

1 **The genome of the crustacean *Parhyale***
2 ***hawaiensis*: a model for animal**
3 **development, regeneration, immunity**
4 **and lignocellulose digestion**

5 **Damian Kao¹, Alvina G. Lai¹, Evangelia Stamatakis², Silvana Rosic^{3,4},**
6 **Nikolaos Konstantinides⁵, Erin Jarvis⁶, Alessia Di Donfrancesco¹, Natalia**
7 **Pouchkina-Stantcheva¹, Marie Sèmon⁵, Marco Grillo⁵, Heather Bruce⁶,**
8 **Suyash Kumar², Igor Siwanowicz², Andy Le², Andrew Lemire²,**
9 **Cassandra Extavour⁷, William Browne⁸, Carsten Wolff⁹, Michalis Averof⁵,**
10 **Nipam H. Patel⁶, Peter Sarkies^{3,4}, Anastasios Pavlopoulos², and A. Aziz**
11 **Aboobaker¹**

12 ¹**Department of Zoology, University of Oxford, South Parks Road, Oxford OX1 3PS,**
13 **United Kingdom**

14 ²**Howard Hughes Medical Institute, Janelia Research Campus, 19700 Helix Drive,**
15 **Ashburn, Virginia 20147, USA**

16 ³**MRC Clinical Sciences Centre, Imperial College London, Hammersmith Hospital**
17 **Campus, Du Cane Road, London W12 0NN**

18 ⁴**Institute for Clinical Sciences, Imperial College London, Hammersmith Hospital**
19 **Campus, Du Cane Road, London, W12 0NN**

20 ⁵**Institut de Génomique Fonctionnelle de Lyon, Lyon, France**

21 ⁶**University of California, Berkeley, Dept. of Molecular and Cell Biology, 519A LSA 3200**
22 **Berkeley, CA 94720-3200**

23 ⁷**Department of Organismic and Evolutionary Biology, Harvard University 16 Divinity**
24 **Avenue, BioLabs Building 4109-4111 Cambridge, MA 02138**

25 ⁸**Cox Science Center, 1301 Memorial Drive, Coral Gables, FL 33146, USA**

26 ⁹**Humboldt-Universität zu Berlin, Institut für Biologie, Vergleichende Zoologie,**
27 **Philippstr. 13, Haus 2, 10115 Berlin, Germany**

28 **Author information**

29 These authors contributed equally to this work: Damian Kao, Alvina G. Lai,
30 Evangelia Stamataki.

31 These authors also contributed equally: Anastasios Pavlopoulos, A. Aziz
32 Aboobaker

33 **Author contributions**

34 AP and AAA conceived, designed and managed project. All authors acquired,
35 analysed and interpreted data. DK, AGL, ES, AP and AAA drafted the
36 manuscript and revised it with input from all other authors.

37 **For correspondence**

38 Aziz.Aboobaker@zoo.ox.ac.uk (AAA)

39 pavlopoulosa@janelia.hhmi.org (AP)

40

41 **Competing interests**

42 The authors declare no competing interests

58

43

44 **Funding**

45 AAA and co-workers are funded by the Biotechnology and Biological Sciences
46 Research Council (BBSRC grant number BB/K007564/1), the Medical
47 Research Council (MRC grant number MR/M000133/1) and the John Fell
48 Fund Oxford University Press (OUP). AGL receives funding from the Human
49 Frontier Science Program postdoctoral fellowship and the Elizabeth Hannah
50 Jenkinson Research Fund. NHP and co-workers are funded by NSF grant
51 IOS-1257379. AP and co-workers are funded by the Howard Hughes Medical
52 Institute. PS and co-workers are funded by the Medical Research council
53 (MRC MC-A652- 5PZ80) and an Imperial College Research Fellowship too
54 PS. MA and colleagues received funding from the Agence Nationale de la
55 Recherche (France), grant ANR-12-CHEX-0001-01. The funding bodies had
56 no role in study design, data collection and interpretation, or the decision to
57 submit the work for publication.

59 ABSTRACT

60 *Parhyale hawaiiensis* is a blossoming model system for studies of developmental mechanisms and
61 more recently adult regeneration. We have sequenced the genome allowing annotation of all key
62 signaling pathways, small non-coding RNAs and transcription factors that will enhance ongoing functional
63 studies. *Parhyale* is a member of the malacostraca, which includes crustacean food crop species. We
64 analysed the immunity related genes of *Parhyale* as an important comparative system for these species,
65 where immunity related aquaculture problems have increased as farming has intensified. We also find
66 that *Parhyale* and other species within multicrustacea contain the enzyme sets necessary to perform
67 lignocellulose digestion (“wood eating”), suggesting this ability may predate the diversification of this
68 lineage. Our data provide an essential resource for further development of the *Parhyale model*. The first
69 malacostracan genome sequence will underpin ongoing comparative work in important food crop species
70 and research investigating lignocellulose as energy source.

71 INTRODUCTION

72 Very few members of the Animal Kingdom hold the esteemed position of major model system for
73 understanding living systems. Inventions in molecular and cellular biology increasingly facilitate the
74 emergence of new experimental systems for developmental genetic studies. The morphological and
75 ecological diversity of the phylum arthropoda makes them an ideal group of animals for comparative
76 studies encompassing embryology, adaptation of adult body plans and life history evolution [1–4]. While
77 the most widely studied group are hexapods, reflected by over a hundred sequencing projects available in
78 the NCBI genome database, genomic data in the other three sub-phyla in arthropoda are still relatively
79 sparse.

80 Recent molecular and morphological studies have placed crustaceans along with hexapods into a
81 Pancrustacean clade (Figure 1A), revealing that crustaceans are paraphyletic [5–9]. Previously, the only
82 available fully sequenced crustacean genome was that of the water flea *D. pulex* which is a member of
83 the branchiopoda [10]. A growing number of transcriptomes for larger phylogenetic analyses have led to
84 differing hypotheses of the relationships of the major Pancrustacean groups (Figure 1B) [11–14]. The
85 *Parhyale* genome addresses the paucity of high quality non-hexapod genomes among the pancrustacean
86 group, and will help to resolve relationships within this group. Crucially, genome sequence data is also
87 necessary to further advance research with *Parhyale*, currently the most tractable crustacean model system.
88 This is particular true for the application of powerful functional genomic approaches, such as genome
89 editing [15–20].

90 *Parhyale* is a member of the diverse malacostraca clade with thousands of extant species including
91 economically and nutritionally important groups such as shrimps, crabs, crayfish and lobsters, as well as
92 common garden animals like woodlice. They are found in all marine, fresh water, and higher humidity
93 terrestrial environments. Apart from attracting research interest as an economically important food

94 crop species, this group of animals has been used to study developmental biology and the evolution of
95 morphological diversity (for example with respect to *Hox* genes) [17, 21–23], stem cell biology [24, 25],
96 innate immunity processes [26, 27] and recently the cellular mechanisms of limb regeneration [24, 28, 29].
97 In addition, members of the malacostraca, specifically both amphipods and isopods, are thought to be
98 capable of “wood eating” or lignocellulose digestion and to have microbiota-free digestive systems
99 [30–33].

100 The life history of *Parhyale* makes it a versatile model organism amenable to experimental manip-
101 ulations (Figure 1C)[34]. Gravid females lay eggs every 2 weeks upon reaching sexual maturity and
102 hundreds of eggs can be easily collected at all stages of embryogenesis. Embryogenesis takes about
103 10 days at 26°C and has been described in detail with an accurate staging system [35]. Early embryos
104 display an invariant cell lineage with blastomere becoming committed to a single germ layer at the
105 8-cell stage (Figure 1D)[35, 36]. Embryonic and post-embryonic stages are amenable to experimental
106 manipulations and direct observation *in vivo* [36–47]. This can be combined with transgenic approaches
107 [23, 45, 48, 49], RNA interference (RNAi) [22] and morpholino-mediated gene knockdown [50], and
108 transgene-based lineage tracing [24]. Most recently the utility of the clustered regularly interspaced short
109 palindromic repeats (CRISPR)/CRISPR-associated (Cas) system for targeted genome editing has been
110 elegantly demonstrated during the systematic study of *Parhyale* Hox genes [16, 17]. This arsenal of
111 experimental tools (Table 1) has already established *Parhyale* as an attractive model system for modern
112 research.

113 So far, work in *Parhyale* has been constrained by the lack of of a reference genome and other
114 standardized genome-wide resources. To address this limitation, we have sequenced, assembled and
115 annotated the genome. At an estimated size of 3.6 Gb, this genome represents one of the largest animal
116 genomes tackled to date. The large size has not been the only challenge of the *Parhyale* genome that also
117 exhibited some of the highest levels of sequence repetitiveness, heterozygosity and polymorphism reported
118 among published genomes. We provide information in our assembly regarding polymorphism to facilitate
119 functional genomic approaches sensitive to levels of sequence similarity, particularly homology-dependent
120 genome editing approaches. We analysed a number of key features of the genome as foundations for
121 new areas of research in *Parhyale*, including innate immunity in crustaceans, lignocellulose digestion,
122 non-coding RNA biology, and epigenetic control of the genome. Our data bring *Parhyale* to the forefront
123 of developing model systems for a broad swathe of important bioscience research questions.

124 **RESULTS AND DISCUSSION**

125 **Genome assembly, annotation, and validation**

126 The *Parhyale* genome contains 23 pairs ($2n=46$) of chromosomes (Figure 2) and with an estimated size of
127 3.6 Gb, it is the second largest reported arthropod genome after the locust genome [51, 52]. Sequencing
128 was performed on genomic DNA isolated from a single adult male. We performed k-mer analyses of the
129 trimmed reads to assess the impact of repeats and polymorphism on the assembly process. We analyzed

130 k-mer frequencies (Figure 3A) and compared k-mer representation between our different sequencing
131 libraries. We observed a 93% intersection of unique k-mers among sequencing libraries, indicating that
132 the informational content was consistent between libraries (Supplemental HTML:assembly). Notably, we
133 observed k-mer frequency peaks at 60x and 120x coverage. While lowering k-mer length reduced the
134 number of k-mers at around 60x coverage, this peak was still apparent down to a k-mer length of 20. This
135 suggested a very high level of heterozygosity in the single male we sequenced.

136 In order to quantify global heterozygosity and repetitiveness of the genome we assessed the de-Bruijn
137 graphs generated from the trimmed reads to observe the frequency of both variant and repeat branches
138 [53] (Figure 3B and C). We found that the frequency of the variant branches was 10x higher than that
139 observed in the human genome and very similar to levels in the highly polymorphic genome of the oyster
140 *Crassostrea gigas* [54]. We also observed a frequency of repeat branches approximately 4x higher than
141 those observed in both the human and oyster genomes (Figure 3C), suggesting that the large size of the
142 *Parhyale* genome can be partly attributed to the expansion of repetitive sequences.

143 These metrics suggested that both contig assembly and scaffolding with mate pair reads were likely
144 to be challenging due to high heterozygosity and repeat content. After an initial contig assembly we
145 remapped reads to assess coverage of each contig. We observed a major peak centered around 75 x
146 coverage and a smaller peak at 150x coverage, reflecting high levels of heterozygosity. This resulted in
147 independent assembly of haplotypes for much of the genome (Figure 3D).

148 One of the prime goals in sequencing the *Parhyale* genome was to achieve an assembly that could
149 assist functional genetic and genomic approaches in this species. Therefore, we aimed for an assembly
150 representative of different haplotypes, allowing manipulations to be targeted to different allelic variants in
151 the assembly. This could be particularly important for homology dependent strategies that are likely to be
152 sensitive to polymorphism. However, the presence of alternative alleles could lead to poor scaffolding
153 as many mate-pair reads may not have uniquely mapping locations to distinguish between alleles in the
154 assembly. To alleviate this problem we conservatively identified pairs of allelic contigs and proceeded
155 to use only one in the scaffolding process. First, we estimated levels of similarity (both identity and
156 alignment length) between all assembled contigs to identify independently assembled allelic regions
157 (Figure 3E). We then kept the longer contig of each pair for scaffolding using our mate-pair libraries
158 (Figure 3F), after which we added back the shorter allelic contigs to produce the final genome assembly
159 (Figure 4A).

160 RepeatModeler and RepeatMasker were used on the final assembly to find repetitive regions, which
161 were subsequently classified into families of transposable elements or short tandem repeats (Supplemental
162 HTML:repeat). We found 1,473 different repeat element sequences representing 57% of the assembly
163 (Supplemental Table:repeatClassification). The *Parhyale* assembly comprises of 133,035 scaffolds (90%
164 of assembly), 259,343 unplaced contigs (4% of assembly), and 584,392 shorter, potentially allelic contigs
165 (6% of assembly), with a total length of 4.02 Gb (Table 2). The N50 length of the scaffolds is 81,190bp.
166 The final genome assembly was annotated with Augustus trained with high confidence gene models

167 derived from assembled transcriptomes, gene homology, and *ab initio* predictions. This resulted in 28,155
168 final gene models (Figure 4B; Supplemental HTML :annotation) across 14,805 genic scaffolds and 357
169 unplaced contigs with an N50 of 161,819, bp and an N90 of 52,952 bp.

170 *Parhyale* has a mean coding gene size (introns and ORFs) of 20kb (median of 7.2kb), which is longer
171 than *D. pulex* (mean: 2kb, median: 1.2kb), while shorter than genes in *Homo sapiens* (mean: 52.9kb,
172 median: 18.5kb). This difference in gene length was consistent across reciprocal blast pairs where ratios of
173 gene lengths revealed *Parhyale* genes were longer than *Caenorhabditis elegans*, *D. pulex*, and *Drosophila*
174 *melanogaster* and similar to *H. sapiens*. (Figure 5A). The mean intron size in *Parhyale* is 5.4kb, similar to
175 intron size in *H. sapiens* (5.9kb) but dramatically longer than introns in *D. pulex* (0.3kb), *D. melanogaster*
176 (0.3kb) and *C. elegans* (1kb) (Figure 5B).

177 For downstream analyses of *Parhyale* protein coding content, a final proteome consisting of 28,666
178 proteins was generated by combining candidate coding sequences identified with TransDecoder [55] from
179 mixed stage transcriptomes. We also included additional high confidence gene predictions that were
180 not found in the transcriptome (Figure 4C). The canonical proteome dataset was annotated with both
181 Pfam, KEGG, and BLAST against Uniprot. Assembly quality was further evaluated by alignment to core
182 eukaryotic genes defined by the Core Eukaryotic Genes Mapping Approach (CEGMA) database [56].
183 We identified 244/248 CEGMA orthology groups from the assembled genome alone and 247/248 with a
184 combination of genome and mapped transcriptome data (Supplemental Figure:cegma). Additionally, 96%
185 of over 280,000 identified transcripts, most of which are fragments that do not contain a large ORF, also
186 mapped to the assembled genome. Together these data suggest that our assembly is close to complete
187 with respect to protein coding genes and transcribed regions that are captured by deep RNA sequencing.

188 **High levels of heterozygosity and polymorphism in the *Parhyale* genome**

189 To estimate the heterozygosity rate in coding regions we first calculated the coverage of genomic reads
190 and rate of heterozygosity for each gene (Figure 6A; Supplemental HTML:variant). This led to most
191 genes falling either into a low coverage or high coverage group of mapped genomic DNA reads. Genes
192 that fell within the higher read coverage group generally had a lower mean heterozygosity rate (mean
193 1.09% of bases displaying polymorphism) than genes that fall within the lower read coverage group
194 (2.68%) (Figure 6B). This is consistent with genes achieving higher mapped genomic read coverage due
195 to having less divergent alleles.

196 The assembled *Parhyale* transcriptome was derived from various laboratory populations, hence we
197 expected to see additional polymorphisms beyond the founder haplotypes of the isofemale Chicago-F
198 strain used for sequencing the genome. Analysing all genes using the transcriptome we found additional
199 variations not found from the genomic reads. We observed that additional variations are not substantially
200 different between genes from the higher (0.88%) versus lower coverage group genes (0.73%; Figure 6C),
201 suggesting that heterozygosity and population variance are independent of each other. We also performed
202 an assessment of polymorphism on previously cloned *Parhyale* developmental genes, and found startling

203 levels of variation. (Supplemental Table:devGeneVariant). For example, we found that the cDNAs of the
204 germ line determinants, *nanos* (78 SNPS, 34 non-synonymous substitutions and one 6bp indel) and *vasa*
205 (37 SNPs, 7 non-synonymous substitutions and a one 6bp indel) can be more distant between *Parhyale*
206 populations than might be observed for orthologs between closely related species.

207 To further evaluate the extent of polymorphism across the genome, we mapped the genomic reads to a
208 set of previously Sanger-sequenced BAC clones of the *Parhyale* HOX cluster from the same isofemale
209 line used for sequencing [16]. We detected SNPs at a rate of 1.3 to 2.5% among the BACs (Table 3)
210 and also additional sequence differences between the BACs and genomic reads, again confirming that
211 additional haplotypes exist in the isofemale line beyond the one or two recovered from the sequenced
212 individual.

213 Overlapping regions of the contiguous BACs gave us the opportunity to directly compare Chicago-F
214 haplotypes and accurately observe polynucleotide polymorphisms (difficult to assess with short reads).
215 (Figure 7A). Since the BAC clones were generated from a pool of Chicago-F animals, we expected
216 each sequenced BAC to be representative of one haplotype. Overlapping regions between clones could
217 potentially represent one or two haplotypes. We found that the genomic reads supported the SNPs
218 observed between the overlapping BAC regions and in many cases display differences supporting the
219 existence of a third allele. This analysis revealed many insertion/deletion (indels) with some cases of
220 indels larger than 100 bases (Figure 7B). The finding that polynucleotide polymorphisms are prevalent
221 between the haplotypes of the sequenced Chicago-F strain explains the extensive independent assembly of
222 haplotypes, and means that special attention will have to be given to those functional genomic approaches
223 that are dependent on homology, such as CRISPR/Cas9 based knock in strategies.

224 **A comparative genomic analysis of the *Parhyale* genome**

225 Assessment of conservation of the proteome using BLAST against a selection of metazoan proteomes was
226 congruent with broad phylogenetic expectations. These analyses included crustacean proteomes likely
227 to be incomplete as they come from limited transcriptome datasets, but nonetheless highlighted genes
228 likely to be specific to the malacostraca (Figure 5C). To better understand global gene content evolution
229 we generated clusters of orthologous and paralogous gene families comparing the *Parhyale* proteome
230 with other complete proteomes across the metazoa using Orthofinder [57] (Figure 5D; Supplemental
231 HTML:orthology). We identified orthologous and paralogous protein groups across 16 species with
232 2,900 and 2,532 orthologous groups containing proteins found only in panarthropoda and arthropoda
233 respectively. We identified 855 orthologous groups that were shared exclusively by mandibulata, 772
234 shared by pancrustacea and 135 shared by crustacea. There were 9,877 *Parhyale* proteins that could not
235 be assigned to an orthologous group, potentially representing rapidly evolving or lineage specific proteins.

236 Our analysis of shared orthologous groups was equivocal with regard to alternative hypotheses on
237 the relationships among pancrustacean subgroups: 44 shared groups of orthologous proteins supported
238 a multicrustacea clade (uniting the malacostraca, copepoda and thecostraca), 37 groups supported the

239 allocarida (branchiopoda and hexapoda) and 49 groups supported the vericrustacea (branchiopoda and
240 multicrustacea)(Supplemental Zip:cladeOrthoGroups).

241 To further analyse the evolution of the *Parhyale* proteome we examined protein families that appeared
242 to be expanded (z-score >2), compared to other taxa (Supplemental Figure:expansion, Supplemental
243 HTML:orthology, Supplemental Txt:orthoGroups). We conservatively identified 29 gene families that
244 are expanded in *Parhyale*. Gene family expansions include the Sidestep (55 genes) and Lachesin (42)
245 immunoglobulin superfamily proteins as well as nephrins (33 genes) and neurotrimins (44 genes), which
246 are thought to be involved in immunity, neural cell adhesion, permeability barriers and axon guidance
247 [58–60]. Other *Parhyale* gene expansions include APN (aminopeptidase N) (38 genes) and cathepsin-like
248 genes (30 genes), involved in proteolytic digestion [61].

249 **Major signaling pathways and transcription factors in *Parhyale***

250 Components of all common metazoan cell-signalling pathways are largely conserved in *Parhyale*. At least
251 13 *Wnt* subfamilies were present in the cnidarian-bilaterian ancestor. *Wnt3* has been lost in protostomes,
252 while lophotrochozoans retain 12 *Wnt* genes [62, 63]. Some sampled ecdysozoans have undergone
253 significant *Wnt* gene loss, for example *C. elegans* has only 5 *Wnt* genes [64]. At most 9 *Wnt* genes
254 are present in any individual hexapod species [65], with *wnt2* and *wnt4* potentially lost before hexapod
255 radiation. The *Parhyale* genome encodes 6 of the 13 *Wnt* subfamily genes; *wnt1*, *wnt4*, *wnt5*, *wnt10*,
256 *wnt11* and *wnt16* (Figure 8). *Wnt* genes are known to have been ancestrally clustered [66]. We observed
257 that *wnt1* and *wnt10* are linked in a single scaffold (phaw_30.0003199), which given *Wnt6* and *Wnt9* loss,
258 this may be the remnant of the ancient *wnt9-1-6-10* cluster conserved in some protostomes.

259 We could identify 2 Fibroblast Growth Factor (*FGF*) genes and only a single FGF receptor (*FGFR*) in
260 the *Parhyale* genome, suggesting one *FGFR* has been lost in the malacostracan lineage (Supplemental
261 Figure:fgf). Within the Transforming Growth Factor beta (*TGF-β*) signaling pathway we found 2 genes
262 from the activin subfamily (an activin receptor and a myostatin), 7 genes from the Bone Morphogen
263 Protein (*BMP*) subfamily and 2 genes from the inhibin subfamily. Of the *BMP* genes, *Parhyale* has a
264 single decapentaplegic homologue (Supplemental Table:geneClassification). Other components of the
265 *TGF-β* pathway were identified such as the neuroblastoma suppressor of tumorigenicity (present in *Aedes*
266 *aegypti* and *Tribolium castaneum* but absent in *D. melanogaster* and *D. pulex*) and TGFβ-induced factor
267 homeobox 1 (*TGFI1*) which is a Smad2-binding protein within the pathway present in arthropods but
268 absent in nematodes (*C. elegans* and *Brugia malayi*; Supplemental Table:geneClassification). We identified
269 homologues of *PITX2*, a downstream target of the *TGF-β* pathway involved in endoderm and mesoderm
270 formation present [67] in vertebrates and crustaceans (*Parhyale* and *D. pulex*) but not in insects and
271 nematodes. With the exception of *SMAD7* and *SMAD8/9*, all other *SMADs* (*SMAD1*, *SMAD2/3*, *SMAD4*,
272 *SMAD6*) are found in arthropods sampled, including *Parhyale*. Components of other pathways interacting
273 with *TGF-β* signaling like the *JNK*, *Par6*, *ROCK1/RhoA*, *p38* and *Akt* pathways were also recovered and
274 annotated in the *Parhyale* genome (Supplemental Table:geneClassification). We identified major Notch

275 signaling components including Notch, Delta, Deltex, Fringe and modulators of the Notch pathway such
276 as *Dvl* and *Numb*. Members of the gamma-secretase complex (Nicastrin, Presenillin, and *APH1*) were
277 also present (Supplemental Table:keggSignallingPathways) as well as to other co-repressors of the Notch
278 pathway such as Groucho and *CtBP* [68].

279 A genome wide survey to annotate all potential transcription factor (TF) discovered a total of 1,143
280 proteins with DNA binding domains that belonged to all the major families previously identified. Importantly,
281 we observed a large expansion of TFs containing the zinc-finger (ZF)-C2H2 domain. *Parhyale* has
282 699 ZF-C2H2-containing genes [69], which is comparable to the number found in *H. sapiens* [70], but
283 significantly expanded compared to other arthropod species like *D. melanogaster* encoding 326 members
284 (Supplemental Table:tfDomain).

285 The *Parhyale* genome contains 126 homeobox-containing genes (Figure 9; Supplemental Table
286 :geneClassification), which is higher than the numbers reported for other arthropods (104 genes in *D.*
287 *melanogaster*, 93 genes in the honey bee *Apis mellifera*, and 113 in the centipede *Strigamia maritima*)
288 [71]. We identified a *Parhyale* specific expansion in the Ceramide Synthase (*CERS*) homeobox proteins,
289 which include members with divergent homeodomains [72]. *H. sapiens* have six *CERS* genes, but only
290 five with homeodomains [73]. We observed an expansion to 12 *CERS* genes in *Parhyale*, compared to
291 1-4 genes found in other arthropods [74] (Supplemental Figure:CERS). In phylogenetic analyses all 12
292 *CERS* genes in *Parhyale* clustered together with a *CERS* from another amphipod *E. veneris* (Supplemental
293 Figure:CERS), suggesting that this is recent expansion in the amphipod lineage.

