SSE3 -m`
612    `MULTIGAMMA -K GTR -n fig2a -s aln.phy -a weights -g`
613    `fig2a.tre`). We then generated per site log likelihoods for these trees (`raxmlHPC-`
614    `SSE3 -f g -m MULTIGAMMA -K GTR -n 6trees -s aln.phy -a`
615    `weights -z 6trees.tre`). Next, we used the per site log likelihoods as input to
616    CONSEL version 1.20 (*95*) to generate p-values for each of the alternative hypotheses
617    (`seqmt --puzzle RAxML_perSiteLLs.6trees; makermt`
618    `RAxML_perSiteLLs; consel RAxML_perSiteLLs; catpv`
619    `RAxML_perSiteLLs`). Under the Approximately Unbiased (AU) test, Kishino-
620    Hasegawa (KH) test, and the Bootstrap Probability (BP) test, the relationship of *M. leidyi*
621    as the sister group to all other animals (Ct,) is significantly (0.05) better than all of the
622    other hypotheses presented in Figure 2 (Table S12). Under the Shimodaira-Hasegawa
623    (SH) test the relationship of *M. leidyi* as sister to all other animals except *A.*
624    *queenslandica* (Po,) cannot be rejected; however it is well-known that the SH test is very
625    conservative and includes more trees in the confidence set than is appropriate (*96*).
626
627    Hypothesis testing of gene content data using SOWH
628
629    We performed a Swofford–Olsen–Waddell–Hillis (SOWH) test on the gene content data.
630    Under the SOWH test, the relationship of *M. leidyi* as the sister group to all other animals
631    (Ct,) is significantly ($P < 0.05$) better than all of the other hypotheses presented in Figure
632    2. We followed the procedure as described by Goldman and coworkers, 2000 (*97*).
633    Briefly, we performed the following: (1) performed a maximum likelihood analysis on
634    the original alignment using the weights file generated as described in S9.1 with and
635    without a constraint tree representing the hypothesis being tested, (2) generated 100 data
636    sets with seq-gen using the alpha and branch lengths obtained from the ML tree produced
637    with the constraint (NOTE: since seq-gen does not produce binary matrices, the
638    frequency of zeroes and ones were assigned to A and T respectively, while the frequency
639    of C and G were set to 0. Since ones and zeroes are required by RAxML for a two-state
640    matrix, occurrences of A and T were substituted with zeroes and ones in the seq-gen
641    output), (3) performed maximum likelihood analyses on each of these 100 simulated
642    datasets with and without the constraint tree being tested, and (4) compared the

643 distribution of the differences in likelihoods with the difference in likelihood between the
644 best tree and the tree generated with the constraint on the real data set with the pnorm
645 function in R.
646
647 <u>Phylogenetic analysis of Near Intron Pairs</u>
648
649 Lehmann and coworkers assembled and analyzed a large set of near intron pairs (i.e.,
650 mutually exclusive introns situated within close proximity) from 48 species using
651 maximum parsimony (*98*). The *M. leidyi* genome assembly was used for this original
652 analysis, but gene models were not publicly available at the time. In their original study
653 *M. leidyi* was placed sister to insects with very low support using unweighted maximum
654 parsimony. Jörg Lehmann supplied us with a version of the near intron pair matrix that
655 used the ML2.2 gene models to increase the number of orthologs considered from 1,847
656 to 2,955 for the 4,405-gene dataset (Jörg Lehmann, University of Leipzig, personal
657 communication). We conducted a weighted likelihood analysis on this new matrix with
658 RAxML (raxmlHPC-PTHREADS-SSE3 version 7.7.8) and constrained known
659 relationships.
660
661 We constructed a tree of the taxa in the matrix used in OrthoMCL with known bilaterian
662 relationships, leaving the non-bilaterian groups in a polytomy. Using this tree as input to
663 the RAxML option (-g) we generated a weight matrix in RAxML (raxmlHPC-
664 PTHREADS-SSE3 -T 4 -f u -g known_relationships.tre -m MULTIGAMMA -n
665 01.int.weights -s concatenatedAll28.REAL_NIPS.phy) to up-weight congruent and
666 down-weight incongruent columns. We used these weights to infer the most likely tree
667 (raxmlHPC-PTHREADS-SSE3 -T 4 -m MULTIGAMMA -K GTR -n
668 03.restrained_w_known -s concatenatedAll28.REAL_NIPS.phy -a
669 RAxML_weights.01.int.weights -g known_relationships.tre). Bootstraps were generated
670 using (raxmlHPC-PTHREADS-SSE3 -T 4 -m MULTIGAMMA -K GTR -n
671 03.restrained_w_known.boot -s concatenatedAll28.REAL_NIPS.phy -a
672 RAxML_weights.01.int.weights -b 4321 -N autoMRE -g known_relationships.tre) and
673 then applied with (raxmlHPC -f b -z RAxML_bootstrap.03.restrained_w_known.boot -t
674 RAxML_bestTree.03.restrained_w_known -s concatenatedAll28.REAL_NIPS.phy -m
675 MULTIGAMMA -n 03.restrained_w_known_boots_applied).
676
677 This analyses resulted in a topology with *M. leidyi* as the sister group to the rest of
678 animals with low support (Fig. S2).
679
680 <u>Presence and absence of *M. leidyi* neural components</u>
681
682 The presence and absence of neural components in *M. leidyi* (Tables S15-17) was
683 determined using the following methods: A neural synaptic gene set of protein sequences
684 of *Homo sapiens* was collected from the NCBI database. This set of neural genes was
685 used as reciprocal best BLASTP and TBLASTN queries against the *M. leidyi* predicted
686 proteome and genome respectively. The domain architecture of candidate protein
687 sequences was predicted using the SMART domain prediction database (*99*) utilizing the
688 outlier homologues and PFAM domains settings. The domain architecture of each protein
689 was then compared to the human proteins, noting when diagnostic domains were absent.

690    The presence or absence of neural components in *N. vectensis* (*Nv*), *H. magnipapillata*
691    (*Hm*), *T. adhaerens* (*Ta*), *A. queenslandica* (*Aq*), *M. leidyi* (*Ml*), and *M. brevicolis* (*Mb*)
692    were verified these candidates using the same approach with queries from the datasets of
693    two other studies (*100, 101*).
694
695    <u>Missing synaptic scaffolding components are also absent from the transcriptomes of</u>
696    <u>seven other ctenophore species</u>
697
698    We searched for neuroligin and glutamate receptors in the transcriptomes of seven other
699    ctenophore species. We found no neuroligin in any of the transcriptomes. We found
700    candidate glutamate receptors in these other species, but our phylogenetic analysis (Fig.
701    S3) shows that these sequences are related equally to AMPA, NMDA, kainate-type, and
702    delta2-like glutamate receptors, suggesting that the ctenophore sequences descended from
703    an ancestral sequence that went on to differentiate into these four classes of glutamate
704    receptors. The top five sequences for neuroligin and glutamate receptors queries (i.e.,
705    those with TBLASTN hits to proteins of interest with E-values less than 0.001) have been
706    submitted to GenBank and accession numbers are included in Table S20.
707
708    <u>Phylogenetic distribution of neural components</u>
709
710    Behavioral sensitivity to several neurotransmitters (acetylcholine, serotonin, epinephrine,
711    and norepinephrine) has been shown in ctenophores (*102*), and immunoreactivity to the
712    neuropeptide FMRFamide, and preprohormone vasopressin has also been reported (*103*).
713    The neurotransmitter synthesis genes acetylcholinesterase and glutamate decarboxylase
714    are present in the *M. leidyi* genome; however, the monoamine enzymes DOPA
715    decarboxylase and dopamine-beta hydroxylase (Table S17) required for the production of
716    the chatelcholamines epinephrine, norepinephrine, and dopamine, are absent, and
717    immunological investigations have failed to detect the presence of serotonin in the
718    ctenophore *Bolinopsis infundibulum* (*104*). Metabotropic glutamate receptors (mGluR
719    and GABAb) are present in the *M. leidyi* genome. Seven-transmembrane receptors with
720    high homology to serotonin, dopamine receptors, and a large number of G-protein
721    coupled receptors are also present, suggesting that these receptors target alternative
722    neuropetides in ctenophores. Neurotransmitter transport genes found in all other
723    metazoans such as SNAP, syntaxin, synaptotagmin, and synaptobrevin are present, but
724    the gene synapsin is missing (also absent in placozoans and sponges). Gap junction genes
725    (innexins), found in cnidarians and bilaterians, are also present within the *M. leidyi*
726    genome (12 predicted paralogs), but are absent in placozoans and sponges.
727
728    <u>Presence and absence of M. leidyi mesoderm components</u>
729
730    We determined the presence and absence of mesoderm components in *M. leidyi* (Tables
731    S18-19) using the same methods used to determine neural components above. The
732    following references detail the role of these genes in mesoderm: twist and snail (*105*),
733    tinman, bagpipe, and ladybird (*106, 107*), nk 2.1, Pax3, and Myf5 (*108*), Eomesodermin
734    (*109*), myoD and myogenin (*110*).
735

736    <u>Trancriptome data from other ctenophore species</u>
737

738    Many genes involved in mesodermal cell types that are present in cnidarians and
739    bilaterians are absent from *M. leidyi*. Since these genes are also missing from non-
740    metazoans and sponges, it is parsimonious to assume that these absences are primary
741    (i.e., these genes arose after ctenophores and sponges split from Parahoxozoa) rather than
742    secondary losses. To reduce the possibility that these were secondarily lost in the lineage
743    leading to *M. leidyi*, we searched for these missing genes in deeply sequenced
744    transcriptomic data from seven ctenophore species: *Bathyctena chuni*, *Beroe forskalii,*
745    *Charistephane fugiens*, *Euplokamis dunlapae, Hormiphora californensis*, *Lampea lactea*,
746    and *Thalassocalyce inconstans*. The top five sequences for each query (i.e., those with
747    TBLASTN hits to proteins of interest with E-values less than 0.001) have been submitted
748    to GenBank and accession numbers are included in Table S20.
749

750    RNA was isolated from adult ctenophores that were collected in either the Gulf of
751    California or Monterey Bay. Total RNA was extracted using RNeasy kit (Qiagen).
752    Preparation of RNA-seq libraries was done using Illumina TruSeq kit for paired end
753    reads. Total RNA was sequenced at the University of Utah. Sequencing was done using
754    the Illumina HiSeq2000 platform on a paired-end protocol with 100 cycles. Sequence
755    data was assembled as described in Francis and coworkers, 2013 (*111*). The following
756    number of unique transcripts were produced: *B. chuni*=399,516, *B. forskalii*=194,616, *C.*
757    *fugiens*=285,056, *E. dunlapae*=276,866, *H. californensis*=528,793, *L. lactea*=198,516,
758    and *T. inconstans*=194,483. The high numbers reflect minor assembly differences at
759    single loci.
760

761    <u>General pipeline for analyzing additional ctenophore transcripts</u>
762

763    One or more human RefSeq proteins were used as a TBLASTN query against the seven
764    ctenophore transcriptomes and a Trinity assembly of the *M. leidyi* RNA-seq data. All hits
765    with E-values ≤ 0.001 were considered candidates. These candidate transcripts were used
766    as BLASTX queries against the Human RefSeq database (edited to include only a single
767    isoform per Entrez Gene ID). Any human sequence that among the top 10 hits (with E-
768    values ≤ 0.001) for any of the ctenophore transcripts was marked for inclusion in our
769    downstream phylogenetic analyses. All candidate ctenophore transcripts were translated
770    using the r2012-08-15 version of TransDecoder (http://transdecoder.sourceforge.net/)
771    with the following command: transcripts_to_best_scoring_ORFs.pl -t candidates.fa --
772    search_pfam Pfam-A.hmm. All translations of ctenophore candidates and marked human
773    proteins were aligned using MAFFT v6.864b (*85*) with the following command: mafft –
774    auto FASTAFILE. The resulting alignment was trimmed with Gblocks version 0.91b (*86*)
775    using the following command: Gblocks -b2=Z -b3=10 -b4=5 -b5=a, where Z = 0.65 x the
776    number of sequences. Sequences where occupancy was less than 50% were removed. We
777    used RAxML version 7.2.8 (*112*) to generate a maximum likelihood tree and 10
778    bootstraps with the following command: raxmlHPC -f a -x RANDOMSEED -m
779    PROTGAMMALG -p RANDOMSEED -N 10 -n NAME -s ALIGNMENT.phy.
780

781   If there is weak evidence uniting any ctenophore candidate with the human gene being
782   tested, the pipeline is rerun on the subset of best candidates. These reruns often lead to
783   longer alignments and better resolution in the resulting trees.
784
785   A small percentage of ctenophore hits were not translated by transdecode. Of the 263
786   total hits that were in the top five for each query, 12 were missed by transdecode. We
787   performed manual analysis on these 12.
788
789   <u>Search for Neuroligin (synaptic scaffolding) in other ctenophore transcriptomes</u>
790
791   Human NLGN1 (accession= NP_055747) was used as a TBLASTN query against the
792   seven ctenophore transcriptomes and the Trinity assembly of the *M. leidyi* RNA-seq data.
793   We generated an alignment of 92 unique ctenophore sequences and 20 human proteins
794   (NCBI Gene ids = AADAC, AADACL2, AADACL3, AADACL4, ACHE, BCHE, CEL,
795   CES1, CES2, CES3, CES4A, CES5A, LIPE, NLGN1, NLGN2, NLGN3, NLGN4X,
796   RPL3, RPL3L, and TG) and a maximum likelihood tree using the general pipeline
797   described in S13.3. The human NLGN formed a polytomy with a small subset of the
798   ctenophore sequences as well as LIPE and the AADAC proteins (zero bootstrap support).
799   The majority of ctenophore sequences formed a large clade that was most closely related
800   to the CES proteins.
801
802   We next removed this large clade of ctenophores and reran the pipeline from S13.3. In
803   this tree, a very long branch *L. lactea* sequence is sister to NLGN3 with a bootstrap value
804   of 20; the highest bootstrap supporting this sequence and the NLGN sequences is 50. An
805   *E. dunlapae* sequence is sister to all the NLGN proteins plus the *L. lactea* sequence with
806   a support value of 10.
807
808   We next reran the pipeline, this time including only the human TG, NLGN and RPL3
809   proteins, as well as the *L. lactea* and *E. dunlapae* sequences. With the subtraction of most
810   of the ctenophore sequence, the number of columns in the alignment expanded from 59 to
811   532. Very few of the *L. lactea* sequence aligned – upon examination the 199 amino acids
812   of this sequence included seven stop codons. This sequence was removed from this
813   analysis based on it having less than 50% occupancy. The resulting tree had the *E.*
814   *dunlapae* sequence as sister to the RPL proteins (or sister to TG + NLGN proteins) with
815   80% bootstrap support. We were unable to find additional candidates using *Nematostella*
816   *vectensis* sequences as queries.
817
818   None of these analyses show any reasonable support for a ctenophore neuroligin and the
819   final analysis shows fairly strong support that the best candidate is not a neuroligin.
820
821   <u>Search for Ionotropic glutamate receptors (synaptic scaffolding) in other ctenophore</u>
822   <u>transcriptomes</u>
823
824   We used human GRIA2 (accession= NP_000817), GRIN1 (accession= NP_000823),
825   GRIK2 (accession= NP_001159719), and GRID2 (accession= NP_001501) as
826   TBLASTN queries against the seven ctenophore transcriptomes and the Trinity assembly

827    of the *M. leidyi* RNA-seq data. GRIA2 is an AMPA-type glutamate receptor, GRIN1 is
828    an NMDA-type glutamate receptor, GRIK2 is a kainate-type glutamate receptor, and
829    GRID2 is a delta2-like glutamate receptor.
830
831    We generated an alignment starting with 136 ctenophore sequences (including Trinity
832    isoforms) and 28 human proteins (NCBI Gene ids = CASR, GRIA1, GRIA2, GRIA3,
833    GRID1, GRID2, GRIK1, GRIK2, GRIK3, GRIK4, GRIK5, GRIN1, GRIN2A, GRIN2B,
834    GRIN2C, GRIN2D, GRIN3A, GRIN3B, GRM1, GRM2, GRM3, GRM4, GRM5, GRM6,
835    GRM7, GRM8, and TAS1R3). Our pipeline produced an alignment of only four columns.
836    We next aligned each ctenophore sequence separately against the human set (using the
837    pipeline alignment procedure) and kept only those in which the resulting alignment
838    included less than 50% gaps in the ctenophore sequence. This resulted in 34 ctenophore
839    sequences (including two *M. leidyi* sequences and one published *Pleurobrachia* sequence
840    with NCBI accession ADV31314). The resulting tree led to a clade of human ionotropic
841    glutamate receptor (iGluR) sequences and a clade of ctenophore sequences that were
842    bisected by the non-iGluR human outgroup sequences. Given the length of the outgroup
843    branch and the fact that the top BLAST hits of all the ctenophore sequences were the
844    iGluR human sequences, we surmise that the attachment of the outgroup branch is likely
845    unreliable. When we run this tree without outgroups the ctenophore sequences form a
846    monophyletic group with 100% bootstrap support. We were unable to find additional
847    candidates using *Capitella telata* sequences as queries.
848
849    These analyses support the ctenophore iGluR candidates being descendants of an
850    ancestral sequence that would later give rise to AMPA, NMDA, kainate-type, and GRID
851    glutamate receptors after bilaterians split from ctenophores. As such, ctenophores do not
852    have bona fide AMPA, NMDA, kainate-type, and GRID iGluRs (Fig. S3).
853
854    <u>Search for Twist (myogenic genes) in other ctenophore transcriptomes</u>
855
856    We used human TWIST1 (accession=NP_000465) as a TBLASTN query against the
857    seven ctenophore transcriptomes and the Trinity assembly of the *M. leidyi* RNA-seq data.
858    We generated an alignment of 30 ctenophore and 20 human proteins (NCBI Gene ids =
859    TWIST1, TWIST2, HAND2, HAND1, TCF15, SCXB, SCXA, TCF21, MSC, ATOH7,
860    NEUROD1, NHLH2, NHLH1, TAL1, TAL2, and LYL1) and a maximum likelihood tree
861    using the general pipeline described above. The original analysis produced a clade that
862    included the human Twist genes, a *M. leidyi* sequence, and two *E. dunlapae* sequences
863    with a support value of 60%. In this same tree there was a large clade of ctenophore
864    sequences that grouped sister to all the human genes. To see if we could get better
865    support for the Twist clade, we removed all ctenophore sequences except for the three
866    that formed a clade with human Twist and reran the analysis. An unrooted view of this
867    tree results in a clade of human LYL1 and TAL proteins sister to the ctenophore
868    sequences with a support value of 40%.
869
870    We ran a third analysis this time running the standard pipeline but running BLAST
871    without SEG filtering. This yielded more ctenophore sequences and a larger alignment:
872    80 columns from 97 sequences as opposed to 66 columns from 46 sequences in the first

873 analysis. Again we recover the majority of ctenophore sequences in a cluster separate
874 from the human sequences. The only exceptions are (1) a very long branch of *B. forskalii*
875 sequences that is sister to a clade that includes human FERD3L and human PTF1A, and
876 (2) a longish *E. dunlapae* branch that is sister to human ID4 sequence. In this tree, the
877 human TWIST1 and TWIST2 sequences form a clade with 18 of the other 26 human
878 proteins. We were unable to find additional candidates using *Capitella telata* sequences
879 as queries.
880
881 From these analyses, it seems unlikely that a true Twist gene exists in ctenophores.
882
883 <u>Search for Snail (myogenic genes) in other ctenophore transcriptomes</u>
884
885 We used human SNAI1 (accession=NP_005976) as a TBLASTN query against the seven
886 ctenophore transcriptomes and the Trinity assembly of the *M. leidyi* RNA-seq data. We
887 generated an alignment of 1,431 ctenophore and 547 human proteins. The alignment that
888 resulted from MAFFT and Gblocks consisted of only seven columns and any tree
889 resulting from this alignment would not be informative, consequently we must rely on
890 reciprocal best BLAST. We identified one *T. inconstans* sequence and two *E. dunlapae*
891 sequences, which are reciprocal best BLAST hits with human Snail genes. This result
892 suggests that Snail was present in the last common ctenophore ancestor, but was lost in
893 the lineage leading to *M. leidyi*. We were unable to find additional candidates using
894 *Capitella telata* sequences as queries.
895
896 <u>Search for Lbx/ladybird, NK4/tinman, NK3/bagpipe, and NK2/vnd (myogenic genes) in</u>
897 <u>other ctenophore transcriptomes</u>
898
899 We fetched the human ANTP class homeodomains from HomeoDB (*113*) and used these
900 sequences as a TBLASTN query against the seven ctenophore transcriptomes and the
901 Trinity assembly of the *M. leidyi* RNA-seq data. After translating transcripts that were
902 hit, we filtered hits based on the presence of the following: 1) complete homeodomains
903 with "WFQN," which is present in all ANTP homeodomains in HomeoDB, 2) absence of
904 LIM, Homez, Pou, Homeobox_KN, and/or PBC domains (these are not associated with
905 ANTP homeodomains). This list produced 822 candidate ANTP homeodomains, of
906 which 124 were unique. We aligned these ctenophore homeoboxes to the ANTP
907 alignment included in the additional file 1 of Ryan and coworkers, 2010 (*114*). We then
908 ran RAxML with 100 fast bootstraps. We found no ctenophore sequences in the
909 Lbx/ladybird, NK4/tinman, NK3/bagpipe, or NK2/vnd clades. Furthermore, we found no
910 evidence of ctenophore Hox or ParaHox genes in this analysis.
911
912 <u>Search for Myf5, Mrf4, Myogenin, and MyoD (myogenic genes) in other ctenophore</u>
913 <u>transcriptomes</u>
914
915 We used human Myf5 (accession=NP_005584), Mrf4 (accession=NP_002460),
916 Myogenin (accession=NP_002470), and MyoD (accession=NP_002469) as a TBLASTN
917 query against the seven ctenophore transcriptomes and the Trinity assembly of the *M.*
918 *leidyi* RNA-seq data. We ran TBLASTN both with and without SEG filtering and found