294 *Parhyale* contains a complement of 9 canonical Hox genes that exhibit both spatial and temporal
295 colinearity in their expression along the anterior-posterior body axis [16]. Chromosome walking experi-
296 ments had shown that the Hox genes labial (*lab*) and proboscipedia (*pb*) are linked and that Deformed
297 (*Dfd*), Sex combs reduced (*Scr*), Antennapedia (*Antp*) and Ultrabithorax (*Ubx*) are also contiguous in
298 a cluster [16]. Previous experiments in *D. melanogaster* had shown that the proximity of nascent tran-
299 scripts in RNA fluorescent *in situ* hybridizations (FISH) coincide with the position of the corresponding
300 genes in the genomic DNA [75, 76]. Thus, we obtained additional information on Hox gene linkage by
301 examining nascent Hox transcripts in cells where Hox genes are co-expressed. We first validated this
302 methodology in *Parhyale* embryos by confirming with FISH, the known linkage of *Dfd* with *Scr* in the
303 first maxillary segment where they are co-expressed (Figure 10A-A“). As a negative control, we detected
304 no linkage between engrailed1 (*en1*) and *Ubx* or *abd-A* transcripts (Figure 10B - B“ and C - C“). We
305 then demonstrated the tightly coupled transcripts of *lab* with *Dfd* (co-expressed in the second antennal
306 segment, Figure (Figure 10D - D“), *Ubx* and *abd-A* (co-expressed in the posterior thoracic segments,
307 (Figure 10E - E“), and *abd-A* with *Abd-B* (co-expressed in the anterior abdominal segments, (Figure 10F
308 - F“). Collectively, all evidence supports the linkage of all analysed Hox genes into a single cluster as
309 shown in (Figure 10G - G“). The relative orientation and distance between certain Hox genes still needs
310 to be worked out. So far, we have not been able to confirm that *Hox3* is also part of the cluster due to
311 the difficulty in visualizing nascent transcripts for *Hox3* together with *pb* or *Dfd*. Despite these caveats,

312 *Parhyale* provides an excellent arthropod model system to understand these still enigmatic phenomena of
313 Hox gene clustering and spatio-temporal colinearity, and compare the underlying mechanisms to other
314 well-studied vertebrate and invertebrate models [77].

315 The Para Hox and *NK* gene clusters encode other *ANTP* class homeobox genes closely related to Hox
316 genes [78]. In *Parhyale*, we found 2 caudal (*Cdx*) and 1 *Gsx* ParaHox genes. Compared to hexapods, we
317 identified expansions in some NK-like genes, including 5 Bar homeobox genes (*BarH1/2*), 2 developing
318 brain homeobox genes (*DBX*) and 6 muscle segment homeobox genes (*MSX/Drop*). Evidence from several
319 bilaterian genomes suggests that *NK* genes are clustered together [79–82]. In the current assembly of the
320 *Parhyale* genome, we identified an *NK2-3* gene and an *NK3* gene on the same scaffold (phaw_30.0004720)
321 and the tandem duplication of an *NK2* gene on another scaffold (phaw_30.0004663). Within the *ANTP*
322 class, we also observed 1 mesenchyme homeobox (*Meox*), 1 motor neuron homeobox (*MNX/Exex*) and 3
323 even-skipped homeobox (*Evx*) genes.

324 **The *Parhyale* genome encodes glycosyl hydrolase enzymes consistent with lignocellu-** 325 **lose digestion (“wood eating”)**

326 Lignocellulosic (plant) biomass is the most abundant raw material on our planet and holds great promise
327 as a source for the production of bio-fuels [83]. Understanding how some animals and their
328 symbionts achieve lignocellulose digestion is a promising research avenue for exploiting lignocellulose-
329 rich material [84, 85]. Amongst metazoans, research into the ability to depolymerize plant biomass
330 into useful catabolites is largely restricted to terrestrial species such as ruminants, termites and beetles.
331 These animals rely on mutualistic associations with microbial endosymbionts that provide cellulolytic
332 enzymes known as glycosyl hydrolases (GHs) [86, 87] (Figure 11). Much less studied is lignocellulose
333 digestion in aquatic animals despite the fact that lignocellulose represents a major energy source in
334 aquatic environments, particularly for benthic invertebrates [88]. Recently, it has been suggested that the
335 marine wood-boring isopod *Limnoria quadripunctata* and the amphipod *Chelura terebrans* may have
336 sterile microbe-free digestive systems and they produce all required enzymes for lignocellulose digestion
337 [30, 31, 89]. Significantly these species have been shown to have endogenous GH7 family enzymes with
338 cellobiohydrolase (beta-1,4-exoglucanase) activity, previously thought to be absent from animal genomes.
339 From an evolutionary perspective, it is likely that GH7 coding genes were acquired by these species via
340 horizontal gene transfer from a protist symbiont.

341 *Parhyale* is a detritivore that can be sustained on a diet of carrots (Figure 11C), suggesting that they
342 too may be able to depolymerize lignocellulose for energy (Figure 11A and B). We searched for GH
343 family genes in *Parhyale* using the classification system of the CAZy (Carbohydrate-Active enZYmes)
344 database [90] and the annotation of protein domains in predicted genes with PFAM [91]. We identified
345 73 GH genes with complete GH catalytic domains that were classified into 17 families (Supplemental
346 Table:geneClassification) including 3 members of the GH7 family. Phylogenetic analysis of *Parhyale*
347 GH7s show high sequence similarity to the known GH7 genes in *L. quadripunctata* and the amphipod

348 *C. terebrans* [31] (Figure 12A; Supplemental Figure:ghAlignment). GH7 family genes were also iden-
349 tified in the transcriptomes of three more species spanning the multicrustacea clade: *Echinogammarus*
350 *veneris* (amphipod), *Eucyclops serrulatus* (copepod) and *Calanus finmarchicus* (copepod) (Supplemental
351 Table:geneClassification). As previously reported [92], we also discovered a closely related GH7 gene
352 in the branchiopod *Daphnia* (Figure 12A). This finding supports the grouping of branchiopoda with
353 multicrustacea (rather than with hexapoda) and the acquisition of a GH7 gene by a vericrustacean ancestor.
354 Alternatively, this suggests an even earlier acquisition of a GH7 gene by a crustacean ancestor with
355 subsequent loss of the GH7 family gene in the lineage leading to insects.

356 GH families 5,9,10, and 45 encode beta-1,4-endoglucanases which are also required for lignocellulose
357 digestion and are commonly found across metazoa. We found 3 GH9 family genes with complete catalytic
358 domains in the *Parhyale* genome as well as in the other three multicrustacean species (Figure 12B).
359 These GH9 enzymes exhibited a high sequence similarity to their homologues in the isopod *Limnoria*
360 and in a number of termites. Beta-glucosidases are the third class of enzyme required for digestion of
361 lignocellulose. They have been classified into a number of GH families: 1, 3, 5, 9 and 30, with GH1
362 representing the largest group [90]. In *Parhyale*, we found 7 beta-glucosidases from the GH30 family and
363 3 from the GH9 family, but none from the GH1 family.

364 Understanding lignocellulose digestion in animals using complex mutualistic interactions with cel-
365 loulolytic microbes has proven a difficult task. The study of “wood-eating” *Parhyale* can offer new
366 insights into lignocellulose digestion in the absence of gut microbes, and the unique opportunity to apply
367 molecular genetic approaches to understand the activity of glycosyl hydrolases in the digestive system.
368 Lignocellulose digestion may also have implications for gut immunity in some crustaceans, since these
369 reactions have been reported to take place in a sterile gut [32, 33].

370 **Characterisation of the innate immune system in a Malacostracan**

371 Immunity research in malacostracans has attracted interest due to the rapid rise in aquaculture related
372 problems [26, 27, 93]. malacostracan food crops represent a huge global industry (>\$40 Billion at point
373 of first sale), and reliance on this crop as a source of animal protein is likely to increase in line with human
374 population growth [93]. Here we provide an overview of immune-related genes in *Parhyale* that were
375 identified by mapping proteins to the ImmunoDB database [94] (Supplemental Table:geneClassification).
376 The ability of the innate immune system to identify pathogen-derived molecules is mediated by pattern
377 recognition receptors (PRRs) [95]. Several groups of invertebrate PRRs have been characterized, i.e.
378 thioester-containing proteins (*TEP*), Toll-like receptors (*TLR*), peptidoglycan recognition proteins (*PGRP*),
379 C-type lectins, galectins, fibrinogen-related proteins (*FREP*), gram-negative binding proteins (*GNBP*),
380 Down Syndrome Cell Adhesion Molecules (*Dscam*) and lipopolysaccharides and beta-1, 3-glucan binding
381 proteins (*LGBP*).

382 The functions of *PGRPs* have been described in detail in insects like *D. melanogaster* [96] and the
383 *PGRP* family has also been reported in vertebrates, molluscs and echinoderms [97, 98]. Surprisingly,

384 we found no PGRP genes in the *Parhyale* genome. *PGRPs* were also not found in other sequence
385 datasets from branchiopoda, copepoda and malacostraca (Figure 13A), further supporting their close
386 phylogenetic relationship (like the GH7 genes). In the absence of *PGRPs*, the freshwater crayfish
387 *Pacifastacus leniusculus* relies on a Lysine-type peptidoglycan and serine proteinases, *SPH1* and *SPH2*
388 that forms a complex with *LGBP* during immune response [99]. In *Parhyale*, we found one *LGBP* gene
389 and two serine proteinases with high sequence identity to *SPH1/2* in *Pacifastacus*. The *D. pulex* genome
390 has also an expanded set of Gram-negative binding proteins (proteins similar to *LGBP*) suggesting a
391 compensatory mechanism for the lost *PGRPs* [100]. Interestingly, we found a putative *PGRP* in the
392 Remipede *Speleonectes tulumensis* (Figure 13A) providing further support for sister group relationship of
393 remipedia and Hexapoda [14].

394 Innate immunity in insects is transduced by three major signaling pathways: the Immune Deficiency
395 (*Imd*), Toll and Janus kinase/signal transducer and activator of transcription (*JAK/STAT*) pathways
396 [101, 102]. We found 16 members of the Toll pathway in *Parhyale* including 10 Toll-like receptors
397 proteins (TLRs) (Figure 13B). Some TLRs have been also implicated in embryonic tissue morphogenesis
398 in *Parhyale* and other arthropods [103]. Additionally, we identified 7 *Imd* and 25 *JAK/STAT* pathway
399 members including two negative regulators: suppressor of cytokine signaling (*SOCS*), and protein inhibitor
400 of activated *STAT* (*PIAS*) [104].

401 The blood of arthropods (hemolymph) contains hemocyanin which is a copper-binding protein involved
402 in the transport of oxygen and circulating blood cells called hemocytes for the phagocytosis of pathogens.
403 Phagocytosis by hemocytes is facilitated by the evolutionarily conserved gene family, the thioester-
404 containing proteins (*TEPs*) [105]. Previously sequenced Pancrustacean species contained between 2 to
405 52 *TEPs*. We find 5 *TEPs* in the *Parhyale* genome. Arthropod hemocyanins themselves are structurally
406 related to phenoloxidases (PO; [106]) and can be converted into POs by conformational changes under
407 specific conditions [107]. POs are involved in several biological processes (like melanization immune
408 response, wound healing, cuticle sclerotization) and we identified 7 PO genes in *Parhyale*. Interestingly,
409 hemocyanins and PO activity have been shown to be highly abundant together with glycosyl hydrolases in
410 the digestive system of isopods and amphipods, raising a potential mechanistic link between gut sterility
411 and degradation of lignocellulose [30, 33].

412 Another well-studied transmembrane protein essential for neuronal wiring and adaptive immune
413 responses in insects is the immunoglobulin (*Ig*)-superfamily receptor Down syndrome cell adhesion
414 molecule (*Dscam*) [108, 109]. Alternative splicing of *Dscam* transcripts can result in thousands of
415 different isoforms that have a common architecture but have sequence variations encoded by blocks
416 of alternative spliced exons. The *D. melanogaster* *Dscam* locus encodes 12 alternative forms of exon
417 4 (encoding the N-terminal half of Ig2), 48 alternative forms of exon 6 (encoding the N-terminal half
418 of Ig3), 33 alternative forms of exon 9 (encoding Ig7), and 2 alternative forms of exon 17 (encoding
419 transmembrane domains) resulting in a total of 38,016 possible combinations. The *Dscam* locus in
420 *Parhyale* (and in other crustaceans analysed) have a similar organization to insects; tandem arrays of

421 multiple exons encode the N-terminal halves of Ig2 (exon 4 array with at least 13 variants) and Ig3 (exon
422 6 array with at least 20 variants) and the entire Ig7 domain (exon 14 array with at least 13 variants)
423 resulting in at least 3,380 possible combinations (Figure 13C-E). The alternative splicing of hypervariable
424 exons in *Parhyale* was confirmed by sequencing of cDNA clones amplified with Dscam-specific primers.
425 Almost the entire *Dscam* gene is represented in a single genomic scaffold and exhibits high amino-acid
426 sequence conservation with other crustacean *Dscams* (Supplemental Figure:dscamVariant). The number
427 of *Dscam* isoforms predicted in *Parhyale* is similar to that predicted for *Daphnia* species [110]. It remains
428 an open question whether the higher number of isoforms observed in insects coincides with the evolution
429 of additional Dscam functions compared to crustaceans.

430 From a functional genomics perspective, the *Parhyale* immune system appears to be a good represen-
431 tative of the malacostracan or even multicrustacean clade that can be studied in detail with existing tools
432 and resources. Interestingly, the loss of *PGRPs* and presence of *GH7* genes in branchiopoda, similar to the
433 presence of *GH7* genes, supports their close relationship with the multicrustacea rather than the hexapoda.

434 **Non-coding RNAs and associated proteins in the *Parhyale* genome**

435 Non-coding RNAs are a central, but still a relatively poorly understood part of eukaryotic genomes. In
436 animal genomes, different classes of small RNAs are key for genome surveillance, host defense against
437 viruses and parasitic elements in the genome, and regulation of gene expression through transcriptional,
438 post-transcriptional and epigenetic control mechanisms [111–119]. The nature of these non-coding
439 RNAs, as well as the proteins involved in their biogenesis and function, can vary between animals. For
440 example, some nematodes have Piwi-interacting short RNAs (piRNAs), while others have replaced these
441 by alternate small RNA based mechanisms to compensate for their loss [120].

442 As first step, we surveyed the *Parhyale* genome for known conserved protein components of the small
443 interfering RNA (siRNA/RNAi) and the piRNA pathways (Table 4). We found key components of all major
444 small RNA pathways, including 4 argonaute family members, 2 PIWI family members, and orthologs
445 of *D. melanogaster* *Dicer-1* and *Dicer-2*, *drosha* and *loquacious*, (Supplemental Figure:dicerPiwiTree).
446 Among Argonaute genes, *Parhyale* has 1 *AGO-1* ortholog and 3 *AGO-2* orthologs, which is presumably
447 a malacostraca-specific expansion. While *Parhyale* only has 2 PIWI family members, other crustacean
448 lineages have clearly undergone independent expansions of this protein family (Supplemental Figure:).
449 Unlike in *C. elegans*, many mammals, fish and insects (but not *D. melanogaster*), we did not find any
450 evidence in the *Parhyale* genome for the *SID-1* (systemic inter-ference defective) transmembrane protein
451 that is essential for systemic RNAi [121–123]. Species without a *SID-1* ortholog can silence genes only
452 in a cell-autonomous manner [124]. This feature has important implications for future design of RNAi
453 experiments in *Parhyale*.

454 We also assessed the miRNA and putative long non-coding RNAs (lncRNA) content of *Parhyale*
455 using both MiRPara and Rfam [125, 126]. We annotated 1405 homologues of known non-coding RNAs
456 using Rfam. This includes 980 predicted tRNAs, 45 rRNA of the large ribosomal subunit, 10 rRNA of

457 the small ribosomal subunit, 175 snRNA components of the major spliceosome (U1, U2, U4, U5 and
458 U6), 5 snRNA components of the minor spliceosome (U11, U12, U4atac and U6atac), 43 ribozymes, 38
459 snoRNAs, 71 conserved cis-regulatory element derived RNAs and 42 highly conserved miRNA genes
460 (Supplemental Table:RFAM; Supplemental HTML:rna). *Parhyale* long non-coding RNAs (lncRNAs)
461 were identified from the transcriptome using a series of filters to remove coding transcripts producing a
462 list of 220,284 putative lncRNAs (32,223 of which are multi-exonic). Only one *Parhyale* lncRNA has
463 clear homology to another annotated lncRNA, the sphinx lncRNA from *D. melanogaster* [127].

464 We then performed a more exhaustive search for miRNAs using MiRPara (Supplemental HTML:rna)
465 and a previously published *Parhyale* small RNA read dataset [128]. We identified 1,403 potential miRNA
466 precursors represented by 100 or more reads. Combining MiRPara and Rfam results, we annotated 31 out
467 of the 34 miRNA families found in all bilateria, 12 miRNAs specific to protostomia, 4 miRNAs specific
468 to arthropoda and 5 miRNAs previously found to be specific to mandibulata (Figure 14). We did not
469 identify *mir-125*, *mir-283* and *mir-1993* in the *Parhyale* genome. The absence of *mir-1993* is consistent
470 with reports that this miRNA was lost during Arthropod evolution [129]. While we did not identify
471 *mir-125*, we observed that *mir-100* and *let-7* occurred in a cluster on the same scaffold (Supplemental
472 Figure:mirnaCluster), where *mir-125* is also present in other animals. The absence of *mir-125* has been
473 also reported for the centipede genome [71]. *mir-100* is one of the most primitive miRNAs shared
474 by bilateria and cnidaria [129, 130]. The distance between *mir-100* and *let-7* genes within the cluster
475 can vary substantially between different species. Both genes in *Parhyale* are localized within a 9.3kb
476 region (Supplemental Figure:mirnaClusterA) as compared to 3.8kb in the mosquito *Anopheles gambiae*
477 and 100bp in the beetle *Tribolium* [131]. Similar to *D. melanogaster* and the polychaete *Platynereis*
478 *dumerilii*, we found that *Parhyale mir-100* and *let-7* are co-transcribed as a single, polycistronic lncRNA.
479 We also found another cluster with *miR-71* and *mir-2* family members which is conserved across many
480 invertebrates [132] (Supplemental Figure:mirnaClusterB).

481 Conserved linkages have also been observed between miRNAs and Hox genes in bilateria [133–137].
482 For example, the phylogenetically conserved *mir-10* is present within both vertebrate and invertebrate
483 Hox clusters between *Hoxb4/Dfd* and *Hoxb5/Scr* [138]. In the *Parhyale* genome and Hox BAC sequences,
484 we found that *mir-10* is also located between *Dfd* and *Src* on BAC clone PA179-K23 and scaffold
485 phaw_30.0001203 (Supplemental Figure:mirnaClusterC,D). However, we could not detect *mir-iab-4* near
486 the *Ubx* and *AbdA* genes in *Parhyale*, the location where it is found in other arthropods/insects [139].

487 Preliminary evidence regarding the presence of PIWI proteins and other piRNA pathway proteins also
488 suggests that the piRNA pathway is likely active in *Parhyale*, although piRNAs themselves await to be
489 surveyed. The opportunity to study these piRNA, miRNA and siRNA pathways in a genetically tractable
490 crustacean system will shed further light into the regulation and evolution of these pathways and their
491 contribution to morphological diversity.

492 **Methylome analysis of the *Parhyale* genome**

493 Methylation of cytosine residues (m5C) in CpG dinucleotides in animal genomes is regulated by a
494 conserved multi-family group of DNA methyltransferases (DNMTs) with diverse roles in the epigenetic
495 control of gene expression, genome stability and chromosome dynamics [140–142]. The phylogenetic
496 distribution of DNMTs in Metazoa suggests that the bilaterian ancestor had at least one member of the
497 Dnmt1 and Dnmt3 families (involved in *de novo* methylation and maintenance of DNA methylation)
498 and the Dnmt2 family (involved in tRNA methylation), as well as additional RNA methyltransferases
499 [143, 144]. Many animal groups have lost some of these DNA methyltransferases, for example *DNMT1*
500 and 3 are absent from *D. melanogaster* and flatworms [145, 146], while *DNMT2* is absent from nematodes
501 *C. elegans* and *C. briggsae*. The *Parhyale* genome encodes members of all 3 families *DNMT1*, *DNMT3*
502 and *DNMT2*, as well as 2 orthologs of conserved methyl-CpG-binding proteins and a single orthologue of
503 *Tet2*, an enzyme involved in DNA demethylation [147] (Figure 15A).

504 We used genome wide bisulfite sequencing to confirm the presence and also assess the distribution of
505 CpG dinucleotide methylation. Our results indicated that 20-30% of *Parhyale* DNA is methylated at CpG
506 dinucleotides (Figure 15B). The *Parhyale* methylation pattern is similar to that observed in vertebrates,
507 with high levels of methylation detected in transposable elements and other repetitive elements, in
508 promoters and gene bodies (Figure 15C). A particular class of rolling-circle transposons are very highly
509 methylated in the genome, potentially implicating methylation in silencing these elements. For comparison,
510 about 1% or less of CpG-associated cytosines are methylated in insects like *Drosophila*, *Apis*, *Bombyx*
511 and *Tribolium*. [140, 148, 149]. These data represent the first documentation of a crustacean methylome.
512 Considering the utility of *Parhyale* for genetic and genomic research, we anticipate future investigations to
513 shed light on the functional importance and spatiotemporal dynamics of epigenetic modifications during
514 normal development and regeneration, as well as their relevance to equivalent processes in vertebrate
515 systems.

516 ***Parhyale* genome editing using homology-independent approaches**

517 *Parhyale* has already emerged as a powerful model for developmental genetic research where the ex-
518 pression and function of genes can be studied in the context of stereotyped cellular processes and with a
519 single-cell resolution. Several experimental approaches and standardized resources have been established
520 to study coding and non-coding sequences (Table 1). These functional studies will be enhanced by
521 the availability of the assembled and annotated genome presented here. As a first application of these
522 resources, we tested the efficiency of CRISPR/Cas system for targeted genome editing in *Parhyale*
523 [15–20]. In these studies, we targeted the *Distal-less* patterning gene (called *PhDII-e*) [22] that has a
524 widely-conserved and highly-specific role in animal limb development [150].

525 We first genotyped our wild-type laboratory culture and found two *PhDII-e* alleles with 23 SNPs
526 and 1 indel in their coding sequences and untranslated regions. For *PhDII-e* knock-out, two sgRNAs
527 targeting both alleles in their coding sequences downstream of the start codon and upstream of the DNA-

528 binding homeodomain were injected individually into 1-cell-stage embryos (F0 generation) together with
529 a transient source of Cas9 (Supplemental Figure:funcConstruct A-B). Both sgRNAs gave rise to animals
530 with truncated limbs (Figure 16A and B); the first sgRNA at a relatively low percentage around 9% and the
531 second one at very high frequencies ranging between 53% and 76% (Supplemental Figure:funcConstruct).
532 Genotyping experiments revealed that injected embryos carried *PhDII-e* alleles modified at the site
533 targeted by each sgRNA (Supplemental Figure:funcConstruct B-D). The number of modified *PhDII-e*
534 alleles recovered from F0s varied from two, in cases of early bi-allelic editing at the 1-cell-stage, to
535 three or more, in cases of later-stage modifications by Cas9 (Supplemental Figure:funcConstruct C). We
536 isolated indels of varying length that were either disrupting the open reading frame, likely producing
537 loss-of-function alleles or were introducing in-frame mutations potentially representing functional alleles
538 (Supplemental Figure:funcConstruct C-D). In one experiment with the most efficient sgRNA, we raised
539 the injected animals to adulthood and set pairwise crosses between 17 fertile F0s (10 male and 7 female):
540 88% (15/17) of these founders gave rise to F1 offspring with truncated limbs, presumably by transmitting
541 *PhDII-e* alleles modified by Cas9 in their germlines. We tested this by genotyping individual F1s from two
542 of these crosses and found that embryos bearing truncated limbs were homozygous for loss-of-function
543 alleles with out-of-frame deletions, while their wild-type siblings carried one loss-of-function allele and
544 one functional allele with an in-frame deletion (Supplemental Figure:funcConstruct D).

545 The non-homologous end joining (NHEJ) repair mechanism operating in the injected cells can be
546 exploited not only for gene knock-out experiments described above, but also for CRISPR knock-in
547 approaches where an exogenous DNA molecule is inserted into the targeted locus in a homology-
548 independent manner. This homology-independent approach could be particularly useful for *Parhyale*
549 that exhibits high levels of heterozygosity and polymorphisms in the targeted laboratory populations,
550 especially in introns and intergenic regions. To this end, we co-injected into 1-cell-stage embryos the
551 Cas9 protein together with the strongest sgRNA and a tagging plasmid. The plasmid was designed in
552 such a way that upon its linearization by the same sgRNA and Cas9 and its integration into the *PhDII-e*
553 locus in the appropriate orientation and open reading frame, it would restore the endogenous *PhDII-e*
554 coding sequence in a bicistronic mRNA also expressing a nuclear fluorescent reporter. Among injected
555 F0s, about 7% exhibited a nuclear fluorescence signal in the distal (telopodite and exopodite) parts of
556 developing appendages (Figure 16C and Supplemental Figure:funcConstruct E), which are the limb
557 segments that were missing in the knock-out experiments (Figure 16B). Genotyping of one of these
558 embryos demonstrated that the tagged *PhDII-e* locus was indeed encoding a functional *PhDII-e* protein
559 with a small in-frame deletion around the targeted region (Supplemental Figure:funcConstruct F).

560 These results, together with the other recent applications of the CRISPR/Cas system to study Hox
561 genes in *Parhyale* [16, 17], demonstrate that the ability to manipulate the fertilized eggs together with the
562 slow tempo of early cleavages can result in very high targeting frequencies and low levels of mosaicism
563 for both knock-out and knock-in approaches. Considering the usefulness of the genome-wide resources
564 described in this report, we anticipate that the *Parhyale* embryo will prove an extremely powerful system

565 for fast and reliable F0 screens of gene expression and function.

566 CONCLUSION

567 In this article we described the first complete genome of a Malacostracan crustacean species, the genome
568 of the marine amphipod *Parhyale hawaiiensis*. With the same chromosome count ($2n=46$) as the human
569 genome and an estimated size of 3.6 Gb, it is among the largest genomes submitted to NCBI. The *Parhyale*
570 genome exhibits high levels of polymorphism, heterozygosity and abundance of repetitive sequence. Our
571 comparative bioinformatics analyses suggest that the expansion of repetitive sequences and the increases
572 in gene size due to an expansion of intron size have contributed to the large size of the genome. Despite
573 these challenges, the *Parhyale* genome and associated transcriptomic resources reported here provide a
574 useful assembly of most genic regions in the genome and a comprehensive description of the *Parhyale*
575 transcriptome and proteome.

576 *Parhyale* has emerged since the early 2000's as an attractive animal model for developmental genetic
577 and molecular cell biology research. It fulfills several desirable biological and technical requirements
578 satisfied by major animal models, including a relatively short life-cycle, year-round breeding under
579 standardized laboratory conditions, availability of thousands of eggs for experimentation on a daily
580 basis, and amenability to various embryological, cellular, molecular genetic and genomic approaches.
581 In addition, it combines some unique features and strengths, like stereotyped cell lineages and cell
582 behaviors, a direct mode of development, a remarkable appendage (limb) diversity and the capacity to
583 regenerate limbs post-embryonically. These qualities can be utilized to address fundamental long-standing
584 questions in developmental biology, like cell fate specification, nervous system development, organ
585 morphogenesis and regeneration [151]. All these *Parhyale* research fields will benefit enormously from
586 the standardized genome-wide resources reported here. Forward and reverse genetic analyses using both
587 unbiased screens and candidate gene approaches have already been devised successfully in *Parhyale*
588 (Table 1). The availability of coding and non-coding sequences for all identified signaling pathway
589 components, transcription factors and various classes of non-coding RNAs will dramatically accelerate
590 the study of the expression and function of genes implicated in the aforementioned processes.

591 Equally importantly, our analyses highlighted additional areas where *Parhyale* could serve as a new
592 experimental model to address other questions of broad biomedical interest. From a functional genomics
593 perspective, the *Parhyale* immune system appears to be a good representative of the malacostracan or
594 even the multicrustacean clade that can be studied in detail with existing tools and resources. Besides
595 the evolutionary implications and the characterization of alternative strategies used by arthropods to
596 defend against pathogens, a deeper mechanistic understanding of the *Parhyale* immune system will be
597 relevant to aquaculture. Some of the greatest setbacks in the crustacean farming industry were caused by
598 severe disease outbreaks. *Parhyale* is closely related to farmed crustaceans (primarily shrimps, prawns
599 and crayfish) and the knowledge acquired from studying its innate immunity could help enhance the
600 sustainability of this industry by preventing or controlling infectious diseases [93, 152–155].

601 An immune-related problem that will be also interesting to explore in *Parhyale* concerns the possibility
602 of a sterile digestive tract similar to that proposed for limnoriid isopods [30]. *Parhyale*, like limnoriid
603 isopods, encodes and expresses all enzymes required for lignocellulose digestion (King et al., 2010),
604 suggesting that it is able to “digest wood” by itself without symbiotic microbial partners. Of course, a lot
605 of work will required to be invested in the characterization of the cellulolytic system in *Parhyale* before
606 any comparisons can be made with other well-established symbiotic digestion systems of lignocellulose.
607 Nevertheless, the possibility of an experimentally tractable animal model that serves as a living bioreactor
608 to convert lignocellulose into simpler metabolites, suggests that future research in *Parhyale* may also have
609 a strong biotechnological potential, especially for the production of biofuels from the most abundant and
610 cheapest raw material, plant biomass.

611 Several of our observations from analysing the *Parhyale* genome and other crustacean data sets also
612 throw light on the relationships among crustacean groups. We and others have observed the absence
613 of *PGRPs* in representatives of branchipoda, copepoda, and malacostraca[100, 156] (Supplementary
614 table 10). Either *PGRPs* were lost independently in multicrustacea and branchiopoda during arthropod
615 evolution or branchiopoda are not a sister taxa of insects but are more closely related to the multicrustacea
616 taxa. We and others also identified a glycosyl hydrolase (GH) family 7 gene in multicrustacean and
617 branchipod genomes, further supporting the close relationship between these groups [30]. Parsimonius
618 interpretation of these data suggest that branchiopoda is a sister group to multicrustacea rather than the
619 hexpoda.

620 Finally, *Parhyale* was introduced recently as a new model for limb regeneration [24]. In many
621 respects, including the segmented body plan, the presence of a blood system and the contribution of
622 lineage-committed adult stem cells to newly formed tissues, the *Parhyale* regenerative process resembles
623 the processes in vertebrates more than other established invertebrate models (e.g. planarians, hydra).
624 Regenerative research in *Parhyale* has been founded on transgenic approaches to label specific populations
625 of cells and will be further assisted by the resources presented here. Likewise, we expect that the new
626 genomic information and CRISPR-based genome editing methodologies together with all other facets of
627 *Parhyale* biology will open other new research avenues not yet imagined.

628 **ACKNOWLEDGMENTS**

629 We are grateful to Serge Picard for sequencing the genome libraries, and Frantisek Marec and Peer Martin
630 for useful advice on *Parhyale* karyotyping.

631 **MATERIALS AND METHODS**

632 A list of software and external datasets used are provided in Supplemental Table:externalDataSoftware.
633 Detailed methodology and codes for each section are provided as supplementary IPython notebooks in
634 HTML format viewable with a web browser.
635 All supplemental data including IPython notebook can be downloaded from this figshare link:

636 [https://figshare.com/articles/supplemental_data_for_Parhyale_hawaniensis_](https://figshare.com/articles/supplemental_data_for_Parhyale_hawaniensis_genome/3498104)
637 [genome/3498104](https://figshare.com/articles/supplemental_data_for_Parhyale_hawaniensis_genome/3498104)

638 Alternatively, the IPython notebooks can also be viewed at the following github repository:

639 https://github.com/damiankao/phaw_genome

640 **Genome library preparation and sequencing**

641 About 10 µg of genomic DNA were isolated from a single adult male from the Chicago-F isofemale line
642 established in 2001 (a.k.a. Iso2) [51]. The animal was starved for one week and treated for 3 days with
643 penicillin-streptomycin (100x, Gibco/Thermo Fisher Scientific), tetracycline hydrochloride (20 µg/ml,
644 Sigma-Aldrich) and amphotericin B (200x, Gibco/Thermo Fisher Scientific). It was then flash frozen in
645 liquid nitrogen, homogenized manually with a pestle in a 1.5 ml microtube (Kimble Kontes) in 600 µl of
646 Lysis buffer (100 mM Tris-HCl pH 8, 100 mM NaCl, 50 mM EDTA, 0.5% SDS, 200 µg/ml Proteinase
647 K, 20 µg/ml RNase A). The lysate was incubated for 3 hours at 37°C, followed by phenol/chloroform
648 extractions and ethanol precipitation. The condensed genomic DNA was fished out with a Pasteur pipette,
649 washed in 70% ethanol, air-dried, resuspended in nuclease-free water and analysed on a Qubit fluorometer
650 (Thermo Fisher Scientific) and on a Bioanalyzer (Agilent Technologies). All genome libraries were
651 prepared from this sample: 1 µg of genomic DNA was used to generate the shotgun libraries using the
652 TruSeq DNA Sample Prep kit (Illumina) combined with size-selection on a LabChip XT fractionation
653 system (Caliper Life Sciences Inc) to yield 2 shotgun libraries with average fragment sizes 421 bp and
654 800 bp, respectively; 4 µg of genomic DNA were used to generate 4 mate-pair libraries with average
655 fragment sizes 5.5 kb, 7.3 kb, 9.3 kb and 13.8 kb using the Nextera Mate Pair Sample Preparation kit
656 (Illumina) combined with agarose size selection. All libraries were sequenced on a HiSeq 2500 instrument
657 (Illumina) using paired-end 150 nt reads.

658 **Karyotyping**

659 For chromosome spreads, tissue was obtained from embryos at stages 14-18 [35]. Eggs were taken from
660 the mother and incubated for 1–2 h in isotonic colchicine solution (0.05% colchicine, artificial sea water).
661 After colchicine incubation, the embryonic tissue was dissected from the egg and placed in hypotonic
662 solution (0.075 M KCl) for 25 min. For tissue fixation, we replaced the hypotonic solution with freshly
663 prepared ice-chilled Carnoy's fixative (six parts ethanol, three parts methanol and one part anhydrous
664 acetic acid) for 25 min. The fixed tissue was minced with a pair of fine tungsten needles in Carnoy's
665 solution and the resulting cell suspension was dropped with a siliconized Pasteur pipette from a height
666 of about 5 cm onto a carefully cleaned ice-chilled microscopic slide. After partial evaporation of the
667 Carnoy's fixative the slides were exposed few times briefly to hot water vapors to rehydrate the tissue.
668 The slides were then dried on a 75°C metal block in a water bath. Finally, the slides with prepared
669 chromosomes were aged overnight at 60°C. After DNA staining either with Hoechst (H33342, Molecular
670 Probes) or with DAPI (Invitrogen), chromosomes were counted on a Zeiss Axioplan II Imaging equipped
671 with C-Apochromat 63x/1.2 NA objective and a PCO pixelfly camera. FIJI was used to improve image

672 quality (contrast and brightness) and FIJI plugin ‘Cell Counter’ was used to determine the number of
673 chromosomes.

674 **Genome assembly and k-mer analyses of polymorphisms repetiveness**

675 The *Parhyale* raw data and assembled data are available on the NCBI website (project accession
676 SRP066767). Genome assembly was done with Abyss [157] at two different k-mer settings (70, 120) and
677 merged with GAM-NGS. Scaffolding was performed with SSPACE [158]. We chose a cut offs of >95%
678 overlap length and >95% identity when removing shorter allelic contigs before scaffolding as these
679 gave better scaffolding results as assessed by assembly metrics. Transcriptome assembly was performed
680 with Trinity [55]. The completeness of the genome and transcriptome was assessed by blasting against
681 CEGMA genes [56] and visualized by plotting the orthologue hit ratio versus e-value. K-mer analysis
682 of variant and repetitive branching was performed with String Graph Asssembler’s preqc module [53].
683 K-mer intersection analysis was performed using jellyfish2 [159]. An in-depth description of the assembly
684 process is detailed in Supplemental HTML:assembly.

685 **Transcriptome library preparation, sequencing and assembly**

686 *Parhyale* transcriptome assembly was generated from Illumina reads collected from diverse embryonic
687 stages (Stages 19, 20, 22, 23, 25, and 28), and adult thoracic limbs and regenerating thoracic limbs (3 and
688 6 days post amputation). For the embryonic samples, RNA was extracted using Trizol; PolyA+ libraries
689 were prepared with the Truseq V1 kit (Illumina), starting with 0.6 - 3.5ug of total mRNA, and sequenced
690 on the Illumina Hiseq 2000 as paired-end 100 base reads, at the QB3 Vincent J. Coates Genomics Sequenc-
691 ing Laboratory. For the limb samples, RNA was extracted using Trizol; PolyA+ libraries were prepared
692 with the Truseq V2 kit (Illumina), starting with 1ug of total mRNA, and sequenced on the Illumina Hiseq
693 2500 as paired-end 100 base reads, at the IGBMC Microarray and Sequencing platform. 260 million
694 reads from embryos and 180 million reads from limbs were used for the transcriptome assembly. Prior to
695 the assembly we trimmed adapter and index sequences using cutadapt [160]. We also removed spliced
696 leader sequences: GAATTTTCACTGTTCCCTTTACCACGTTTTACTG, TTACCAATCACCCCTTTAC-
697 CAAGCGTTTACTG, CCCTTTACCAACTCTTAACTG, CCCTTTACCAACTTTACTG using cutadapt
698 with 0.2 error allowance to remove all potential variants. To assemble the transcriptome we used Trinity
699 (version trinityrnaseq_r20140413) [55] with settings: -min_kmer_cov 2, -path_reinforcement_distance 50.

700 **Gene model prediction and canonical proteome dataset generation**

701 Gene prediction was done with a combination of Evidence Modeler [161] and Augustus [162]. The
702 transcriptome was first mapped to the genome using GMAP [163]. A secondary transcriptome reference
703 assembly was performed with STAR/Cufflinks [164, 165]. The transcriptome mapping and Cufflinks
704 assembly was processed through the PASA pipeline [161] to consolidate the annotations. The PASA
705 dataset, a set of Exonerate [166] mapped Uniprot proteins, and Ab initio GeneMark [167] predictions
706 were consolidated with Evidence Modeler to produce a set of gene annotations. A high confidence set

707 of gene models from Evidence Modeler containing evidence from all three sources was used to train
708 Augustus. Evidence from RepeatMasker [168], PASA and Exonerate was then used to generate Augustus
709 gene predictions. A final list of genes for down-stream analysis was generated using both transcriptome
710 and gene predictions (canonical proteome dataset). Detailed methods are described in Supplemental
711 HTML:annotations.

712 **Polymorphism analysis on genic regions and BAC clones**

713 For variant analysis on the BAC clones, the short shot-gun library genomic reads were mapped to the
714 BAC clones individually. GATK was then used to call variants. For variant analysis on the genic regions,
715 transcript sequences from the canonical proteome dataset were first aligned to the genome assembly.
716 Genome alignments less than 30 bases were discarded. The possible genome alignments were sorted based
717 on number of mismatches with the top alignment having the least amount of mismatches. For each base
718 of the transcript, the top two genome aligned bases were recorded as the potential variants. Bases where
719 there were more than five genomic mapping loci were discarded as potentially highly conserved domains
720 or repetitive region. Detailed methods of this process are described in Supplemental HTML:variant.

721 **Polymorphisms in *Parhyale* developmental genes**

722 *Parhyale* genes (nucleotide sequences) were downloaded from GenBank. Each gene was used as a query
723 for blastn against the *Parhyale* genome using the Geneious software [169]. In each case two reference con-
724 tig hits were observed where both had E values of close to zero. A new sequence called geneX_snp was cre-
725 ated and this sequence was annotated with the snps and/or indels present in the alternative genomic contigs.
726 To determine the occurrence of synonymous and non-synonymous substitution, the original query and the
727 newly created sequence (with polymorphisms annotated) were in silico translated into protein sequences
728 followed by pairwise alignment. Regions showing amino acid changes were annotated as non-synonymous
729 substitutions. Five random genes from the catalogue were selected for PCR, cloning and Sanger sequenc-
730 ing to confirm genomic polymorphisms and assess further polymorphism in the lab population. Primers
731 for genomic PCR designed to capture exon regions are listed as the following: dachshund (PH1F = 5'-
732 GGTGCGCTAAATTGAAGAAATTACG-3' and PH1R = 5'- ACTCAGAGGGTAATAGTAACAGAA-3'),
733 distalless exon 2 (PH2F = 5'-CACGGCCCCGGCACTA ACTATCTC-3' and PH2R = 5'-GTAATATATCTTACAACAACGA
734 3'), distalless exon 3 (PH3F = 5'-GGTGAACGGGCCGGAGTCTC-3' and PH3R = 5'-GCTGTGGGTGCTGTGGGT-
735 3'), homothorax (PH4F = 5'-TCGGGGTGTA AAAAGGACTCTG-3' and PH4R = 5'-AACATAGGAACTCACCTGGTG
736 3'), orthodenticle (PH5F = 5'-TTTGCCACTAACACATATTTGAAA-3' and PH5R = 5'-TCCCAAGTAGATGATCCCT
737 3') and prospero (PH6F = 5'-TACACTGCAACATCCGATGACTTA-3' and PH6R = 5'-CGTGTTATGTTCTCTCGTGGC
738 3').

739 **Evolutionary analyses of orthologous groups**

740 Evolutionary analyses and comparative genomics were performed with 16 species (*D. melanogaster*, *A.*
741 *gambiae*, *D. pulex*, *L. salmonis*, *S. maritima*, *S. mimosarum*, *M. martensii*, *I. scapularis*, *H. dujardini*, *C.*

742 *elegans*, *B. malayi*, *T. spiralis*, *M. musculus*, *H. sapiens*, and *B. floridae*. For orthologous group analyses,
743 gene families were identified using OrthoFinder [57]. The canonical proteome was used as a query in
744 BlastP against proteomes from 16 species to generate a distance matrix for OrthoFinder to normalize
745 and then cluster with MCL. Detailed methods are described in Supplemental HTML:orthology. For
746 the comparative BLAST analysis, five additional transcriptome datasets were used from the following
747 crustacean species: *Litopenaeus vannamei*, *Echinogammarus veneris*, *Eucyclops serrulatus*, *Calanus*
748 *finmarchicus*, *Speleonectes tulumensis*

749 **Fluorescence in situ hybridization detection of Hox genes**

750 Embryo fixation and in-situ hybridization was performed according to [170]. To enhance the nascent nu-
751 clear signal over mature cytoplasmic transcript, we used either early germband embryos (Stages 11 – 15) in
752 which expression of *lab*, *Dfd*, and *Scr* are just starting [16], or probes that contain almost exclusively intron
753 sequence (*Ubx*, *abd-A*, *Abd-B*, and *en1*). *Lab*, *Dfd*, and *Scr* probes are described in [16]. Template for the
754 intron-spanning probes were amplified using the following primers: *en1*-Intron1, AAGACACGACGAG-
755 CATCCTG and CTGTGTATGGCTACCCGTCC; *Ubx*-Intron1, GGTATGACAGCCGTCCAACA and
756 AGAGTGCCAAGGATACCCGA; *abd-A*, CGATATACCCAGTCCGGTGC and TCATCAGCGAGGGCA-
757 CAATT; *Abd-B*, GCTGCAGGATATCCACACGA and TGCAGTTGCCGCCATAGTAA. A T7-adapter
758 was appended to the 5' end of each reverse primer to enable direct transcription from PCR product. Probes
759 were labeled with either Digoxigenin (DIG) or Dinitrophenol (DNP) conjugated UTPs, and visualized
760 using sheep α -DIG (Roche) and donkey α -Sheep AlexaFluor 555 (Thermo Fischer Scientific), or Rabbit
761 α -DNP (Thermo Fischer Scientific) and Donkey α -Rabbit AlexaFluor 488 (Jackson ImmunoResearch),
762 respectively following the procedure of Ronshaugen and Levine (2004). Preparations were imaged on an
763 LSM 780 scanning laser confocal (Zeiss), and processed using Volocity software (Perkin-Elmer).

764 Cross species identification of GH family genes and immune-related genes. The identification of GH
765 family genes was done by obtaining Pfam annotations [91] for the *Parhyale* canonical proteome. Pfam
766 domains were classified into different GH families based on the CAZy database [90]. For immune-related
767 genes, best-reciprocal blast was performed with ImmunoDB genes [94].

768 **Phylogenetic tree construction**

769 Multiple sequence alignments of protein sequences for gene families of *FGF*, *FGFR*, *CERS*, *GH7*,
770 *GH9*, *PGRP*, Toll-like receptors, *DICER*, Piwi and Argonaute were performed using MUSCLE [171].
771 Phylogenetic tree construction was performed with RAxML [172] using the WAG+G model from
772 MUSCLE multiple alignments.

773 **Bisulfite sequencing**

774 Libraries for DNA methylation analysis by bisulfite sequencing were constructed from 100ng of genomic
775 DNA extracted from one *Parhyale* male individual, using the Illumina Truseq DNA methylation kit
776 according to manufacturers instructions. Alignments to the *Parhyale* genome were generated using the

777 core Bismark module from the program Bismark [173], having first artificially joined the *Parhyale* contigs
778 to generate 10 pseudo-contigs as the program is limited as to the number of separate contigs it can analyse.
779 We then generated genome-wide cytosine coverage maps using the `bismark_methylation_extraction`
780 module with the parameter `-CX` specified to generate annotations of CG, CHH and CHG sites. In order
781 to analyse genome-wide methylation patterns, cytosines with more than 10 read depth coverage were
782 selected. Overall methylation levels at CG, CHH and CHG sites were generated using a custom Perl
783 script. To analyse which regions were methylated we mapped back from the joined contigs to the original
784 contigs and assigned these to functional regions based on RepeatMasker [168] and transcript annotations
785 of repeats and genes respectively. To generate overall plots of methylation levels in different features we
786 averaged over all sites mapping to particular features, focusing on CG methylation and measuring the
787 %methylation at each site as the number of reads showing methylation divided by the total number of
788 reads covering the site. Meta gene plots over particular features were generated similarly except that sites
789 mapping within a series of 100bp wide bins from 1000bp upstream of the feature start site onwards were
790 collated.

791 **Identification and cloning of *Dscam* alternative spliced variants**

792 For the identification of *Dscam* in the *Parhyale*, we used the *Dscam* protein sequence from crustaceans *D.*
793 *pulex* [110] and *L. vannamei* [174] as queries to probe the assembled genome using tBlastN. A 300kb
794 region on scaffold phaw_30.0003392 was found corresponding to the *Parhyale Dscam* extending from
795 IG1 to FN6 exons. This sequence was annotated using transcriptome data together with manual searches
796 for open reading frames to identify IG, FN exons and exon-intron boundaries (Figure 10). Hypervariable
797 regions of IG2, IG3 and IG7 were also annotated accordingly on the scaffold (Figure 8). This region
798 represents a bona fide *Dscam* paralog as it matches the canonical extracellular *Dscam* domain structure
799 of nine IGs – four FNs – one IG and two FNs. *Parhyale* mRNA extractions were performed using
800 the Zymo Research Direct-zol RNA MiniPrep kit according to manufacturer's instructions. Total RNA
801 extract was used for cDNA synthesis using the Qiagen QuantiTect Reverse Transcription Kit according to
802 manufacturer's instructions. To identify and confirm potential hypervariable regions from the *Parhyale*
803 *Dscam* (PhDscam) transcript, three regions of PhDscam was corresponding to IG2, IG3 and IG7 exons
804 respectively were amplified using the following primer pairs. IG2 region:

805 DF1 = 5'-CCCTCGTGTTCCCGCCCTTCAAC-3'

806 DR1 = 5'-GCGATGTGCAGCTCTCCAGAGGG-3'

807 IG3 region:

808 DF2 = 5'-TCTGGAGAGCTGCACATCGCTAAT-3'

809 DR2 = 5'-GTGGTCATTGCGTACGAAGCACTG-3'

810 IG7 region:

811 DF3 = 5'-CGGATACCCCATCGACTCCATCG-3'

812 DR3 = 5'-GAAGCCGTCAGCCTTGCATTCAA-3'

813 PCR of each region was performed using Phusion High-fidelity polymerase from Thermo Fisher Scientific
814 and thermal cycling was done as the following: 98°C 30s, followed by 30 cycles of 98°C 10s, 67°C 30s,
815 72°C 1m30s, and then 72°C 5m. PCR products were cloned into pGEMT-Easy vector and a total of 81
816 clones were selected and Sanger sequenced and in silico translated in the correct reading frame using
817 Geneious (R7; [169]) for multiple sequence alignment.

818 **Identification of non-protein-coding RNAs**

819 *Parhyale* non-protein-coding RNAs were identified using two independent approaches. Infernal 1.1.1
820 [175] was used with the RFAM 12.0 database [126] to scan the genome to identify potential non-protein-
821 coding RNAs according. Additionally, MiRPara [125] was used to scan the genome for potential miRNA
822 precursors. These potential precursors were further filtered using small RNA read mapping and miRBase
823 mapping [176]. Putative lncRNAs were identified from the transcriptome by applying filtering criteria
824 including removal of known coding proteins and removal of predicted proteins. Detailed methods are
825 available in Supp_rna.

826 **CRISPR/Cas genome editing**

827 To genotype our wild-type population, extraction of total RNA and preparation of cDNA from embryos
828 were carried out as previously described [23]. The PhDII-e cDNA was amplified with primers PhDIIe_2For
829 (5'-TTTGTCAGGGATCTGCCATT-3') and PhDIIe_1852Rev (5'-TAGCGGCTGACGGTTGTTAC-3'),
830 purified with the DNA Clean and Concentrator kit (Zymo Research), cloned with the Zero Blunt
831 TOPO PCR Cloning Kit (Thermo Fisher Scientific) and sequenced with primers M13 forward (5'-
832 GTAAACGACGGCCAG-3') and M13 reverse (5'-CAGGAAACAGCTATGAC-3').

833 Each template for sgRNA synthesis was prepared by annealing and PCR amplification of the sgRNA-
834 specific forward primer DII1: (18 nt PhDII-e-targeted sequence underlined)

835 5'-GAAATTAATACGACTCACTATA

836 AGAGTTGTTACCAAAGAAGTTTTAGAGCTAGAAATAGC-3'

837 or DII2: (20 nt PhDII-e-targeted sequence underlined)

838 5'-GAAATTAATACGACTCACTAT

839 AGGCTTCCCCGCCGCCATGTAGTTTTAGAGCTAGAAATAGC-3'

840 together with the universal reverse primer:

841 5'-AAAAGCACCGACTCGGTGCCACTTTTTCAAGTTGATAA

842 CGGACTAGCCTTATTTAACTTGCTATTTCTAGCTCTAAAAC-3'

843 using the Phusion DNA polymerase (New England Biolabs).