919  no hits with E-values below 0.001. These human genes are identifiable by TBLASTN to
920  protostome transcriptomes at much lower E-values. We were unable to find additional
921  candidates using *Capitella telata* sequences as queries.
922
923  <u>Search for Noggin (myogenic genes) in other ctenophore transcriptomes</u>
924
925  We used human NOG (accession=NP_005441) as a TBLASTN query against the seven
926  ctenophore transcriptomes and the Trinity assembly of the *M. leidyi* RNA-seq data. The
927  search produced no hits with E-values < 0.2. We ran the searches with and without SEG
928  filtering. This human gene is identifiable by TBLASTN to protostome transcriptomes at
929  much lower E-values. We were unable to find additional candidates using *Capitella*
930  *telata* sequences as queries.
931
932  <u>Search for Eomesodermin (myogenic genes) in other ctenophore transcriptomes</u>
933
934  We used human Eomes (accession=NP_005433) as a TBLASTN query against the seven
935  ctenophore transcriptomes and the Trinity assembly of the *M. leidyi* RNA-seq data. We
936  generated an alignment of 126 ctenophore (58 unique) and 21 human proteins. The
937  alignment that resulted from MAFFT and Gblocks consisted of 177 columns. The
938  RAxML analysis produced a clade with the human proteins Tbr1 and Eomes,
939  (bootstrap=90%). Sister to this clade was human TBX21 (bootstrap=30%). These results
940  suggest that that ctenophore sequences are not direct orthologs to human Eomes. We
941  were unable to find additional candidates using *Danio rerio* sequences as queries.
942
943  <u>Search for GATA (myogenic genes) in other ctenophore transcriptomes</u>
944
945  We used human GATA1 (accession=NP_002040) as a TBLASTN query against the
946  seven ctenophore transcriptomes and the Trinity assembly of the *M. leidyi* RNA-seq data.
947  We generated an alignment of 60 candidate transcripts, of which 33 were unique. There
948  were 31 columns in the alignment. The resulting tree produced a clade that included
949  ctenophore sequences from six species of ctenophores along with the human GATA
950  genes and the human TRPS1 and ZGLP1 genes (bootstrap=60). Like GATA proteins in
951  other animals, these ctenophore sequences contain two GATA zinc finger domains
952  (PFAM accession=PF00320). From this analysis we can say that GATA is present in *B.*
953  *chuni*, *B. forskalii, E. dunlapae, H. californensis*, and *T. inconstans*, and was lost in *M.*
954  *leidyi*.
955
956  <u>Search for Troponin (myogenic genes) in other ctenophore transcriptomes</u>
957
958  We used human TNNI1 (accession=NP_003272) as a TBLASTN query against the seven
959  ctenophore transcriptomes and the Trinity assembly of the *M. leidyi* RNA-seq data. The
960  search produced no hits with E-Value < 0.001. We ran the searches with and without
961  SEG filtering. This human gene is identifiable by TBLASTN to protostome
962  transcriptomes at much lower E-values.
963
964  <u>Search for FGF (myogenic genes) in other ctenophore transcriptomes</u>

965
966    We used human FGF1 (accession=NP_001138364) as a TBLASTN query against the
967    seven ctenophore transcriptomes and the Trinity assembly of the *M. leidyi* RNA-seq data.
968    The search produced no hits with E-values < 0.001. We ran the searches with and without
969    SEG filtering. This human gene is identifiable by TBLASTN to protostome
970    transcriptomes at much lower E-values.
971
972    <u>Search for Nodal (myogenic genes) in other ctenophore transcriptomes</u>
973
974    We used human Nodal (accession=NP_060525) as a TBLASTN query against the seven
975    ctenophore transcriptomes and the Trinity assembly of the *M. leidyi* RNA-seq data. We
976    generated an alignment with 145 ctenophore sequences (50 unique) and 31 human
977    sequences. The RAxML analysis produced a clade that included human Nodal with seven
978    ctenophore sequences from four species. Support for this clade was 10%. Included in this
979    clade is a *M. leidyi* sequence (accession=AEP16385) that was analyzed as part of a
980    genome analysis of *M. leidyi* TGF-beta and was named Tgf2 (*115*). In this detailed study,
981    it was determined that no *M. leidyi* sequence (including Tgf2) was orthologous to Nodal.
982
983    To test whether the low support of the clade of Nodal and ctenophore sequences was due
984    to long ctenophore branches, we reran this analysis, pruning all ctenophore sequences
985    besides those that appeared in the Nodal clade. In this analysis only two of the ctenophore
986    sequences formed a clade with human Nodal. The bootstrap support for this clade was
987    20%. The other ctenophore sequences formed poorly supported clades with other human
988    sequences. To test whether the low support of this Nodal clade may have been due to the
989    longer human branches, we ran another maximimum-likelihood analysis on a smaller
990    subset of human sequences. In this tree, the two ctenophore NODAL proteins formed a
991    sister clade to a clade of Nodal and GDF5 with 20% bootstrap support. From this analysis
992    and consistent with our previous work (*115*), we conclude that there is not convincing
993    evidence that a Nodal ortholog is present in either *M. leidyi* or the other seven ctenophore
994    transcriptomes that we have investigated.
995
996    <u>Search for Shh/hh (myogenic genes) in other ctenophore transcriptomes</u>
997
998    We used human Shh (accession= NP_000184) as a TBLASTN query against against the
999    seven ctenophore transcriptomes and the Trinity assembly of the *M. leidyi* RNA-seq data.
1000   We recovered two hits to partial *B. chuni* transcripts. One coded 266 amino acids and the
1001   other coded 340 amino acids. The latter was incomplete at the C-terminal end. Like the
1002   full length *M. leidyi* candidate (ML073718a) and unlike the hedgehog proteins from
1003   cnidarians and bilaterians, a "Hog" domain (Pfam accession=PF01079) was present in
1004   these sequences, but there no "Hedge" domain (Pfam accession=PF01085) was present
1005   (Fig. S10). Turning off SEG filtering for the same TBLASTN search, we identified an
1006   additional two *B. chuni* candidates. The two were 970 and 981 amino acids long. Both
1007   contained "Hog" domains, but no "Hedge" domains. These two also contained
1008   MMR_HSR1 and MACPF domains. As such, we see no proper hedgehog gene in any of
1009   the ctenophore transcriptomes.
1010

1011  <u>Genome browser</u>
1012
1013  As part of this project, we have launched the *Mnemiopsis leidyi* genome Project Portal.
1014  At this site, users can perform BLAST searches against the genome sequence, proteome,
1015  and transcriptome. An interface allows users to retrieve individual sequences or scaffolds,
1016  and assemble custom data sets. We have implemented JBrowse, a JavaScript-based
1017  genome browser (*116*), for viewing the *M. leidyi* genome assembly, gene predictions,
1018  RNA-seq data, and public EST and mRNA sequences. Gene pages are editable wiki
1019  pages and users are encouraged to annotate sequences of interest. The *Mnemiopsis* web
1020  site is freely accessible at http://research.nhgri.nih.gov/mnemiopsis/
1021
1022  Genomic scaffolds are named using a six-character convention (e.g. MLXXXX); the ML
1023  designates the species (*Mnemiopsis leidyi*), and the individual scaffolds are numbered
1024  from 0001 to 5100 (e.g., ML4323). Gene identifiers (e.g., ML103316a) are prefixed with
1025  the scaffold on which the gene is located (in this example, "ML1033"), followed a non-
1026  padded integer that is unique in combination with the scaffold identifier (in this case
1027  "16"), followed by a lower-case letter that corresponds to the genes isoform (in this case
1028  "a"). A gene identifier is typically (but not always) ordered by its most 5' position on the
1029  scaffold. Newly added genes are assigned the next integer, regardless of its position on
1030  the scaffold.
1031
1032
1033

1034 **Supplementary Table S1: Summary of phylogenetic placements of ctenophores**
1035
1036 Below is a table of results from a series of phylogenetic analyses that include
1037 ctenophores. The result is left blank if the reference refers only to the relationship of
1038 ctenophores rather than the branching pattern of the five earliest branching lineages.
1039 Abreviations are as follows: Porifera(Po), Ctenophora(Ct), Placozoa(Tr), Cnidaria(Cn)
1040 and Bilateria(Bi). Multiple "Po" entries indicate that a study included data from multiple
1041 poriferans and recovered a paraphyletic Porifera. Parahoxozoa indicates a monophyletic
1042 group consisting of Cnidaria, Bilateria and Placozoa (if included). Taxa in square
1043 brackets indicates a sister relationship to all other animals based on a rooted tree.
1044 Modeled after Table 1 of Wallberg and coworkers, 2004 (*117*) and Ryan and Baxevanis
1045 2007 (*118*).
1046

| Authors | Year | Data/Meth. | Result | Hypothesis |
|---|---|---|---|---|
| Lang (*119*) | 1884 | Morph-MP | | Ct->Platyhelminthes |
| Hyman (*120*) | 1940 | Morph | (Po,((Cn,Ct),Bi)) | Coelenterata |
| Hadzi (*121*) | 1953 | Morph | | Ct->Protostomia |
| Brusca & Brusca (*122*) | 1990 | Morph | | Coelenterata |
| Ehlers (*123*) | 1993 | Morph | | (Ct,Bi) |
| Ruppert & Barnes (*124*) | 1994 | Morph | | Coelenterata |
| Nielsen (*125*) | 1995 | Morph | (Po,(Tr,(Cn,(Bi,Ct,Bi)))) | Ct->Deuterostomia |
| Nielsen (*126*) | 1996 | Morph | (Po,(Tr,(Cn,(Ct,Bi)))) | (Ct,Bi) |
| Ax (*127*) | 1996 | Morph | | (Ct,Bi) |
| Margulis & Schwartz (*128*) | 1998 | Morph | | Coelenterata |
| Nielsen (*129*) | 2001 | Morph | | (Ct,Bi) |
| Wainright et al. (*130*) | 1993 | Ribo-ML | (Po,(Ct,((Tr,Cn),Bi))) | Parahoxozoa [Po] |
| Katayama et al. (*131*) | 1995 | Ribo-DI | ((((Po,Ct),Tr),Cn),Bi) | Diploblastica |
| "" | | Ribo-MP | ((((Po,Ct),Tr),Cn),Bi) | Diploblastica |
| "" | | Ribo-ML | ((Po,(Ct,Tr)),(Cn,Bi)) | |

| | | | | |
|---|---|---|---|---|
| Hanelt et al. (*132*) | 1996 | Ribo-DI | ((((Po,Ct),Tr),Cn),Bi) | Diploblastica [Bi] |
| Van de Peer & Wachter (*133*) | 1997 | Ribo-DI | (((Po,(Po,Ct)),(Tr,Cn)),Bi) | Parahoxozoa (Po,Ct) [Bi] |
| Abouheif et al. (*134*) | 1998 | Ribo-MP | (Po,(Ct,(Tr,(Cn,Bi)))) | Parahoxozoa [Po] |
| Collins (*135*) | 1998 | Ribo-MP | (Po,(Po,Ct,(Tr,(Cn,Bi))) | Parahoxozoa [Po] |
| "" | | Ribo-ML | (Po,((Po,Ct),(Tr,(Cn,Bi)))) | Parahoxozoa (Po,Ct) [Po] |
| Halanych (*136*) | 1998 | Ribo-MP | (Po,(Tr,((Cn,(Ct,Cn)),Bi))) ; | Coelenterata [Po] |
| "" | | Ribo-MP | (Po,(Tr,(Ct,Cn,Bi)))) | |
| Lipscomb et al. (*137*) | 1998 | Ribo-MP | (Po,Po,Ct,((Tr,Cn),Bi)); | Parahoxozoa [Po] |
| Winnepenninckx et al. (*138*) | 1998 | Ribo-DI | ((Po,(Po,Ct)),(Cn,(Tr,Bi))) | Parahoxozoa [Po] |
| Zrzavý et al. (*139*) | 1998 | Ribo-MP | ((Po,(Po,Ct)),(Tr,(Cn,(Cn, Bi)))) | Parahoxozoa [Po] |
| Kim et al. (*140*) | 1999 | Ribo-DI | (Po,(Po,(Ct,(Tr,(Cn,Bi))))) | Parahoxozoa [Po] |
| "" | | Ribo-ML | (Po,(Po,(Ct,(Tr,(Cn,Bi))))) | Parahoxozoa [Po] |
| Giribet et al. (*141*) | 1999 | Ribo-MP | (Po,(Ct,(Cn,(Tr,Bi)))) | Parahoxozoa [Po] |
| Siddall & Whiting (*142*) | 1999 | Ribo-MP | ((Po,Ct),(Cn,(Tr,Bi))) | Parahoxozoa (Po,Ct) |
| Medina et al. (*143*) | 2001 | Ribo-ML | (Po,(Po,(Ct,(Cn,Bi)))) | Parahoxozoa [Po] |
| "" | | Ribo-ML | ((Cn,Ct),(Po,(Po,Bi))) | Coelenterata |
| "" | | Ribo-ML | (Po,(Ct,(Cn,Bi)) | Parahoxozoa [Po] |
| "" | | Ribo-MP | (Po,Po,(Ct,(Cn,Bi))) | Parahoxozoa [Po] |
| Peterson & Eernisse (*144*) | 2001 | Ribo-MP | (Po,(Po,Ct,(Tr,(Cn,Bi)))) | Parahoxozoa [Po] |
| Podar et al. (*145*) | 2001 | Ribo-ML | (Po,(Ct,(Tr,(Cn,Bi)))) | Parahoxozoa [Po] |

| | | | | |
|---|---|---|---|---|
| Collins et al. (*146*) | 2002 | Ribo-MP | (Po,(Po,(Ct,(Cn,(Tr,Bi))))) | Parahoxozoa [Po] |
| Jondelius et al. (*135*) | 2002 | Ribo-ML | (Po,(Po,(Ct,(Cn,(Tr,Bi))))) | Parahoxozoa [Po] |
| Martinelli & Spring (*147*) | 2003 | Ribo-ML | (Po,(((Ct,Tr),Cn),Bi)) | |
| Zrzavý & Hypša (*148*) | 2003 | Ribo-MP | (Po,(Ct,((Tr,Cn),Bi))) | Parahoxozoa [Po] |
| "" | | Ribo-MP | (Po,(Ct,(Tr,(Cn,Bi)))) | Parahoxozoa [Po] |
| "" | | Ribo-MP | (Po,((Ct,Tr),(Cn,Bi))) | |
| Wallberg et al. (*117*) | 2004 | Ribo-MP | (Po,Po,(Ct,(Cn,(Tr,Bi)))) | Parahoxozoa [Po] |
| Dunn et al. (*149*) | 2008 | EST-ML | (Ct,((Po,Cn),Bi)) | [Ct] |
| Philippe et al. (*150*) | 2009 | EST-Ba | (Po,(Tr,((Cn,Ct),Bi))) | Coelenterata [Po] |
| Hejnol et al. (*88*) | 2009 | EST-ML | (Ct,(Po,(Tr,(Po,(Cn,Bi))))) | [Ct] |
| Schierwater et al. (*151*) | 2009 | MT,EST, Morph | (Bi,(Tr,(Po,(Cn,Ct)))) | [Bi] |
| Pick et al. (*152*) | 2010 | EST-Ba | (Po,(Ct,(Cn,(Pl,Bi)))) | Parahoxozoa [Po] |
| Srivastava et al. (*80*) | 2010 | Geno-ML | (Ct,(Po,(Pl,(Cn,Bi)))) | Parahoxozoa [Ct] |
| Mallat et al. (*153*) | 2012 | rRNA-ML | ((Po,Ct),(Cn,Tr),Bi) | Parahoxozoa (Po,Ct) |

1047
1048
1049 **Supplementary Table S2: Nested intronic genes in *Mnemiopsis leidyi* and other**
1050 **animal genomes**
1051 Numbers for bilaterian species are from Kumar, 2009 and references therein (*154*).
1052 Genome assembly versions from which numbers were generated are in parentheses after
1053 gene number. "% Nested" indicates the percent of genes in the genome that nested genes
1054 comprise.
1055

| Species | Number of Nested Intronic Genes | Total Number of Genes | % Nested |
|---|---|---|---|
| *M. leidyi* | 1,323 | 16,554 | ~8.0 |
| *D. melanogaster* | 792 | 14,601 (r5.1) | ~5.4 |

| | | | |
|---|---|---|---|
| *C. elegans* | 429 | 20,061 (WSI176) | ~2.1 |
| *C. briggsae* | 233 | 19,500 | ~1.2 |
| *H. sapiens* | 158 | 28,755 (r36.2) | ~0.5 |

1056

1057 **Supplementary Table S3: Genomic repeat content as detected by VMatch**

1058

| Organism | Total bases | Masked bases | % Masked |
|---|---|---|---|
| *M. leidyi* | 155,865,547 | 15,693,038 | 10.07 |
| *A. queenslandica* | 147,463,102 | 15,964,649 | 10.83 |
| *M. brevicollis* | 41,633,360 | 5,135,212 | 12.33 |
| *N. vectensis* | 356,613,585 | 88,039,756 | 24.69 |
| *T. adhaerens* | 105,632,827 | 2,220,249 | 2.10 |
| *S. rosetta* | 55,440,309 | 4,644,212 | 8.38 |

1059

1060 **Supplementary Table S4: Known genomic repeats (RepeatMasker)**
1061 *Ml=Mnemiopsis leidyi, Aq=Amphimedon queenslandica, Mb=Monosiga brevicolis,*
1062 *Ta=Trichoplax adhaerens, Sr=Salpingoeca rosetta*

1063

| | *Ml* | *Aq* | *Mb* | *Ta* | *Sr* |
|---|---|---|---|---|---|
| Sequences | 5100 | 3579 | 218 | 1415 | 154 |
| Total length (bp) | 155,865,547 | 147,463,102 | 41,633,360 | 105,632,827 | 55,440,309 |
| GC level (%) | 38.86% | 35.87% | 54.89% | 32.74% | 56.01% |
| bases masked (bp) | 3,874,208 | 5,678,584 | 1,258,976 | 1,541,813 | 9,862,436 |
| bases masked (%) | 2.49% | 3.85% | 3.02% | 1.46% | 17.79% |
| **Retroelements** | 0.49% | 0.93% | 0.23% | 0.06% | 1.68% |
| **SINEs** | 0.01% | 0.00% | 0.01% | 0.00% | 0.00% |
| **Penelope** | 0.01% | 0.05% | 0.03% | 0.00% | 0.02% |
| **LINEs** | 0.38% | 0.25% | 0.09% | 0.02% | 0.60% |
| CRE/SLACS | 0.00% | 0.00% | 0.00% | 0.00% | 0.00% |
| L2/CR1/Rex | 0.31% | 0.18% | 0.00% | 0.01% | 0.11% |
| R1/LOA/Jockey | 0.01% | 0.01% | 0.03% | 0.00% | 0.25% |
| R2/R4/NeSL | 0.00% | 0.00% | 0.00% | 0.00% | 0.00% |
| RTE/Bov-B | 0.02% | 0.00% | 0.00% | 0.00% | 0.00% |
| L1/CIN4 | 0.01% | 0.01% | 0.01% | 0.00% | 0.01% |
| **LTR elements** | 0.11% | 0.67% | 0.13% | 0.04% | 1.08% |
| BEL/Pao | 0.01% | 0.06% | 0.01% | 0.00% | 0.01% |
| Ty1/Copia | 0.00% | 0.09% | 0.04% | 0.00% | 0.04% |
| Gypsy/DIRS1 | 0.05% | 0.48% | 0.03% | 0.03% | 0.30% |

| | | | | | |
|---|---|---|---|---|---|
| Retroviral | 0.02% | 0.02% | 0.03% | 0.01% | 0.08% |
| **DNA transposons** | 0.26% | 0.41% | 0.17% | 0.04% | 3.88% |
| hobo-Activator | 0.05% | 0.12% | 0.08% | 0.01% | 0.35% |
| Tc1-IS630-Pogo | 0.00% | 0.06% | 0.01% | 0.00% | 0.09% |
| En-Spm | 0.02% | 0.03% | 0.02% | 0.01% | 0.55% |
| MuDR-IS905 | 0.00% | 0.00% | 0.00% | 0.00% | 0.00% |
| PiggyBac | 0.00% | 0.01% | 0.00% | 0.00% | 0.01% |
| Tourist/Harbinger | 0.00% | 0.04% | 0.00% | 0.00% | 0.00% |
| Other (Mirage, P-element, Transib) | 0.01% | 0.01% | 0.00% | 0.00% | 0.00% |
| **Rolling-circles** | 0.00% | 0.00% | 0.00% | 0.00% | 0.00% |
| **Unclassified** | 0.02% | 0.04% | 0.00% | 0.01% | 0.01% |
| **Total interspersed repeats** | 0.77% | 1.38% | 0.40% | 0.11% | 5.57% |
| **Small RNA** | 0.03% | 0.01% | 0.04% | 0.03% | 0.10% |
| **Satellites** | 0.05% | 0.01% | 0.01% | 0.01% | 0.04% |
| **Simple repeats** | 1.16% | 1.05% | 1.97% | 0.68% | 11.58% |
| **Low complexity** | 0.49% | 1.40% | 0.61% | 0.64% | 0.51% |

1064
1065 **Supplementary Table S5: Scaffold size frequencies**
1066

| Scaffold Size | Frequency | Total base pairs |
|---|---|---|
| **> 2 kbp** | 2342 | 152300948 |
| **> 5 kbp** | 1517 | 149739537 |
| **> 10 kbp** | 1253 | 147912957 |
| **> 20 kbp** | 1076 | 145345415 |
| **> 50 kbp** | 794 | 135886036 |
| **> 100 kbp** | 509 | 115294090 |
| **> 120 kbp** | 429 | 106527368 |
| **> 200 kbp** | 215 | 72812389 |
| **> 500 kbp** | 27 | 16400572 |
| **> 1 Mbp** | 1 | 1222598 |

1067
1068 **Supplementary Table S6: Scaffold gap frequencies**
1069

| Gap Size | Frequency | Total gaps |
|---|---|---|

| | | |
|---|---|---|
| **>= 10 bp** | 1882 | 5527158 |
| **>= 100 bp** | 1508 | 5520825 |
| **>= 1 kbp** | 1274 | 5378496 |
| **>= 5 kbp** | 350 | 3040555 |
| **>= 10 kbp** | 90 | 1252099 |

1070

**Supplementary Table S7: GC content (% of genome)**
Genome-wide GC content was calculated for *M. leidyi* and several other species. The GC content of the *M. leidyi* genome in the same range as *A. queenslandica*, *N. vectensis*, and *H. sapiens*. *Hydra magnipapillata (Hm), Trichoplax adhaerens (Ta), Amphimedon queenslandica (Aq), Mnemiopsis leidyi (Ml), Nematostella vectensis (Nv), Homo sapiens (Hs), Monosiga brevicolis (Mb), Salpingoeca rosetta (Sr)*

| | *Hm* | *Ta* | *Aq* | *Ml* | *Nv* | *Hs* | *Dm* | *Mb* | *Sr* |
|---|---|---|---|---|---|---|---|---|---|
| **GC** | 28 | 32 | 36 | 38 | 40 | 40 | 42 | 54 | 56 |

1078

**Supplementary Table S8: Dinucleotide odds ratios**
The quantity [XpY]/[X][Y] was calculated over all dinucleotides for *M. leidyi* and other selected species. The mean value of each dinucleotide ratio for these species was also calculated. The dinucleotide odds ratios for *M. leidyi* show little notable divergence from the calculated mean values of each dinucleotide of the other sampled genomes.
*Mnemiopsis leidyi* (*Ml*), *Amphimedon queenslandica* (*Aq*), *Hydra magnipapillata (Hm)*, *Monosiga brevicolis* (*Mb*), *Nematostella vectensis* (*Nv*), *Trichoplax adhaerens* (*Ta*), *Salpingoeca rosetta* (*Sr*), *Drosophila melanogaster* (*Dm*), *Homo sapiens* (*Hs*)

1087

| **Dinucleotide** | *Ml* | *Aq* | *Hm* | *Mb* | *Nv* | *Ta* | *Sr* | *Dm* | *Hs* | μ |
|---|---|---|---|---|---|---|---|---|---|---|
| **AA** | 1.04 | 0.98 | 1.16 | 1.13 | 1.11 | 0.95 | 0.83 | 1.19 | 1.11 | 1.06 |
| **AC** | 1.02 | 1.04 | 0.99 | 0.97 | 1.00 | 0.92 | 1.30 | 0.82 | 0.83 | 0.99 |
| **AG** | 1.02 | 1.04 | 0.99 | 0.97 | 1.00 | 1.10 | 0.81 | 0.82 | 1.17 | 0.99 |
| **AT** | 0.83 | 0.98 | 0.93 | 0.95 | 0.89 | 1.04 | 0.83 | 0.95 | 0.89 | 0.92 |
| **CA** | 1.19 | 1.22 | 1.19 | 1.29 | 1.17 | 1.10 | 1.79 | 1.15 | 1.17 | 1.25 |
| **CC** | 1.11 | 0.93 | 1.02 | 0.96 | 1.25 | 0.78 | 0.77 | 1.13 | 1.25 | 1.02 |
| **CG** | 0.83 | 0.31 | 0.51 | 0.82 | 0.75 | 0.78 | 0.89 | 0.91 | 0.25 | 0.67 |
| **CT** | 1.02 | 1.04 | 0.99 | 0.97 | 1.00 | 1.10 | 0.81 | 0.82 | 1.17 | 0.99 |
| **GA** | 1.02 | 1.04 | 0.79 | 0.97 | 1.00 | 0.92 | 0.81 | 0.99 | 1.00 | 0.95 |
| **GC** | 1.11 | 0.93 | 1.02 | 1.23 | 1.00 | 1.17 | 1.15 | 1.36 | 1.00 | 1.11 |

| | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|
| **GG** | 1.11 | 0.93 | 1.02 | 0.96 | 1.25 | 0.78 | 0.77 | 1.13 | 1.25 | 1.02 |
| **GT** | 1.02 | 1.04 | 0.99 | 0.97 | 1.00 | 0.92 | 1.30 | 0.82 | 0.83 | 0.99 |
| **TA** | 0.83 | 0.88 | 0.85 | 0.38 | 0.78 | 0.95 | 0.41 | 0.71 | 0.78 | 0.73 |
| **TC** | 1.02 | 1.04 | 0.79 | 0.97 | 1.00 | 0.92 | 0.81 | 0.99 | 1.00 | 0.95 |
| **TG** | 1.19 | 1.22 | 1.19 | 1.29 | 1.17 | 1.10 | 1.79 | 1.15 | 1.17 | 1.25 |
| **TT** | 1.04 | 0.98 | 1.16 | 1.13 | 1.11 | 0.95 | 0.83 | 1.19 | 1.11 | 1.06 |

1088
1089
1090 **Supplementary Table S9: Summary statistics for final gene models**
1091

| | Mean (bp) | Mode (bp) |
|---|---|---|
| **Intron Length** | 898 | 297 |
| **Length of Intergenic Regions** | 2463 | 45492 |
| **Exon Length** | 314 | 132 |
| **Length of Predicted Transcript** | 5799 | 315 |
| **Protein Length** | 463 | 132 |
| **Scaffold Length** | 30562 | 1012 |

1092
1093
1094
1095
1096 **Supplementary Table S10: Total coding percentage of *Mnemiopsis leidyi* genome**
1097

| | Percent |
|---|---|
| **Coding** | 14.76 |
| **Non-coding** | 85.24 |

1098
1099 **Supplementary Table S11: Non-metazoan outgroups used in the analyses of the**
1100 **Genome and EST sets**
1101

| Name | Genome Set | EST Set |
|---|---|---|

| Opisthokonta | *Sphaeroforma arctica*<br>*Saccharomyces cerevisiae*<br>*Spizellomyces punctatus*<br>*Capsaspora owczarzaki*<br>*Monosiga brevicollis*<br>*Salpingoeca rosetta* | *Spizellomyces punctatus*<br>*Batrachochytrium dendrobatidis*<br>*Cryptococcus neoformans*<br>*Saccharomyces cerevisiae*<br>*Phycomyces blakesleeanus*<br>*Rhizopus orizae*<br>*Amoebidium parasiticum*<br>*Sphaeroforma arctica*<br>*Capsaspora owczarzaki*<br>*Monosiga ovata*<br>*Monosiga brevicollis*<br>*Salpingoeca rosetta* |
|---|---|---|
| Holozoa | *Capsaspora owczarzaki*<br>*Monosiga brevicollis*<br>*Salpingoeca rosetta* | *Capsaspora owczarzaki*<br>*Monosiga ovata*<br>*Monosiga brevicollis*<br>*Salpingoeca rosetta* |
| Choanimalia | *Monosiga brevicollis*<br>*Salpingoeca rosetta* | *Monosiga ovata*<br>*Monosiga brevicollis*<br>*Salpingoeca rosetta* |
| Animalia | None | None |

1102
1103
1104 **Table S12: Estimated branch lengths for each Genome Set tree**
1105 For each tree, we calculated the branch length from root to tip for each taxa. Table is
1106 sorted based on average branch lengths across all analyses (longest to shortest). Species
1107 are in the first column and the latin names are abbreviated to first letter of the first part
1108 and first three of the second. Opist.=Opisthokonta (all opisthokontan outgroups
1109 included), Holo=Holozoa (only holozoan outgroups included), Cho=Choanimalia (only
1110 choanoflagellate outgroups included), ML=maximum-likelihood, Bayes=Bayesian,
1111 Avg=Average.
1112

|  | Opist.<br>ML | Holo.<br>ML | Cho.<br>ML | Opist.<br>Bayes | Holo.<br>Bayes | Cho.<br>Bayes |
|---|---|---|---|---|---|---|
| *Cele* | 1.04 | 1.06 | 1.23 | 2.03 | 2.30 | 2.76 |
| *Ppac* | 1.04 | 1.06 | 1.23 | 2.02 | 2.29 | 2.76 |
| *Mlei* | 0.67 | 0.69 | 0.73 | 1.26 | 1.37 | 1.75 |
| *Dmel* | 0.68 | 0.69 | 0.85 | 1.06 | 1.14 | 1.59 |
| *Hrob* | 0.63 | 0.64 | 0.80 | 0.90 | 0.97 | 1.41 |
| *Aque* | 0.52 | 0.53 | 0.66 | 0.88 | 0.94 | 1.31 |
| *Tadh* | 0.51 | 0.53 | 0.68 | 0.70 | 0.75 | 1.16 |

| | | | | | |
|---|---|---|---|---|---|
| *Ctel* | 0.48 | 0.49 | 0.65 | 0.65 | 0.69 | 1.14 |
| *Lgig* | 0.46 | 0.47 | 0.62 | 0.62 | 0.66 | 1.11 |
| *Spur* | 0.45 | 0.46 | 0.62 | 0.58 | 0.62 | 1.06 |
| *Hsap* | 0.43 | 0.44 | 0.60 | 0.57 | 0.62 | 1.06 |
| *Bflo* | 0.40 | 0.41 | 0.56 | 0.50 | 0.54 | 0.98 |
| *Nvec* | 0.37 | 0.38 | 0.54 | 0.45 | 0.48 | 0.92 |
| **Median** | 0.51 | 0.53 | 0.66 | 0.7 | 0.75 | 1.16 |
| **STDEV** | 0.22 | 0.23 | 0.23 | 0.53 | 0.62 | 0.62 |

1113

**Table S13: Estimated branch lengths for each EST Set tree**

For each tree, we calculated the branch length from root to tip for each taxon. Since each
of the two independent Bayesian analyses did not converge for any of the analyses of the
EST Sets, we report both trees in the main text. However, for the branch lengths analyses
we calculate a single tree using bpcomp for each analysis. Table is sorted based on
average branch lengths across all analyses (longest to shortest). Species are in the first
column and the Latin names are abbreviated to first letter of the first part and first three
(or four) of the second. Opist.=Opisthokonta (all opisthokontan outgroups included),
Holo=Holozoa (only holozoan outgroups included), Cho=Choanimalia (only
choanoflagellate outgroups included), ML=maximum-likelihood, Bayes=Bayesian,
Avg=Average.

1125

| | Opist. ML | Holo. ML | Cho. ML | Opist. Bayes | Holo. Bayes | Cho. Bayes |
|---|---|---|---|---|---|---|
| *Smed* | 0.93 | 0.96 | 1.03 | 1.47 | 1.83 | 2.11 |
| *Clon* | 0.82 | 0.84 | 0.91 | 1.37 | 1.67 | 1.92 |
| *Ipul* | 0.82 | 0.85 | 0.92 | 1.35 | 1.66 | 1.93 |
| *Sros* | 0.80 | 0.82 | 0.89 | 1.27 | 1.54 | 1.45 |
| *Msti* | 0.75 | 0.75 | 0.82 | 1.17 | 1.37 | 1.62 |
| *Omin* | 0.69 | 0.71 | 0.77 | 1.14 | 1.34 | 1.60 |
| *Xind* | 0.71 | 0.72 | 0.79 | 0.99 | 1.16 | 1.41 |
| *Poli* | 0.71 | 0.72 | 0.79 | 0.98 | 1.14 | 1.40 |
| *Ppil* | 0.57 | 0.59 | 0.63 | 1.00 | 1.28 | 1.52 |

| | | | | | |
|---|---|---|---|---|---|
| *Nwes* | 0.65 | 0.66 | 0.73 | 0.95 | 1.12 | 1.36 |
| *Dmel* | 0.68 | 0.70 | 0.76 | 0.92 | 1.08 | 1.33 |
| *Mlei* | 0.56 | 0.58 | 0.61 | 0.96 | 1.22 | 1.47 |
| *Msp* | 0.56 | 0.58 | 0.61 | 0.94 | 1.21 | 1.45 |
| *Hrob* | 0.62 | 0.64 | 0.70 | 0.82 | 0.95 | 1.20 |
| *Hror* | 0.63 | 0.64 | 0.71 | 0.82 | 0.96 | 1.19 |
| *Cint* | 0.62 | 0.64 | 0.71 | 0.81 | 0.95 | 1.18 |
| *Dpul* | 0.58 | 0.59 | 0.66 | 0.77 | 0.90 | 1.14 |
| *Ehor* | 0.58 | 0.59 | 0.66 | 0.74 | 0.86 | 1.10 |
| *Srap* | 0.53 | 0.54 | 0.60 | 0.78 | 0.89 | 1.12 |
| *Chem* | 0.55 | 0.56 | 0.63 | 0.76 | 0.87 | 1.11 |
| *Bmic* | 0.56 | 0.57 | 0.64 | 0.73 | 0.85 | 1.09 |
| *Aque* | 0.51 | 0.52 | 0.58 | 0.76 | 0.88 | 1.11 |
| *Lbai* | 0.52 | 0.53 | 0.59 | 0.75 | 0.84 | 1.07 |
| *Cfol* | 0.52 | 0.53 | 0.59 | 0.73 | 0.83 | 1.06 |
| *Hmag* | 0.54 | 0.55 | 0.62 | 0.70 | 0.81 | 1.05 |
| *Esco* | 0.55 | 0.56 | 0.63 | 0.69 | 0.81 | 1.05 |
| *Pmar* | 0.54 | 0.55 | 0.61 | 0.66 | 0.77 | 1.00 |
| *Emue* | 0.49 | 0.50 | 0.55 | 0.71 | 0.80 | 1.04 |
| *Hech* | 0.52 | 0.53 | 0.60 | 0.67 | 0.77 | 1.01 |
| *Cvir* | 0.53 | 0.53 | 0.60 | 0.65 | 0.76 | 1.01 |
| *Lcha* | 0.49 | 0.50 | 0.56 | 0.70 | 0.79 | 1.02 |
| *Aero* | 0.51 | 0.52 | 0.59 | 0.65 | 0.77 | 1.01 |
| *Sdom* | 0.47 | 0.48 | 0.54 | 0.68 | 0.76 | 1.00 |
| *Apall* | 0.53 | 0.55 | 0.62 | 0.64 | 0.72 | 0.96 |
| *Ocar* | 0.48 | 0.49 | 0.55 | 0.68 | 0.76 | 0.99 |

| | | | | | | |
|---|---|---|---|---|---|---|
| **Lgig** | 0.51 | 0.52 | 0.59 | 0.64 | 0.75 | 0.99 |
| **Ctel** | 0.51 | 0.52 | 0.59 | 0.64 | 0.75 | 0.99 |
| **Pcar** | 0.51 | 0.52 | 0.59 | 0.64 | 0.74 | 0.98 |
| **Xboc** | 0.48 | 0.49 | 0.55 | 0.65 | 0.75 | 1.00 |
| **Tadh** | 0.46 | 0.47 | 0.54 | 0.65 | 0.77 | 1.01 |
| **Spur** | 0.50 | 0.51 | 0.58 | 0.60 | 0.71 | 0.96 |
| **Ttra** | 0.48 | 0.49 | 0.56 | 0.60 | 0.70 | 0.95 |
| **Olob** | 0.46 | 0.46 | 0.53 | 0.64 | 0.72 | 0.95 |
| **Clac** | 0.48 | 0.49 | 0.56 | 0. 60 | 0.70 | 0.94 |
| **Ccap** | 0.48 | 0.49 | 0.55 | 0.60 | 0.69 | 0.93 |
| **Apec** | 0.48 | 0.49 | 0.56 | 0.58 | 0.68 | 0.92 |
| **Ekan** | 0.46 | 0.47 | 0.54 | 0.59 | 0.69 | 0.93 |
| **Ggal** | 0.46 | 0.47 | 0.53 | 0.57 | 0.66 | 0.89 |
| **Bflo** | 0.45 | 0.46 | 0.53 | 0.55 | 0.64 | 0.87 |
| **Avir** | 0.44 | 0.44 | 0.51 | 0.53 | 0.59 | 0.83 |
| **Skow** | 0.43 | 0.44 | 0.50 | 0.50 | 0.60 | 0.84 |
| **Apalm** | 0.43 | 0.44 | 0.51 | 0.52 | 0.59 | 0.83 |
| **Msen** | 0.43 | 0.44 | 0.51 | 0.52 | 0.59 | 0.83 |
| **Mfav** | 0.43 | 0.44 | 0.51 | 0.52 | 0.59 | 0.83 |
| **Amil** | 0.42 | 0.42 | 0.49 | 0.51 | 0.57 | 0.82 |
| **Pfla** | 0.41 | 0.42 | 0.48 | 0.49 | 0.58 | 0.82 |
| **Past** | 0.41 | 0.42 | 0.49 | 0.48 | 0.57 | 0.80 |
| **Nvec** | 0.40 | 0.41 | 0.48 | 0.49 | 0.55 | 0.79 |
| **Median** | 0.515 | 0.525 | 0.59 | 0.68 | 0.77 | 1.01 |
| **STDEV** | 0.12 | 0.12 | 0.12 | 0.23 | 0.30 | 0.29 |

1126
1127

1128 **Supplementary Table S14: Hypotheses comparisons of likelihood analysis of gene**
1129 **content**
1130 Values in blue indicate values larger than the significance level, 0.05, and indicate
1131 inclusion within the confidence set. See Figure 2 for graphical representation of
1132 hypotheses. Hypotheses were determined *a priori*. AU = the p-value from the
1133 approximately unbiased test, NP = bootstrap probability of the selection, BP = same as np
1134 but calculated directly from the replicates with $r_k = 1$, PP = Bayesian posterior probability
1135 calculated by the BIC approximation, KH = the p-value from the Kishino-Hasegawa test,
1136 SH = the p-value from the Shimodaira-Hasegawa test, WKH = the p-value from the
1137 weighted Kishino-Hasegawa test, WSH = the p-value from the weighted Shimodaira-
1138 Hasegawa test. SOWH = the p-value from the Swofford–Olsen–Waddell–Hillis test.
1139 Asterisk indicates that the SH and WSH tests tend to include more trees in the confidence
1140 set than is necessary (*96*).
1141

| Hypothesis (figure) | AU | NP | BP | PP | KH | SH | WKH | WSH | SOWH |
|---|---|---|---|---|---|---|---|---|---|
| (Ct,) Fig. 2d | 0.973 | 0.973 | 0.973 | 1 | 0.972 | 0.991 | 0.972 | 0.998 | 1 |
| (Po,) Fig. 2c | 0.027 | 0.027 | 0.027 | $2 \times 10^{-14}$ | 0.028 | 0.32* | 0.028 | 0.086* | 0 |
| (Bi,) Fig. 2f | $1 \times 10^{-7}$ | $3 \times 10^{-07}$ | 0 | $2 \times 10^{-220}$ | 0 | 0 | 0 | 0 | 0 |
| (Cn,Ct) Fig. 2a | $3 \times 10^{-41}$ | $3 \times 10^{-15}$ | 0 | $8 \times 10^{-129}$ | 0 | 0 | 0 | 0 | 0 |
| (Ct,Bi) Fig. 2b | $8 \times 10^{-45}$ | $5 \times 10^{-16}$ | 0 | $1 \times 10^{-129}$ | 0 | 0 | 0 | 0 | 0 |
| (Tr,) Fig. 2e | $4 \times 10^{-53}$ | $2 \times 10^{-17}$ | 0 | $1 \times 10^{-93}$ | 0 | 0 | 0 | 0 | 0 |

1142
1143
1144 **Supplementary Table S15: Presence and absence of notch pathway components**
1145 A check mark in a green box indicates the presence of the gene. A check mark in a
1146 yellow box indicates uncertainty because of missing diagnostic domains. A grey box with
1147 a minus sign (-) indicates the absence of the gene. See Table S14 for *M. leidyi* gene
1148 identifiers. *Hs=Homo sapiens*, *Nv=Nematostella vectensis*, *Hm=Hydra magnipapillata,*
1149 *Ta=Trichoplax adhaerens*,  *Aq=Amphimedon queenslandica*, *Ml=Mnemiopsis leidyi*,
1150 *Mb=Monosiga brevicolis*, Choano=Choanoflagellata.

1151

| | Bilateria | Cnidaria | | Placozoa | Porifera | Ctenophora | Choano |
|---|---|---|---|---|---|---|---|
| | *Hs* | *Nv* | *Hm* | *Ta* | *Aq* | *Ml* | *Mb* |
| Notch | √ | √ | √ | √ | √ | √ | - |
| Delta | √ | √ | √ | √ | √ | √ | - |
| O-fut | √ | √ | √ | √ | √ | - | √ |
| Fringe | √ | √ | - | - | √ | - | - |
| furin | √ | √ | √ | - | - | √ | - |
| Tace=Adam17 | √ | √ | √ | √ | √ | √ | √ |
| Kuzbanian=Adam10 | √ | √ | √ | √ | √ | √ | √ |
| Presenillin | √ | √ | √ | √ | √ | √ | √ |
| Nicastrin | √ | √ | √ | √ | √ | √ | √ |
| APH1 | √ | √ | √ | √ | √ | √ | √ |
| Pen2 | √ | √ | √ | √ | √ | √ | √ |
| Su(H) | √ | √ | √ | √ | √ | √ | √ |
| Mastermind (Co-A) | √ | √ | - | - | - | - | - |
| SMRT (Co-R) | √ | - | - | - | - | - | - |
| Numb | √ | √ | √ | √ | - | √ | - |
| Hes/Hey | √ | √ | √ | √ | √ | √ | - |
| Strawberry notch | √ | √ | √ | √ | √ | - | √ |
| Neuralized | √ | √ | - | √ | √ | - | - |
| Mindbomb | √ | √ | √ | - | √ | √ | - |
| Deltex | √ | √ | - | √ | √ | √ | √ |
| Nedd4/sudx | √ | √ | √ | √ | √ | √ | √ |
| notchless | √ | √ | √ | √ | √ | √ | √ |
| **Total** | 22 | 16-21 | 8-17 | 14-17 | 17-18 | 12-16 | 10-12 |

1152
1153 **Supplementary Table S16:** *M. leidyi* **identifiers corresponding to genes in**
1154 **Supplementary Table S15**
1155 A yellow box indicates uncertainty because of missing diagnostic domains.
1156

| | |
|---|---|
| Notch | ML128617a |
| Delta | ML21438a |
| furin | ML07022a |
| Tace=Adam17 | ML17408a |
| Kuzbanian=Adam10 | ML03054a |
| Presenillin | ML01594a |
| Nicastrin | ML102219a |
| APH1 | ML305514a |
| Pen2 | ML0708a, MLRB070835a |
| Su(H) | ML141212a |
| Numb | ML00718a |

| | | ML065313a |
|---|---|---|
| Hes/Hey | | ML065313a |
| Mindbomb | | ML26791a |
| Deltex | | ML030223a |
| Nedd4/sudx | | ML044111a |
| notchless | | ML45849a |

1157

**Supplementary Table S17: Presence and absence of post-synaptic genes**

1159 A check mark indicates the presence of the gene and a grey box with a minus sign (-)
1160 indicates the absence of the gene. See Table S16 for *M. leidyi* gene identifiers. *Hs=Homo*
1161 *sapiens*, *Nv=Nematostella vectensis*, *Hm=Hydra magnipapillata, Ta=Trichoplax*
1162 *adhaerens*, *Aq=Amphimedon queenslandica, Ml=Mnemiopsis leidyi, Mb=Monosiga*
1163 *brevicolis*, Choano=Choanoflagellata