844 Each PCR product was gel-purified with the Zymoclean DNA recovery kit (Zymo Research) and 150 ng of
845 DNA were used as template in an in vitro transcription reaction with the Megashortscript T7 kit (Thermo
846 Fisher Scientific). A 4-hour incubation at 37°C was followed by DNase digestion, phenol/chloroform
847 extraction, ethanol precipitation and storage in ethanol at -20°C according to the manufacturer's instructions.
848 Before microinjection, a small aliquot of the sgRNA was centrifuged, the pellet was washed with 70%

849 ethanol, resuspended in nuclease-free water and quantified on a Nanodrop spectrophotometer (Thermo
850 Scientific). The Cas9 was provided either as in vitro synthesized capped mRNA or as recombinant protein.
851 Cas9 mRNA synthesis was carried out as previously described [45] using plasmid T7-Cas9 (a gift from
852 David Stern and Justin Crocker) linearized with EcoRI digestion. The lyophilized Cas9 protein (PNA
853 Bio Inc) was resuspended in nuclease-free water at a concentration of 1.25 µg/µl and small aliquots were
854 stored at -80°C. For microinjections, we mixed 400 ng/µl of Cas9 protein with 40-200 ng/µl sgRNA,
855 incubated at 37°C for 5 min, transferred on ice, added the inert dye phenol red (5x from Sigma-Aldrich)
856 and, for knock-in experiments, the tagging plasmid at a concentration of 10 ng/µl. The injection mix was
857 centrifuged for 20 min at 4°C and the cleared solution was microinjected into 1-cell-stage embryos as
858 previously described [45].

859 In the knock-out experiments, embryos were scored for phenotypes under a bright-field stereomicro-
860 scope 7-8 days after injection (stage S25-S27) when organogenesis is almost complete and the limbs are
861 clearly visible through the transparent egg shell. To image the cuticle, anaesthetized hatchlings were fixed
862 in 2% paraformaldehyde in 1xPBS for 24 hours at room temperature. The samples were then washed in
863 PTx (1xPBS containing 1% TritonX-100) and stained with 1 mg/ml Congo Red (Sigma-Aldrich) in PTx
864 at room temperature with agitation for 24 hours. Stained samples were washed in PTx and mounted in
865 70% glycerol for imaging. Serial optical sections were obtained at 2 µm intervals with the 562 nm laser
866 line on a Zeiss 710 confocal microscope using the Plan-Apochromat 10x/0.45 NA objective. Images were
867 processed with Fiji (<http://fiji.sc>) and Photoshop (Adobe Systems Inc).

868 This methodology enabled us to also extract genomic DNA for genotyping from the same imaged
869 specimen. Each specimen was disrupted with a disposable pestle in a 1.5 ml microtube (Kimble Kontes)
870 in 50 µl of Squishing buffer (10 mM Tris-HCl pH 8, 1 mM EDTA, 25 mM NaCl, 200 µg/ml Proteinase
871 K). The lysate was incubated at 37°C for a minimum of 2 hours, followed by heat inactivation of the
872 Proteinase K for 5 min at 95°C, centrifugation at full speed for 5 min and transferring of the cleared
873 lysate to a new tube. To recover the sequences in the PhDII-e locus targeted by the DII1 and DII2 sgRNAs,
874 5 µl of the lysate were used as template in a 50 µl PCR reaction with the Phusion DNA polymerase
875 (New England Biolabs) and primers 313For (5'-TGGTTTTAGCAACAGTGAAGTGA-3') and 557Rev
876 (5'-GACTGGGAGCGTGAGGGTA-3'). The amplified products were purified with the DNA Clean and
877 Concentrator kit (Zymo Research), cloned with the Zero Blunt TOPO PCR Cloning Kit (Thermo Fisher
878 Scientific) and sequenced with the M13 forward primer.

879 For the knock-in experiments, we constructed the tagging plasmid pCRISPR-NHEJ-KI-DII-T2A-H2B-
880 Ruby2 that contained the PhDII-e coding sequence fused in-frame with the T2A self-cleaving peptide,
881 the *Parhyale histone* H2B and the Ruby 2 monomeric red fluorescent protein, followed by the PhDII-e
882 3'UTR and the pGEM-T Easy vector backbone (Promega). This tagging plasmid has a modular design
883 with unique restriction sites for easy exchange of any desired part. More details are available upon request.
884 Embryos co-injected with the Cas9 protein, the DII2 sgRNA and the pCRISPR-NHEJ-KI-DII-T2A-H2B-
885 Ruby2 tagging plasmid were screened for nuclear fluorescence in the developing appendages under an

886 Olympus MVX10 epi-fluorescence stereomicroscope. To image expression, live embryos at stage S22
887 were mounted in 0.5% SeaPlaque low-melting agarose (Lonza) in glass bottom microwell dishes (MatTek
888 Corporation) and scanned as described above acquiring both the fluorescence and transmitted light on an
889 inverted Zeiss 880 confocal microscope. To recover the chromosome-plasmid junctions, genomic DNA
890 was extracted from transgenic siblings with fluorescent limbs and used as template in PCR reaction as
891 described above with primer pair 313For and H2BRev (5'-TTACTTAGAAGAAGTGTACTTTG-3') for
892 the left junction and primer pair M13 forward and 557Rev for the right junction. Amplified products were
893 purified and cloned as described above and sequenced with the M13 forward and M13 reverse primers.

894 FIGURES AND TABLES

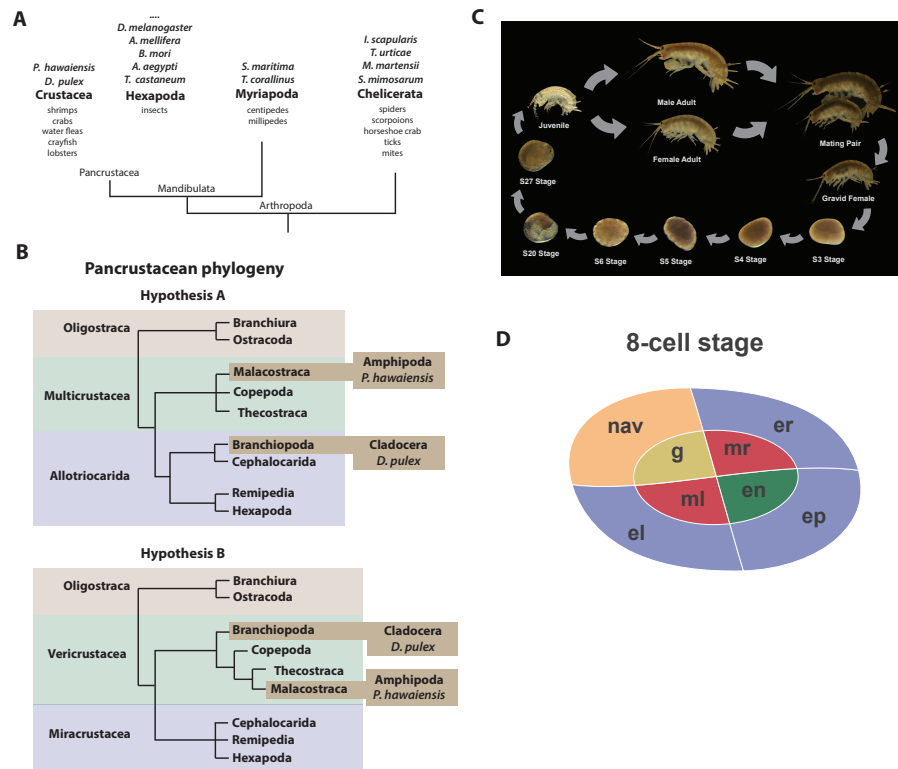


Figure 1. Introduction. (A) Phylogenetic relationship of Arthropods showing the Chelicerata as an outgroup to Mandibulata and the Pancrustacea clade which includes crustaceans and insects. Species listed for each clade have ongoing or complete genomes. Species for Crustacea include: *Parhyale hawaiiensis*, *D. pulex*; Hexapoda: *Drosophila melanogaster*, *Apis mellifera*, *Bombyx mori*, *Aedes aegypti*, *Tribolium castaneum*; Myriapoda: *Strigamia maritima*, *Trigoniulus corallinus*; Chelicerata: *Ixodes scapularis*, *Tetranychus urticae*, *Mesobuthus martensii*, *Stegodyphus mimosarum*. (B) One of the unresolved issues concerns the placement of the Branchiopoda either together with the Cephalocarida, Remipedia and Hexapoda (Allotriocarida hypothesis A) or with the Copepoda, Thecostraca and Malacostraca (Vericrustacea hypothesis B). (C) Life cycle of *Parhyale* that takes about two months at 26°C. *Parhyale* is a direct developer and a sexually dimorphic species. The fertilized egg undergoes stereotyped total cleavages and each blastomere becomes committed to a particular germ layer already at the 8-cell stage depicted in (D). The three macromeres Er, El, and Ep give rise to the anterior right, anterior left, and posterior ectoderm, respectively, while the fourth macromere Mav gives rise to the visceral mesoderm and anterior head somatic mesoderm. Among the 4 micromeres, the mr and ml micromeres give rise to the right and left somatic trunk mesoderm, en gives rise to the endoderm, and g gives rise to the germline.

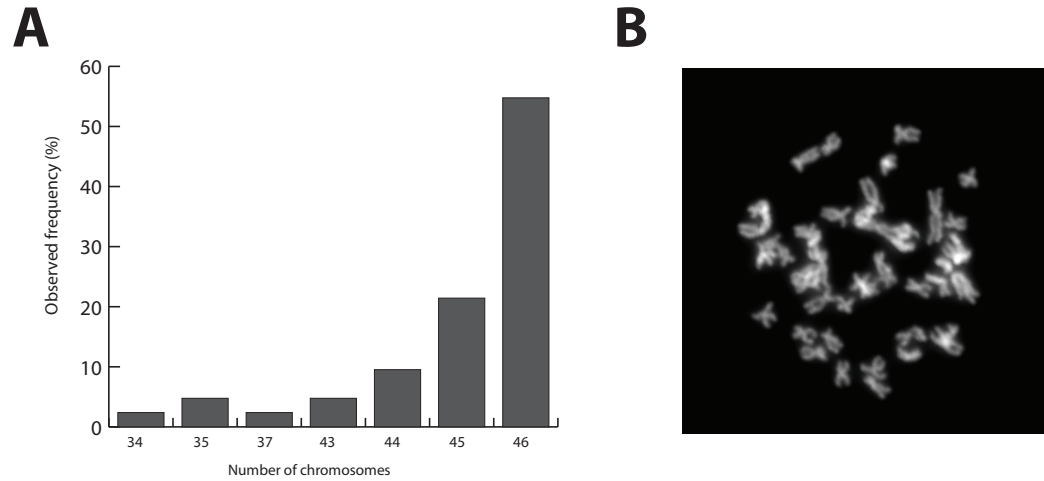


Figure 2. *Parhyale* karyotype. (A) Frequency of the number of chromosomes observed in 42 mitotic spreads. Forty-six chromosomes were observed in more than half preparations. (B) Representative image of Hoechst-stained chromosomes.

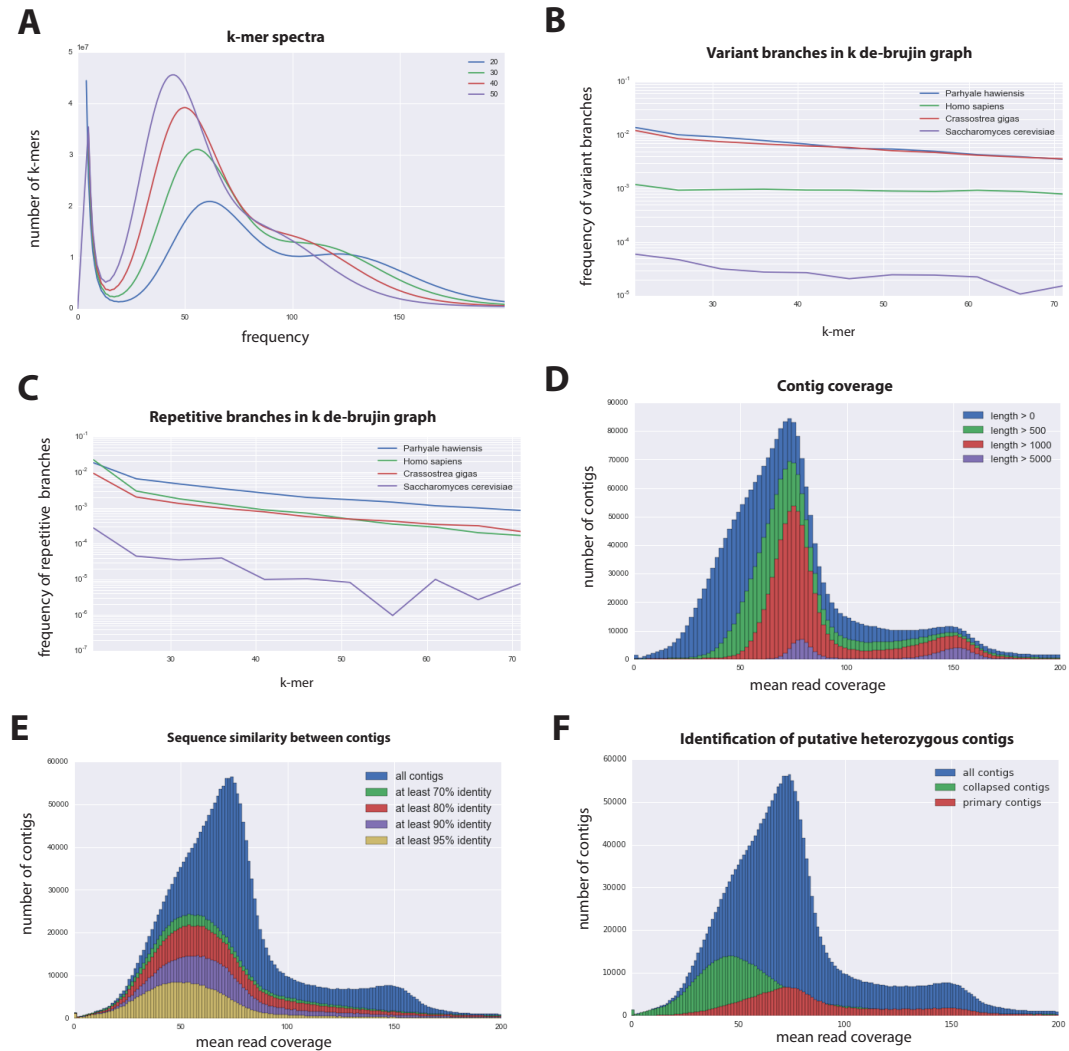


Figure 3. *Parhyale* genome assembly metrics. (A) K-mer frequency spectra of all reads for k-length from 20 to 50. (B) K-mer branching analysis showing the frequency of k-mer branches classified as variants compared to *Homo sapiens* (human), *Crassostrea gigas* (oyster), and *Saccharomyces cerevisiae* (yeast). (C) K-mer branching analysis showing the frequency of k-mer branches classified as repetitive compared to *H. sapiens*, *C. gigas* and *S. cerevisiae*. (D) Histogram of read coverage of assembled contigs. (E) The number of contigs with an identity ranging from 70-95% to another contig in the set of assembled contigs. (F) Collapsed contigs (green) are contigs with at least 95% identity with a longer primary contig (red). These contigs were removed prior to scaffolding and added back as potential heterozygous contigs after scaffolding.

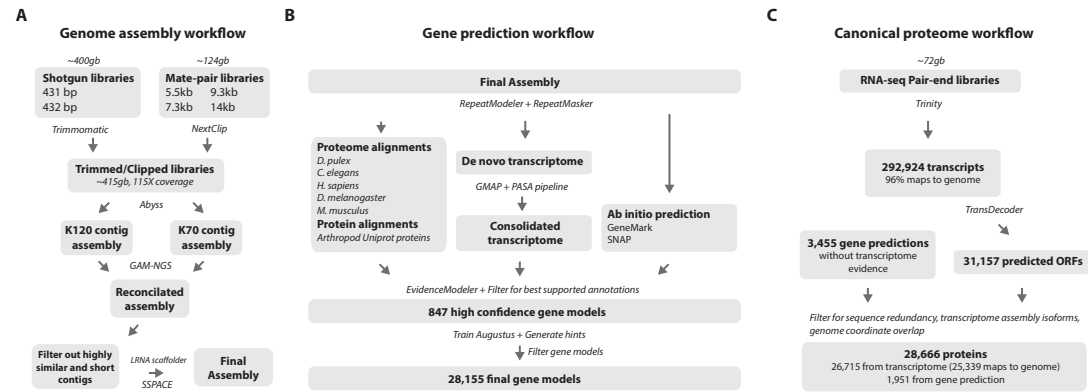


Figure 4. Workflows of assembly, annotation, and proteome generation. (A) Flowchart of the genome assembly. Two shotgun libraries and four mate-pair libraries with the indicated average sizes were prepared from a single male animal and sequenced at a 115x coverage after read filtering. Contigs were assembled at two different k-mers with Abyss and the two assemblies were merged with GAM-NGS. Filtered contigs were scaffolded with SSPACE. (B) The final scaffolded assembly was annotated with a combination of Evidence Modeler to generate 847 high quality gene models and Augustus for the final set of 28,155 predictions. These protein-coding gene models were generated based on a *Parhyale* transcriptome consolidated from multiple developmental stages and condition, their homology to the species indicated, and ab initio predictions with GeneMark and SNAP. (C) The *Parhyale* proteome contains 28,666 entries based on the consolidated transcriptome and gene predictions. The transcriptome contains 292,924 coding and non-coding RNAs, 96% of which could be mapped to the assembled genome.

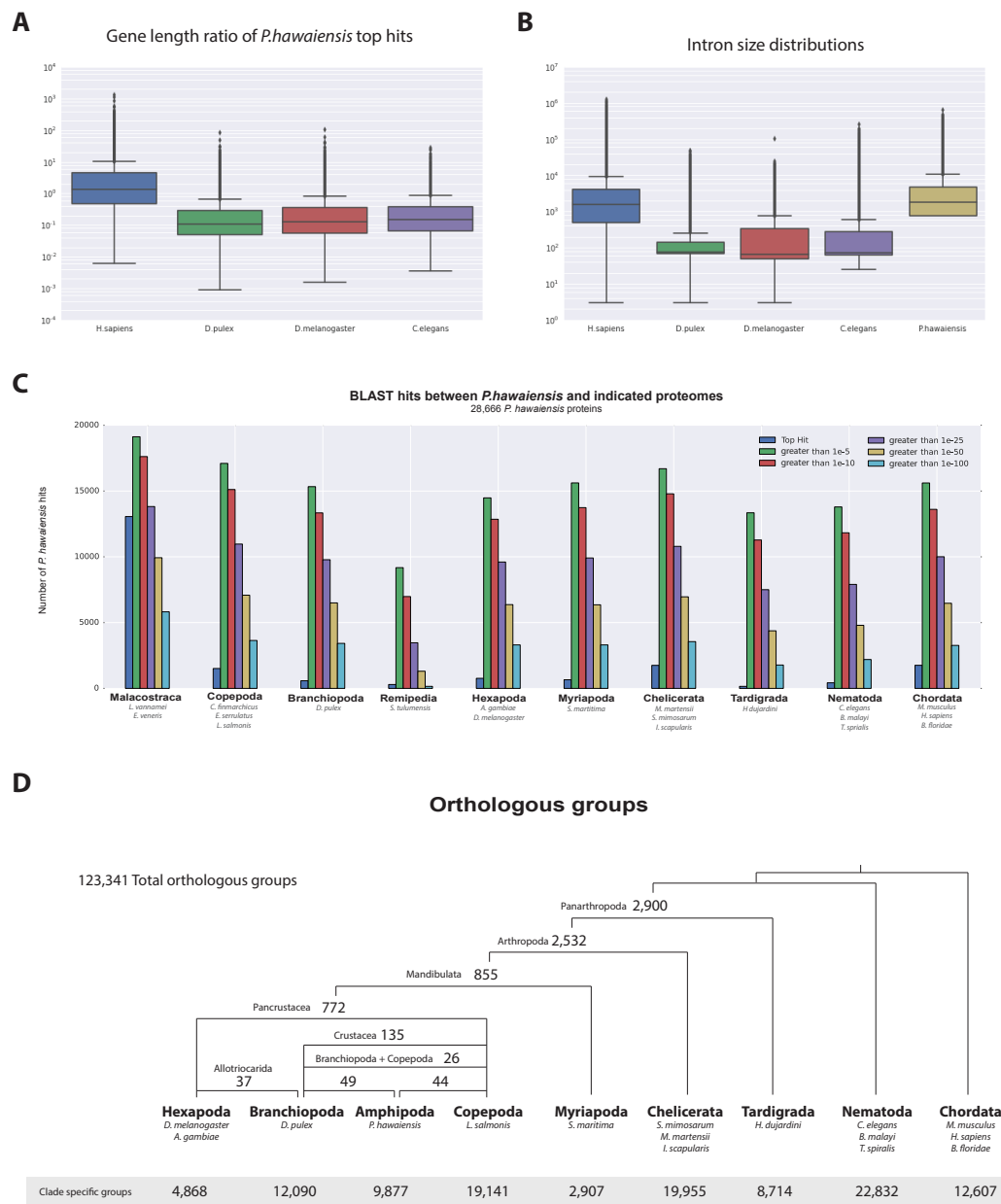


Figure 5. *Parhyale* genome comparisons. (A) Box plots comparing gene size between *Parhyale* and humans (*H. sapiens*), water fleas (*D. pulex*), flies (*D. melanogaster*) and nematodes (*C. elegans*). Ratios were calculated by dividing the size of the top blast hits in each species with the corresponding *Parhyale* gene size. (B) Box plots showing the distribution of intron size in the same species used in A. (C) Comparison between *Parhyale* and representative proteomes from the indicated animal taxa. Colored bars indicate the number of blast hits recovered across various thresholds of E-values. The top hit value represents the number of proteins with a top hit corresponding to the respective species. (D) Cladogram showing the number of shared orthologous protein groups at various taxonomic levels, as well as the number of clade-specific groups. A total of 123,341 orthogroups were identified with Orthofinder across the 16 genomes used in this analysis. Within Pancrustacea, 37 orthogroups were shared between Branchiopoda with Hexapoda (supporting the Allotriocarida hypothesis) and 49 orthogroups were shared between Branchiopoda and Amphipoda (supporting the Vericrustacea hypothesis).

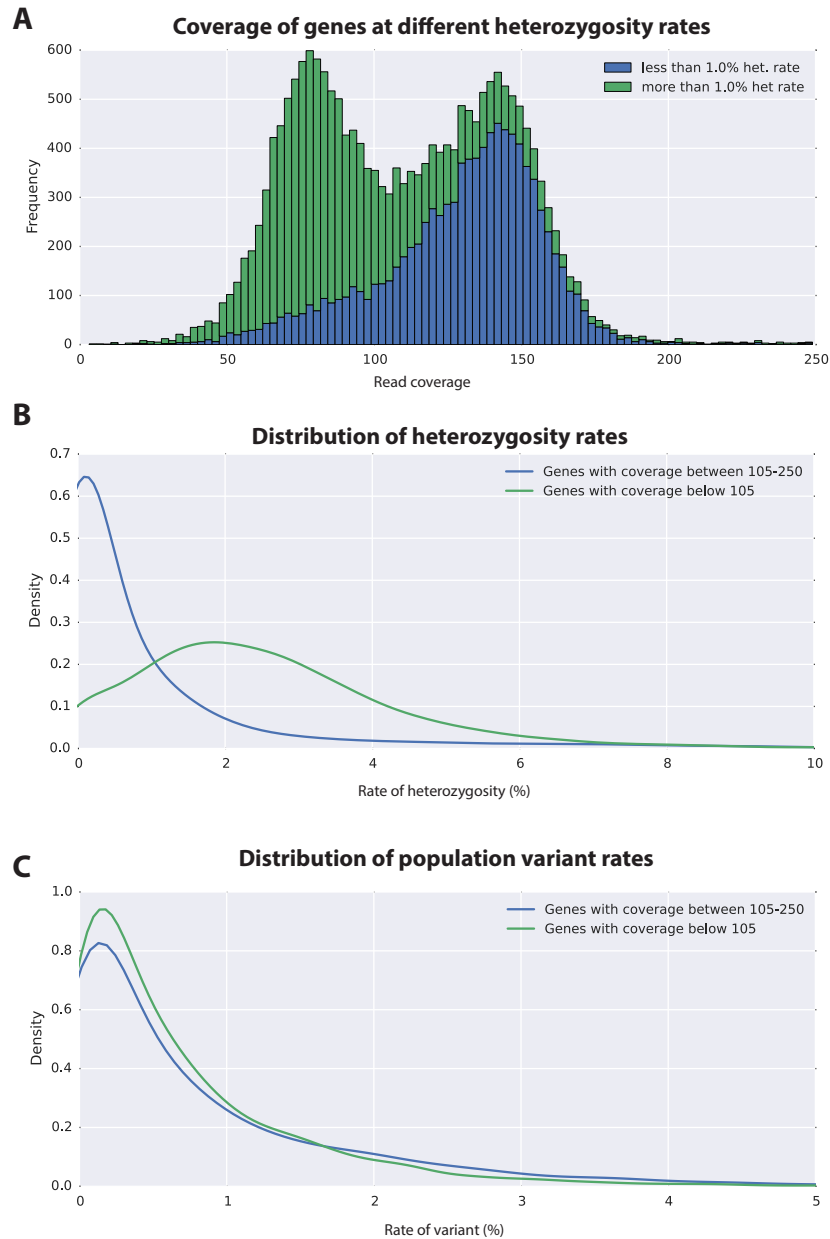


Figure 6. Variation analyses predicted genes. (A) A read coverage histogram of predicted genes. Reads were first mapped to the genome, then coverage were calculated for each defined locus. (B) Distribution plot shows that genes in the lower coverage region (<105 coverage) have a higher heterozygosity rate than genes in the higher coverage region (>105 coverage). (C) Distribution plot indicates that mean population variant rates are similar for genes in the higher and lower coverage regions.

A Variation in contiguous BAC sequences

	PA264-B19		PA40-O15		PA272-M04		PA284-I07		PA76-H18	
	% identity according to BAC % identity according to reads	100% ident. 98% ident.	99% ident. 96% ident.	97% ident. 94% ident.	96% ident. 96% ident.	100% ident. 96% ident.	100% ident. 93% ident.	99% ident. 97% ident.	98% ident. 98% ident.	
	PA179-K23		PA81-D11		PA92-D22		PA221-A05			
overlap length	19,846	3,135	16,536	20,707	32,587	3,155	24,345	24,892		
BAC supported SNPs	1	89	543	842	8	2	122	395		
Genomic reads supported SNPs	425	121	902	854	1,269	206	633	541		
BAC + Genomic reads supported SNPs	0	88	539	841	0	0	120	395		
Third allele	0	1	13	1	0	0	2	10		
Number of INDELS	64	17	106	115	127	24	88	85		
Number of INDELS >= 100	2	1	5	1	1	0	0	6		

B Position and length of indels > 1bp in overlapping BAC regions

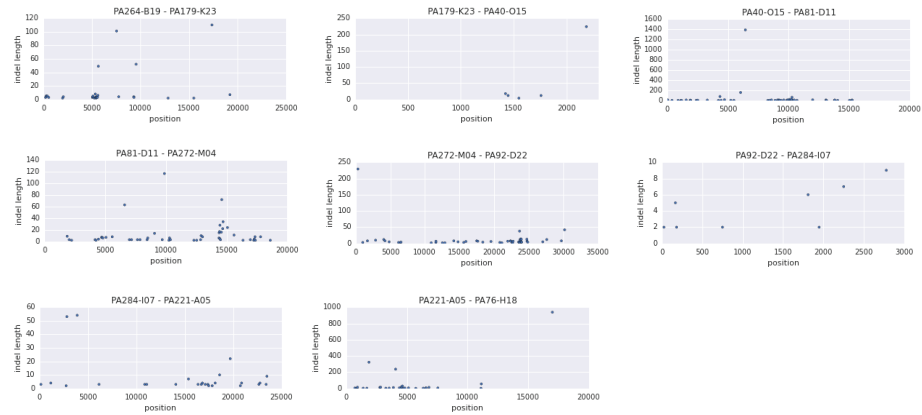


Figure 7. Variation observed in contiguous BAC sequences. (A) Schematic diagram of the contiguous BAC clones tiling across the HOX cluster and their % sequence identities. “Overlap length” refers to the lengths (bp) of the overlapping regions between two BAC clones. “BAC supported single nucleotide polymorphisms (SNPs)” refer to the number of SNPs found in the overlapping regions by pairwise alignment. “Genomic reads supported SNPs” refer to the number of SNPs identified in the overlapping regions by mapping all reads to the BAC clones and performing variant calling with GATK. “BAC + Genomic reads supported SNPs” refer to the number of SNPs identified from the overlapping regions by pairwise alignment that are supported by reads. “Third allele” refers to presence of an additional polymorphism not detected by genomic reads. “Number of INDELS” are the number of all insertion or deletions found in the contiguous region. “Number of INDELS >100” are insertion or deletions greater than or equal to 100. (B) Position versus indel lengths across each overlapping BAC region.

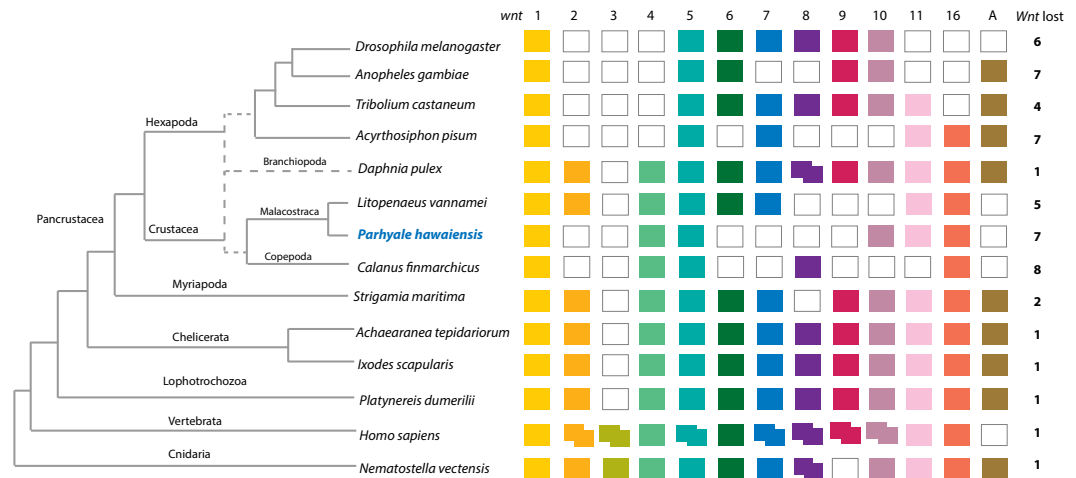


Figure 8. Comparison of Wnt family members across Metazoa. Comparison of Wnt genes across Metazoa. Tree on the left illustrates the phylogenetic relationships of species used. Dotted lines in the phylogenetic tree illustrate the alternative hypothesis of Branchiopoda + Hexapoda versus Branchiopoda + Multicrustacea. Colour boxes indicate the presence of certain Wnt subfamily members (wnt1 to wnt11, wnt16 and wntA) in each species. Empty boxes indicate the loss of particular Wnt genes. Two overlapping colour boxes represent duplicated Wnt genes.

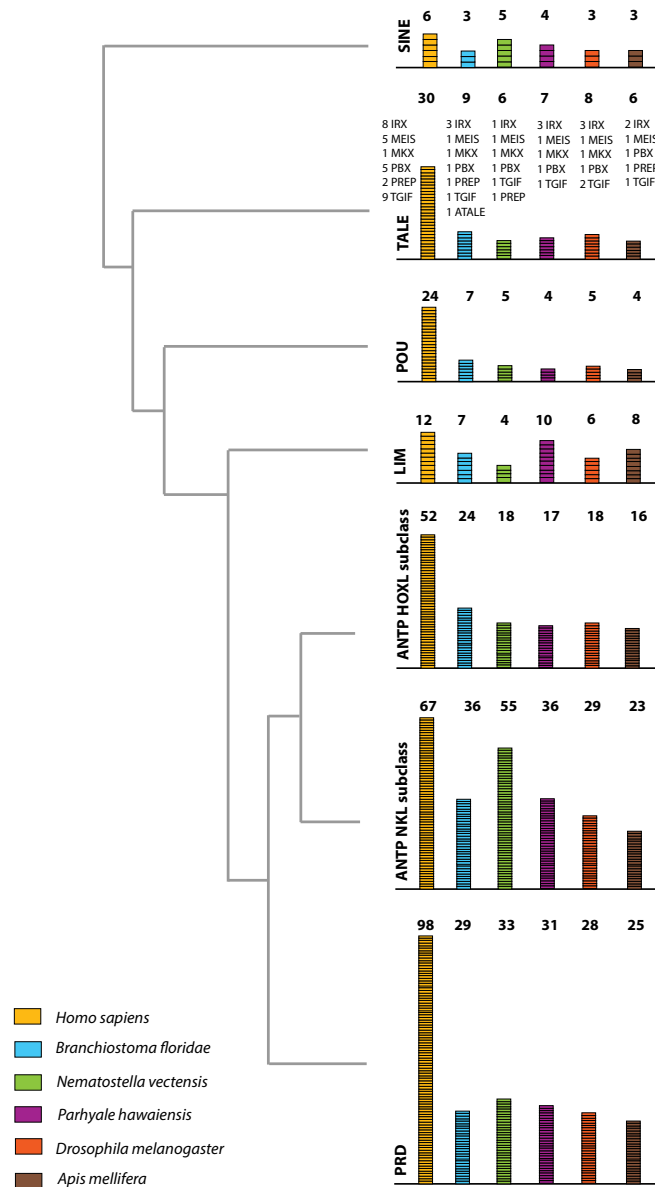


Figure 9. Homeodomain protein family tree. The overview of homeodomain radiation and phylogenetic relationships among homeodomain proteins from Arthropoda (*P. hawaiensis*, *D. melanogaster* and *A. mellifera*), Chordata (*H. sapiens* and *B. floridae*) Cnidaria (*N. vectensis*). Six major homeodomain classes are illustrated (SINE, TALE, POU, LIM, ANTP and PRD) with histograms indicating the number of genes in each species belonging to a given class.

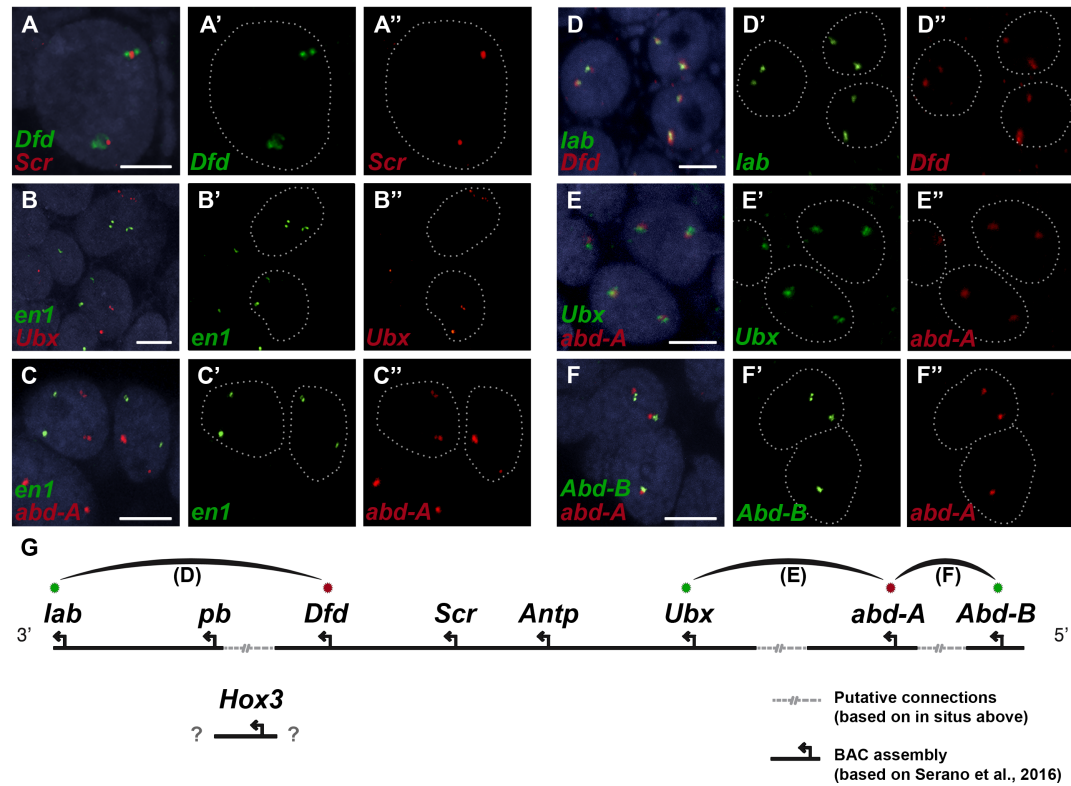


Figure 10. Variation observed in contiguous BAC sequences. (A-F'') Double fluorescent in situ hybridizations (FISH) for nascent transcripts of genes. (A-A'') Deformed (*Dfd*) and Sex combs reduced (*Scr*), (B-B'') engrailed 1 (*en1*) and Ultrabithorax (*Ubx*), (C-C'') *en1* and abdominal-A (*abd-A*), (D-D'') labial (*lab*) and *Dfd*, (E-E'') *Ubx* and *abd-A*, and (F-F'') Abdominal-B (*Abd-B*) and *abd-A*. Cell nuclei are stained with DAPI (blue) in panels A-F and outlined with white dotted lines in panels A'-F' and A''-F''. Co-localization of nascent transcript dots in A, D, E and F suggest the proximity of the corresponding Hox genes in the genomic DNA. As negative controls, the *en1* nascent transcripts in B and C do not co-localize with those of Hox genes *Ubx* or *abd-A*. (G) Schematic representation of the predicted configuration of the Hox cluster in Parhyale. Previously identified genomic linkages are indicated with solid black lines, whereas linkages established by FISH are shown with dotted gray lines. The arcs connecting the green and red dots represent the linkages identified in D, E and F, respectively. The position of the *Hox3* gene is still uncertain. Scale bars are 5 μ m.

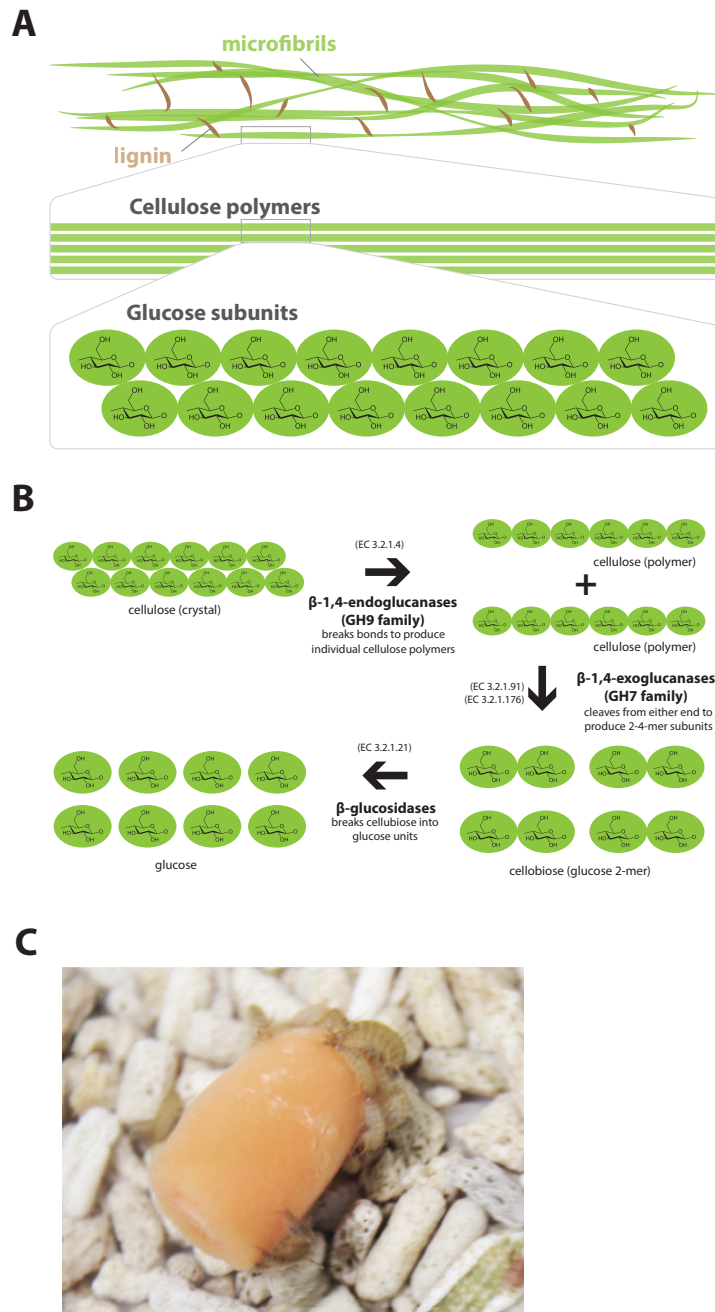


Figure 11. Lignocellulose digestion overview. (A) Simplified drawing of lignocellulose structure. The main component of lignocellulose is cellulose, which is a β -1,4-linked chain of glucose monosaccharides. Cellulose and lignin are organized in structures called microfibrils, which in turn form macrofibrils. (B) Summary of cellulolytic enzymes and reactions involved in the breakdown of cellulose into glucose. β -1,4-endoglucanases of the GH9 family catalyze the hydrolysis of crystalline cellulose into cellulose chains. β -1,4-exoglucanases of the GH7 family break down cellulose chains into cellobiose (glucose disaccharide) that can be converted to glucose by β -glucosidases. (C) Adult Parhyale feeding on a slice of carrot.

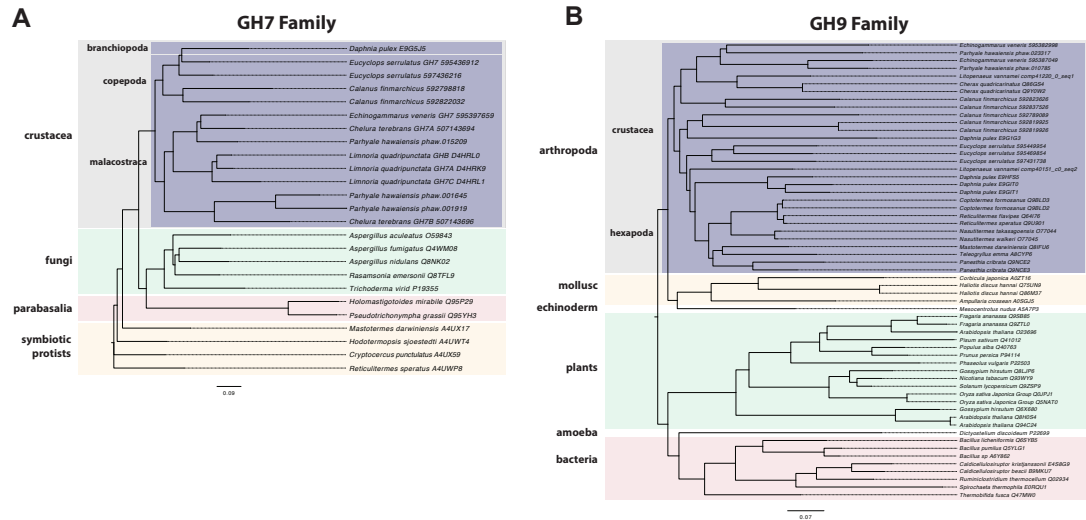


Figure 12. Phylogenetic analysis of GH7 and GH9 family proteins. (A) Phylogenetic tree showing the relationship between GH7 family proteins of *Parhyale*, other crustaceans from Vericrustacea (Malacostraca, Branchiopoda, Copepoda), fungi and symbiotic protists (root). UniProt and GenBank accessions are listed next to the species names. **(B)** Phylogenetic tree showing the relationship between GH9 family proteins of *Parhyale*, crustaceans, insects, molluscs, echinoderms, amoeba, bacteria and plants (root). UniProt and GenBank accessions are listed next to the species names. Both trees were constructed with RAxML using the WAG+G model from multiple alignments of protein sequences created with MUSCLE.

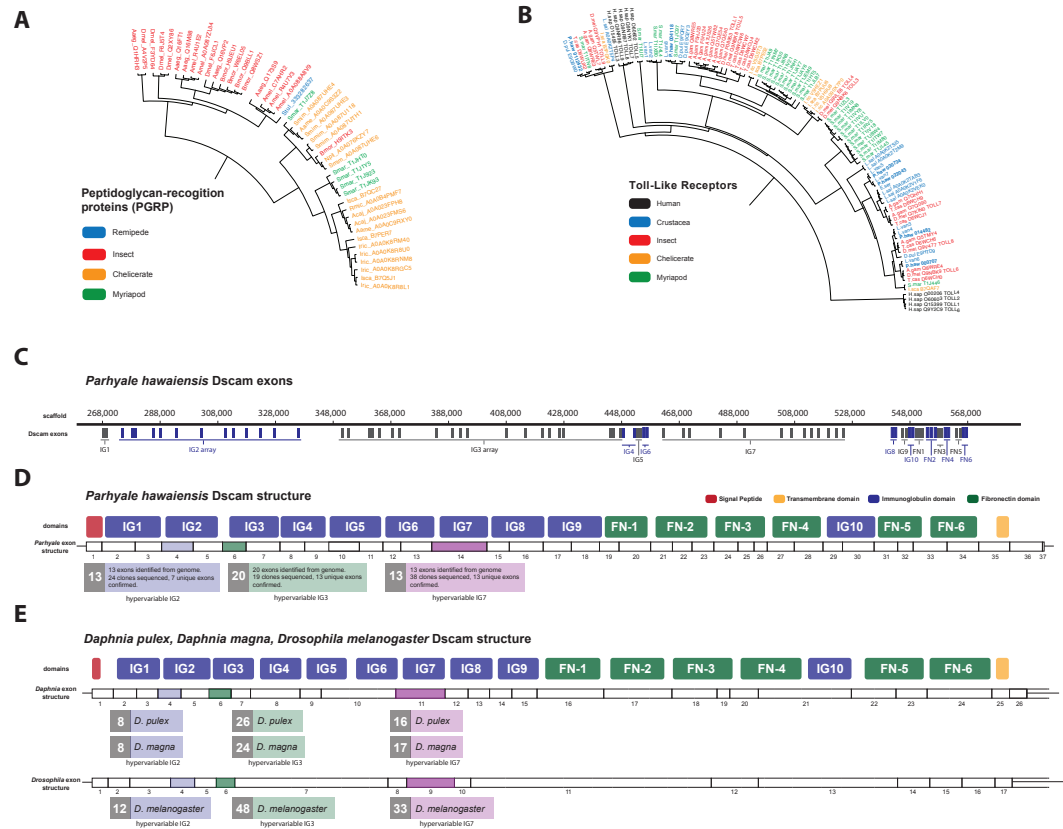


Figure 13. Peptidoglycan recognition proteins (PGRPs) and Toll-like receptors (TLRs) phylogeny. (A) Phylogenetic tree of peptidoglycan recognition proteins (PGRPs). With the exception of remipedes, PGRPs were not found in crustaceans. PGRPs have been found in the rest arthropods, including insects, myriapods and chelicerates. (B) Phylogenetic tree of Toll-like receptors (TLRs) generated from five crustaceans, three hexapods, two chelicerates, one myriapod and one vertebrate species. (C) Genomic organization of the *Parhyale* Dscam locus showing the individual exons and exon arrays encoding the immunoglobulin (IG) and fibronectin (FN) domains of the protein. (D) Structure of the *Parhyale* Dscam locus and comparison with the (E) Dscam loci from *Daphnia pulex*, *Daphnia magna* and *Drosophila melanogaster*. The white boxes represent the number of predicted exons in each species encoding the signal peptide (red), the IGs (blue), the FNs and transmembrane (yellow) domains of the protein. The number of alternative spliced exons in the arrays encoding the hypervariable regions IG2 (exon 4 in all species), IG3 (exon 6 in all species) and IG7 (exon 14 in *Parhyale*, 11 in *D. pulex* and 9 in *Drosophila*) are indicated under each species schematic in the purple, green and magenta boxes, respectively. Abbreviations of species used: *Parhyale hawaiiensis* (Phaw), *Bombyx mori* (Bmor), *Aedes aegypti* (Aaeg), *Drosophila melanogaster* (Dmel), *Apis mellifera* (Amel), *Speleonectes tulumensis* (Stul), *Strigamia maritima* (Smar), *Stegodyphus mimosarum* (Smim), *Ixodes scapularis* (Isca), *Amblyomma americanum* (Aame), *Nephila pilipes* (Npil), *Rhipicephalus microplus* (Rmic), *Ixodes ricinus* (Iric), *Amblyomma cajennense* (Acaj), *Anopheles gambiae* (Agam), *Daphnia pulex* (Apul), *Tribolium castaneum* (Tcas), *Litopenaeus vannamei* (Lvan), *Lepeophtheirus salmonis* (Lsal), *Eucyclops serrulatus* (Eser), *Homo sapiens* (H.sap). Both trees were constructed with RAXML using the WAG+G model from multiple alignments of protein sequences created with MUSCLE.