1164

| | Bilateria | Cnidaria | | Placozoa | Porifera | Ctenophora | Choano |
|---|---|---|---|---|---|---|---|
| | *Hs* | *Nv* | *Hm* | *Ta* | *Aq* | *Ml* | *Mb* |
| Alpha Catenin | √ | √ | - | √ | √ | √ | - |
| AMPA iGluR | √ | √ | √ | - | - | - | - |
| Beta Catenin | √ | √ | √ | √ | √ | √ | √ |
| CamKII | √ | √ | √ | √ | √ | √ | √ |
| CASK | √ | √ | √ | √ | - | - | - |
| Citron | √ | √ | √ | √ | √ | √ | √ |
| Classical Cadherin | √ | √ | √ | √ | √ | √ | - |
| Cortactin | √ | √ | √ | √ | √ | √ | √ |
| CRIPT | √ | √ | √ | √ | √ | √ | √ |
| p120/δ Catenin | √ | √ | √ | √ | √ | √ | - |
| DLG | √ | √ | √ | √ | √ | √ | √ |
| Ephrin Receptor | √ | √ | √ | √ | √ | √ | √ |
| ErbB Receptor | √ | - | - | - | √ | - | - |
| Erbin | √ | √ | √ | - | - | - | - |
| GKAP | √ | √ | √ | √ | √ | √ | - |
| GRIP | √ | - | - | - | √ | - | - |
| Homer | √ | √ | √ | √ | √ | √ | √ |
| IP3R | √ | √ | √ | √ | √ | √ | √ |
| K+ channel shaker | √ | √ | √ | √ | √ | √ | √ |
| LIN-7 | √ | √ | √ | √ | √ | √ | - |
| MAGI | √ | √ | √ | √ | √ | √ | - |
| mGluR | √ | √ | √ | √ | √ | √ | - |
| Neuroligin | √ | √ | - | - | - | - | - |
| NMDA iGluR | √ | √ | √ | - | - | - | - |
| NOS | √ | √ | √ | √ | √ | √ | - |

| | | | | | | | |
|---|---|---|---|---|---|---|---|
| PICK1 | √ | √ | √ | √ | √ | √ | - |
| PKC | √ | √ | √ | √ | √ | √ | √ |
| PMCA | √ | √ | √ | √ | √ | √ | √ |
| Shank | √ | √ | √ | - | √ | √ | √ |
| SPAR | √ | √ | √ | √ | √ | √ | - |
| stargazin | √ | - | - | - | - | - | - |
| SynGAP | √ | √ | √ | - | √ | √ | - |
| **Total** | 32 | 29 | 27 | 23 | 26 | 24 | 13 |

1165

**Supplementary Table S18: *M. leidyi* identifiers corresponding to genes in Supplementary Table S17**

1168

| | |
|---|---|
| **Alpha Catenin** | ML257617a |
| **Beta Catenin** | ML073715a |
| **CamKII** | ML070269a |
| **Citron** | ML154113a |
| **Classical Cadherin** | ML00359a |
| **Cortactin** | ML045237a |
| **CRIPT** | ML04051a |
| **Delta Catenin** | ML002622a |
| **DLG** | ML01744a |
| **Ephrin Receptor** | ML35913a |
| **GKAP** | ML03326a |
| **Homer** | ML06361a |
| **IP3R** | ML00401a |
| **K+ channel shaker** | ML18152a |
| **LIN-7** | ML05296a |
| **MAGI** | ML02503a |
| **mGluR** | ML17995a |
| **NOS** | ML074215a |
| **PICK1** | ML19124a |
| **PKC** | ML13931a |
| **PMCA** | ML11054a |
| **Shank** | ML01578a |
| **SPAR** | ML11651a |
| **SynGAP** | ML038810a |

1169

**Supplementary Table S19: Presence and absence of Dopamine / Norepinephrine / Epinephrine Biosynthetic Pathway components**
A check mark indicates the presence of the gene and a grey box with a minus sign (-) indicates the absence of the gene. The identifier for the Qdpr gene in *M. leidyi* is

1174 ML08064. Qdpr = quinoid dihydropteridine reductase, Th = tyrosine hydroxylase,
1175 Slc18A2 = Homo sapiens solute carrier family 18 member 2,  Ddc = dopa decarboxylase,
1176 Dbh = dopamine beta-hydroxylase, Pnmt = phenylethanolamine N-methyltransferase,
1177 *Hs=Homo sapiens*, *Nv=Nematostella vectensis*,  *Hm=Hydra magnipapillata*,
1178 *Ta=Trichoplax adhaerens*,  *Aq=Amphimedon queenslandica*, *Ml=Mnemiopsis leidyi*,
1179 *Mb=Monosiga brevicolis*, *Cowc=Capsaspora owczarzaki,* Choano=Choanoflagellata
1180

|  | Bilateria | Cnidaria | | Placozoa | Porifera | Ctenophora | Choano |
|---|---|---|---|---|---|---|---|
|  | *Hs* | *Nv* | *Hm* | *Ta* | *Aq* | *Ml* | *Mb* |
| Qdpr | √ | √ | √ | √ | √ | √ | (in *Cowc*) |
| Th | √ | - | - | - | - | - | - |
| Slc18A2 | √ | - | - | - | - | - | - |
| Ddc | √ | √ | √ | √ | - | - | - |
| Dbh | √ | √ | - | √ | √ | - | - |
| Pnmt | √ | - | - | - | - | - | - |

1181
1182 **Supplementary Table S20: Presence and absence of mesoderm components in model**
1183 **genomes**
1184 A check mark indicates the presence of the gene and a grey box with a minus sign (-)
1185 indicates its absence. See Table S19 for *M. leidyi* gene identifiers. *Hs=Homo sapiens*,
1186 *Nv=Nematostella vectensis*, *Hm=Hydra magnipapillata, Ta=Trichoplax adhaerens*,
1187 *Aq=Amphimedon queenslandica*, *Ml=Mnemiopsis leidyi*, *Mb=Monosiga brevicolis*,
1188 Choano=Choanoflagellata
1189

|  | Bilateria | Cnidaria | | Placozoa | Porifera | Ctenophora | Choano |
|---|---|---|---|---|---|---|---|
|  | *Hs* | *Nv* | *Hm* | *Ta* | *Aq* | *Ml* | *Mb* |
| Twist | √ | √ | - | - | - | - | - |
| MyoD family myogenin/Myf5 | √ | √ | - | - | - | - | - |
| Brachyury | √ | √ | √ | √ | (a) | √ | (b) |
| Snail | √ | √ | √ | √ | - | (c) | - |
| Eomesodermin/TBR2 | √ | - | - | - | - | - | - |
| gli/glis | √ | √ | √ | - | √ | √ | - |
| Tinman/Nkx-2.6 | √ | √ | √ | √ | √ | - | - |
| nk2.1/NKX2-1 | √ | √ | √ | √ | √ | - | - |
| Bagpipe/NKX3-2 | √ | √ | √ | - | - | - | - |
| Ladybird | √ | √ | - | - | - | - | - |
| Pax3 | √ | √ | - | - | - | - | - |
| Gsc | √ | √ | √ | √ | - | - | - |
| Forkhead/HNF3 FoxA/group 1 | √ | √ | √(d) | √ | - | - | - |
| Mef2 | √ | √ | √ | √ | √ | √ | (b) |
| GATA | √ | √ | √ | √ | √ | (e) | (f) |

| Muscle LIM/CSRP3 | √ | √ | √ | - | √ | √ | √ |
|---|---|---|---|---|---|---|---|
| Troponin T | √ | - | - | - | - | - | - |
| Troponin I | √ | - | - | - | - | - | - |
| Troponin C | √ | - | - | - | - | - | - |

(a) Brachyury in *A. queenslandica* but in 3 other sponges(*155, 156*)
(b) Brachyury and Muscle LIM are present in *Capsaspora owczarzaki*
(c) Snail is present in the ctenophores *Euplokamis dunlapae* and *Thalassocalyce inconstans*
(d) Forkhead is a pseudogene in *Hydra magnipapillata* but it is present in *Hydra vulgaris (157)*
(e) GATA is found in several other ctenophores
(f) Genes with single Znf GATA domains are present in *C. owczarzaki*

1190

1191 **Supplementary Table S21: *M. leidyi* identifiers corresponding to genes in**
1192 **Supplementary Table S20**

1193

| **Brachyury** | ML174736a |
|---|---|
| **gli/glis** | ML145833a |
| **Mef2** | ML07781a |
| **Muscle LIM/CSRP3** | ML02959a |

1194

1195 **Supplementary Table S22: Accession numbers of sequences from other ctenophore**
1196 **transcriptomes that were top TBLASTN hits to myogenic and synaptic scaffolding**
1197 **genes of interest. The row headers are the human NCBI Gene IDs used as the**
1198 **TBLASTN queries. In all cases except for ACCESSIONs (in red), these genes were**
1199 **determined not to be orthologous to the gene of interest.**

1200

| | *B. chuni* | *B. forskalii* | *C. fugiens* | *E. dunlapae* | *H. californensis* | *L. lactea* | *T. inconstans* |
|---|---|---|---|---|---|---|---|
| GRIA2, GRIN1, GRIK2, GRID2 | KF317296 KF317294 KF317293 KF317292 KF317295 KF317321 KF317322 KF317300 | KF317345 KF317346 KF317347 KF317348 KF317344 KF317358 | KF317388 KF317389 KF317387 KF317390 KF317391 | KF317429 KF317428 KF317426 KF317427 KF317425 KF317440 | KF317468 KF317466 KF317465 KF317467 KF317469 KF317476 KF317475 | KF317491 KF317490 KF317493 KF317497 KF317492 KF317494 KF317496 | KF317522 KF317524 KF317521 KF317523 KF317520 KF317531 |
| TWIST1 | KF317289 KF317290 KF317287 KF317302 KF317301 | KF317332 KF317328 KF317330 KF317331 KF317329 | KF317374 KF317375 KF317373 KF317372 KF317371 | KF317408 KF317407 KF317411 KF317410 KF317409 | KF317452 KF317454 KF317453 | | KF317505 KF317506 KF317507 KF317504 KF317503 |
| SNAI1 | KF317279 KF317277 KF317276 KF317278 KF317280 | KF317323 KF317325 KF317324 KF317326 KF317327 | KF317369 KF317366 KF317367 KF317370 KF317368 | KF317402 KF317403 KF317404 KF317405 KF317406 | KF317449 KF317447 KF317450 KF317448 | KF317484 KF317482 KF317483 KF317485 KF317481 | KF317498 KF317501 KF317500 KF317499 KF317502 |

| | | | | | | | |
|---|---|---|---|---|---|---|---|
| NLGN1 | KF317281<br>KF317284<br>KF317282<br>KF317283<br>KF317285 | KF317340<br>KF317338<br>KF317339<br>KF317342<br>KF317341 | KF317383<br>KF317386<br>KF317382<br>KF317384<br>KF317385 | KF317418<br>KF317417<br>KF317419<br>KF317420<br>KF317421 | KF317462<br>KF317463<br>KF317460<br>KF317459<br>KF317461 | KF317489<br>KF317488 | KF317517<br>KF317514<br>KF317513<br>KF317515<br>KF317516 |
| MYF5,<br>MYF6,<br>MYOG,<br>MYOD1 | KF317297<br>KF317286<br>KF317288<br>KF317316 | KF317349<br>KF317350<br>KF317351<br>KF317353<br>KF317352<br>KF317360<br>KF317359 | | KF317422<br>KF317423 | KF317464 | | KF317534<br>KF317535 |
| NOG | KF317291 | | | KF317424 | | | |
| EOMES | KF317309<br>KF317307<br>KF317310<br>KF317306<br>KF317308 | KF317337<br>KF317336<br>KF317334<br>KF317335<br>KF317333 | KF317376 | KF317412<br>KF317413<br>KF317416<br>KF317415<br>KF317414 | KF317458<br>KF317457<br>KF317456<br>KF317455<br>KF317451 | KF317486 | KF317508<br>KF317511<br>KF317512<br>KF317509<br>KF317510 |
| GATA1 | KF317319<br>KF317318<br>KF317320<br>KF317317 | KF317365<br>KF317361<br>KF317362<br>KF317364<br>KF317363 | KF317399<br>KF317397<br>KF317400<br>KF317401<br>KF317398 | KF317446<br>KF317443<br>KF317444<br>KF317442<br>KF317445 | KF317479<br>KF317478<br>KF317477<br>KF317480 | KF317487 | KF317536<br>KF317537<br>KF317525<br>KF317526<br>KF317527 |
| TNNC1,<br>TNNT1,<br>TNNI1 | KF317533<br>KF317532 | | | KF317441 | | | |
| NODAL | KF317313<br>KF317312<br>KF317314<br>KF317311<br>KF317315 | KF317343<br>KF317354<br>KF317356<br>KF317355<br>KF317357 | KF317392<br>KF317395<br>KF317396<br>KF317393<br>KF317378<br>KF317394 | KF317435<br>KF317437<br>KF317436<br>KF317439<br>KF317438 | KF317472<br>KF317470<br>KF317471<br>KF317473<br>KF317474 | KF317495 | KF317519<br>KF317518<br>KF317529<br>KF317530 |
| FGF | KF317305<br>KF317303<br>KF317304 | | KF317377<br>KF317380<br>KF317381<br>KF317379 | | | | |
| SHH | KF317299<br>KF317298 | | | KF317431<br>KF317434<br>KF317433<br>KF317432<br>KF317430 | | | KF317528 |

1201
1202

**Table S23: ANTP class homeodomains from other ctenophore transcriptomes**

```
B.chuni.2 PRKDRTIFSKIQLFRLERRFHDQKYLSSYEKALIAHSLQLTQQQVQTWFQNRRAKAKREE
B.chuni.3 EKRVRTIFSISQLFRLERRFNAQKYLTASERARLAYSLQLTETQVKIWFQNRRAKWKREM
B.chuni.5 KKKTRTTFSRAQVFGLERKFATQKYLSTHDRILLAAALQMTECQVKIWFQNRRTKWRREK
B.chuni.6 RRGPRTTIKSKQLEILKDAFNSTPKPTRHIREQLATKTGLNMRVIQVWFQNRRSKERRMK
B.chuni.7 QRRNRTQFSTYQLDQLEAEFEKSHYPDVLTREELANGLELTEARVQVWFQNRRAKWRKKQ
B.chuni.8 ERKARTVFTNDQLQLLENRFKSQPRLTSLEREELAEQIQLSATQIRVWFQNRRAKMKRDK
B.chuni.11 PRKTQTVFSTHQLSNLERRYNSNNYITTEERKCIADSLHMSVPQVKNWFQNKRNKERKLG
B.chuni.12 QRRNRITFTAAQLDLGEKSFQETHYPDVFGREEIAYSLQLTEQRVQVWFQNRRAKWRKRE
B.chuni.16 KRKPRTCFSVGQMLVLENTFQQTKYVSITDRGKLAGSLGLTDSQIKVWFQNRRSKWRKTM
B.chuni.18 KWRNRHTFTATQLEELEQVFNSTHYPDIFTREELAMKHNLTESRVQVWFQNRRAKHRKNE
B.chuni.20 SRKPRTVFTWEQLKQMEETFKQKKYLCTEDRMILAQKLDLTDEQIKVWFQNRRQKWKKGG
B.chuni.22 QRRNRTNYSAIQLNELEIVFSKSKYPDIFTREELALRLGVPEARIQVWFQNRRAKWRKQG
B.chuni.32 IKRPRTTITAKQLETLKSAYENSPKPARHVREQLSTETGLDMRVVQVWFQNRRAKEKRMK
B.chuni.36 GTNDRTIYTPVQAVRLEELFTERPHITKEQRDTLSGELTIHPDRIKVWFQNRRAKQRRED
B.chuni.39 PKYSRKNFSPEQIEGLEAAFSQHRFVKKDLRKQLAKQLSLSERSISYWFQNKRARSKPPI
B.chuni.40 KRNERHGIDDHQATTLKDWFQRYTYLTIENRKIVSIETGLPEKTVMYWFQNQRRKIKRQQ
B.chuni.41 PRRPRTIFSASQLLELEERFRYQKYLSTAERSCLAFTLGLSEEQVKVWFQNRRSKWKKGE
B.chuni.49 QRRNRTTFSSSVQLHELERAFQQSHYPDVFTREELAMRLDLTEARVQVWFQNRRAKWRKRE
B.chuni.60 KRRKRTTISSNSKEILEQYYNTNPLPSTDEIGNLSNNLALDKRVIRIWFQNRRAIGKRLS
B.chuni.64 PRKPRTHFTDSQIEDLEKVFEDKKYLSATERQIIANDLNLQEEQVKNWFQNRRSRWRKDC
B.chuni.68 SRRCRTKIEPDMLDLLETKYQESHFISPFERKSLSETLGITERAVIYWFQNRRRKDIKNM
B.chuni.83 PRKPRTIFTAAQLLELELEQKFRHQKYLSTSERSCLAFMLGLSEEQIKVWFQNRRSKWKKGG
B.chuni.111 TMRHRTRFTSSQLEKLEDAFSDTQYPDLSSRESISRDIGLSESCVQVWFQNRRARWRKSL
C.fugiens.1 KKKTRTTFSRAQVFGLERKFATQKYLSTHDRILLAAALHMTECQVKIWFQNRRTKWRREK
C.fugiens.4 QRRNRTNYSAAQLNELEIVFQRTRYPDIFTREELSLRLGIPEARIQVWFQNRRAKWRKRT
C.fugiens.8 PRKPRTHFTNGQVEELEKIFEDKKYLDAKEREVVAVDLNLAEEQVKNWFQNRRSKWRKDI
C.fugiens.10 PRKTQTVFSSFQVQHLEQKYNSANYISTEERQKIASRLQMSVPQVKNWFQNKRNKERKMG
C.fugiens.11 QRRNRTNYSAIQLNELEMVFGKSKYPDIFTREELALRLGVPEARIQVWFQNRRAKWRKQS
C.fugiens.12 RRKARTVFSDDQLQGLERKFKMQKYLSVPERMELAGMLSLSETQVKTWFQNRRMKWKKQG
C.fugiens.13 KRRRRTVFTERQLQGLEEAYSKSQYLDRESRLELCKKLSLSLHTVVYWFQNKRAISRRRG
C.fugiens.16 KRRRRTEYNEYQVAYLEMAFTENHYPSIAMREELADHIQIPEARIQIWFQNRRNKFRNHG
C.fugiens.20 KWRNRHTFTQAQLDELEKVFSKTHYPDIFTREDLAMKHNLTESRVQVWFQNRRAKHRKSE
C.fugiens.22 IKRPRTTITAKQLETLKNAYENSPKPARHVREQLSTETGLDMRVVQVWFQNRRAKEKRMK
C.fugiens.30 KRKPRTSFTNLQLFELERKFYHKKYLASSERKKLAQLLNLTDIQVKTWFQNRRNHKRSK
C.fugiens.33 QRRNRTQFTTYQLDQLEAEFDKSHYPDVLTREELANGLDLTEARVQVWFQNRRAKWRKKQ
C.fugiens.38 CSKRTKFSNEQLRVLEHYYHQVNKYVVGANKTALCHVTGLELNTILMWFQNKRAREKKQR
C.fugiens.52 RCRPRTSFTNFQLAHLEVAFSQTHYADIHTREELAKRLQLHESRIQVWFQNRRAKFRKAG
C.fugiens.55 PKYTRKNFSPEQIEGLEAAFQEHRFVKKDLRRELARKLNLSERSISYWFQNKRARSKPPI
C.fugiens.79 PRKPRTIFSAAQLNELEERFKYQKYLSTSERSCLAYSLGLTEEQIKVWFQNRRSKWKKGD
C.fugiens.80 RRKARTVFTDDQLKGLETQFGTQKYLSVPERMELAVSLRLSETQVKTWFQNRRMKWKKQV
C.fugiens.81 EKRVRTIFSISQLFRLERRFNAQKYLSASERARLAYSLQLTETQVKIWFQNRRAKWKREM
C.fugiens.89 RRKARTVFSDEQLTGLEDKFRVQKYLSVPERVELAVSLDLSETQVKTWFQNRRMKWKKGQ
C.fugiens.99 ERKARTVFTNEQLQLLENRFKAQPRLTSLEREELAQQIKLSATQIRVWFQNKRAKVKRDR
C.fugiens.123 RRKARTVFSDVQLEGLEKKFRSQKYLSVPERLDVATGLGLSETQVKTWFQNRRMKWKKQV
C.fugiens.128 GKKPRTIFSRDQVQKLEEAYLSKRYLTRRERKDLATAAQLSHTQVKIWFQNRRAKAKLKD
C.fugiens.164 RKRTRTTFSSAQVYELEKKFQRSQYLSAVDRLNLAAALSMSDVQVKRWFQNRRCKERHRA
C.fugiens.178 QRRNRITFTAAQLEGLEKSFQETHYPDVYGREEIAYSLQLTEQRVQVWFQNRRAKWRKRQ
L.lactea.2 QRRNRTNYSAVQLNELEIVFSKSKYPDIFTREELALRLGVPEARIQVWFQNRRAKWRKQG
L.lactea.3 SRKNRTSFSDKQLSVLEGVFRYKMYVSVTDRTMLSSRLELSDMQIKTWFQNRRTKWKKDN
L.lactea.8 PRKPRTHFTCSQVGDLEKVFEDKKYLSASERQTIAIDLNLQEEQVKNWFQNRRSRWRKDC
L.lactea.9 QRRNRTNYSAAQLNELELVFQRTRYPDIFTREELSLRLGIPEARIQVWFQNRRAKWRKRA
L.lactea.30 RRKARTVFSDIQLEGLERKFRNQKYLSVPERLDIATGLGLSETQVKTWFQNRRMKWKKQV
L.lactea.44 RRGPRTTIKSKQLEILKDAFNKTPKPTRHIREQLATKTGLNMRVIQVWFQNRRSKERRMK
L.lactea.46 HVNSRTMFTVHQVAAMEERYLKSSTILKPEREKFGMGLGLSETAIRTWFQNRRARQKRQE
T.inconstans.4 PRKPRTHFTDTQIEDLEKVFEDKKYLSANERQIIANDLSLHEEQVKNWFQNRRSRWRKDC
T.inconstans.8 KRRKRTTISSNSKEILEQYYQTNPLPSTDEIGNLSNNLSLDKRVIRIWFQNRRAIGKRLS
T.inconstans.14 PRKPRTIFSAAQLNELEERFKYQKYLSTTERSCLAYSLGLTEEQIKVWFQNRRSKWKKGD
T.inconstans.21 ERKPRTHFTDTQIEDLEKVFEDKKYLSANERQIIANDLSLHEEQVKNWFQNRRSRWRKDC
T.inconstans.23 PRKTQTVFTSHQLNNLEAKYRRNNYITTDERVSIAETLHMSVPQVKNWFQNKRNKERKMG
T.inconstans.37 QRRNRTQFTTCQLDQLEAEFDRSHYPDVLTREELAKCLGLTEARVQVWFQNRRAKWRKKQ
T.inconstans.51 YAGGARKQRRNRTNELEIVFSKSKYPDIFTREELALRLGVPEARIQVWFQNRRAKWRKRQ
T.inconstans.71 QRRNRITFTAAQLEGLEKSFQETHYPDVFGREEIAYSLQLTEQRVQVWFQNRRAKWRKRQ
T.inconstans.84 ERKARTVFTNDQLQLLENRFRSQPRLTSLEREELAEQIQLSATQIRVWFQNKRAKMKRDR
T.inconstans.99 KWRNRHTFTQAQIDELEQVFATTHYPDIFTREELANKHKLTEARVQVWFQNRRAKHRKNE
B.forskalii.5 ERKARTVFTNDQLQLLEDRFKSQPRLTSLEREELAEQIQLSATQIRVWFQNKRAKMKRDK
B.forskalii.7 PRKPRTHFTDSLIEDLEKVFEDKKYLSANERIIIANDLGLQEEQVKNWFQNRRSRWRKDC
B.forskalii.8 RHRTRFNDDQVAAMERYYNHVSRYARPDNGLPELIRETGLTHDTIMLWFQNKRARDKRKG
B.forskalii.12 AKYARRNFSATQIRGLEAAFREQRFIKKEVREQLAKDLGLSERSISYWFQNKRARSKGPI
B.forskalii.21 PRKPRTHFTDSQIEDLEKVFEDKKYLSANERIIIANDLGLQEEQVKNWFQNRRSRWRKDC
B.forskalii.25 VKRPRTVLSSVQRKVFKEAFDRTPRPCRKEREKLSSQTGLSVRVVQVWFQNQRAKVKKLA
B.forskalii.27 TKRFRSSFSTEQLRTMEDTFRHRPYLSTAQVEELAGNLQLSSRQVKIWFQNRRTKLKKQV
B.forskalii.31 QRRNRTQFSTYQLDQLEAEFVKSHYPDVLTREELANGLDLTEARVQVWFQNRRAKWRKKE
B.forskalii.33 QRRNRITFTAAQLEGLEKSFQETHYPDVFGREEIAYSLQLTEQRVQVWFQNRRAKWRKRE
B.forskalii.40 SRRFRTTFTSCQLQALEGAFRQTHYPDMYMREELAMRIDLTEARVQVWFQNRRAKWRKRE
B.forskalii.41 KRRRRTVFTERQLQGLEEAYSRSQYLDRESRLELCKKLSLSLHTVVYWFQNKRAISRRRG
B.forskalii.46 SKKPRTIFSREQVQKLEDAYQTKRYLTRRERKELAADAEISHTQVKIWFQNRRAKAKLKD
B.forskalii.50 IKRPRTTITAKQLETLKTAYENSPKPARHVREQLSTETGLDMRVVQVWFQNRRAKEKRMK
B.forskalii.56 PRKTQTVFTRQQLSELEAEYNVNNYISTEDRNEIARRLNMTVPQVKNWFQNKRNKERKMG
B.forskalii.59 PRKPRTIFSATQLNELEERFKYQKYLSTTERSCLAYSLGLTEEQIKVWFQNRRSKWKKGD
B.forskalii.68 KRRKRTTISSNSKEILEQYYQTNPLPSTEEIGNLSNNLSLDKRVIRIWFQNRRAIGKRLS
B.forskalii.81 EKRNRTKFTWEQLDRLEAEYQREQFAVADRKDALARELGVPARTISLWFQNRRAKQRREQ
B.forskalii.84 KKKTRTTFSRAQVFGLERKFATQKYLSTHDRILLAAALQMTECQVKIWFQNRRAKAKLKD
B.forskalii.86 PRKDRTIFSKIQLFRLERRFHAQKYLSSYEKALIAHSLQLTQQQVQTWFQNRRAKAKREE
E.dunlapae.1 RRKARTVFSDDQLNGLNGKFKDQKYLSVPERVELAVSLELSETQVKTWFQNRRMKWKKGQ
E.dunlapae.2 QRRNRTQFTTYQLDQLEGEFDKSHYPDVLTREELAHSLELTEARVQVWFQNRRAKWRKKQ
E.dunlapae.5 PRKDRTIFTKIQLFRLERRFHDQKYLSSYEKALIAHSLHLTQQQVQTWFQNRRAKAKREE
E.dunlapae.6 RRKARTVFTDDQLQGLESQFGAQKYLSVPERMELAVSLRLSETQVKTWFQNRRMKWKKQV
```

```
E.dunlapae.7   RRKARTVFSDVQLEGLERKFRSQKSLSVPERMDIASGLGLSETQVKTWFQNRRMKWKKQI
E.dunlapae.9   ERKARTVFTNAQLQLLEDRFKAQPRLTSLEREELAEQMKLSATQIRVWFQNKRAKMKRDK
E.dunlapae.10  KWRNRHTFTAAQLEELEHVFNTTHYPDIFTREELAMKHNLTEARVQVWFQNRRAKFRKTL
E.dunlapae.11  QRRNRTNYSAAQLNELEMVFQRTRYPDIFTREELSLRLGIPEARIQVWFQNRRAKWRKRA
E.dunlapae.14  RRKARTVFSDVQLEGLERKFRSQKYLSVPERMDIASGLGLSETQVKTWFQNRRMKWKKQI
E.dunlapae.16  TTRVRTVLNDRQLRILRTCYNNNPRPDALMKEQMTKLTGLSARVIRVWFQNKRCKDKKKA
E.dunlapae.17  PRKPRTIFTASQLLELEQKFRYQKYLSTSERSCLAFGLGLSEEQIKVWFQNRRSKWKKGG
E.dunlapae.25  PRKPRTIFSASQLNELEERFKYQKYLSTTERSCLAYSLGLTEEQIKVWFQNRRSKWKKGD
E.dunlapae.27  GECDRHTFTAAQLEELEHVFNTTHYPDIFTREELAMKHNLTEARVQVWFQNRRAKFRKTE
E.dunlapae.30  VMRHRTRFNNEQLFVLENSFNDSQYPDLAARENIAGLVGLSENCVQVWFQNRRARWRKIV
E.dunlapae.33  KRRKRTTISSNSKEILEQYYQTNPLPSTEEIGNLSNNLALDKRVIRIWFQNRRAIGKRLS
E.dunlapae.36  SRKPRTVFTWEQLKQMEETFKLKKYLCTEDRMLLAQKLDLTDEQIKVWFQNRRQKWKKGG
E.dunlapae.38  GKKPRTIFTREQVQRLEEAYTQKRYLTRRERKELAKEADISHTQVKIWFQNRRAKAKLKD
E.dunlapae.41  PRKPRTHFTISQIDDLEKVFVDKKYLSVTERSTLAYNLGLQEEQVKNWFQNRRSRWRKDC
E.dunlapae.48  TVRHRTRFTVSQLDQLEQVFNKTHYPDLSLREHLAMRTSLTEACVQIWFQNRRARWRKAV
E.dunlapae.50  RHRTRFNDDQIIAMEQYYNNVSHYARPDNGLPDLIQSTGLSHDTIMLWFQNKRARDKRKV
E.dunlapae.54  MKRPRTTITAKQLETLKTAYEKSPKPARHVREQLSTETGLDMRVVQVWFQNRRAKEKRLK
E.dunlapae.55  QRRNRTNYSAAQLNELEIIFQKSRYPDIFTREEMSLRLGIPEARIQVWFQNRRAKFRKQV
E.dunlapae.57  EEGTLANGKKPRTIFTREQVQRLEEAYTQKRYLTRRERKELHTQVKIWFQNRRAKAKLKD
E.dunlapae.59  RRKARTVFSDDQLQGLERKFFKMQKYLSVPERMELANMLALSETQVKTWFQNRRMKWKKQG
E.dunlapae.71  KWRNRHTFTAAQLEELEHVFNTTHYPDIFTREELAMKHNLTEARVQVWFQNRRAKFRKTE
E.dunlapae.79  QRRNRTTFSSIQLHELERAFQQSHYPDVFTREELAMRLDLTEARVQVWFQNRRAKWRKRE
E.dunlapae.88  QRRNRTNYSAAQLNELEMVFQRTRYPDIFTREELSLRLGIPEARIQVWFQNCHVGVMSPC
E.dunlapae.91  QRRNRTNYSAAQLNELEMVFQRTRYPDIFTREELSLRLGIPEARIQVWFQNRRATRRRPQ
E.dunlapae.109 TKISRRNFSVNQIEGLEAVYLKHRFVKRDMRRELAKRLNLSERSVSYWFQNKRARSKDPV
E.dunlapae.113 IKRYRSSFSTDQLRRMEDTFRHRPYLSTSQVEELAGSLKLSNRQVKIWFQNRRTKLKKQV
E.dunlapae.115 QRRNRTNYSAAQLNELEMVFQRTRYPDIFTREELSLRLGIPEARIQVWFQNRRAKWRKRT
E.dunlapae.116 RKRTRTTFSSAQVYELEKKFQRSQYLSAVDRLNLASALNMSDVQVKRWFQNRRCKERHRA
E.dunlapae.123 SRRFRTKIEPHMLSHLEEMYQRKHFISCEERKELADLLGMTERAVVYWFQNRRKDMRNY
E.dunlapae.142 RRGPRTTIKSKQLEILKDAFSSTPKPTRHIREQLAAKTGLNMRVIQVWFQNRRSKERRMK
E.dunlapae.144 VKRPRTVLSSAQRRVFKEAFDRSPRPCRKEREKLSSQTGLSVRVVQVWFQNQRAKVKKLA
E.dunlapae.170 WRRRRTVFTQEELSVLESVYSQNKFLNPDLKAEILSKVNVPGNVIVMWFQNRRAKDRSAG
E.dunlapae.185 QRRNRTNYSAAQLNELEMVFQRTRYPDIFTREELSLRLGIPEARIQVWFQNVMSPCSNEK
E.dunlapae.204 SNEQRHTFTAAQLEELEHVFNTTHYPDIFTREELAMKHNLTEARVQVWFQNRRAKFRKTE
E.dunlapae.219 PRRPRTIFSASQLLELEERFRYQKYLSTAERSCLAFTLGLSEEQVKVWFQNRRSKERRMK
E.dunlapae.225 GRPLRRRFTDDQLQGLESQFGAQKYLSVPERMELAVSLRLSETQVKTWFQNRRMKWKKQG
```