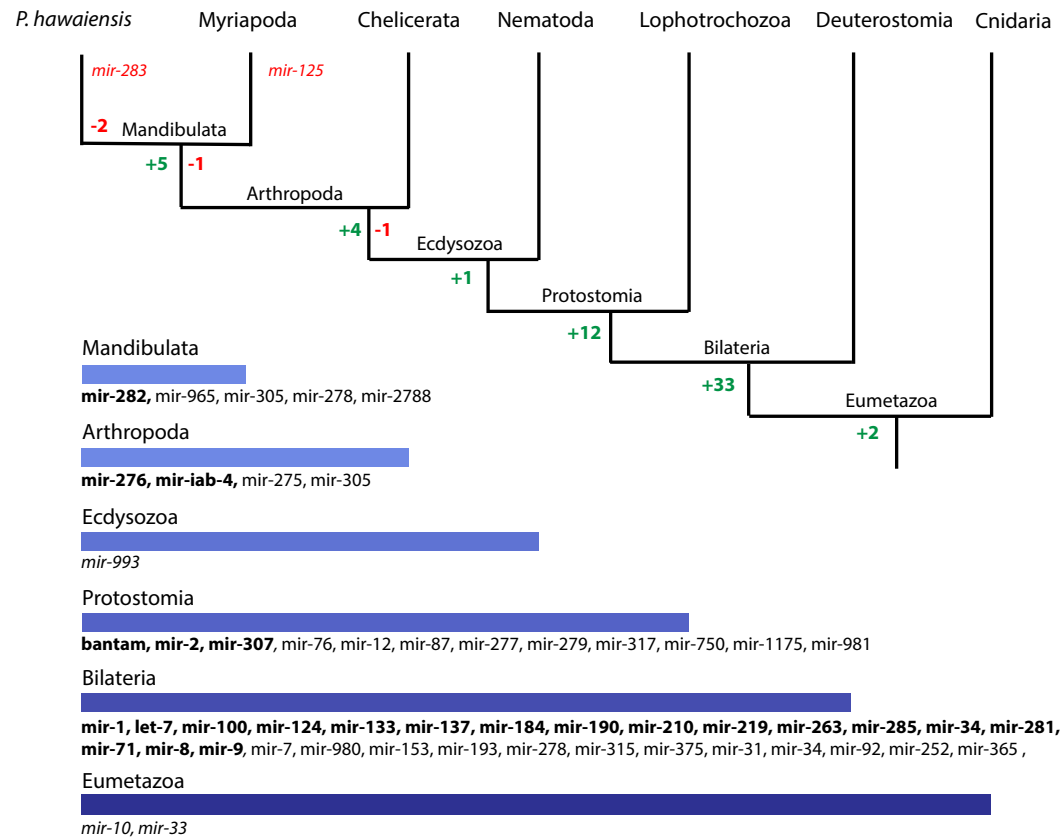


Figure 14. Evolution of miRNA families in Eumetazoans. Phylogenetic tree showing the gains (in green) and losses (in red) of miRNA families at various taxonomic levels of the Eumetazoan tree leading to Parhyale. miRNAs marked with plain characters were identified by MirPara with small RNA sequencing read support. miRNAs marked with bold characters were identified by Rfam and MirPara with small RNA sequencing read support.

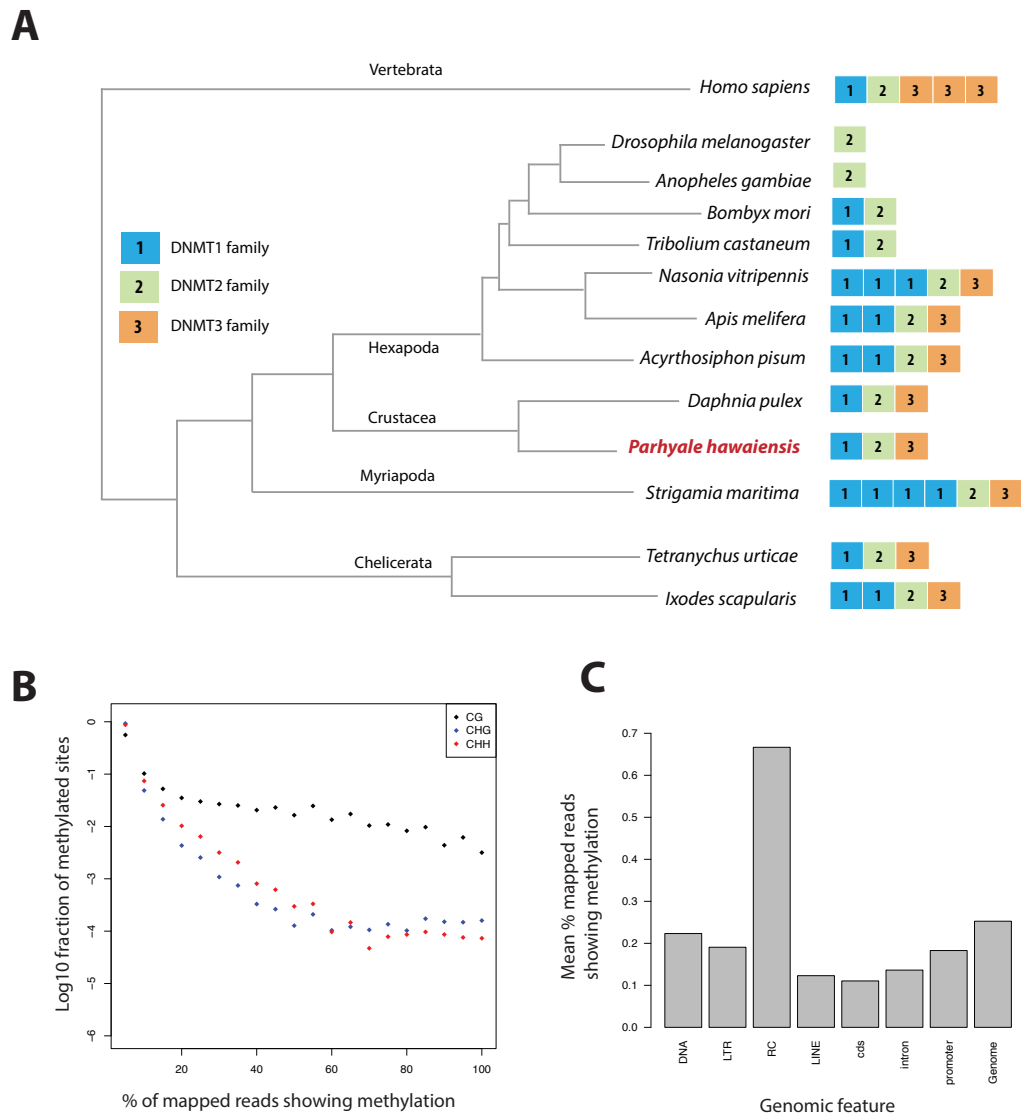


Figure 15. Analysis of *Parhyale* genome methylation. (A) Phylogenetic tree showing the families and numbers of DNA methyltransferases (DNMTs) present in the genomes of indicated species. *Parhyale* has one copy from each DNMT family. (B) Amounts of methylation detected in the *Parhyale* genome. Amount of methylation is presented as percentage of reads showing methylation in bisulfite sequencing data. DNA methylation was analyzed in all sequence contexts (CG shown in dark, CHG in blue and CHH in red) and was detected preferentially in CpG sites. (C) Histograms showing mean percentages of methylation in different fractions of the genome: DNA transposons (DNA), long terminal repeat transposable elements (LTR), rolling circle transposable elements (RC), long interspersed elements (LINE), coding sequences (cds), introns, promoters, and the rest of the genome.

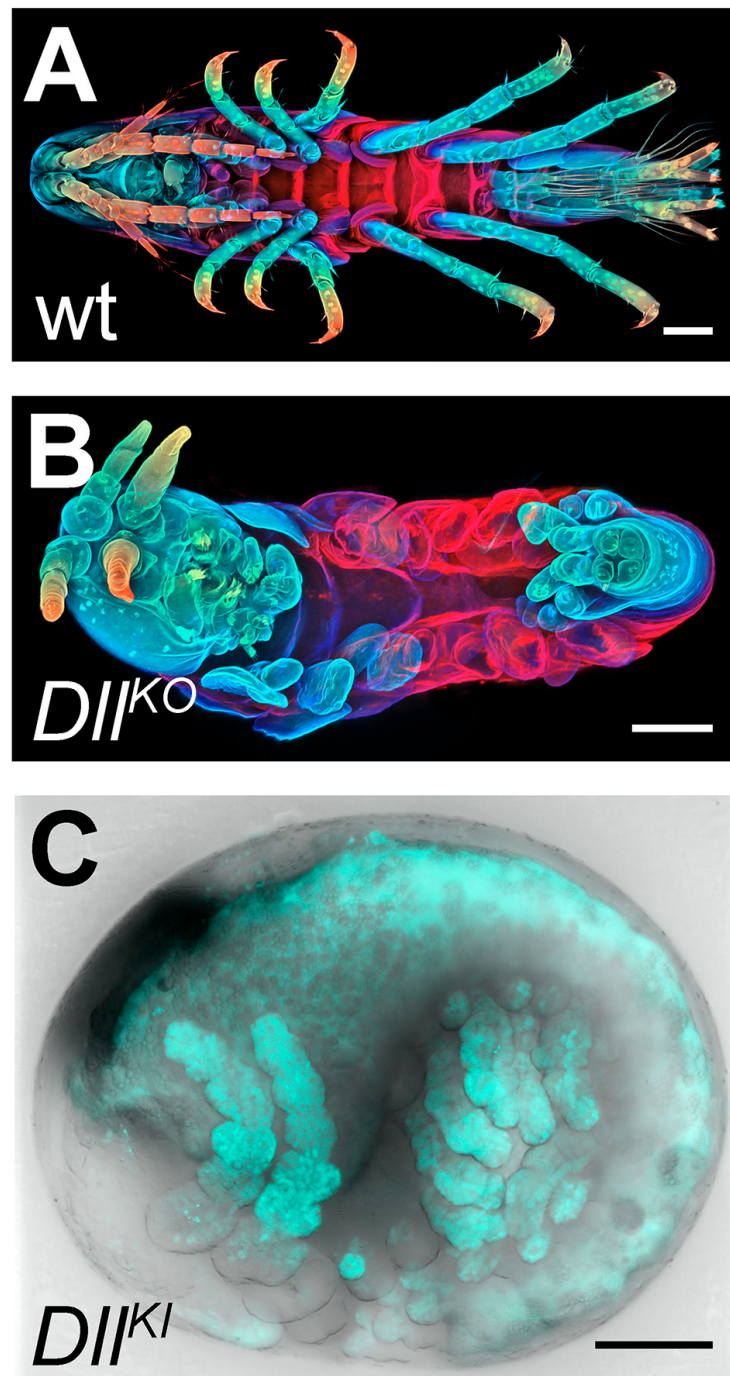


Figure 16. CRISPR/Cas9-based genome editing in *Parhyale*. (A) Wild-type morphology. (B) Mutant *Parhyale* with truncated limbs after CRISPR-mediated knock-out (DIIKO) of the limb patterning gene *Distal-less* (*PhDII-e*). Panels show ventral views of juveniles stained for cuticle and color-coded by depth with anterior to the left. (C) Fluorescent tagging of *PhDII-e* expressed in most limbs (shown in cyan) by CRISPR-mediated knock-in (DIIKI) using the non-homologous-end-joining repair mechanism. Panel shows a lateral view with anterior to the left and dorsal to the top of a live embryo (stage S22) with merged bright-field and fluorescence channels. Yolk autofluorescence produces a dorsal crescent of fluorescence in the gut. Scale bars are 100 μ m.

Table 1. Experimenta resources. Available experimental resources in *Parhyale* and corresponding references.

Experimental Resources	References
Embryological manipulations Cell microinjection, isolation, ablation	[36–38, 41–46]
Gene expression studies In situ hybridization, antibody staining	[39, 40]
Gene knock-down RNA interference, morpholinos	[22, 50]
Transgenesis Transposon-based, integrase-based	[45, 48, 49]
Gene trapping Exon/enhancer trapping, iTRAC (trap conversion)	[49]
Gene misexpression Heat-inducible	[23]
Gene knock-out CRISPR/Cas	[17]
Gene knock-in CRISPR/Cas homology-dependent or homology-independent	[16]
Live imaging Bright-field, confocal, light-sheet microscopy	[43, 44, 47]

Table 2. Assembly statistics. Length metrics of assembled scaffolds and contigs.

	# sequences	N90	N50	N10	Sum Length	Max Length	# Ns
scaffolds	133,035	14,799	81,190	289,705	3.63GB	1,285,385	1.10GB
unplaced contigs	259,343	304	627	1,779	146MB	40,222	23,431
hetero. contigs	584,392	265	402	1,038	240MB	24,461	627
genic scaffolds	15,160	52,952	161,819	433,836	1.49GB	1,285,385	323MB

Table 3. BAC variant statistics. Rate of heterozygosity of each BAC sequence determined by mapping genomic reads to each BAC individually. Population variance rate represent additional alleles found (more than 2 alleles) from genomic reads.

BAC ID	Length	Heterozygosity	Pop.Variance
PA81-D11	140,264	1.654	0.568
PA40-O15	129,957	2.446	0.647
PA76-H18	141,844	1.824	0.199
PA120-H17	126,766	2.673	1.120
PA222-D11	128,542	1.344	1.404
PA31-H15	140,143	2.793	0.051
PA284-I07	141,390	2.046	0.450
PA221-A05	148,703	1.862	1.427
PA93-L04	139,955	2.177	0.742
PA272-M04	134,744	1.925	0.982
PA179-K23	137,239	2.671	0.990
PA92-D22	126,848	2.650	0.802
PA268-E13	135,334	1.678	1.322
PA264-B19	108,571	1.575	0.157
PA24-C06	141,446	1.946	1.488

Table 4. Small RNA processing pathway members. The *Parhyale* orthologs of small RNA processing pathway members.

Gene	Counts	Gene ID
Armitage	2	phaw_30_tra_m.006391 phaw_30_tra_m.007425
Spindle_E	3	phaw_30_tra_m.000091 phaw_30_tra_m.020806 phaw_30_tra_m.018110
rm62	7	phaw_30_tra_m.014329 phaw_30_tra_m.012297 phaw_30_tra_m.004444 phaw_30_tra_m.012605 phaw_30_tra_m.001849 phaw_30_tra_m.006468 phaw_30_tra_m.023485
Piwi/aubergine	2	phaw_30_tra_m.011247 phaw_30_tra_m.016012
Dicer 1	1	phaw_30_tra_m.001257
Dicer 2	1	phaw_30_tra_m.021619
argonaute 1	1	phaw_30_tra_m.006642
argonaute 2	3	phaw_30_tra_m.021514 phaw_30_tra_m.018276 phaw_30_tra_m.012367
Loquacious	2	phaw_30_tra_m.006389 phaw_30_tra_m.000074
Drosha	1	phaw_30_tra_m.015433

895 REFERENCES

- 896 [1] M Akam. Arthropods: Developmental diversity within a (super) phylum. *Proceedings of the*
897 *National Academy of Sciences of the United States of America*, 97(9):1–4, April 2000.
- 898 [2] Graham E Budd and Maximilian J Telford. The origin and evolution of arthropods. *Nature*,
899 457(7231):812–817, February 2009.
- 900 [3] Andrew D Peel, Ariel D Chipman, and Michael Akam. Arthropod Segmentation: beyond the
901 *Drosophila* paradigm. *Nature reviews. Genetics*, 6(12):905–916, November 2005.
- 902 [4] G Scholtz and C Wolff. Arthropod embryology: cleavage and germ band development. *Arthropod*
903 *Biology and Evolution*, 2013.
- 904 [5] Jon M Mallatt, James R Garey, and Jeffrey W Shultz. Ecdysozoan phylogeny and Bayesian inference:
905 first use of nearly complete 28S and 18S rRNA gene sequences to classify the arthropods and their
906 kin. *Molecular Phylogenetics and Evolution*, 31(1):178–191, April 2004.
- 907 [6] C E Cook, Q Yue, and M Akam. Mitochondrial genomes suggest that hexapods and crustaceans are
908 mutually paraphyletic. *Proceedings. Biological sciences / The Royal Society*, 272(1569):1295–1304,
909 June 2005.
- 910 [7] Jerome C Regier, Jeffrey W Shultz, and Robert E Kambic. Pancrustacean phylogeny: hexapods are
911 terrestrial crustaceans and maxillopods are not monophyletic. *Proceedings. Biological sciences /*
912 *The Royal Society*, 272(1561):395–401, February 2005.
- 913 [8] B Ertas, B M von Reumont, J W Wagele, B Misof, and T Burmester. Hemocyanin Suggests a Close
914 Relationship of Remipedia and Hexapoda. *Molecular biology and evolution*, 26(12):2711–2718,
915 November 2009.
- 916 [9] S Richter. The Tetraconata concept: hexapod-crustacean relationships and the phylogeny of
917 Crustacea. *Organisms Diversity & Evolution*, 2(3):217–237, 2002.
- 918 [10] John K Colbourne, Michael E Pfrender, Donald Gilbert, W Kelley Thomas, Abraham Tucker, Todd H
919 Oakley, Shinichi Tokishita, Andrea Aerts, Georg J Arnold, Malay Kumar Basu, Darren J Bauer,
920 Carla E Caceres, Liran Carmel, Claudio Casola, Jeong-Hyeon Choi, John C Detter, Qunfeng Dong,
921 Serge Dusheyko, Brian D Eads, Thomas Froehlich, Kerry A Geiler-Samerotte, Daniel Gerlach, Phil
922 Hatcher, Sanjuro Jogdeo, Jeroen Krijgsveld, Evgenia V Kriventseva, Dietmar Kueltz, Christian
923 Laforsch, Erika Lindquist, Jacqueline Lopez, J Robert Manak, Jean Muller, Jasmyn Pangilinan,
924 Rupali P Patwardhan, Samuel Pitluck, Ellen J Pritham, Andreas Rechtsteiner, Mina Rho, Igor B
925 Rogozin, Onur Sakarya, Asaf Salamov, Sarah Schaack, Harris Shapiro, Yasuhiro Shiga, Courtney
926 Skalitzky, Zachary Smith, Alexander Souvorov, Way Sung, Zuojian Tang, Dai Tsuchiya, Hank Tu,
927 Harmjan Vos, Mei Wang, Yuri I Wolf, Hideo Yamagata, Takuji Yamada, Yuzhen Ye, Joseph R Shaw,

- 928 Justen Andrews, Teresa J Crease, Haixu Tang, Susan M Lucas, Hugh M Robertson, Peer Bork,
929 Eugene V Koonin, Evgeny M Zdobnov, Igor V Grigoriev, Michael Lynch, and Jeffrey L Boore. The
930 Ecoresponsive Genome of *Daphnia pulex*. *Science*, 331(6017):555–561, 2011.
- 931 [11] K Meusemann, B M von Reumont, S Simon, F Roeding, S Strauss, P Kuck, I Ebersberger, M Walz,
932 G Pass, S Breuers, V Achter, A von Haeseler, T Burmester, H Hadrys, J W Wagele, and B Misof. A
933 Phylogenomic Approach to Resolve the Arthropod Tree of Life. *Molecular biology and evolution*,
934 27(11):2451–2464, October 2010.
- 935 [12] Jerome C Regier, Jeffrey W Shultz, Andreas Zwick, April Hussey, Bernard Ball, Regina Wetzer,
936 Joel W Martin, and Clifford W Cunningham. Arthropod relationships revealed by phylogenomic
937 analysis of nuclear protein-coding sequences. *Nature*, 463(7284):1079–1083, February 2010.
- 938 [13] T H Oakley, J M Wolfe, A R Lindgren, and A K Zaharoff. Phylotranscriptomics to Bring the Under-
939 studied into the Fold: Monophyletic Ostracoda, Fossil Placement, and Pancrustacean Phylogeny.
940 *Molecular biology and evolution*, 30(1):215–233, December 2012.
- 941 [14] Bjoern M von Reumont, Ronald A Jenner, Matthew A Wills, Emiliano Dell’ampio, Günther Pass,
942 Ingo Ebersberger, Benjamin Meyer, Stefan Koenemann, Thomas M Iliffe, Alexandros Stamatakis,
943 Oliver Niehuis, Karen Meusemann, and Bernhard Misof. Pancrustacean phylogeny in the light of
944 new phylogenomic data: support for Remipedia as the possible sister group of Hexapoda. *Molecular
945 biology and evolution*, 29(3):1031–1045, March 2012.
- 946 [15] Le Cong, F Ann Ran, David Cox, Shuailiang Lin, Robert Barretto, Naomi Habib, Patrick D Hsu,
947 Xuebing Wu, Wenyan Jiang, Luciano A Marraffini, and Feng Zhang. Multiplex genome engineering
948 using CRISPR/Cas systems. *Science*, 339(6121):819–823, February 2013.
- 949 [16] Julia M Serano, Arnaud Martin, Danielle M Liubicich, Erin Jarvis, Heather S Bruce, Konnor La,
950 William E Browne, Jane Grimwood, and Nipam H Patel. Comprehensive analysis of Hox gene
951 expression in the amphipod crustacean *Parhyale hawaiiensis*. *Developmental Biology*, pages 1–13,
952 November 2015.
- 953 [17] Arnaud Martin, Julia M Serano, Erin Jarvis, Heather S Bruce, Jennifer Wang, Shagnik Ray, Carryn A
954 Barker, Liam C O’Connell, and Nipam H Patel. CRISPR/Cas9 Mutagenesis Reveals Versatile Roles
955 of Hox Genes in Crustacean Limb Specification and Evolution. *Current biology : CB*, December
956 2015.
- 957 [18] Prashant Mali, Luhan Yang, Kevin M Esvelt, John Aach, Marc Guell, James E DiCarlo, Julie E
958 Norville, and George M Church. RNA-guided human genome engineering via Cas9. *Science*,
959 339(6121):823–826, February 2013.

- 960 [19] Martin Jinek, Krzysztof Chylinski, Ines Fonfara, Michael Hauer, Jennifer A Doudna, and Em-
961 manuelle Charpentier. A programmable dual-RNA-guided DNA endonuclease in adaptive bacterial
962 immunity. *Science*, 337(6096):816–821, August 2012.
- 963 [20] Anna F Gilles and Michalis Averof. Functional genetics for all: engineered nucleases, CRISPR and
964 the gene editing revolution. *EvoDevo*, 5(1):43–13, 2014.
- 965 [21] M Averof and N H Patel. Crustacean appendage evolution associated with changes in Hox gene
966 expression. *Nature*, 388(6643):682–686, 1997.
- 967 [22] Danielle M Liubicich, Julia M Serano, Anastasios Pavlopoulos, Zacharias Kontarakis, Meredith E
968 Protas, Elaine Kwan, Sandip Chatterjee, Khoa D Tran, Michalis Averof, and Nipam H Patel.
969 Knockdown of Parhyale Ultrabithorax recapitulates evolutionary changes in crustacean appendage
970 morphology. *Proceedings of the National Academy of Sciences of the United States of America*,
971 106(33):13892–13896, August 2009.
- 972 [23] Anastasios Pavlopoulos, Zacharias Kontarakis, Danielle M Liubicich, Julia M Serano, Michael
973 Akam, Nipam H Patel, and Michalis Averof. Probing the evolution of appendage specialization by
974 Hox gene misexpression in an emerging model crustacean. *Proceedings of the National Academy of
975 Sciences of the United States of America*, 106(33):13897–13902, August 2009.
- 976 [24] Nikolaos Konstantinides and Michalis Averof. A common cellular basis for muscle regeneration in
977 arthropods and vertebrates. *Science*, 343(6172):788–791, February 2014.
- 978 [25] Jeanne L. Benton, Rachel Kery, Jingjing Li, Chadanat Noonin, Irene Söderhäll, and Barbara S. Beltz.
979 Cells from the Immune System Generate Adult-Born Neurons in Crayfish. *Developmental Cell*,
980 30(3):322–333, August 2014.
- 981 [26] L Vazquez, J Alpuche, G Maldonado, C Agundis, A Pereyra-Morales, and E Zenteno. Review:
982 Immunity mechanisms in crustaceans. *Innate Immunity*, 15(3):179–188, May 2009.
- 983 [27] Chris Hauton. The scope of the crustacean immune system for disease control. *Journal of Inverte-
984 brate Pathology*, 110(2):251–260, June 2012.
- 985 [28] T L Maginnis. The costs of autotomy and regeneration in animals: a review and framework for
986 future research. *Behavioral Ecology*, 17(5):857–872, June 2006.
- 987 [29] Sunetra Das and David S Durica. Ecdysteroid receptor signaling disruption obstructs blastemal cell
988 proliferation during limb regeneration in the fiddler crab, *Uca pugilator*. *Molecular and cellular
989 endocrinology*, 365(2):249–259, January 2013.
- 990 [30] Andrew J King, Simon M Cragg, Yi Li, Jo Dymond, Matthew J Guille, Dianna J Bowles, Neil C
991 Bruce, Ian A Graham, and Simon J McQueen-Mason. Molecular insight into lignocellulose digestion

- 992 by a marine isopod in the absence of gut microbes. *Proceedings of the National Academy of Sciences*,
993 107(12):5345–5350, March 2010.
- 994 [31] Marcelo Kern, John E. McGeehan, Simon D. Streeter, Richard N. A. Martin, Katrin Besser, Luisa
995 Elias, Will Eborall, Graham P. Malyon, Christina M. Payne, Michael E. Himmel, Kirk Schnorr,
996 Gregg T. Beckham, Simon M. Cragg, Neil C. Bruce, and Simon J. McQueen-Mason. Structural
997 characterization of a unique marine animal family 7 cellobiohydrolase suggests a mechanism of
998 cellulase salt tolerance. volume 110, pages 10189–10194, June 2013.
- 999 [32] P J Boyle and R Mitchell. Absence of Microorganisms in Crustacean Digestive Tracts. *Science*,
1000 200(4346):1157–1159, 1978.
- 1001 [33] M Zimmer, J Danko, S Pennings, A Danford, and T Carefoot. Cellulose digestion and phenol
1002 oxidation in coastal isopods (Crustacea: Isopoda). *Marine Biology*, 2002.
- 1003 [34] Carsten Wolff and Matthias Gerberding. “Crustacea”: Comparative Aspects of Early Development.
1004 In *Evolutionary Developmental Biology of Invertebrates 4*, pages 39–61. Springer Vienna, Vienna,
1005 2015.
- 1006 [35] William E Browne, Alivia L Price, Matthias Gerberding, and Nipam H Patel. Stages of embryonic
1007 development in the amphipod crustacean, *Parhyale hawaiiensis*. *Genesis (New York, N.Y. : 2000)*,
1008 42(3):124–149, July 2005.
- 1009 [36] Matthias Gerberding, William E Browne, and Nipam H Patel. Cell lineage analysis of the amphipod
1010 crustacean *Parhyale hawaiiensis* reveals an early restriction of cell fates. *Development*, 129(24):5789–
1011 5801, December 2002.
- 1012 [37] Cassandra G Extavour. The fate of isolated blastomeres with respect to germ cell formation in the
1013 amphipod crustacean *Parhyale hawaiiensis*. *Developmental Biology*, 277(2):387–402, January 2005.
- 1014 [38] E J Rehm, R L Hannibal, R C Chaw, M A Vargas-Vila, and N H Patel. Fixation and Dissection of
1015 *Parhyale hawaiiensis* Embryos. *Cold Spring Harbor Protocols*, 2009(1):pdb.prot5127–pdb.prot5127,
1016 January 2009.
- 1017 [39] E J Rehm, R L Hannibal, R C Chaw, M A Vargas-Vila, and N H Patel. Antibody Staining of *Parhyale*
1018 *hawaiiensis* Embryos. *Cold Spring Harbor Protocols*, 2009(1):pdb.prot5129–pdb.prot5129, January
1019 2009.
- 1020 [40] E Jay Rehm, Roberta L Hannibal, R Crystal Chaw, Mario A Vargas-Vila, and Nipam H Patel. In situ
1021 hybridization of labeled RNA probes to fixed *Parhyale hawaiiensis* embryos. *Cold Spring Harbor*
1022 *Protocols*, 2009(1):pdb.prot5130–pdb.prot5130, January 2009.