1207
## Table S24: Total genes predicted by various gene-modeling programs
1209

|  | FGENESH | AUGUSTUS | HMMGene | GenomeScan |
|---|---|---|---|---|
| Gene models | 16,367 | 29,359 | 13,948 | 6,443 |

1210
## Table S25: EVidenceModeler (EVM) weight assignments
1212

|  | 1st Run | 2nd Run | 3rd Run | Final |
|---|---|---|---|---|
| Transcripts | 10 | 10 | 10 | 10 |
| RNA-seq | 5 | 5 | 5 | 5 |
| FGENESH | 3 | 5 | 5 | 5 |
| AUGUSTUS | 2 | 2 | -- | -- |
| GenomeScan | 2 | -- | -- | -- |
| HMMgene | 1 | -- | -- | -- |
| RACE | -- | -- | -- | 10 |

1213
## Table S26: Manual comparison of EVM gene predictions with various settings vs. genes in which the structure has been determined with RACE PCR
1216

| Gene Name | Scaffold | 1st Run | 2nd Run | 3rd run | Final |
|---|---|---|---|---|---|
| Innexin1 | ML2599 | √ | √ | √ | √ |
| Un119 | ML3221 | √ | √ | √ | √ |

| | | | | | |
|---|---|---|---|---|---|
| Mago | ML2679 | √ | √ | √ | √ |
| Nova1 | ML2207 | √ | √ | √ | √ |
| Nova2 | ML1054 | √ | √ | √ | √ |
| Syn | ML3422 | √ | √ | √ | √ |
| Cornichon | ML4816 | √ | Miss | Join | √ |
| Tropomyosin | ML0772 | √ | √ | √ | √ |
| Tubby | ML0823 | √ | √ | √ | √ |
| Kv1/4 | ML0780 | Miss | Miss | Miss | √ |
| Titinlike | ML1985 | Miss | √ | √ | √ |
| Neb | ML2528 | Join | Join | Join | √ |
| Eph | ML3591 | √ | √ | √ | √ |
| Kv3 | ML0252 | Join | Join | Join | Join |
| Chibby | ML0412 | √ | √ | √ | √ |
| DIXD | ML0866 | Partial | Partial | Miss | √ |
| PaxA | ML0693 | Split | Split | Split | √ |
| PaxB | ML0788 | Join | Join | Join | √ |
| SixB | ML0649 | Partial | Partial | Partial | √ |
| SixC | ML2235 | Join | Join | √ | √ |
| SixD | ML2340 | Miss | Miss | Miss | √ |
| SixE | ML0120 | Join | Join | √ | Join |
| SixA | ML0146 | Miss | Miss | Miss | √ |
| Six | ML0179 | √ | √ | √ | √ |
| Lhx3 | ML0681 | Join | Join | Join | √ |
| Evx | ML4500 | √ | Split | Split | √ |
| Ml4195 | ML0547 | √ | √ | √ | √ |
| M4121 | ML0144 | √ | √ | √ | √ |
| Msxlike | ML0001 | Miss | Miss | Miss | √ |
| MN3 | ML0500 | Join | Join | Join | √ |
| DL2 | ML1029 | Miss | √ | √ | √ |
| Mleng | ML2705 | √ | √ | √ | √ |
| MINK1 | ML0633 | √ | √ | √ | √ |
| Pbx | ML0476 | √ | √ | √ | √ |
| Lim/Isl | ML0530 | Join | Join | √ | √ |
| Lhx1 | ML1325 | Join | Join | Join | Join |
| Smad4 | ML0219 | Split | Join | Join | √ |
| Smad1/5a | ML0930 | Miss | Miss | Miss | √ |
| Smad1/5b | ML0120 | Join | Join | √ | √ |
| Smad6 | ML1970 | √ | √ | √ | √ |
| Smad2/3 | ML0117 | Join | Join | √ | √ |
| TgfR | ML0859 | √ | √ | √ | Join |
| TgfR | ML0821 | √ | √ | √ | √ |
| TgfR | ML1311 | Join | Join | √ | √ |
| TgfR | ML0465 | √ | √ | √ | √ |
| SMURF | ML2068 | √ | √ | √ | √ |
| MlBmp | ML2188 | Join | Join | √ | √ |

| | | | | | |
|---|---|---|---|---|---|
| MlInh | ML1022 | Join | Join | √ | √ |
| MlMst | ML2002 | Partial | Partial | Partial | √ |
| MlMst2 | ML2002 | √ | √ | √ | √ |
| MITG2 | ML3487 | Miss | Miss | Miss | Miss |
| MITGF | ML0482 | Join | Join | √ | √ |
| MITG3 | ML1932 | Partial | Join | Join | Partial |
| MITG4 | ML3588 | Miss | Miss | Miss | √ |
| MITG5 | ML3689 | Miss | Miss | Join | √ |
| Piwi4 | ML0091 | √ | √ | √ | √ |
| HINT1 | ML0737 | √ | √ | √ | √ |
| HINT2 | ML3209 | √ | Miss | Join | Join |
| Opsin | ML1305 | √ | Join | Join | √ |
| Piwi2 | ML1497 | Miss | Miss | Miss | √ |
| Piwi | ML0742 | Join | Join | Join | √ |
| Pygopus | ML2207 | √ | √ | √ | √ |
| Tcf | ML1122 | √ | √ | √ | √ |
| Dishevelled | ML0053 | Partial | Partial | Partial | √ |
| Beta-catenin | ML0737 | √ | √ | √ | √ |
| Sec. frizzled | ML2235 | Join | Join | Join | Join |
| FrizzledB | ML0346 | Partial | Partial | Partial | Join |
| FrizzledA | ML0032 | Join | Join | √ | √ |
| Wnt9 | ML1010 | √ | √ | √ | √ |
| WntX | ML0105 | Miss | √ | √ | √ |
| Wnt6 | ML0602 | Partial | Partial | Partial | √ |
| WntA | ML0752 | Join | Join | √ | √ |
| COE | ML0447 | √ | √ | √ | √ |
| Bar | ML0315 | √ | √ | √ | √ |
| NKL1 | ML0585 | √ | √ | √ | √ |
| Tlx-like | ML0143 | √ | √ | √ | √ |
| Prd1 | ML0211 | √ | √ | √ | √ |
| Prd2 | ML1011 | √ | √ | √ | √ |
| Prd3 | ML0015 | Join | Join | Partial | Partial |
| Glis | ML1458 | √ | √ | √ | √ |
| Nanos1 | ML1302 | Split | Split | Split | √ |
| Nanos2 | ML2208 | √ | √ | Split | √ |
| Vasa | ML0471 | √ | √ | Miss | √ |
| Zn finger | ML0317 | Join | Join | Partial | √ |
| Opsinlike | ML1790 | Partial | Partial | Partial | √ |
| Mef2 | ML1232 | √ | √ | √ | √ |
| Syntaxin | ML3422 | √ | √ | √ | √ |
| Rag1 | ML0821 | Join | Join | Join | √ |
| Paraxis | ML0324 | Join | Join | Join | Join |
| Atonal | ML1043 | Miss | Miss | √ | √ |
| Musashi | ML0617 | √ | √ | √ | √ |
| Opsin | ML1204 | Partial | Partial | Partial | √ |

**Table S27: EVidenceModeler (EVM) consensus gene predictions**

|          | Genes | Exons | Scaffolds |
|----------|-------|-------|-----------|
| 1st Run  | 14537 | 85446 | 2037      |
| 2nd Run  | 14835 | 86712 | 1998      |
| 3rd Run  | 16845 | 89564 | 1915      |
| Final    | 16545 | 91482 | 1748      |

**Table S28: Species tree used for gene clustering algorithm (Newick)**

(((((((((((((*Homo_sapiens,Gallus_gallus*),*Xenopus_tropicalis*),*Danio_rerio*),*Ciona_intestinalis*),*Branchiostoma_floridae*),*Strongylocentrotus_purpuratus*),(((*Lottia_gigantea*,(*Capitella_teleta,Helobdella_robusta*)),*Schistosoma_mansoni*),((*Pristionchus_pacificus,Caenorhabditis_elegans*),((*Drosophila_melanogaster,Daphnia_pulex*),*Ixodes_scapularis*)))),(*Nematostella_vectensis,Hydra_magnipapillata*))),*Trichoplax_adhaerens*),*Amphimedon_queenslandica*),*Mnemiopsis_leidyi*);

**Table S29: Pfam domains - patterns of presence-absence**
Columns 2-6 indicate the presence absence pattern that defines column 1. The "TOTAL" column indicates how many Pfam domains fit this pattern. Co=*Capsaspora owczarzaki*, Mb=*Monosiga brevicolis*, Sr=*Salpingoeca rosetta*, Ml=*Mnemiopsis leidyi*, Aq=*Amphimedon queenslandica*, Ta=*Trichoplax adhaerens*, Nv=*Nematostella vectensis*, Bf=*Branchiostoma floridae*, Ce=Caenorhabditis *elegans*, Ci=*Ciona intestinalis*, Ct=*Capitella teleta*, Dm=*Drosophila melanogaster*, Dr=*Danio rerio*, Gg=*Gallus gallus*, Hs=*Homo sapiens*, Lg=*Lottia gigantea*

|  | Non-metazoan | Ctenophora | Porifera | Placozoa | Cnidaria | Bilateria | |
|---|---|---|---|---|---|---|---|
|  | Co,Mb,Sr | *Ml* | *Aq* | *Ta* | *Nv* | *Bf, Ce, Ci, Ct, Dm, Dr, Gg, Hs, Lg* | TOTAL |
| **Lost in Ml** | √ | X | √ | √ | √ | √ | 190 |
| **Lost in Aq** | √ | √ | X | √ | √ | √ | 122 |
| **Bilaterian** | X | X | X | X | X | √ | 1786 |
| **Parahoxozoa** | X | X | X | √ | √ | √ | 50 |
| **Coelenterata** | X | √ | X | X | √ | √ | 56 |
| **All animals but Porifera** | X | √ | X | √ | √ | √ | 58 |
| **All animals but Ctenophora** | X | X | √ | √ | √ | √ | 62 |

**Table S30: Domains absent from *M. leidyi* but present in *A. queenslandica*, *T.*

1235 *adhaerens*, *N. vectensis*, **at least one bilaterian and at least one non-metazoan**
1236

| Pfam name | Pfam ID | Pfam description |
|---|---|---|
| 3-HAO | PF06052.5 | 3-hydroxyanthranilic acid dioxygenase |
| 4HBT | PF03061.15 | Thioesterase superfamily |
| AICARFT_IMPCHas | PF01808.11 | AICARFT/IMPCHase bienzyme |
| AIRC | PF00731.13 | AIR carboxylase |
| AOX | PF01786.10 | Alternative oxidase |
| Adeno_IVa2 | PF02456.8 | Adenovirus IVa2 protein |
| Ald_Xan_dh_C | PF01315.15 | Aldehyde oxidase and xanthine dehydrogenase, a/b hammerhead domain |
| Ald_Xan_dh_C2 | PF02738.11 | Molybdopterin-binding domain of aldehyde dehydrogenase |
| Alpha_L_fucos | PF01120.10 | Alpha-L-fucosidase |
| Arginosuc_synth | PF00764.12 | Arginosuccinate synthase |
| BAAT_C | PF08840.4 | BAAT / Acyl-CoA thioester hydrolase C terminal |
| BRE | PF06113.5 | Brain and reproductive organ-expressed protein (BRE) |
| BSMAP | PF12280.1 | Brain specific membrane anchored protein |
| Bac_rhamnosid | PF05592.4 | Bacterial alpha-L-rhamnosidase |
| Baculo_LEF5_C | PF11792.1 | Baculoviridae late expression factor 5 C-terminal domain |
| Branch | PF02485.14 | Core-2/I-Branching enzyme |
| BtpA | PF03437.8 | BtpA family |
| CDKN3 | PF05706.5 | Cyclin-dependent kinase inhibitor 3 (CDKN3) |
| CHGN | PF05679.9 | Chondroitin N-acetylgalactosaminyltransferase |
| COLFI | PF01410.11 | Fibrillar collagen C-terminal domain |
| CO_deh_flav_C | PF03450.10 | CO dehydrogenase flavoprotein C-terminal domain |
| CR6_interact | PF10147.2 | Growth arrest and DNA-damage-inducible proteins-interacting protein 1 |
| CTP_transf_3 | PF02348.12 | Cytidylyltransferase |
| CYTH | PF01928.14 | CYTH domain |
| Catalase | PF00199.12 | Catalase |
| Catalase-rel | PF06628.5 | Catalase-related immune-responsive |
| Caveolin | PF01146.10 | Caveolin |
| Chitin_bind_3 | PF03067.8 | Chitin binding domain |
| Churchill | PF06573.4 | Churchill protein |
| Cob_adeno_trans | PF01923.11 | Cobalamin adenosyltransferase |
| CorA | PF01544.11 | CorA-like Mg2+ transporter protein |
| DAHP_synth_1 | PF00793.13 | DAHP synthetase I family |
| DNA_alkylation | PF08713.4 | DNA alkylation repair enzyme |
| DNA_binding_1 | PF01035.13 | 6-O-methylguanine DNA methyltransferase, DNA binding domain |
| DREV | PF05219.5 | DREV methyltransferase |

| Pfam name | Pfam ID | Pfam description |
| --- | --- | --- |
| DUF1075 | PF06388.4 | Protein of unknown function (DUF1075) |
| DUF1103 | PF06513.4 | Repeat of unknown function (DUF1103) |
| DUF1183 | PF06682.5 | Protein of unknown function (DUF1183) |
| DUF1228 | PF06779.7 | Protein of unknown function (DUF1228) |
| DUF1242 | PF06842.5 | Protein of unknown function (DUF1242) |
| DUF1258 | PF06869.5 | Protein of unknown function (DUF1258) |
| DUF143 | PF02410.8 | Domain of unknown function DUF143 |
| DUF1493 | PF07377.5 | Protein of unknown function (DUF1493) |
| DUF1604 | PF07713.6 | Protein of unknown function (DUF1604) |
| DUF1647 | PF07801.4 | Protein of unknown function (DUF1647) |
| DUF1713 | PF08213.4 | Mitochondrial domain of unknown function (DUF1713) |
| DUF1736 | PF08409.4 | Domain of unknown function (DUF1736) |
| DUF1855 | PF08910.3 | Protein of unknown function (DUF1855) |
| DUF1903 | PF08991.3 | Domain of unknown function (DUF1903) |
| DUF2054 | PF10218.2 | Uncharacterized conserved protein (DUF2054) |
| DUF2201 | PF09967.2 | Predicted metallopeptidase (DUF2201) |
| DUF2209 | PF09974.2 | Uncharacterized protein conserved in archaea (DUF2209) |
| DUF2215 | PF10225.2 | Uncharacterized conserved protein (DUF2215) |
| DUF2225 | PF09986.2 | Uncharacterized protein conserved in bacteria (DUF2225) |
| DUF2315 | PF10231.2 | Uncharacterised conserved protein (DUF2315) |
| DUF2348 | PF09807.2 | Uncharacterized conserved protein (DUF2348) |
| DUF2373 | PF10180.2 | Uncharacterised conserved protein (DUF2373) |
| DUF2721 | PF11026.1 | Protein of unknown function (DUF2721) |
| DUF2962 | PF11176.1 | Protein of unknown function (DUF2962) |
| DUF3074 | PF11274.1 | Protein of unknown function (DUF3074) |
| DUF3128 | PF11326.1 | Protein of unknown function (DUF3128) |
| DUF3184 | PF11380.1 | Protein of unknown function (DUF3184) |
| DUF3377 | PF11857.1 | Domain of unknown function (DUF3377) |
| DUF3639 | PF12341.1 | Protein of unknown function (DUF3639) |
| DUF3657 | PF12394.1 | Protein of unknown function (DUF3657) |
| DUF3752 | PF12572.1 | Protein of unknown function (DUF3752) |
| DUF3754 | PF12576.1 | Protein of unknown function (DUF3754) |
| DUF543 | PF04418.5 | Domain of unknown function (DUF543) |
| DUF563 | PF04577.7 | Protein of unknown function (DUF563) |
| DUF608 | PF04685.6 | Protein of unknown function, DUF608 |
| DUF803 | PF05653.7 | Protein of unknown function (DUF803) |
| DUF818 | PF05677.5 | Chlamydia CHLPS protein (DUF818) |
| DUF872 | PF05915.5 | Eukaryotic protein of unknown function (DUF872) |

| Pfam name | Pfam ID | Pfam description |
|---|---|---|
| DUF971 | PF06155.5 | Protein of unknown function (DUF971) |
| EBP | PF05241.5 | Emopamil binding protein |
| EF-1_beta_acid | PF10587.2 | Eukaryotic elongation factor 1 beta central acidic region |
| EPSP_synthase | PF00275.13 | EPSP synthase (3-phoshoshikimate 1-carboxyvinyltransferase) |
| ERG2_Sigma1R | PF04622.5 | ERG2 and Sigma1 receptor like protein |
| ETC_C1_NDUFA4 | PF04800.5 | ETC complex I subunit conserved region |
| Ependymin | PF00811.11 | Ependymin |
| Erv26 | PF04148.6 | Transmembrane adaptor Erv26 |
| Exostosin | PF03016.8 | Exostosin family |
| FAD_binding_5 | PF00941.14 | FAD binding domain in molybdopterin dehydrogenase |
| FA_FANCE | PF11510.1 | Fanconi Anaemia group E protein FANCE |
| Fe_dep_repress | PF01325.12 | Iron dependent repressor, N-terminal DNA binding domain |
| Fer2_2 | PF01799.13 | [2Fe-2S] binding domain |
| Ferritin | PF00210.17 | Ferritin-like domain |
| Fibrinogen_C | PF00147.11 | Fibrinogen beta and gamma chains, C-terminal globular domain |
| Flavodoxin_2 | PF02525.10 | Flavodoxin-like fold |
| FliL | PF03748.7 | Flagellar basal body-associated protein FliL |
| Folate_rec | PF03024.7 | Folate receptor family |
| Frataxin_Cyay | PF01491.9 | Frataxin-like domain |
| FtsH_ext | PF06480.8 | FtsH Extracellular |
| GBA2_N | PF12215.1 | beta-Glucocerebrosidase 2 N terminal |
| Galactosyl_T_2 | PF02709.7 | Galactosyltransferase |
| Glyco_hydro_2_N | PF02837.11 | Glycosyl hydrolases family 2, sugar binding domain |
| Glyco_hydro_30 | PF02055.9 | O-Glycosyl hydrolase family 30 |
| Glyco_transf_10 | PF00852.12 | Glycosyltransferase family 10 (fucosyltransferase) |
| Glyco_transf_54 | PF04666.6 | N-Acetylglucosaminyltransferase-IV (GnT-IV) conserved region |
| Glyco_transf_64 | PF09258.3 | Glycosyl transferase family 64 domain |
| Glycophorin_A | PF01102.11 | Glycophorin A |
| HemN_C | PF06969.9 | HemN C-terminal region |
| HisKA | PF00512.18 | His Kinase A (phosphoacceptor) domain |
| IF3_C | PF00707.15 | Translation initiation factor IF-3, C-terminal domain |
| IF3_N | PF05198.9 | Translation initiation factor IF-3, N-terminal domain |
| IIGP | PF05049.6 | Interferon-inducible GTPase (IIGP) |
| IRK | PF01007.13 | Inward rectifier potassium channel |
| Interfer-bind | PF09294.3 | Interferon-alpha/beta receptor, fibronectin type III |
| IspD | PF01128.12 | Uncharacterized protein family UPF0007 |
| Kinetochor_Ybp2 | PF08568.3 | Central kinetochore-associated |
| LMWPc | PF01451.14 | Low molecular weight phosphotyrosine protein phosphatase |

| Pfam name | Pfam ID | Pfam description |
| --- | --- | --- |
| LuxC | PF05893.7 | Acyl-CoA reductase (LuxC) |
| MMtag | PF10159.2 | Kinase phosphorylation protein |
| MRP-S33 | PF08293.4 | Mitochondrial ribosomal subunit S27 |
| MRP-S35 | PF10246.2 | Mitochondrial ribosomal protein MRP-S35 |
| MaoC_dehydratas | PF01575.12 | MaoC like domain |
| Menin | PF05053.6 | Menin |
| Methyltransf_5 | PF01795.12 | MraW methylase family |
| MitoNEET_N | PF10660.2 | Iron-containing outer mitochondrial membrane protein N-terminus |
| Myb_DNA-bind_2 | PF08914.4 | Rap1 Myb domain |
| NUC194 | PF08163.5 | NUC194 domain |
| Neur_chan_memb | PF02932.9 | Neurotransmitter-gated ion-channel transmembrane region |
| Nnf1 | PF03980.7 | Nnf1 |
| Nuf2 | PF03800.7 | Nuf2 family |
| O-FucT | PF10250.2 | GDP-fucose protein O-fucosyltransferase |
| Ocular_alb | PF02101.8 | Ocular albinism type 1 protein |
| PAC2 | PF09754.2 | PAC2 family |
| PAF-AH_p_II | PF03403.6 | isoform II |
| PAPA-1 | PF04795.5 | PAPA-1-like conserved region |
| PD40 | PF07676.5 | WD40-like Beta Propeller Repeat |
| PEP-utilizers_C | PF02896.11 | PEP-utilising enzyme, TIM barrel domain |
| PNPOx_C | PF10590.2 | Pyridoxine 5'-phosphate oxidase C-terminal dimerisation region |
| PPR | PF01535.13 | PPR repeat |
| Peptidase_C15 | PF01470.10 | Pyroglutamyl peptidase |
| Peptidase_M19 | PF01244.14 | Membrane dipeptidase (Peptidase family M19) |
| Peptidase_M49 | PF03571.8 | Peptidase family M49 |
| Peptidase_S37 | PF05576.4 | PS-10 peptidase S37 |
| Peptidase_S49 | PF01343.11 | Peptidase family S49 |
| Peroxin-3 | PF04882.5 | Peroxin-3 |
| Phlebovirus_NSM | PF07246.4 | Phlebovirus nonstructural protein NS-M |
| PhzC-PhzF | PF02567.9 | Phenazine biosynthesis-like protein |
| PigN | PF04987.7 | Phosphatidylinositolglycan class N (PIG-N) |
| Pox_A32 | PF04665.5 | Poxvirus A32 protein |
| Pox_A_type_inc | PF04508.5 | Viral A-type inclusion protein repeat |
| Pr_beta_C | PF12465.1 | Proteasome beta subunits C terminal |
| Pyridox_oxidase | PF01243.13 | Pyridoxamine 5'-phosphate oxidase |
| RRN7 | PF11781.1 | RNA polymerase I-specific transcription initiation factor Rrn7 |
| Rib_hydrolayse | PF02267.10 | ADP-ribosyl cyclase |
| Ribosomal_L17 | PF01196.12 | Ribosomal protein L17 |

| Pfam name | Pfam ID | Pfam description |
|---|---|---|
| Ribosomal_L28 | PF00830.12 | Ribosomal L28 family |
| Ribosomal_L9_N | PF01281.12 | Ribosomal protein L9, N-terminal domain |
| Ribosomal_S21 | PF01165.13 | Ribosomal protein S21 |
| Ribosomal_S6 | PF01250.10 | Ribosomal protein S6 |
| SAICAR_synt | PF01259.11 | SAICAR synthetase |
| SE | PF08491.3 | Squalene epoxidase |
| SKI | PF01202.15 | Shikimate kinase |
| SRA1 | PF07304.4 | Steroid receptor RNA activator (SRA1) |
| SUV3_C | PF12513.1 | Mitochondrial degradasome RNA helicase subunit C terminal |
| Scramblase | PF03803.8 | Scramblase |
| Sec2p | PF06428.4 | GDP/GTP exchange factor Sec2p |
| Sec39 | PF08314.4 | Secretory pathway protein Sec39 |
| SecA_DEAD | PF07517.7 | SecA DEAD-like domain |
| Seipin | PF06775.7 | Putative adipose-regulatory protein (Seipin) |
| Selenoprotein_S | PF06936.4 | Selenoprotein S (SelS) |
| Sigma70_ner | PF04546.6 | Sigma-70, non-essential region |
| TMP-TENI | PF02581.10 | Thiamine monophosphate synthase/TENI |
| TPPII | PF12580.1 | Tripeptidyl peptidase II |
| TPX2 | PF06886.4 | Targeting protein for Xklp2 (TPX2) |
| TRP | PF06011.5 | Transient receptor potential (TRP) ion channel |
| TYW3 | PF02676.7 | Methyltransferase TYW3 |
| Tmemb_14 | PF03647.6 | Transmembrane proteins 14C |
| Tmemb_40 | PF10160.2 | Predicted membrane protein |
| Translin | PF01997.9 | Translin family |
| TrmB | PF01978.12 | Sugar-specific transcriptional regulator TrmB |
| Trp_halogenase | PF04820.7 | Tryptophan halogenase |
| Trp_syntA | PF00290.13 | Tryptophan synthase alpha chain |
| UCR_TM | PF02921.7 | Ubiquinol cytochrome reductase transmembrane region |
| UDPG_MGDP_dh | PF00984.12 | UDP-glucose/GDP-mannose dehydrogenase family, central domain |
| UDPG_MGDP_dh _C | PF03720.8 | UDP-glucose/GDP-mannose dehydrogenase family, UDP binding domain |
| UPF0041 | PF03650.6 | Uncharacterised protein family (UPF0041) |
| UPF0054 | PF02130.10 | Uncharacterized protein family UPF0054 |
| VPS11_C | PF12451.1 | Vacuolar protein sorting protein 11 C terminal |
| Vitellogenin_N | PF01347.15 | Lipoprotein amino terminal region |
| Vps55 | PF04133.7 | Vacuolar protein sorting 55 |
| Wyosine_form | PF08608.5 | Wyosine base formation |
| XLF | PF09302.4 | XLF (XRCC4-like factor) |
| Xylo_C | PF12529.1 | Xylosyltransferase C terminal |

| Pfam name | Pfam ID | Pfam description |
|---|---|---|
| dCMP_cyt_deam_2 | PF08211.4 | Cytidine and deoxycytidylate deaminase zinc-binding region |
| tRNA_Me_trans | PF03054.9 | tRNA methyl transferase |
| zf-C4H2 | PF10146.2 | Zinc finger-containing protein |

1237

**Table S31: PhyloBayes stats**

We ran two instances of PhyloBayes for each data matrix (genome and EST) with varied outgroups (Opisthokonta, Holozoa, Choanimalia, Animalia). We used the burn-in below with bpcomp, which produced the reported maxdiff values. EST trees did not converge.

| Dataset | Total cycles.01 | Total cycles.02 | Burn-in | Maxdiff |
|---|---|---|---|---|
| Genome.Opisthokonta | 7964 | 8504 | 2500 | 0.05 |
| Genome.Holozoa | 12190 | 12895 | 4000 | 0 |
| Genome.Choanimalia | 8653 | 8643 | 2500 | 0.01 |
| Genome.Animalia | 14409 | 14451 | 4500 | 0 |
| EST.Opisthokonta | 10321 | 11814 | 3440 | 1 |
| EST.Holozoa | 12569 | 11124 | 3708 | 0.749153 |
| EST.Choanimalia | 15592 | 15747 | 5197 | 0.978199 |
| EST.Animalia | 9998 | 15801 | 3332 | 0.866308 |

1243
1244

**Figure S1: Phylogenetic results of concatenated amino acid analyses**

**a) RaxML Genome.Opisthokonta dataset**

1250
1251 **b) RaxML Genome.Holozoa dataset**
1252

```
                        ┌──────── Cowc
                        │                        ┌──── Mbre
                        │                  ┌─100─┤
                        │                  │      └──── Sros
                        │            ┌──────┤
                        └────────────┤      │  69 ┌──── Aque
                                     │      └─────┤
                                     │            └──── Mlei
                                     │ 100 ┌─73 ┌──── Nvec
                                     └─────┤    └──── Tadh
                                           │ 90      ┌──── Spur
                                           │   ┌100 ┤
                                           │   │    │ 94 ┌── Hsap
                                           │   │    └────┤
                                           └100┤         └── Bflo
                                               │    ┌100 ┌──── Dmel
                                               │    │    │    ┌──── Ppac
                                               └100 ┤    └100─┤
                                                    │         └──── Cele
                                                    │ 100 ┌──── Lgig
                                                    └─────┤ 100 ┌── Ctel
                                                          └─────┤ 100
                                                                └── Hrob
```

├──────── 0.2

1253
1254

**c) RaxML Genome.Choanimalia dataset**

1257 **d) RaxML Genome.Animalia dataset**
1258
1259

Mlei

Aque

Tadh

Nvec

100

Spur

98

Hsap

83

Bflo

92

100

Dmel

100

Ppac

100

Cele

100

Lgig

100

Ctel

100

Hrob

0.1

1260

1261

**e) PhyloBayes Genome.Opisthokonta dataset**

1263
1264
1265



1266

1267

1268 **f) PhyloBayes Genome.Holozoa dataset**

1269



1270
1271

1272

**g) PhyloBayes Genome.Choanimalia dataset**

1273

1274

1275



1276

1277

1278 **h) PhyloBayes Genome.Animalia dataset**

1279



Mlei

Aque

Tadh

1 Nvec

1 Spur

1 Hsap

Bflo

1 Lgig

1 Hrob

1 Ctel

1 Dmel

1 Ppac

Cele

0.3

1280

1281

1282 **i) RaxML EST.Opisthokonta dataset**

1283



Monosiga_ovata
Salpingoeca_rosetta
Monosiga_brevicollis
Capsaspora_owczarzaki_ATCC_30864
Amoebidium_parasiticum
Sphaeroforma_arctica
Spizellomyces_punctatus
Batrachochytrium_dendrobatidis
Cryptococcus_neoformans
Saccharomyces_cerevisiae
Phycomyces_blakesleeanus
Rhizopus_orizae
Pleurobrachia_pileus
Mnemiopsis
Mertensiid_sp
Leucetta_chagosensis
Sycon_raphanus
Oscarella_carmela
Oscarella_lobularis
Oopsacas_minuta
Carteriospongia_foliascens
Amphimedon_queenslandica
Suberites_domuncula
Lubomirskia_baicalensis
Ephydatia_muelleri
Trichoplax_adhaerens
Cyanea_capillata
Clytia_hemisphaerica
Hydra_magnipapillata
Podocoryna_carnea
Hydractinia_echinata
Nematostella_vectensis
Anemonia_viridis
Aiptasia_pallida
Metridium_senile
Acropora_palmata
Acropora_millepora
Porites_astreoides
Montastraea_faveolata
Xenoturbella_bocki
Nemertoderma_westbladi
Meara_stichopi
Isodiametra_pulchra
Symsagittifera_roscoffensis
Convolutriloba_longifissura
Saccoglossus_kowalevskii
Ptychodera_flava
Strongylocentrotus_purpuratus
Asterina_pectinifera
Branchiostoma_floridae
Petromyzon_marinus
Gallus
Ciona_intestinalis
Halocynthia_roretzi
Echinoderes_horni
Xiphinema_index
Euperipatoides_kanangrensis
Anoplodactylus_eroticus
Boophilus_microplus
Daphnia_pulex
Drosophila_melanogaster
Schmidtea_mediterranea
Paraplanocera_oligoglena
Capitella_telata
Helobdella_robusta
Cerebratulus_lacteus
Terebratalia_transversa
Euprymna_scolopes
Lottia_gigantea
Crassostrea_virginica

0.2

1284

## 1285    j) RaxML EST.Holozoa dataset
1286



Capsaspora_owczarzaki_ATCC_30864
Amoebidium_parasiticum
Sphaeroforma_arctica
Monosiga_ovata
Monosiga_brevicollis
Salpingoeca_rosetta
Pleurobrachia_pileus
Mnemiopsis
Mertensiid_sp
Oscarella_carmela
Oscarella_lobularis
Leucetta_chagosensis
Sycon_raphanus
Oopsacas_minuta
Carteriospongia_foliascens
Amphimedon_queenslandica
Suberites_domuncula
Lubomirskia_baicalensis
Ephydatia_muelleri
Trichoplax_adhaerens
Cyanea_capillata
Clytia_hemisphaerica
Hydra_magnipapillata
Hydractinia_echinata
Podocoryna_carnea
Nematostella_vectensis
Anemonia_viridis
Metridium_senile
Aiptasia_pallida
Porites_astreoides
Montastraea_faveolata
Acropora_palmata
Acropora_millepora
Xenoturbella_bocki
Nemertoderma_westbladi
Meara_stichopi
Isodiametra_pulchra
Symsagittifera_roscoffensis
Convolutriloba_longifissura
Saccoglossus_kowalevskii
Ptychodera_flava
Asterina_pectinifera
Strongylocentrotus_purpuratus
Branchiostoma_floridae
Halocynthia_roretzi
Ciona_intestinalis
Gallus
Petromyzon_marinus
Echinoderes_horni
Xiphinema_index
Euperipatoides_kanangrensis
Drosophila_melanogaster
Daphnia_pulex
Boophilus_microplus
Anoplodactylus_eroticus
Paraplanocera_oligoglena
Schmidtea_mediterranea
Euprymna_scolopes
Crassostrea_virginica
Lottia_gigantea
Cerebratulus_lacteus
Terebratalia_transversa
Capitella_telata
Helobdella_robusta

0.2

1287

1288
1289 **k) RaxML EST.Choanimalia dataset**
1290



- 100 — Monosiga_ovata
- 100 — Salpingoeca_rosetta
- Monosiga_brevicollis
- 100 — Pleurobrachia_pileus
- 100 — Mertensiid_sp
- Mnemiopsis
- 100 — Sycon_raphanus
- 92 — Leucetta_chagosensis
- 100 — Oscarella_carmela
- 54 — Oscarella_lobularis
- Oopsacas_minuta
- 93 — Carteriospongia_foliascens
- 92 — Amphimedon_queenslandica
- 100 — Suberites_domuncula
- 95 — Lubomirskia_baicalensis
- 100 — Ephydatia_muelleri
- 93 — Trichoplax_adhaerens
- Cyanea_capillata
- 100 — Clytia_hemisphaerica
- 100 — Hydra_magnipapillata
- 95 — Hydractinia_echinata
- 100 — Podocoryna_carnea
- 100 — Acropora_millepora
- 100 — Acropora_palmata
- 100 — Porites_astreoides
- 100 — Montastraea_faveolata
- Nematostella_vectensis
- 100 — Anemonia_viridis
- 100 — Aiptasia_pallida
- 100 — Metridium_senile
- Xenoturbella_bocki
- 50 — Nemertoderma_westbladi
- 98 — Meara_stichopi
- 99 — Isodiametra_pulchra
- 100 — Convolutriloba_longifissura
- 100 — Symsagittifera_roscoffensis
- Strongylocentrotus_purpuratus
- 100 — Asterina_pectinifera
- 100 — Ptychodera_flava
- 100 — Saccoglossus_kowalevskii
- 50 — Branchiostoma_floridae
- 84 — Halocynthia_roretzi
- 100 — Ciona_intestinalis
- 100 — Gallus
- 100 — Petromyzon_marinus
- Xiphinema_index
- 87 — Echinoderes_horni
- 97 — Euperipatoides_kanangrensis
- 98 — Daphnia_pulex
- 100 — Drosophila_melanogaster
- 100 — Anoplodactylus_eroticus
- 99 — Boophilus_microplus
- Paraplanocera_oligoglena
- 100 — Schmidtea_mediterranea
- 99 — Euprymna_scolopes
- 100 — Lottia_gigantea
- 100 — Crassostrea_virginica
- 85 — Helobdella_robusta
- 100 — Capitella_telata
- 41 — Cerebratulus_lacteus
- 86 — Terebratalia_transversa