- 1023 [41] E Jay Rehm, Roberta L Hannibal, R Crystal Chaw, Mario A Vargas-Vila, and Nipam H Patel.
1024 Injection of Parhyale hawaiiensis blastomeres with fluorescently labeled tracers. *Cold Spring Harbor*
1025 *Protocols*, 2009(1):pdb.prot5128–pdb.prot5128, January 2009.
- 1026 [42] Alivia L Price, Melinda S Modrell, Roberta L Hannibal, and Nipam H Patel. Mesoderm and
1027 ectoderm lineages in the crustacean Parhyale hawaiiensis display intra-germ layer compensation.
1028 *Developmental Biology*, 341(1):256–266, May 2010.
- 1029 [43] Frederike Alwes, Billy Hinchin, and Cassandra G Extavour. Patterns of cell lineage, movement,
1030 and migration from germ layer specification to gastrulation in the amphipod crustacean Parhyale
1031 hawaiiensis. *Developmental Biology*, 359(1):110–123, November 2011.
- 1032 [44] Roberta L Hannibal, Alivia L Price, and Nipam H Patel. The functional relationship between
1033 ectodermal and mesodermal segmentation in the crustacean, Parhyale hawaiiensis. *Developmental*
1034 *Biology*, 361(2):427–438, January 2012.
- 1035 [45] Zacharias Kontarakis and Anastasios Pavlopoulos. Transgenesis in Non-model Organisms: The
1036 Case of Parhyale. In *Molecular Methods for Evolutionary Genetics*, pages 145–181. Springer New
1037 York, New York, NY, July 2014.
- 1038 [46] Anastasia R Nast and Cassandra G Extavour. Ablation of a Single Cell From Eight-cell Embryos
1039 of the Amphipod Crustacean Parhyale hawaiiensis. *Journal of visualized experiments : JoVE*, (85),
1040 2014.
- 1041 [47] R Crystal Chaw and Nipam H Patel. Independent migration of cell populations in the early
1042 gastrulation of the amphipod crustacean Parhyale hawaiiensis. *Developmental Biology*, 371(1):94–
1043 109, November 2012.
- 1044 [48] Anastasios Pavlopoulos and Michalis Averof. Establishing genetic transformation for comparative
1045 developmental studies in the crustacean Parhyale hawaiiensis. *Proceedings of the National Academy*
1046 *of Sciences of the United States of America*, 102(22):7888–7893, May 2005.
- 1047 [49] Zacharias Kontarakis, Anastasios Pavlopoulos, Alexandros Kiupakis, Nikolaos Konstantinides,
1048 Vassilis Douris, and Michalis Averof. A versatile strategy for gene trapping and trap conversion in
1049 emerging model organisms. *Development*, 138(12):2625–2630, June 2011.
- 1050 [50] Günes Özhan-Kizil, Johanna Havemann, and Matthias Gerberding. Germ cells in the crustacean
1051 Parhyale hawaiiensis depend on Vasa protein for their maintenance but not for their formation.
1052 *Developmental Biology*, 327(1):230–239, March 2009.
- 1053 [51] Ronald J Parchem, Francis Poulin, Andrew B Stuart, Chris T Amemiya, and Nipam H Patel. BAC
1054 library for the amphipod crustacean, Parhyale hawaiiensis. *Genomics*, 95(5):261–267, May 2010.

- 1055 [52] Xianhui Wang, Xiaodong Fang, Pengcheng Yang, Xuating Jiang, Feng Jiang, Dejian Zhao, Bolei
1056 Li, Feng Cui, Jianing Wei, Chuan Ma, Yundan Wang, Jing He, Yuan Luo, Zhifeng Wang, Xiaojiao
1057 Guo, Wei Guo, Xuesong Wang, Yi Zhang, Meiling Yang, Shuguang Hao, Bing Chen, Zongyuan
1058 Ma, Dan Yu, Zhiqiang Xiong, Yabing Zhu, Dingding Fan, Lijuan Han, Bo Wang, Yuanxin Chen,
1059 Junwen Wang, Lan Yang, Wei Zhao, Yue Feng, Guanxing Chen, Jinmin Lian, Qiye Li, Zhiyong
1060 Huang, Xiaoming Yao, Na Lv, Guojie Zhang, Yingrui Li, Jian Wang, Jun Wang, Baoli Zhu, and
1061 Le Kang. The locust genome provides insight into swarm formation and long-distance flight. *Nature*
1062 *communications*, 5:2957–2959, 2014.
- 1063 [53] Jared T Simpson. Exploring genome characteristics and sequence quality without a reference.
1064 *Bioinformatics*, 30(9):1228–1235, May 2014.
- 1065 [54] Guofan Zhang, Xiaodong Fang, Ximing Guo, Li Li, Ruibang Luo, Fei Xu, Pengcheng Yang, Linlin
1066 Zhang, Xiaotong Wang, Haigang Qi, Zhiqiang Xiong, Huayong Que, Yinlong Xie, Peter W H
1067 Holland, Jordi Paps, Yabing Zhu, Fucun Wu, Yuanxin Chen, Jiafeng Wang, Chunfang Peng, Jie
1068 Meng, Lan Yang, Jun Liu, Bo Wen, Na Zhang, Zhiyong Huang, Qihui Zhu, Yue Feng, Andrew
1069 Mount, Dennis Hedgecock, Zhe Xu, Yunjie Liu, Tomislav Domazet-Lošo, Yishuai Du, Xiaoqing
1070 Sun, Shoudu Zhang, Binghang Liu, Peizhou Cheng, Xuating Jiang, Juan Li, Dingding Fan, Wei
1071 Wang, Wenjing Fu, Tong Wang, Bo Wang, Jibiao Zhang, Zhiyu Peng, Yingxiang Li, Na Li, Jinpeng
1072 Wang, Maoshan Chen, Yan He, Fengji Tan, Xiaorui Song, Qiumei Zheng, Ronglian Huang, Hailong
1073 Yang, Xuedi Du, Li Chen, Mei Yang, Patrick M Gaffney, Shan Wang, Longhai Luo, Zhicai She,
1074 Yao Ming, Wen Huang, Shu Zhang, Baoyu Huang, Yong Zhang, Tao Qu, Peixiang Ni, Guoying
1075 Miao, Junyi Wang, Qiang Wang, Christian E W Steinberg, Haiyan Wang, Ning Li, Lumin Qian,
1076 Guojie Zhang, Yingrui Li, Huanming Yang, Xiao Liu, Jian Wang, Ye Yin, and Jun Wang. The oyster
1077 genome reveals stress adaptation and complexity of shell formation. *Nature*, 490(7418):49–54,
1078 September 2012.
- 1079 [55] Brian J Haas, Alexie Papanicolaou, Moran Yassour, Manfred Grabherr, Philip D Blood, Joshua
1080 Bowden, Matthew Brian Couger, David Eccles, Bo Li, Matthias Lieber, Matthew D MacManes,
1081 Michael Ott, Joshua Orvis, Nathalie Pochet, Francesco Strozzi, Nathan Weeks, Rick Westerman,
1082 Thomas William, Colin N Dewey, Robert Henschel, Richard D LeDuc, Nir Friedman, and Aviv
1083 Regev. De novo transcript sequence reconstruction from RNA-seq using the Trinity platform for
1084 reference generation and analysis. *Nature Protocols*, 8(8):1494–1512, July 2013.
- 1085 [56] G Parra, K Bradnam, and I Korf. CEGMA: a pipeline to accurately annotate core genes in eukaryotic
1086 genomes. *Bioinformatics*, 23(9):1061–1067, May 2007.
- 1087 [57] David M Emms and Steven Kelly. OrthoFinder: solving fundamental biases in whole genome
1088 comparisons dramatically improves orthogroup inference accuracy. *Genome biology*, 16:157, 2015.

- 1089 [58] Maura Strigini, Rafael Cantera, Xavier Morin, Michael J Bastiani, Michael Bate, and Domna
1090 Karagozeos. The IgLON protein Lachesin is required for the blood-brain barrier in *Drosophila*.
1091 *Molecular and cellular neurosciences*, 32(1-2):91–101, May 2006.
- 1092 [59] Lindsey S Garver, Zhiyong Xi, and George Dimopoulos. Immunoglobulin superfamily members
1093 play an important role in the mosquito immune system. *Developmental & Comparative Immunology*,
1094 32(5):519–531, 2008.
- 1095 [60] Matthias Siebert, Daniel Banovic, Bernd Goellner, and Hermann Aberle. *Drosophila* motor axons
1096 recognize and follow a Sidestep-labeled substrate pathway to reach their target fields. *Genes &*
1097 *development*, 23(9):1052–1062, May 2009.
- 1098 [61] C Deraison, I Darboux, L Duportets, T Gorjankina, Y Rahbe, and L Jouanin. Cloning and
1099 characterization of a gut-specific cathepsin L from the aphid *Aphis gossypii*. *Insect Molecular*
1100 *Biology*, 13(2):165–177, April 2004.
- 1101 [62] B Prud'homme, N Lartillot, G Balavoine, and A Adoutte. Phylogenetic analysis of the Wnt gene
1102 family: insights from lophotrochozoan members. *Current Biology*, 12(16):1395–1400, 2002.
- 1103 [63] Sung-Jin Cho, Yvonne Vallès, Vincent C Giani, Elaine C Seaver, and David A Weisblat. Evolutionary
1104 dynamics of the wnt gene family: a lophotrochozoan perspective. *Molecular biology and evolution*,
1105 27(7):1645–1658, July 2010.
- 1106 [64] Massimo A Hilliard and Cornelia I Bargmann. Wnt Signals and Frizzled Activity Orient Anterior-
1107 Posterior Axon Outgrowth in *C. elegans*. *Developmental Cell*, 10(3):379–390, March 2006.
- 1108 [65] Renata Bolognesi, Laila Farzana, Tamara D Fischer, and Susan J Brown. Multiple Wnt Genes Are
1109 Required for Segmentation in the Short-Germ Embryo of *Tribolium castaneum*. *Current Biology*,
1110 18(20):1624–1629, October 2008.
- 1111 [66] Thomas W. Holstein. The evolution of the wnt pathway. *Cold Spring Harbor Perspectives in*
1112 *Biology*, 4(7), 2012.
- 1113 [67] A K Ryan, B Blumberg, C Rodriguez-Esteban, S Yonei-Tamura, K Tamura, T Tsukui, J de la Pena,
1114 W Sabbagh, J Greenwald, S Choe, D P Norris, E J Robertson, R M Evans, M G Rosenfeld, and
1115 JCI Belmonte. Pitx2 determines left-right asymmetry of internal organs in vertebrates. *Nature*,
1116 394(6693):545–551, 1998.
- 1117 [68] Anja C Nagel, Alena Krejci, Gennady Tenin, Alejandro Bravo-Patiño, Sarah Bray, Dieter Maier, and
1118 Anette Preiss. Hairless-mediated repression of notch target genes requires the combined activity of
1119 Groucho and CtBP corepressors. *Molecular and cellular biology*, 25(23):10433–10441, December
1120 2005.

- 1121 [69] Ho-Ryun Chung, Ulrich Schäfer, Herbert Jäckle, and Siegfried Böhm. Genomic expansion and
1122 clustering of ZAD-containing C2H2 zinc-finger genes in *Drosophila*. *EMBO reports*, 3(12):1158–
1123 1162, December 2002.
- 1124 [70] Hamed S Najafabadi, Sanie Mnaimneh, Frank W Schmitges, Michael Garton, Kathy N Lam, Ally
1125 Yang, Mihai Albu, Matthew T Weirauch, Ernest Radovani, Philip M Kim, Jack Greenblatt, Brendan J
1126 Frey, and Timothy R Hughes. C2H2 zinc finger proteins greatly expand the human regulatory lexicon.
1127 *Nature Biotechnology*, 33(5):555–562, February 2015.
- 1128 [71] Ariel D Chipman, David E K Ferrier, Carlo Brena, Jiaxin Qu, Daniel S T Hughes, Reinhard Schröder,
1129 Montserrat Torres-Oliva, Nadia Znassi, Huaiyang Jiang, Francisca C Almeida, Claudio R Alonso,
1130 Zivkos Apostolou, Peshtewani Aqrawi, Wallace Arthur, Jennifer C J Barna, Kerstin P Blankenburg,
1131 Daniela Brites, Salvador Capella-Gutiérrez, Marcus Coyle, Peter K Dearden, Louis Du Pasquier,
1132 Elizabeth J Duncan, Dieter Ebert, Cornelius Eibner, Galina Erikson, Peter D Evans, Cassandra G
1133 Extavour, Liezl Francisco, Toni Gabaldón, William J Gillis, Elizabeth A Goodwin-Horn, Jack E
1134 Green, Sam Griffiths-Jones, Cornelis J P Grimmelikhuijzen, Sai Gubbala, Roderic Guigó, Yi Han,
1135 Frank Hauser, Paul Havlak, Luke Hayden, Sophie Helbing, Michael Holder, Jerome H L Hui, Julia P
1136 Hunn, Vera S Hunnekuhl, LaRonda Jackson, Mehwish Javaid, Shalini N Jhangiani, Francis M
1137 Jiggins, Tamsin E Jones, Tobias S Kaiser, Divya Kalra, Nathan J Kenny, Viktoriya Korchina,
1138 Christie L Kovar, F Bernhard Kraus, François Lapraz, Sandra L Lee, Jie Lv, Christigale Mandapat,
1139 Gerard Manning, Marco Mariotti, Robert Mata, Tittu Mathew, Tobias Neumann, Irene Newsham,
1140 Dinh N Ngo, Maria Ninova, Geoffrey Okwuonu, Fiona Onger, William J Palmer, Shobha Patil,
1141 Pedro Patraquim, Christopher Pham, Ling-Ling Pu, Nicholas H Putman, Catherine Rabouille,
1142 Olivia Mendivil Ramos, Adelaide C Rhodes, Helen E Robertson, Hugh M Robertson, Matthew
1143 Ronshaugen, Julio Rozas, Nehad Saada, Alejandro Sánchez-Gracia, Steven E Scherer, Andrew M
1144 Schurko, Kenneth W Siggins, DeNard Simmons, Anna Stief, Eckart Stolle, Maximilian J Telford,
1145 Kristin Tessmar-Raible, Rebecca Thornton, Maurijn van der Zee, Arndt von Haeseler, James M
1146 Williams, Judith H Willis, Yuanqing Wu, Xiaoyan Zou, Daniel Lawson, Donna M Muzny, Kim C
1147 Worley, Richard A Gibbs, Michael Akam, and Stephen Richards. The First Myriapod Genome
1148 Sequence Reveals Conservative Arthropod Gene Content and Genome Organisation in the Centipede
1149 *Strigamia maritima*. *PLoS biology*, 12(11):e1002005–24, November 2014.
- 1150 [72] Y Pewzner-Jung, S Ben-Dor, and A H Futerman. When Do Lasses (Longevity Assurance Genes)
1151 Become CerS (Ceramide Synthases)? : INSIGHTS INTO THE REGULATION OF CERAMIDE
1152 SYNTHESIS. *Journal of Biological Chemistry*, 281(35):25001–25005, August 2006.
- 1153 [73] Peter WH Holland, H Anne F Booth, and Elspeth A Bruford. Classification and nomenclature of all
1154 human homeobox genes. *BMC biology*, 5(1):47–28, 2007.
- 1155 [74] Ying-fu Zhong and Peter W H Holland. HomeoDB2: functional expansion of a comparative

- 1156 homeobox gene database for evolutionary developmental biology. *Evolution & Development*,
1157 13(6):567–568, November 2011.
- 1158 [75] Dave Kosman, Claudia M Mizutani, Derek Lemons, W Gregory Cox, William McGinnis, and Ethan
1159 Bier. Multiplex detection of RNA expression in *Drosophila* embryos. *Science*, 305(5685):846,
1160 August 2004.
- 1161 [76] Matthew Ronshaugen and Mike Levine. Visualization of trans-Homolog Enhancer-Promoter In-
1162 teractions at the Abd-B Hox Locus in the *Drosophila* Embryo. *Developmental Cell*, 7(6):925–932,
1163 December 2004.
- 1164 [77] József Zákány, Marie Kmita, and Denis Duboule. A dual role for hox genes in limb anterior-posterior
1165 asymmetry. *Science*, 304(5677):1669–1672, 2004.
- 1166 [78] N M Brooke, J Garcia-Fernandez, and PWH Holland. The ParaHox gene cluster is an evolutionary
1167 sister of the Hox gene cluster. *Nature*, 392(6679):920–922, 1998.
- 1168 [79] S L Pollard and P W Holland. Evidence for 14 homeobox gene clusters in human genome ancestry.
1169 *Current Biology*, 10(17):1059–1062, September 2000.
- 1170 [80] K Jagla, M Bellard, and M Frasch. A cluster of *Drosophila* homeobox genes involved in mesoderm
1171 differentiation programs. *BioEssays*, 23(2):125–133, February 2001.
- 1172 [81] G N Luke, L F C Castro, K McLay, C Bird, A Coulson, and P W H Holland. Dispersal of NK
1173 homeobox gene clusters in amphioxus and humans. *Proceedings of the National Academy of
1174 Sciences of the United States of America*, 100(9):1–4, April 2003.
- 1175 [82] L F C Castro and P W H Holland. Chromosomal mapping of ANTP class homeobox genes in
1176 amphioxus: piecing together ancestral genomes. *Evolution & Development*, 5(5):1–7, August 2003.
- 1177 [83] Michael E Himmel, Shi-You Ding, David K Johnson, William S Adney, Mark R Nimlos, John W
1178 Brady, and Thomas D Foust. Biomass recalcitrance: Engineering plants and enzymes for biofuels
1179 production. *Science*, 315(5813):804–807, 2007.
- 1180 [84] David B Wilson. Microbial diversity of cellulose hydrolysis. *Current Opinion in Microbiology*,
1181 14(3):259–263, June 2011.
- 1182 [85] Simon M Cragg, Gregg T Beckham, Neil C Bruce, Timothy DH Bugg, Daniel L Distel, Paul Dupree,
1183 Amaia Green Etxabe, Barry S Goodell, Jody Jellison, John E McGeehan, Simon J McQueen-Mason,
1184 Kirk Schnorr, Paul H Walton, Joy EM Watts, and Martin Zimmer. ScienceDirect Lignocellulose
1185 degradation mechanisms across the Tree of Life. *Current Opinion in Chemical Biology*, 29(C):108–
1186 119, December 2015.

- 1187 [86] C J Duan, L Xian, G C Zhao, Y Feng, H Pang, X L Bai, J L Tang, Q S Ma, and J X Feng. Isolation
1188 and partial characterization of novel genes encoding acidic cellulases from metagenomes of buffalo
1189 rumens. *Journal of Applied Microbiology*, 107(1):245–256, July 2009.
- 1190 [87] Falk Warnecke, Peter Luginbühl, Natalia Ivanova, Majid Ghassemian, Toby H Richardson, Justin T
1191 Stege, Michelle Cayouette, Alice C McHardy, Gordana Djordjevic, Nahla Aboushadi, Rotem
1192 Sorek, Susannah G Tringe, Mircea Podar, Hector Garcia Martin, Victor Kunin, Daniel Dalevi,
1193 Julita Madejska, Edward Kirton, Darren Platt, Ernest Szeto, Asaf Salamov, Kerrie Barry, Natalia
1194 Mikhailova, Nikos C Kyrpides, Eric G Matson, Elizabeth A Ottesen, Xinning Zhang, Myriam
1195 Hernández, Catalina Murillo, Luis G Acosta, Isidore Rigoutsos, Giselle Tamayo, Brian D Green,
1196 Cathy Chang, Edward M Rubin, Eric J Mathur, Dan E Robertson, Philip Hugenholtz, and Jared R
1197 Leadbetter. Metagenomic and functional analysis of hindgut microbiota of a wood-feeding higher
1198 termite. *Nature*, 450(7169):560–565, November 2007.
- 1199 [88] Daniel L Distel, Mehwish Amin, Adam Burgoyne, Eric Linton, Gustaf Mamangkey, Wendy Morrill,
1200 John Nove, Nicole Wood, and Joyce Yang. Molecular phylogeny of Pholadoidea Lamarck, 1809
1201 supports a single origin for xylophagy (wood feeding) and xylophagous bacterial endosymbiosis in
1202 *Bivalvia*. *Molecular Phylogenetics and Evolution*, 61(2):245–254, November 2011.
- 1203 [89] Amaia Green Etxabe. The wood boring amphipod *Chelura* (terebrans). pages 1–254, 2013.
- 1204 [90] B L Cantarel, P M Coutinho, C Rancurel, T Bernard, V Lombard, and B Henrissat. The Carbohydrate-
1205 Active EnZymes database (CAZy): an expert resource for Glycogenomics. *Nucleic Acids Research*,
1206 37(Database):D233–D238, January 2009.
- 1207 [91] Robert D. Finn, Jaina Mistry, Benjamin Schuster-Böckler, Sam Griffiths-Jones, Volker Hollich,
1208 Timo Lassmann, Simon Moxon, Mhairi Marshall, Ajay Khanna, Richard Durbin, Sean R. Eddy,
1209 Erik L. L. Sonnhammer, and Alex Bateman. Pfam: clans, web tools and services. *Nucleic Acids
1210 Research*, 34(Database issue):D247–251, January 2006.
- 1211 [92] Simon M Cragg, Gregg T Beckham, Neil C Bruce, Timothy D H Bugg, Daniel L Distel, Paul Dupree,
1212 Amaia Green Etxabe, Barry S Goodell, Jody Jellison, John E McGeehan, Simon J McQueen-Mason,
1213 Kirk Schnorr, Paul H Walton, Joy E M Watts, and Martin Zimmer. Lignocellulose degradation
1214 mechanisms across the Tree of Life. *Current Opinion in Chemical Biology*, 29:108–119, December
1215 2015.
- 1216 [93] G D Stentiford, D M Neil, E J Peeler, J D Shields, H J Small, T W Flegel, J M Vlak, B Jones,
1217 F Morado, S Moss, J Lotz, L Bartholomay, D C Behringer, C Hauton, and D V Lightner. Disease
1218 will limit future food supply from the global crustacean fishery and aquaculture sectors. *Journal of
1219 Invertebrate Pathology*, 110(2):141–157, June 2012.

- 1220 [94] Robert M Waterhouse, Evgenia V Kriventseva, Stephan Meister, Zhiyong Xi, Kanwal S Alvarez,
1221 Lyric C Bartholomay, Carolina Barillas-Mury, Guowu Bian, Stephanie Blandin, Bruce M Chris-
1222 tensen, Yuemei Dong, Haobo Jiang, Michael R Kanost, Anastasios C Koutsos, Elena A Levashina,
1223 Jianyong Li, Petros Ligoxygakis, Robert M Maccallum, George F Mayhew, Antonio Mendes, Kristin
1224 Michel, Mike A Osta, Susan Paskewitz, Sang Woon Shin, Dina Vlachou, Lihui Wang, Weiqi Wei,
1225 Liangbiao Zheng, Zhen Zou, David W Severson, Alexander S Raikhel, Fotis C Kafatos, George
1226 Dimopoulos, Evgeny M Zdobnov, and George K Christophides. Evolutionary dynamics of immune-
1227 related genes and pathways in disease-vector mosquitoes. *Science*, 316(5832):1738–1743, June
1228 2007.
- 1229 [95] Charles A Janeway and Ruslan Medzhitov. Innate immune recognition. *Annual review of immunol-*
1230 *ogy*, 20:197–216, 2002.
- 1231 [96] T Werner, K Borge-Renberg, P Mellroth, H Steiner, and D Hultmark. Functional Diversity of
1232 the *Drosophila* PGRP-LC Gene Cluster in the Response to Lipopolysaccharide and Peptidoglycan.
1233 *Journal of Biological Chemistry*, 278(29):26319–26322, July 2003.
- 1234 [97] C Liu, Z Xu, D Gupta, and R Dziarski. Peptidoglycan Recognition Proteins: A novel family
1235 of four human innate immunity pattern recognition molecules. *Journal of Biological Chemistry*,
1236 276(37):34686–34694, September 2001.
- 1237 [98] Abdur Rehman, Ping Taishi, Jidong Fang, Jeannine A Majde, and James M Krueger. The cloning
1238 of a rat peptidoglycan recognition protein (PGRP) and its induction in brain by sleep deprivation.
1239 *Cytokine*, 13(1):8–17, January 2001.
- 1240 [99] Haipeng Liu, Chenglin Wu, Yasuyuki Matsuda, Shun-ichiro Kawabata, Bok Luel Lee, Kenneth
1241 Söderhäll, and Irene Söderhäll. Peptidoglycan activation of the proPO-system without a peptidogly-
1242 can receptor protein (PGRP)? *Developmental & Comparative Immunology*, 35(1):51–61, January
1243 2011.
- 1244 [100] Seanna J McTaggart, Claire Conlon, John K Colbourne, Mark L Blaxter, and Tom J Little. The
1245 components of the *Daphnia pulex* immune system as revealed by complete genome sequencing.
1246 *BMC Genomics*, 10(1):175–119, 2009.
- 1247 [101] Catherine Dostert, Emmanuelle Jouanguy, Phil Irving, Laurent Troxler, Delphine Galiana-Arnoux,
1248 Charles Hetru, Jules A Hoffmann, and Jean-Luc Imler. The Jak-STAT signaling pathway is required
1249 but not sufficient for the antiviral response of *drosophila*. *Nature Immunology*, 6(9):946–953, August
1250 2005.
- 1251 [102] T Tanji, X Hu, A N R Weber, and Y T Ip. Toll and IMD Pathways Synergistically Activate an Innate
1252 Immune Response in *Drosophila melanogaster*. *Molecular and cellular biology*, 27(12):4578–4588,
1253 May 2007.