0.2

1291
1292

1293 **l) RaxML EST.Animalia dataset**



Pleurobrachia_pileus
100 Mertensiid_sp
100 Mnemiopsis

Oopsacas_minuta
83 Carteriospongia_foliascens
80 Amphimedon_queenslandica
100 Suberites_domuncula
96 Lubomirskia_baicalensis
100 Ephydatia_muelleri

Trichoplax_adhaerens
50 Oscarella_carmela
100 Oscarella_lobularis
94 Leucetta_chagosensis
100 Sycon_raphanus

Cyanea_capillata
100 Clytia_hemisphaerica
100 Hydra_magnipapillata
95 Podocoryna_carnea
100 Hydractinia_echinata
69 100 Nematostella_vectensis
100 Anemonia_viridis
100 Aiptasia_pallida
100 Metridium_senile
100 Acropora_palmata
100 Acropora_millepora
100 Porites_astreoides
100 Montastraea_faveolata

100 Xenoturbella_bocki
41 Meara_stichopi
100 Nemertoderma_westbladi
98 Isodiametra_pulchra
100 Symsagittifera_roscoffensis
100 Convolutriloba_longifissura

Saccoglossus_kowalevskii
100 Ptychodera_flava
100 Asterina_pectinifera
100 Strongylocentrotus_purpuratus
40 Branchiostoma_floridae
82 Gallus
100 Petromyzon_marinus
99 Ciona_intestinalis
100 Halocynthia_roretzi
38 Xiphinema_index
83 Echinoderes_horni
95 Euperipatoides_kanangrensis
98 Drosophila_melanogaster
100 Daphnia_pulex
100 Boophilus_microplus
98 Anoplodactylus_eroticus
92 Paraplanocera_oligoglena
100 Schmidtea_mediterranea
Euprymna_scolopes
96 100 Crassostrea_virginica
100 Lottia_gigantea
85 Cerebratulus_lacteus
87 Terebratalia_transversa
44 Helobdella_robusta
100 Capitella_telata

0.2

1294

1295    **m) PhyloBayes EST.Opisthokonta dataset run 1**

1296



Spizellomyces_punctatus
Batrachochytrium_dendrobatidis
Rhizopus_orizae
Phycomyces_blakesleeanus
Saccharomyces_cerevisiae
Cryptococcus_neoformans
Sphaeroforma_arctica
Amoebidium_parasiticum
Capsaspora_owczarzaki_ATCC_30864
Monosiga_ovata
Salpingoeca_rosetta
Monosiga_brevicollis
Pleurobrachia_pileus
Mnemiopsis
Mertensiid_sp
Oscarella_lobularis
Oscarella_carmela
Sycon_raphanus
Leucetta_chagosensis
Oopsacas_minuta
Carteriospongia_foliascens
Amphimedon_queenslandica
Suberites_domuncula
Lubomirskia_baicalensis
Ephydatia_muelleri
Trichoplax_adhaerens
Cyanea_capillata
Hydra_magnipapillata
Clytia_hemisphaerica
Podocoryna_carnea
Hydractinia_echinata
Nematostella_vectensis
Anemonia_viridis
Metridium_senile
Aiptasia_pallida
Porites_astreoides
Montastraea_faveolata
Acropora_millepora
Acropora_palmata
Saccoglossus_kowalevskii
Ptychodera_flava
Strongylocentrotus_purpuratus
Asterina_pectinifera
Branchiostoma_floridae
Petromyzon_marinus
Gallus
Halocynthia_roretzi
Ciona_intestinalis
Xenoturbella_bocki
Nemertoderma_westbladi
Meara_stichopi
Isodiametra_pulchra
Symsagittifera_roscoffensis
Convolutriloba_longifissura
Echinoderes_horni
Xiphinema_index
Euperipatoides_kanangrensis
Drosophila_melanogaster
Daphnia_pulex
Boophilus_microplus
Anoplodactylus_eroticus
Terebratalia_transversa
Cerebratulus_lacteus
Helobdella_robusta
Capitella_telata
Schmidtea_mediterranea
Paraplanocera_oligoglena
Euprymna_scolopes
Lottia_gigantea
Crassostrea_virginica

0.4

1297
1298

1299 **n) PhyloBayes EST.Opisthokonta dataset run 2**

1300

Spizellomyces_punctatus
Batrachochytrium_dendrobatidis
Rhizopus_orizae
Phycomyces_blakesleeanus
Saccharomyces_cerevisiae
Cryptococcus_neoformans
Sphaeroforma_arctica
Amoebidium_parasiticum
Capsaspora_owczarzaki_ATCC_30864
Monosiga_ovata
Salpingoeca_rosetta
Monosiga_brevicollis
Pleurobrachia_pileus
Mnemiopsis
Mertensiid_sp
Oscarella_lobularis
Oscarella_carmela
Sycon_raphanus
Leucetta_chagosensis
Oopsacas_minuta
Carteriospongia_foliascens
Amphimedon_queenslandica
Suberites_domuncula
Lubomirskia_baicalensis
Ephydatia_muelleri
Trichoplax_adhaerens
Cyanea_capillata
Hydra_magnipapillata
Clytia_hemisphaerica
Podocoryna_carnea
Hydractinia_echinata
Nematostella_vectensis
Anemonia_viridis
Metridium_senile
Aiptasia_pallida
Porites_astreoides
Montastraea_faveolata
Acropora_millepora
Acropora_palmata
Saccoglossus_kowalevskii
Ptychodera_flava
Strongylocentrotus_purpuratus
Asterina_pectinifera
Branchiostoma_floridae
Petromyzon_marinus
Gallus
Halocynthia_roretzi
Ciona_intestinalis
Xenoturbella_bocki
Nemertoderma_westbladi
Meara_stichopi
Isodiametra_pulchra
Symsagittifera_roscoffensis
Convolutriloba_longifissura
Echinoderes_horni
Xiphinema_index
Euperipatoides_kanangrensis
Drosophila_melanogaster
Daphnia_pulex
Boophilus_microplus
Anoplodactylus_eroticus
Euprymna_scolopes
Lottia_gigantea
Crassostrea_virginica
Schmidtea_mediterranea
Paraplanocera_oligoglena
Terebratalia_transversa
Cerebratulus_lacteus
Helobdella_robusta
Capitella_telata

0.4

1301
1302
1303

**o) PhyloBayes EST.Holozoa dataset run 1**

- Sphaeroforma_arctica
- Amoebidium_parasiticum
- Capsaspora_owczarzaki_ATCC_30864
- Monosiga_ovata
- Salpingoeca_rosetta
- Monosiga_brevicollis
- Oscarella_lobularis
- Oscarella_carmela
- Sycon_raphanus
- Leucetta_chagosensis
- Oopsacas_minuta
- Carteriospongia_foliascens
- Amphimedon_queenslandica
- Suberites_domuncula
- Lubomirskia_baicalensis
- Ephydatia_muelleri
- Pleurobrachia_pileus
- Mnemiopsis
- Mertensiid_sp
- Trichoplax_adhaerens
- Cyanea_capillata
- Hydra_magnipapillata
- Clytia_hemisphaerica
- Podocoryna_carnea
- Hydractinia_echinata
- Nematostella_vectensis
- Anemonia_viridis
- Metridium_senile
- Aiptasia_pallida
- Porites_astreoides
- Montastraea_faveolata
- Acropora_millepora
- Acropora_palmata
- Xenoturbella_bocki
- Saccoglossus_kowalevskii
- Ptychodera_flava
- Strongylocentrotus_purpuratus
- Asterina_pectinifera
- Branchiostoma_floridae
- Petromyzon_marinus
- Gallus
- Halocynthia_roretzi
- Ciona_intestinalis
- Nemertoderma_westbladi
- Meara_stichopi
- Isodiametra_pulchra
- Symsagittifera_roscoffensis
- Convolutriloba_longifissura
- Echinoderes_horni
- Xiphinema_index
- Euperipatoides_kanangrensis
- Drosophila_melanogaster
- Daphnia_pulex
- Boophilus_microplus
- Anoplodactylus_eroticus
- Terebratalia_transversa
- Cerebratulus_lacteus
- Helobdella_robusta
- Capitella_telata
- Schmidtea_mediterranea
- Paraplanocera_oligoglena
- Euprymna_scolopes
- Lottia_gigantea
- Crassostrea_virginica

0.3

1308    **p) PhyloBayes EST.Holozoa dataset run 2**

1309



1310

0.3

**1311    q) PhyloBayes EST.Choanimalia dataset run 1**

1312



- Monosiga_ovata
- Salpingoeca_rosetta
- Monosiga_brevicollis
- Oscarella_lobularis
- Oscarella_carmela
- Sycon_raphanus
- Leucetta_chagosensis
- Oopsacas_minuta
- Carteriospongia_foliascens
- Amphimedon_queenslandica
- Suberites_domuncula
- Lubomirskia_baicalensis
- Ephydatia_muelleri
- Pleurobrachia_pileus
- Mnemiopsis
- Mertensiid_sp
- Trichoplax_adhaerens
- Cyanea_capillata
- Hydra_magnipapillata
- Clytia_hemisphaerica
- Podocoryna_carnea
- Hydractinia_echinata
- Nematostella_vectensis
- Anemonia_viridis
- Metridium_senile
- Aiptasia_pallida
- Porites_astreoides
- Montastraea_faveolata
- Acropora_millepora
- Acropora_palmata
- Xenoturbella_bocki
- Saccoglossus_kowalevskii
- Ptychodera_flava
- Strongylocentrotus_purpuratus
- Asterina_pectinifera
- Branchiostoma_floridae
- Petromyzon_marinus
- Gallus
- Halocynthia_roretzi
- Ciona_intestinalis
- Nemertoderma_westbladi
- Meara_stichopi
- Isodiametra_pulchra
- Symsagittifera_roscoffensis
- Convolutriloba_longifissura
- Echinoderes_horni
- Xiphinema_index
- Euperipatoides_kanangrensis
- Drosophila_melanogaster
- Daphnia_pulex
- Boophilus_microplus
- Anoplodactylus_eroticus
- Terebratalia_transversa
- Cerebratulus_lacteus
- Helobdella_robusta
- Capitella_telata
- Schmidtea_mediterranea
- Paraplanocera_oligoglena
- Euprymna_scolopes
- Lottia_gigantea
- Crassostrea_virginica

0.3

1313

**r) PhyloBayes EST.Choanimalia dataset run 2**

- Monosiga_ovata
- Salpingoeca_rosetta
- Monosiga_brevicollis
- Oscarella_lobularis
- Oscarella_carmela
- Sycon_raphanus
- Leucetta_chagosensis
- Oopsacas_minuta
- Carteriospongia_foliascens
- Amphimedon_queenslandica
- Suberites_domuncula
- Lubomirskia_baicalensis
- Ephydatia_muelleri
- Pleurobrachia_pileus
- Mnemiopsis
- Mertensiid_sp
- Trichoplax_adhaerens
- Cyanea_capillata
- Hydra_magnipapillata
- Clytia_hemisphaerica
- Podocoryna_carnea
- Hydractinia_echinata
- Nematostella_vectensis
- Anemonia_viridis
- Metridium_senile
- Aiptasia_pallida
- Porites_astreoides
- Montastraea_faveolata
- Acropora_millepora
- Acropora_palmata
- Branchiostoma_floridae
- Petromyzon_marinus
- Gallus
- Halocynthia_roretzi
- Ciona_intestinalis
- Xenoturbella_bocki
- Saccoglossus_kowalevskii
- Ptychodera_flava
- Strongylocentrotus_purpuratus
- Asterina_pectinifera
- Nemertoderma_westbladi
- Meara_stichopi
- Isodiametra_pulchra
- Symsagittifera_roscoffensis
- Convolutriloba_longifissura
- Echinoderes_horni
- Xiphinema_index
- Euperipatoides_kanangrensis
- Drosophila_melanogaster
- Daphnia_pulex
- Boophilus_microplus
- Anoplodactylus_eroticus
- Terebratalia_transversa
- Cerebratulus_lacteus
- Helobdella_robusta
- Capitella_telata
- Schmidtea_mediterranea
- Paraplanocera_oligoglena
- Euprymna_scolopes
- Lottia_gigantea
- Crassostrea_virginica

0.3

1318    **s) PhyloBayes EST.Animalia dataset  run 1**

1319



Pleurobrachia_pileus
Mnemiopsis
Mertensiid_sp
Oscarella_lobularis
Oscarella_carmela
Sycon_raphanus
Leucetta_chagosensis
Oopsacas_minuta
Carteriospongia_foliascens
Amphimedon_queenslandica
Suberites_domuncula
Lubomirskia_baicalensis
Ephydatia_muelleri
Trichoplax_adhaerens
Cyanea_capillata
Hydra_magnipapillata
Clytia_hemisphaerica
Podocoryna_carnea
Hydractinia_echinata
Nematostella_vectensis
Anemonia_viridis
Metridium_senile
Aiptasia_pallida
Porites_astreoides
Montastraea_faveolata
Acropora_millepora
Acropora_palmata
Saccoglossus_kowalevskii
Ptychodera_flava
Strongylocentrotus_purpuratus
Asterina_pectinifera
Branchiostoma_floridae
Petromyzon_marinus
Gallus
Halocynthia_roretzi
Ciona_intestinalis
Xenoturbella_bocki
Nemertoderma_westbladi
Meara_stichopi
Isodiametra_pulchra
Symsagittifera_roscoffensis
Convolutriloba_longifissura
Echinoderes_horni
Xiphinema_index
Euperipatoides_kanangrensis
Drosophila_melanogaster
Daphnia_pulex
Boophilus_microplus
Anoplodactylus_eroticus
Terebratalia_transversa
Cerebratulus_lacteus
Helobdella_robusta
Capitella_telata
Schmidtea_mediterranea
Paraplanocera_oligoglena
Euprymna_scolopes
Lottia_gigantea
Crassostrea_virginica

0.3

1320
1321

1322 **t) PhyloBayes EST.Animalia dataset run 2**

1323



Pleurobrachia_pileus
1 Mnemiopsis
0.53 Mertensiid_sp
1 Oscarella_lobularis
Oscarella_carmela
0.99
1 Sycon_raphanus
Leucetta_chagosensis
1
Oopsacas_minuta
0.99 Carteriospongia_foliascens
0.84 Amphimedon_queenslandica
1 Suberites_domuncula
0.67 Lubomirskia_baicalensis
1 Ephydatia_muelleri
1
Trichoplax_adhaerens
Cyanea_capillata
1 Hydra_magnipapillata
1 Clytia_hemisphaerica
0.99 Podocoryna_carnea
1 Hydractinia_echinata
1
0.96
Nematostella_vectensis
1 Anemonia_viridis
1 Metridium_senile
1 Aiptasia_pallida
1
Porites_astreoides
1 Montastraea_faveolata
1 Acropora_millepora
1 Acropora_palmata
1
Saccoglossus_kowalevskii
1 Ptychodera_flava
1 Strongylocentrotus_purpuratus
1 Asterina_pectinifera
Xenoturbella_bocki
0.77 Branchiostoma_floridae
0.98 Petromyzon_marinus
1 Gallus
1 Halocynthia_roretzi
1 Ciona_intestinalis
0.82 Nemertoderma_westbladi
1 Meara_stichopi
1 Isodiametra_pulchra
1 Symsagittifera_roscoffensis
1 Convolutriloba_longifissura
Echinoderes_horni
0.84 Xiphinema_index
0.95 Euperipatoides_kanangrensis
0.99 Drosophila_melanogaster
1 Daphnia_pulex
1 Boophilus_microplus
0.96 Anoplodactylus_eroticus
1
Terebratalia_transversa
0.91 Cerebratulus_lacteus
0.54 Helobdella_robusta
1 Capitella_telata
1
Schmidtea_mediterranea
1 Paraplanocera_oligoglena
0.54 Euprymna_scolopes
1 Lottia_gigantea
1 Crassostrea_virginica

0.3

1324

1325 **Figure S2: Maximum likelihood tree of an near intron pair matrix**
1326 A Hexagon at a node indicates the node was constrained in the analysis.
1327



1328
1329

1330 **Figure S3: Ionotropic glutamate receptor phylogeny of human and ctenophore**
1331 **sequences from 8 ctenophores**
1332 The tree demonstrates that the ionotropic glutamate receptors (iGluR) of ctenophores are
1333 not direct orthologs to AMPA (GRIA), NMDA (GRIN), kainate-type (GRIK), or delta2-
1334 like (GRID) glutamate receptors. NOTE: The *Pleurobrachia bachei* sequence
1335 (ADV31314) is erroneously labeled as a kainate-type in GenBank.
1336



1337

**Figure S4: Maximum-likelihood analysis of gene content data with known**
**relationships constrained**

1342 **Figure S5: Example of bona fide nested genes**
1343 A screenshot from the *Mnemiopsis* Genome Browser shows ML000127a and
1344 ML000128a nested within introns of ML000126a.
1345



1346

1347 **Figure S6: Example of likely spurious prediction of a nested gene**
1348 A screenshot from the *Mnemiopsis* Genome Browser shows ML000317a nested within an
1349 intron of ML000316a. Additional transcriptomic evidence that was not included in the
1350 initial gene prediction pipeline suggests that ML000317a is actually an exon incorporated
1351 in a rare isoform of ML000316a.
1352



1353

1354 **Figure S7: Comparison of lineage-specific genes**
1355 For each species, we count the occurrence of clusters that only include genes from that
1356 species (blue bars). We also count the number of genes within those clusters (red bars).
1357



1358
1359
1360

1361 **Figure S8: Comparison of gene duplications**
1362 For each cluster thought to represent a single gene present in filozoan and/or
1363 choanimalian ancestors where a species contains more than one gene, we count the
1364 number of genes in this cluster minus one and total this among all these clusters.
1365



Duplicates in Filozoa and Choanimalia Clusters

1366
1367

1368 **Figure S9: Comparison of gene losses**
1369 For each species, we count clusters in which a gene from this species is not present, but a
1370 gene from at least one metazoan and one non-metazoan are present.
1371



1372
1373

1374
**Figure S10: Hedgehog and Hedgling**
**a,** Domain structure of Hedgehog and Hedgling gene products. **b,** Phylogenetic
distribution of Hedge domains, Hog domains, Hedgling genes, and Hedgehog genes.
1378



1379
1380
1381

**References**

54. J. C. Mullikin, Z. Ning, The phusion assembler. *Genome Res* **13**, 81-90 (2003).
55. W. J. Kent, BLAT--the BLAST-like alignment tool. *Genome Res* **12**, 656-664 (2002).
56. J. F. Ryan, Baa. pl: A tool to evaluate de novo genome assemblies with RNA transcripts. *arXiv preprint arXiv:1309.2087*, (2013).
57. M. G. Grabherr, B. J. Haas, M. Yassour, J. Z. Levin, D. A. Thompson, I. Amit, X. Adiconis, L. Fan, R. Raychowdhury, Q. Zeng, Full-length transcriptome assembly from RNA-Seq data without a reference genome. *Nature biotechnology* **29**, 644-652 (2011).
58. C. Alkan, S. Sajjadian, E. E. Eichler, Limitations of next-generation genome sequence assembly. *Nature methods* **8**, 61-65 (2010).
59. W. Huang, L. Li, J. R. Myers, G. T. Marth, ART: a next-generation sequencing read simulator. *Bioinformatics* **28**, 593-594 (2012).
60. S. Kurtz, The Vmatch large scale sequence analysis software. *Ref Type: Computer Program*, 4-12 (2003).
61. A. F. A. Smit, R. Hubley, P. Green, RepeatMasker Open-3.0. (2004).
62. J. Jurka, V. V. Kapitonov, A. Pavlicek, P. Klonowski, O. Kohany, J. Walichiewicz, Repbase Update, a database of eukaryotic repetitive elements. *Cytogenetic and genome research* **110**, 462-467 (2005).
63. C. Trapnell, L. Pachter, S. L. Salzberg, TopHat: discovering splice junctions with RNA-Seq. *Bioinformatics* **25**, 1105-1111 (2009).
64. C. Trapnell, B. A. Williams, G. Pertea, A. Mortazavi, G. Kwan, M. J. van Baren, S. L. Salzberg, B. J. Wold, L. Pachter, Transcript assembly and quantification by RNA-Seq reveals unannotated transcripts and isoform switching during cell differentiation. *Nat Biotechnol* **28**, 511-515 (2010).
65. B. J. Haas, A. L. Delcher, S. M. Mount, J. R. Wortman, R. K. Smith Jr, L. I. Hannick, R. Maiti, C. M. Ronning, D. B. Rusch, C. D. Town, Improving the Arabidopsis genome annotation using maximal transcript alignment assemblies. *Nucleic Acids Research* **31**, 5654-5666 (2003).
66. A. A. Salamov, V. V. Solovyev, Ab initio gene finding in Drosophila genomic DNA. *Genome Res* **10**, 516-522 (2000); published online EpubApr (
67. M. Stanke, S. Waack, Gene prediction with a hidden Markov model and a new intron submodel. *Bioinformatics* **19**, ii215-ii225 (2003).
68. A. Krogh, Two methods for improving performance of an HMM and their application for gene finding. *Center for Biological Sequence Analysis. Phone* **45**, 4525 (1997).
69. R. F. Yeh, L. P. Lim, C. B. Burge, Computational inference of homologous gene structures in the human genome. *Genome research* **11**, 803-816 (2001).
70. B. J. Haas, S. L. Salzberg, W. Zhu, M. Pertea, J. E. Allen, J. Orvis, O. White, C. R. Buell, J. R. Wortman, Automated eukaryotic gene structure annotation using EVidenceModeler and the Program to Assemble Spliced Alignments. *Genome Biol* **9**, R7 (2008).
71. T. D. Wu, C. K. Watanabe, GMAP: a genomic mapping and alignment program for mRNA and EST sequences. *Bioinformatics* **21**, 1859-1875 (2005).

1428    72.    B. J. Koch, J. F. Ryan, A. D. Baxevanis, The Diversification of the LIM
1429           Superclass at the Base of the Metazoa Increased Subcellular Complexity and
1430           Promoted Multicellular Specialization. *PLoS One* **7**, e33261 (2012).
1431    73.    C. E. Schnitzler, K. Pang, M. L. Powers, A. M. Reitzel, J. F. Ryan, D. Simmons,
1432           T. Tada, M. Park, J. Gupta, S. Y. Brooks, Genomic organization, evolution, and
1433           expression of photoprotein and opsin genes in Mnemiopsis leidyi: a new view of
1434           ctenophore photocytes. *BMC Biology* **10**, 107 (2012).
1435    74.    D. K. Simmons, K. Pang, M. Q. Martindale, Lim homeobox genes in the
1436           Ctenophore Mnemiopsis leidyi: the evolution of neural cell type specification.
1437           *Evodevo* **3**, 2 (2012).
1438    75.    E. A. Gladyshev, M. Meselson, I. R. Arkhipova, Massive horizontal gene transfer
1439           in bdelloid rotifers. *Science* **320**, 1210-1213 (2008).
1440    76.    S. F. Altschul, T. L. Madden, A. A. Schaffer, J. Zhang, Z. Zhang, W. Miller, D. J.
1441           Lipman, Gapped BLAST and PSI-BLAST: a new generation of protein database
1442           search programs. *Nucleic Acids Res* **25**, 3389-3402 (1997).
1443    77.    N. H. Putnam, M. Srivastava, U. Hellsten, B. Dirks, J. Chapman, A. Salamov, A.
1444           Terry, H. Shapiro, E. Lindquist, V. V. Kapitonov, J. Jurka, G. Genikhovich, I. V.
1445           Grigoriev, S. M. Lucas, R. E. Steele, J. R. Finnerty, U. Technau, M. Q.
1446           Martindale, D. S. Rokhsar, Sea anemone genome reveals ancestral eumetazoan
1447           gene repertoire and genomic organization. *Science* **317**, 86-94 (2007).
1448    78.    N. H. Putnam, T. Butts, D. E. K. Ferrier, R. F. Furlong, U. Hellsten, T.
1449           Kawashima, M. Robinson-Rechavi, E. Shoguchi, A. T. J. K. Yu, The amphioxus
1450           genome and the evolution of the chordate karyotype. *Nature* **453**, 1064-1071
1451           (2008).
1452    79.    M. Srivastava, E. Begovic, J. Chapman, N. H. Putnam, U. Hellsten, T.
1453           Kawashima, A. Kuo, T. Mitros, A. Salamov, M. L. Carpenter, A. Y. Signorovitch,
1454           M. A. Moreno, K. Kamm, J. Grimwood, J. Schmutz, H. Shapiro, I. V. Grigoriev,
1455           L. W. Buss, B. Schierwater, S. L. Dellaporta, D. S. Rokhsar, The Trichoplax
1456           genome and the nature of placozoans. *Nature* **454**, 955-960 (2008).
1457    80.    M. Srivastava, O. Simakov, J. Chapman, B. Fahey, M. E. Gauthier, T. Mitros, G.
1458           S. Richards, C. Conaco, M. Dacre, U. Hellsten, C. Larroux, N. H. Putnam, M.
1459           Stanke, M. Adamska, A. Darling, S. M. Degnan, T. H. Oakley, D. C. Plachetzki,
1460           Y. Zhai, M. Adamski, A. Calcino, S. F. Cummins, D. M. Goodstein, C. Harris, D.
1461           J. Jackson, S. P. Leys, S. Shu, B. J. Woodcroft, M. Vervoort, K. S. Kosik, G.
1462           Manning, B. M. Degnan, D. S. Rokhsar, The Amphimedon queenslandica
1463           genome and the evolution of animal complexity. *Nature* **466**, 720-726 (2010).
1464    81.    O. Simakov, F. Marletaz, S. J. Cho, E. Edsinger-Gonzales, P. Havlak, U. Hellsten,
1465           D. H. Kuo, T. Larsson, J. Lv, D. Arendt, R. Savage, K. Osoegawa, P. de Jong, J.
1466           Grimwood, J. Chapman, H. Shapiro, A. Aerts, R. P. Otillar, A. Y. Terry, J. L.
1467           Boore, I. V. Grigoriev, D. R. Lindberg, E. C. Seaver, D. A. Weisblat, N. Putnam,
1468           D. S. Rokhsar, Insights into bilaterian evolution from three spiralian genomes.
1469           *Nature* **493**, 526-531 (*2012*).
1470    82.    L. Li, C. J. Stoeckert, D. S. Roos, OrthoMCL: identification of ortholog groups
1471           for eukaryotic genomes. *Genome research* **13**, 2178-2189 (2003).
1472    83.    L. S. Johnson, S. R. Eddy, E. Portugaly, Hidden Markov model speed heuristic
1473           and iterative HMM search procedure. *BMC Bioinformatics* **11**, 431 (2010).

1474   84.   M. Punta, P. C. Coggill, R. Y. Eberhardt, J. Mistry, J. Tate, C. Boursnell, N. Pang,
1475         K. Forslund, G. Ceric, J. Clements, A. Heger, L. Holm, E. L. Sonnhammer, S. R.
1476         Eddy, A. Bateman, R. D. Finn, The Pfam protein families database. *Nucleic Acids*
1477         *Res* **40**, D290-301 (2012).
1478   85.   K. Katoh, G. Asimenos, H. Toh, Multiple alignment of DNA sequences with
1479         MAFFT. *Methods Mol Biol* **537**, 39-64 (2009).
1480   86.   G. Talavera, J. Castresana, Improvement of phylogenies after removing divergent
1481         and ambiguously aligned blocks from protein sequence alignments. *Systematic*
1482         *Biology* **56**, 564-577 (2007).
1483   87.   R. M. Waterhouse, E. M. Zdobnov, F. Tegenfeldt, J. Li, E. V. Kriventseva,
1484         OrthoDB: the hierarchical catalog of eukaryotic orthologs in 2011. *Nucleic Acids*
1485         *Res* **39**, D283-288 (2011).
1486   88.   A. Hejnol, M. Obst, A. Stamatakis, M. Ott, G. W. Rouse, G. D. Edgecombe, P.
1487         Martinez, J. Baguna, X. Bailly, U. Jondelius, M. Wiens, W. E. Muller, E. Seaver,
1488         W. C. Wheeler, M. Q. Martindale, G. Giribet, C. W. Dunn, Assessing the root of
1489         bilaterian animals with scalable phylogenomic methods. *Proc Biol Sci* **276**, 4261-
1490         4270 (2009).
1491   89.   D. L. Swofford, PAUP*: phylogenetic analysis using parsimony, version 4.0 b10.
1492         (2003).
1493   90.   A. Rokas, B. L. Williams, N. King, S. B. Carroll, Genome-scale approaches to
1494         resolving incongruence in molecular phylogenies. *Nature* **425**, 798-804 (2003).
1495   91.   A. Rambaut, A. Drummond, FigTree v1. 3.1. *Program distributed by the author.*
1496         *Institute of Evolutionary Biology, University of Edinburgh. Edinburgh, United*
1497         *Kingdom*, (2009).
1498   92.   A. Rambaut, A. J. Drummond, TreeStat v1. 2: tree statistic calculation tool.
1499         *Program distributed by the author. Institute of Evolutionary Biology, University*
1500         *of Edinburgh. Edinburgh, United Kingdom*, (2008).
1501   93.   S. A. Berger, A. Stamatakis, R. Lucking, Morphology-based phylogenetic binning
1502         of the lichen genera Graphis and Allographa (Ascomycota: Graphidaceae) using
1503         molecular site weight calibration. *Taxon* **60**, 1450-1457 (2011).
1504   94.   H. Fang, M. E. Oates, R. B. Pethica, J. M. Greenwood, A. J. Sardar, O. J. L.
1505         Rackham, P. C. J. Donoghue, A. Stamatakis, D. A. de Lima Morais, J. Gough, A
1506         daily-updated tree of (sequenced) life as a reference for genome research.
1507         *Scientific Reports* **3**, (2013).
1508   95.   H. Shimodaira, M. Hasegawa, CONSEL: for assessing the confidence of
1509         phylogenetic tree selection. *Bioinformatics* **17**, 1246-1247 (2001).
1510   96.   H. Shimodaira, An approximately unbiased test of phylogenetic tree selection.
1511         *Systematic Biology* **51**, 492-508 (2002).
1512   97.   N. Goldman, J. P. Anderson, A. G. Rodrigo, Likelihood-based tests of topologies
1513         in phylogenetics. *Systematic Biology* **49**, 652-670 (2000).
1514   98.   J. Lehmann, P. F. Stadler, V. Krauss, Near intron pairs and the metazoan tree.
1515         *Molecular Phylogenetics and Evolution*, (2013).
1516   99.   I. Letunic, T. Doerks, P. Bork, SMART 7: recent updates to the protein domain
1517         annotation resource. *Nucleic Acids Res* **40**, D302-305 (2012).

1518　100. A. Alie, M. Manuel, The backbone of the post-synaptic density originated in a
1519　　　　unicellular ancestor of choanoflagellates and metazoans. *BMC Evol Biol* **10**, 34
1520　　　　(2010).
1521　101. O. Sakarya, K. A. Armstrong, M. Adamska, M. Adamski, I. F. Wang, B. Tidor, B.
1522　　　　M. Degnan, T. H. Oakley, K. S. Kosik, A post-synaptic scaffold at the origin of
1523　　　　the animal kingdom. *PLoS One* **2**, e506 (2007).
1524　103. M. Jager, R. Chiori, A. Alie, C. Dayraud, E. Queinnec, M. Manuel, New insights
1525　　　　on ctenophore neural anatomy: immunofluorescence study in Pleurobrachia pileus
1526　　　　(Muller, 1776). *J Exp Zool B Mol Dev Evol* **316B**, 171-187 (2011).
1527　104. A. Hay-Schmidt, The evolution of the serotonergic nervous system. *P Roy Soc
1528　　　　Lond B Bio* **267**, 1071-1079 (2000).
1529　105. M. Leptin, twist and snail as positive and negative regulators during Drosophila
1530　　　　mesoderm development. *Genes Dev* **5**, 1568-1576 (1991).
1531　106. N. Azpiazu, M. Frasch, tinman and bagpipe: two homeo box genes that determine
1532　　　　cell fates in the dorsal mesoderm of Drosophila. *Genes Dev* **7**, 1325-1340 (1993).
1533　107. K. Jagla, M. Bellard, M. Frasch, A cluster of Drosophila homeobox genes
1534　　　　involved in mesoderm differentiation programs. *Bioessays* **23**, 125-133 (2001).
1535　108. T. Sato, D. Rocancourt, L. Marques, S. Thorsteinsdottir, M. Buckingham, A
1536　　　　Pax3/Dmrt2/Myf5 regulatory cascade functions at the onset of myogenesis. *PLoS
1537　　　　Genet* **6**, e1000897 (2010).
1538　109. K. Ryan, N. Garrett, A. Mitchell, J. B. Gurdon, Eomesodermin, a key early gene
1539　　　　in Xenopus mesoderm differentiation. *Cell* **87**, 989-1000 (1996).
1540　110. P. J. Gianakopoulos, V. Mehta, A. Voronova, Y. Cao, Z. Yao, J. Coutu, X. Wang,
1541　　　　M. S. Waddington, S. J. Tapscott, I. S. Skerjanc, MyoD directly up-regulates
1542　　　　premyogenic mesoderm factors during induction of skeletal myogenesis in stem
1543　　　　cells. *J Biol Chem* **286**, 2517-2525 (2011).
1544　111. W. R. Francis, L. M. Christianson, R. Kiko, M. L. Powers, N. C. Shaner, S. H. D.
1545　　　　Haddock, A comparison across non-model animals suggests an optimal
1546　　　　sequencing depth for de novo transcriptome assembly. *BMC Genomics* **14**, 167
1547　　　　(2013).
1548　112. A. Stamatakis, RAxML-VI-HPC: maximum likelihood-based phylogenetic
1549　　　　analyses with thousands of taxa and mixed models. *Bioinformatics* **22**, 2688-2690
1550　　　　(2006).
1551　113. Y. f. Zhong, P. W. H. Holland, HomeoDB2: functional expansion of a
1552　　　　comparative homeobox gene database for evolutionary developmental biology.
1553　　　　*Evolution & Development* **13**, 567-568 (2011).
1554　114. J. F. Ryan, K. Pang, J. C. Mullikin, M. Q. Martindale, A. D. Baxevanis, The
1555　　　　homeodomain complement of the ctenophore Mnemiopsis leidyi suggests that
1556　　　　Ctenophora and Porifera diverged prior to the ParaHoxozoa. *Evodevo* **1**, 9 (2010).
1557　115. K. Pang, J. F. Ryan, A. D. Baxevanis, M. Q. Martindale, Evolution of the TGF-
1558　　　　beta signaling pathway and its potential role in the ctenophore, Mnemiopsis
1559　　　　leidyi. *PLoS One* **6**, e24152 (2011).
1560　116. M. E. Skinner, A. V. Uzilov, L. D. Stein, C. J. Mungall, I. H. Holmes, JBrowse: A
1561　　　　next-generation genome browser. *Genome research* **19**, 1630-1638 (2009).

1562    117.   A. Wallberg, M. Thollesson, J. S. Farris, U. Jondelius, The phylogenetic position
1563           of the comb jellies (Ctenophora) and the importance of taxonomic sampling.
1564           *Cladistics* **20**, 558-578 (2004).
1565    118.   J. F. Ryan, A. D. Baxevanis, Hox, Wnt, and the evolution of the primary body
1566           axis: insights from the early-divergent phyla. *Biol Direct* **2**, 37 (2007).
1567    119.   A. Lang, Die Polycladen des Golfes von Neapel und der angrenzenden
1568           Meeresabschnitte. *Fauna u. Flora Neapel* **11**, 1-688 (1884).
1569    120.   L. H. Hyman, *The invertebrates*. (McGraw-Hill, New York,, ed. 1st, 1940), pp. v.
1570    121.   H. Hadzi, *Turbellarijska teorija knidarijev (Turbellarien-Theorie der Knidarien)*.
1571           Slovenian Academy of Sciences and Arts (Ljubljana, Slovenia, 1944), vol. 3, pp.
1572           238.
1573    122.   R. C. Brusca, G. J. Brusca, Invertebrates. *Sinauer, Sunderland, Mass* **264**, (1990).
1574    123.   U. Ehlers, Ultrastructure of the spermatozoa of Halammohydra schulzei
1575           (Cnidaria, Hydrozoa): The significance of acrosomal structures for the
1576           systematization of the Eumetazoa. *Microfauna Marina* **8**, 115-130 (1993).
1577    124.   E. E. Ruppert, R. D. Barnes, R. S. Fox, *Invertebrate zoology*. (Saunders College
1578           Pub., 1994), vol. 6.
1579    125.   C. Nielsen, *Animal evolution: interrelationships of the animal phyla*. (Oxford
1580           University Press, Oxford, 1995).
1581    126.   C. Nielsen, Cladistic analyses of the animal kingdom. *Biological Journal of the
1582           Linnean Society* **57**, 385-410 (1996).
1583    127.   P. Ax, *Multicellular animals. A New Approach to the Phylogenetic Order in
1584           Nature, Volume I* (Springer, Berlin ; New York, 1996).
1585    128.   L. Margulis, K. V. Schwartz, *Five Kingdoms: An Illustrated Guide to the Phyla of
1586           Life on Earth*. (Freeman New York, ed. 3, 1998).
1587    129.   C. Nielsen, *Animal evolution : interrelationships of the living phyla*. (Oxford
1588           University Press, Oxford, ed. 2nd, 2001), pp. X, 563 s.
1589    130.   P. O. Wainright, G. Hinkle, M. L. Sogin, S. K. Stickel, Monophyletic origins of
1590           the metazoa: an evolutionary link with fungi. *Science* **260**, 340-342 (1993).
1591    131.   T. Katayama, H. Wada, H. Furuya, N. Satoh, M. Yamamoto, Phylogenetic
1592           position of the dicyemid mesozoa inferred from 18S rDNA sequences. *Biological
1593           Bulletin* **189**, 81-90 (1995).
1594    132.   B. Hanelt, D. Van Schyndel, C. M. Adema, L. A. Lewis, E. S. Loker, The
1595           phylogenetic position of Rhopalura ophiocomae (Orthonectida) based on 18S
1596           ribosomal DNA sequence analysis. *Molecular Biology and Evolution* **13**, 1187-
1597           1191 (1996).
1598    133.   Y. Van De Peer, R. De Wachter, Evolutionary relationships among the eukaryotic
1599           crown taxa taking into account site-to-site rate variation in 18S rRNA. *Journal of
1600           Molecular Evolution* **45**, 619-630 (1997).
1601    134.   E. Abouheif, R. Zardoya, A. Meyer, Limitations of metazoan 18s rRNA sequence
1602           data: Implications for reconstructing a phylogeny of the animal kingdom and
1603           inferring the reality of the cambrian explosion. *Journal of Molecular Evolution*
1604           **47**, 394-405 (1998).
1605    135.   A. G. Collins, Evaluating multiple alternative hypotheses for the origin of
1606           Bilateria: an analysis of 18S rRNA molecular evidence. *Proc. Natl. Acad. Sci.
1607           USA* **95**, 15458-15463 (1998).

1608    136.    K. M. Halanych, Considerations for Reconstructing Metazoan History: Signal,
1609            Resolution, and Hypothesis Testing. *Integrative and Comparative Biology* **38**,
1610            929-941 (1998).
1611    137.    D. L. Lipscomb, J. S. Farris, M. Ka?llersjo, A. Tehler, Support, ribosomal
1612            sequences and the phylogeny of the eukaryotes. *Cladistics* **14**, 303-338 (1998).
1613    138.    B. M. H. Winnepenninckx, Y. D. E. Van Peer, T. Backeljau, Metazoan
1614            Relationships on the Basis of 18S rRNA Sequences: A Few Years Later.
1615            *Integrative and Comparative Biology* **38**, 888-906 (1998).
1616    139.    J. Zrzavy, S. Mihulka, P. Kepka, A. Bezde?k, D. Tietz, Phylogeny of the Metazoa
1617            Based on Morphological and 18S Ribosomal DNA Evidence. *Cladistics* **14**, 249-
1618            285 (1998).
1619    140.    J. Kim, W. Kim, C. W. Cunningham, A new perspective on lower metazoan
1620            relationships from 18S rDNA sequences [2]. *Molecular Biology and Evolution* **16**,
1621            423-427 (1999).
1622    141.    G. Giribet, W. C. Wheeler, The Position of Arthropods in the Animal Kingdom:
1623            Ecdysozoa, Islands, Trees, and the "Parsimony Ratchet". *Molecular Phylogenetics*
1624            *and Evolution* **13**, 619-623 (1999).
1625    142.    M. E. Siddall, M. F. Whiting, Long-branch abstractions. *Cladistics* **15**, 9-24
1626            (1999).
1627    143.    M. Medina, A. G. Collins, J. D. Silberman, M. L. Sogin, Evaluating hypotheses of
1628            basal animal phylogeny using complete sequences of large and small subunit
1629            rRNA. *Proc Natl Acad Sci U S A* **98**, 9707-9712 (2001).
1630    144.    K. J. Peterson, D. J. Eernisse, Animal phylogeny and the ancestry of bilaterians:
1631            Inferences from morphology and 18S rDNA gene sequences. *Evolution and*
1632            *Development* **3**, 170-205 (2001).
1633    145.    M. Podar, S. H. Haddock, M. L. Sogin, G. R. Harbison, A molecular phylogenetic
1634            framework for the phylum Ctenophora using 18S rRNA genes. *Mol Phylogenet*
1635            *Evol* **21**, 218-230 (2001).
1636    146.    A. G. Collins, Phylogeny of Medusozoa and the evolution of cnidarian life cycles.
1637            *Journal of Evolutionary Biology* **15**, 418-432 (2002).
1638    147.    C. Martinelli, J. Spring, Distinct expression patterns of the two T-box homologues
1639            Brachyury and Tbx2/3 in the placozoan Trichoplax adhaerens. *Dev Genes Evol*
1640            **213**, 492-499 (2003).
1641    148.    J. Zrzavy, V. Hypsa, Myxozoa, Polypodium, and the origin of the Bilateria: The
1642            phylogenetic position of "Endocnidozoa" in light of the rediscovery of
1643            Buddenbrockia [1]. *Cladistics* **19**, 164-169 (2003).
1644    149.    C. W. Dunn, A. Hejnol, D. Q. Matus, K. Pang, W. E. Browne, S. A. Smith, E.
1645            Seaver, G. W. Rouse, M. Obst, G. D. Edgecombe, M. V. Sorensen, S. H.
1646            Haddock, A. Schmidt-Rhaesa, A. Okusu, R. M. Kristensen, W. C. Wheeler, M. Q.
1647            Martindale, G. Giribet, Broad phylogenomic sampling improves resolution of the
1648            animal tree of life. *Nature* **452**, 745-749 (2008).
1649    150.    H. Philippe, R. Derelle, P. Lopez, K. Pick, C. Borchiellini, N. Boury-Esnault, J.
1650            Vacelet, E. Renard, E. Houliston, E. Queinnec, C. Da Silva, P. Wincker, H. Le
1651            Guyader, S. Leys, D. J. Jackson, F. Schreiber, D. Erpenbeck, B. Morgenstern, G.
1652            Worheide, M. Manuel, Phylogenomics revives traditional views on deep animal
1653            relationships. *Curr Biol* **19**, 706-712 (2009).

1654 151. B. Schierwater, M. Eitel, W. Jakob, H. J. Osigus, H. Hadrys, S. L. Dellaporta, S.
1655      O. Kolokotronis, R. Desalle, Concatenated analysis sheds light on early metazoan
1656      evolution and fuels a modern "urmetazoon" hypothesis. *PLoS Biol* **7**, e20 (2009).
1657 152. K. S. Pick, H. Philippe, F. Schreiber, D. Erpenbeck, D. J. Jackson, P. Wrede, M.
1658      Wiens, A. Alie, B. Morgenstern, M. Manuel, G. Worheide, Improved
1659      phylogenomic taxon sampling noticeably affects non-bilaterian relationships. *Mol*
1660      *Biol Evol*, (2010).
1661 153. J. Mallatt, C. W. Craig, M. J. Yoder, Nearly complete rRNA genes from 371
1662      Animalia: updated structure-based alignment and detailed phylogenetic analysis.
1663      *Mol Phylogenet Evol* **64**, 603-617 (2012).
1664 154. A. Kumar, An overview of nested genes in eukaryotic genomes. *Eukaryot Cell* **8**,
1665      1321-1329 (2009).
1666 155. T. Adell, V. A. Grebenjuk, M. Wiens, W. E. Muller, Isolation and
1667      characterization of two T-box genes from sponges, the phylogenetically oldest
1668      metazoan taxon. *Dev Genes Evol* **213**, 421-434 (2003).
1669 156. M. Manuel, Y. Le Parco, C. Borchiellini, Comparative analysis of Brachyury T-
1670      domains, with the characterization of two new sponge sequences, from a
1671      hexactinellid and a calcisponge. *Gene* **340**, 291-301 (2004).
1672 157. D. E. Martinez, M. L. Dirksen, P. M. Bode, M. Jamrich, R. E. Steele, H. R. Bode,
1673      Budhead, a fork head/HNF-3 homologue, is expressed during axis formation and
1674      head specification in hydra. *Dev Biol* **192**, 523-536 (1997).
1675
1676