- 1254 [103] Matthew A. Benton, Matthias Pechmann, Nadine Frey, Dominik Stappert, Kai H. Conrads, Yen-
1255 Ta Chen, Evangelia Stamataki, Anastasios Pavlopoulos, and Siegfried Roth. Toll genes have an
1256 ancestral role in axis elongation. *Current Biology*, 26(12):1609–1615, 2016.
- 1257 [104] Natalia I Arbouzova and Martin P Zeidler. JAK/STAT signalling in *Drosophila*: insights into
1258 conserved regulatory and cellular functions. *Development*, 133(14):2605–2616, July 2006.
- 1259 [105] E A Levashina, L F Moita, S Blandin, G Vriend, M Lagueux, and F C Kafatos. Conserved role of
1260 a complement-like protein in phagocytosis revealed by dsRNA knockout in cultured cells of the
1261 mosquito, *Anopheles gambiae*. *Cell*, 104(5):709–718, 2001.
- 1262 [106] H Decker. Recent findings on phenoloxidase activity and antimicrobial activity of hemocyanins.
1263 *Developmental & Comparative Immunology*, 28(7-8):673–687, June 2004.
- 1264 [107] So Young Lee, Bok Luel Lee, and Kenneth Söderhäll. Processing of crayfish hemocyanin subunits
1265 into phenoloxidase. *Biochemical and Biophysical Research Communications*, 322(2):490–496,
1266 September 2004.
- 1267 [108] D Schmucker, J C Clemens, H Shu, C A Worby, J Xiao, M Muda, J E Dixon, and S L Zipursky.
1268 *Drosophila* Dscam is an axon guidance receptor exhibiting extraordinary molecular diversity. *Cell*,
1269 101(6):671–684, June 2000.
- 1270 [109] Fiona L Watson, Roland Püttmann-Holgado, Franziska Thomas, David L Lamar, Michael Hughes,
1271 Masahiro Kondo, Vivienne I Rebel, and Dietmar Schmucker. Extensive diversity of Ig-superfamily
1272 proteins in the immune system of insects. *Science*, 309(5742):1874–1878, September 2005.
- 1273 [110] Daniela Brites, Seanna McTaggart, Krystalynne Morris, Jobriah Anderson, Kelley Thomas, Isabelle
1274 Colson, Thomas Fabbro, Tom J Little, Dieter Ebert, and Louis Du Pasquier. The Dscam homologue
1275 of the crustacean *Daphnia* is diversified by alternative splicing like in insects. *Molecular biology
1276 and evolution*, 25(7):1429–1439, July 2008.
- 1277 [111] Stephane E Castel and Robert A Martienssen. RNA interference in the nucleus: roles for small
1278 RNAs in transcription, epigenetics and beyond. *Nature reviews. Genetics*, 14(2):100–112, February
1279 2013.
- 1280 [112] A. A. Aravin, N. M. Naumova, A. V. Tulin, V. V. Vagin, Y. M. Rozovsky, and V. A. Gvozdev.
1281 Double-stranded RNA-mediated silencing of genomic tandem repeats and transposable elements in
1282 the *D. melanogaster* germline. *Current biology: CB*, 11(13):1017–1027, July 2001.
- 1283 [113] N J Caplen, S Parrish, F Imani, A Fire, and R A Morgan. Specific inhibition of gene expression by
1284 small double-stranded RNAs in invertebrate and vertebrate systems. *Proceedings of the National
1285 Academy of Sciences of the United States of America*, 98(17):1–7, August 2001.

- 1286 [114] Julius Brennecke, Alexei A Aravin, Alexander Stark, Monica Dus, Manolis Kellis, Ravi Sachidanandam, and Gregory J Hannon. Discrete Small RNA-Generating Loci as Master Regulators of
1287 Transposon Activity in *Drosophila*. *Cell*, 128(6):1089–1103, March 2007.
1288
- 1289 [115] Weifeng Gu, Masaki Shirayama, Darryl Conte Jr, Jessica Vasale, Pedro J Batista, Julie M Claycomb,
1290 James J Moresco, Elaine M Youngman, Jennifer Keys, Matthew J Stoltz, Chun-Chieh G Chen,
1291 Daniel A Chaves, Shenghua Duan, Kristin D Kasschau, Noah Fahlgren, John R Yates III, Shohei
1292 Mitani, James C Carrington, and Craig C Mello. Distinct Argonaute-Mediated 22G-RNA Pathways
1293 Direct Genome Surveillance in the *C. elegans* Germline. *Molecular cell*, 36(2):231–244, October
1294 2009.
- 1295 [116] Heng-Chi Lee, Weifeng Gu, Masaki Shirayama, Elaine Youngman, Darryl Conte, and Craig C
1296 Mello. *C. elegans* piRNAs mediate the genome-wide surveillance of germline transcripts. *Cell*,
1297 150(1):78–87, July 2012.
- 1298 [117] Lin He and Gregory J Hannon. MicroRNAs: small RNAs with a big role in gene regulation. *Nature*
1299 *reviews. Genetics*, 5(7):522–531, July 2004.
- 1300 [118] J Michael Thomson, Martin Newman, Joel S Parker, Elizabeth M Morin-Kensicki, Tricia Wright,
1301 and Scott M Hammond. Extensive post-transcriptional regulation of microRNAs and its implications
1302 for cancer. *Genes & development*, 20(16):2202–2207, August 2006.
- 1303 [119] Witold Filipowicz, Suwendra N Bhattacharyya, and Nahum Sonenberg. Mechanisms of post-
1304 transcriptional regulation by microRNAs: are the answers in sight? *Nature reviews. Genetics*,
1305 2008(2):102–114, February 2008.
- 1306 [120] Peter Sarkies, Murray E Selkirk, John T Jones, Vivian Blok, Thomas Boothby, Bob Goldstein,
1307 Ben Hanelt, Alex Ardila-Garcia, Naomi M Fast, Phillip M Schiffer, Christopher Kraus, Mark J
1308 Taylor, Georgios Koutsovoulos, Mark L Blaxter, and Eric A Miska. Ancient and Novel Small RNA
1309 Pathways Compensate for the Loss of piRNAs in Multiple Independent Nematode Lineages. *PLoS*
1310 *biology*, 13(2):e1002061–20, February 2015.
- 1311 [121] Ying Dong and Markus Friedrich. Nymphal RNAi: systemic RNAi mediated gene knockdown in
1312 juvenile grasshopper. *BMC Biotechnology*, 5:25, 2005.
- 1313 [122] George M Weinstock, Gene E Robinson, Richard A Gibbs, George M Weinstock, George M
1314 Weinstock, Gene E Robinson, Kim C Worley, Hugh M Robertson, Daniel B Weaver, Martin Beye,
1315 Peer Bork, Jay D Evans, Klaus Hartfelder, Greg J Hunt, Gene E Robinson, Ryszard Maleszka,
1316 George M Weinstock, Klaus Hartfelder, Gro V Amdam, Mrcia M G Bitondi, Anita M Collins,
1317 Alexandre S Cristino, H Michael, G Lattorff, Carlos H Lobo, Robin F A Moritz, Francis M F Nunes,
1318 Robert E Page, Zil L P Simões, Diana Wheeler, Piero Carninci, Shiro Fukuda, Yoshihide Hayashizaki,
1319 Chikatoshi Kai, Jun Kawai, Naoko Sakazume, Daisuke Sasaki, Michihira Tagami, Gro V Amdam,

1320 Stefan Albert, Geert Baggerman, Kyle T Beggs, Guy Bloch, Giuseppe Cazzamali, Mira Cohen,
1321 Mark David Drapeau, Dorothea Eisenhardt, Christine Emore, Michael A Ewing, Susan E Fahrbach,
1322 Sylvain Foret, Cornelis J P Grimmelikhuijzen, Frank Hauser, Amanda B Hummon, Greg J Hunt,
1323 Jurgen Huybrechts, Andrew K Jones, Noam Kaplan, Gérard Lebouille, Michal Linial, J Troy
1324 Littleton, Alison R Mercer, Robert E Page, Gene E Robinson, Timothy A Richmond, Sandra L
1325 RodriguezZas, Elad B Rubin, David B Sattelle, David Schlipalius, Liliane Schoofs, Yair Shemesh,
1326 Jonathan V Sweedler, Rodrigo Velarde, Peter Verleyen, Evy Vierstraete, Michael R Williamson,
1327 Martin Beye, Seth A Ament, Susan J Brown, Miguel Corona, Peter K Dearden, W Augustine
1328 Dunn, Michelle M Elekonich, Christine G Elsik, Tomoko Fujiyuki, Irene Gattermeier, Tanja Gempe,
1329 Martin Hasselmann, Tatsuhiko Kadowaki, Eriko Kage, Azusa Kamikouchi, Takeo Kubo, Robert
1330 Kucharski, Takekazu Kunieda, Marcé Lorenzen, Natalia V Milshina, Mizue Morioka, Kazuaki
1331 Ohashi, Ross Overbeek, Robert E Page, Gene E Robinson, Christian A Ross, Morten Schioett, Teresa
1332 Shippy, Hideaki Takeuchi, Amy L Toth, Judith H Willis, Megan J Wilson, Evgeny M Zdobnov,
1333 Karl H J Gordon, Ivica Letunic, Kevin Hackett, Jane Peterson, Adam Felsenfeld, Mark Guyer,
1334 Michel Solignac, Richa Agarwala, Jean Marie Cornuet, Christine Emore, Greg J Hunt, Monique
1335 Monnerot, Florence Mougél, Justin T Reese, David Schlipalius, Dominique Vautrin, Daniel B
1336 Weaver, Joseph J Gillespie, Jamie J Cannone, Robin R Gutell, J Spencer Johnston, Michael B
1337 Eisen, Amanda B Hummon, Venky N Iyer, Vivek Iyer, Peter Kosarev, Aaron J Mackey, Timothy A
1338 Richmond, Victor Solovyev, Alexandre Souvorov, George M Weinstock, Michael R Williamson,
1339 Katherine A Aronstein, Katarina Bilikova, Yan Ping Chen, Andrew G Clark, Laura I Decanini,
1340 William M Gelbart, Charles Hetru, Dan Hultmark, Jean-Luc Imler, Haobo Jiang, Michael Kanost,
1341 Kiyoshi Kimura, Brian P Lazzaro, Dawn L Lopez, Jozef Simuth, Graham J Thompson, Zhen Zou,
1342 Pieter De Jong, Erica Sodergren, Miklós Csűrös, Aleksandar Milosavljevic, J Spencer Johnston,
1343 Kazutoyo Osoegawa, Stephen Richards, Chung-Li Shu, George M Weinstock, Laurent Duret, Eran
1344 Elhaik, Dan Graur, Daniel B Weaver, Gro V Amdam, Juan M Anzola, Kathryn S Campbell, Kevin L
1345 Childs, Derek Collinge, Madeline A Crosby, C Michael Dickens, Karl H J Gordon, L Sian Gramates,
1346 Christina M Grozinger, Peter L Jones, Mireia Jorda, Xu Ling, Beverly B Matthews, Jonathan Miller,
1347 Natalia V Milshina, Craig Mizzen, Miguel A Peinado, Jeffrey G Reid, Gene E Robinson, Susan M
1348 Russo, Andrew J Schroeder, Susan E St Pierre, Ying Wang, Pinglei Zhou, Richa Agarwala, Natalia V
1349 Milshina, Daniel B Weaver, Kevin L Childs, C Michael Dickens, William M Gelbart, Huaiyang Jiang,
1350 Paul Kitts, Natalia V Milshina, Barbara Ruef, Susan M Russo, Anand Venkatraman, George M
1351 Weinstock, Lan Zhang, Pinglei Zhou, J Spencer Johnston, Gildardo Aquino-Perez, Jean Marie
1352 Cornuet, Monique Monnerot, Michel Solignac, Dominique Vautrin, Charles W Whitfield, Susanta K
1353 Behura, Stewart H Berlocher, Andrew G Clark, J Spencer Johnston, Walter S Sheppard, Deborah R
1354 Smith, Andrew V Suarez, Neil D Tsutsui, and Daniel B and... Weaver. Insights into social insects
1355 from the genome of the honeybee *Apis mellifera*. *Nature*, 443(7114):931–949, October 2006.

1356 [123] Weina Xu and Zhaojun Han. Cloning and phylogenetic analysis of *sid-1*-like genes from aphids.

- 1357 *Journal of insect science (Online)*, 8(30):1–6, 2008.
- 1358 [124] J Y Roignant, C Carre, R Mugat, D Szymczak, J A Lepesant, and C Antoniewski. Absence of
1359 transitive and systemic pathways allows cell-specific and isoform-specific RNAi in *Drosophila*. *RNA*,
1360 9(3):299–308, March 2003.
- 1361 [125] Yonggan Wu, Bo Wei, Haizhou Liu, Tianxian Li, and Simon Rayner. MiRPara: a SVM-based
1362 software tool for prediction of most probable microRNA coding regions in genome scale sequences.
1363 *BMC bioinformatics*, 12(1):107, 2011.
- 1364 [126] Eric P Nawrocki, Sarah W Burge, Alex Bateman, Jennifer Daub, Ruth Y Eberhardt, Sean R Eddy,
1365 Evan W Floden, Paul P Gardner, Thomas A Jones, John Tate, and Robert D Finn. Rfam 12.0:
1366 updates to the RNA families database. *Nucleic Acids Research*, 43(Database issue):D130–7, January
1367 2015.
- 1368 [127] W Wang, F G Brunet, E Nevo, and M Long. Origin of sphinx, a young chimeric RNA gene in
1369 *Drosophila melanogaster*. *Proceedings of the National Academy of Sciences of the United States of*
1370 *America*, 99(7):4448–4453, 2002.
- 1371 [128] Martin J. Blythe, Sunir Malla, Richard Overall, Yu-Huan H. Shih, Virginie Lemay, Joanna Moreton,
1372 Raymond Wilson, and Aziz A. Aboobaker. High Through-Put sequencing of the parhyale hawaiiensis
1373 mRNAs and microRNAs to aid comparative developmental studies. *PloS one*, 7(3), 2012.
- 1374 [129] Benjamin M Wheeler, Alysha M Heimberg, Vanessa N Moy, Erik A Sperling, Thomas W Holstein,
1375 Steffen Heber, and Kevin J Peterson. The deep evolution of metazoan microRNAs. *Evolution &*
1376 *Development*, 11(1):50–68, January 2009.
- 1377 [130] Andrew Grimson, Mansi Srivastava, Bryony Fahey, Ben J Woodcroft, H Rosaria Chiang, Nicole
1378 King, Bernard M Degan, Daniel S Rokhsar, and David P Bartel. Early origins and evolution of
1379 microRNAs and Piwi-interacting RNAs in animals. *Nature*, 455(7217):1193–1197, October 2008.
- 1380 [131] Susanta K Behura. Insect microRNAs: Structure, function and evolution. *Insect Biochemistry and*
1381 *Molecular Biology*, 37(1):3–9, January 2007.
- 1382 [132] Antonio Marco, Katarzyna Hooks, and Sam Griffiths-Jones. Evolution and function of the extended
1383 miR-2 microRNA family. *RNA Biology*, 9(3):242–248, November 2014.
- 1384 [133] Anton J Enright, Bino John, Ulrike Gaul, Thomas Tuschl, Chris Sander, and Debora S Marks.
1385 MicroRNA targets in *Drosophila*. *Genome biology*, 5(1):R1, 2003.
- 1386 [134] Andrea Tanzer, Chris T Amemiya, Chang-Bae Kim, and Peter F Stadler. Evolution of microR-
1387 NAs located withinHox gene clusters. *Journal of Experimental Zoology Part B: Molecular and*
1388 *Developmental Evolution*, 304B(1):75–85, 2005.

- 1389 [135] Derek Lemons and William McGinnis. Genomic evolution of Hox gene clusters. *Science*,
1390 313(5795):1918–1922, 2006.
- 1391 [136] A Stark, N Bushati, C H Jan, P Kheradpour, E Hodges, J Brennecke, D P Bartel, S M Cohen, and
1392 M Kellis. A single Hox locus in *Drosophila* produces functional microRNAs from opposite DNA
1393 strands. *Genes & development*, 22(1):8–13, January 2008.
- 1394 [137] Teresa D Shippy, Matthew Ronshaugen, Jessica Cande, JianPing He, Richard W Beeman, Michael
1395 Levine, Susan J Brown, and Robin E Denell. Analysis of the *Tribolium* homeotic complex: insights
1396 into mechanisms constraining insect Hox clusters. *Development Genes and Evolution*, 218(3-4):127–
1397 139, April 2008.
- 1398 [138] Anton J Enright, Bino John, Ulrike Gaul, Thomas Tuschl, Chris Sander, and Debora S Marks.
1399 MicroRNA targets in *Drosophila*. *Genome biology*, 5(1):R1–14, 2003.
- 1400 [139] S Cumberledge, A Zaratzian, and S Sakonju. Characterization of two RNAs transcribed from the
1401 cis-regulatory region of the *abd-A* domain within the *Drosophila* bithorax complex. *Proceedings of*
1402 *the National Academy of Sciences of the United States of America*, 87(9):3259–3263, May 1990.
- 1403 [140] Assaf Zemach, Ivy E McDaniel, Pedro Silva, and Daniel Zilberman. Genome-Wide Evolutionary
1404 Analysis of Eukaryotic DNA Methylation. *Science*, 328(5980):916–919, 2010.
- 1405 [141] Julie A Law and Steven E Jacobsen. Establishing, maintaining and modifying DNA methylation
1406 patterns in plants and animals. *Nature reviews. Genetics*, 11(3):204–220, February 2010.
- 1407 [142] Peter A Jones. Functions of DNA methylation: islands, start sites, gene bodies and beyond. *Nature*
1408 *reviews. Genetics*, 13(7):484–492, May 2012.
- 1409 [143] Peter A. Jones and Gangning Liang. Rethinking how DNA methylation patterns are maintained.
1410 *Nature Reviews Genetics*, 10(11):805–811, September 2009.
- 1411 [144] Albert Jeltsch, Ann Ehrenhofer-Murray, Tomasz P. Jurkowski, Frank Lyko, Gunter Reuter, Serge
1412 Ankri, Wolfgang Nellen, Matthias Schaefer, and Mark Helm. Mechanism and biological role of
1413 *dnmt2* in nucleic acid methylation. *RNA Biology*, 0(0):1–16, 0. PMID: 27232191.
- 1414 [145] Mary Grace Goll, Finn Kirpekar, Keith A Maggert, Jeffrey A Yoder, Chih-Lin Hsieh, Xiaoyu Zhang,
1415 Kent G Golic, Steven E Jacobsen, and Timothy H Bestor. Methylation of tRNA^{Asp} by the DNA
1416 methyltransferase homolog *Dnmt2*. *Science*, 311(5759):395–398, January 2006.
- 1417 [146] Farah Jaber-Hijazi, Priscilla J K P Lo, Yuliana Mihaylova, Jeremy M Foster, Jack S Benner,
1418 Belen Tejada Romero, Chen Chen, Sunir Malla, Jordi Solana, Alexey Ruzov, and A Aziza Aboobaker.
1419 Planarian MBD2/3 is required for adult stem cell pluripotency independently of DNA methylation.
1420 *Developmental Biology*, 384(1):141–153, December 2013.

- 1421 [147] Jamie A Hackett, Roopsha Sengupta, Jan J Zylicz, Kazuhiro Murakami, Caroline Lee, Thomas A
1422 Down, and M Azim Surani. Germline DNA Demethylation Dynamics and Imprint Erasure Through
1423 5-Hydroxymethylcytosine. *Science*, 339(6118):448–452, 2013.
- 1424 [148] Suhua Feng, Shawn J. Cokus, Xiaoyu Zhang, Pao-Yang Chen, Magnolia Bostick, Mary G. Goll,
1425 Jonathan Hetzel, Jayati Jain, Steven H. Strauss, Marnie E. Halpern, Chinweike Ukomadu, Kirsten C.
1426 Sadler, Sriharsa Pradhan, Matteo Pellegrini, and Steven E. Jacobsen. Conservation and divergence
1427 of methylation patterning in plants and animals. *Proceedings of the National Academy of Sciences*,
1428 107(19):8689–8694, 2010.
- 1429 [149] Albert Jeltsch. Phylogeny of methylomes. *Science*, 328(5980):837–838, 2010.
- 1430 [150] G Panganiban, S M Irvine, C Lowe, H Roehl, L S Corley, B Sherbon, J K Grenier, J F Fallon,
1431 J Kimble, M Walker, G A Wray, B J Swalla, M Q Martindale, and S B Carroll. The origin and
1432 evolution of animal appendages. *Proceedings of the National Academy of Sciences of the United*
1433 *States of America*, 94(10):5162–5166, 1997.
- 1434 [151] Evangelia Stamataki and Anastasios Pavlopoulos. Non-insect crustacean models in developmental
1435 genetics including an encomium to *Parhyale hawaiiensis*. *Current Opinion in Genetics & Develop-*
1436 *ment*, 39:149–156, August 2016.
- 1437 [152] Karyn N Johnson, Marielle C W van Hulst, and Andrew C Barnes. “Vaccination” of shrimp
1438 against viral pathogens: Phenomenology and underlying mechanisms. *Vaccine*, 26(38):4885–4892,
1439 September 2008.
- 1440 [153] Yanan Lu, Junjun Liu, Liji Jin, Xiaoyu Li, YuHong Zhen, Hongyu Xue, Jiansong You, and Yongping
1441 Xu. Passive protection of shrimp against white spot syndrome virus (WSSV) using specific antibody
1442 from egg yolk of chickens immunized with inactivated virus or a WSSV-DNA vaccine. *Fish and*
1443 *Shellfish Immunology*, 25(5):604–610, November 2008.
- 1444 [154] S Rajesh Kumar, V P Ishaq Ahamed, M Sarathi, A Nazeer Basha, and A S Sahul Hameed. Immuno-
1445 logical responses of *Penaeus monodon* to DNA vaccine and its efficacy to protect shrimp against
1446 white spot syndrome virus (WSSV). *Fish and Shellfish Immunology*, 24(4):467–478, April 2008.
- 1447 [155] Andrew F Rowley and Edward C Pope. Vaccines and crustacean aquaculture—A mechanistic
1448 exploration. *Aquaculture*, 334-337(C):1–11, March 2012.
- 1449 [156] William J Palmer and Francis M Jiggins. Comparative Genomics Reveals the Origins and Diversity
1450 of Arthropod Immune Systems. *Molecular biology and evolution*, 32(8):2111–2129, August 2015.
- 1451 [157] J T Simpson, K Wong, S D Jackman, J E Schein, S J M Jones, and I Birol. ABySS: A parallel
1452 assembler for short read sequence data. *Genome Research*, 19(6):1117–1123, June 2009.

- 1453 [158] M Boetzer, C V Henkel, H J Jansen, D Butler, and W Pirovano. Scaffolding pre-assembled contigs
1454 using SSPACE. *Bioinformatics*, 27(4):578–579, February 2011.
- 1455 [159] Guillaume Marçais and Carl Kingsford. A fast, lock-free approach for efficient parallel counting of
1456 occurrences of k-mers. *Bioinformatics*, 27(6):764–770, March 2011.
- 1457 [160] Marcel Martin. Cutadapt removes adapter sequences from high-throughput sequencing reads.
1458 *EMBnet*, 17(1):10–12, August 2011.
- 1459 [161] Brian J Haas, Steven L Salzberg, Wei Zhu, Mihaela Pertea, Jonathan E Allen, Joshua Orvis, Owen
1460 White, C Robin Buell, and Jennifer R Wortman. Automated eukaryotic gene structure annotation
1461 using EvidenceModeler and the Program to Assemble Spliced Alignments. *Genome biology*,
1462 9(1):R7, 2008.
- 1463 [162] M Stanke and S Waack. Gene prediction with a hidden Markov model and a new intron submodel.
1464 *Bioinformatics*, 19(Suppl 2):ii215–ii225, October 2003.
- 1465 [163] Thomas D Wu and Colin K Watanabe. GMAP: a genomic mapping and alignment program for
1466 mRNA and EST sequences. *Bioinformatics*, 21(9):1859–1875, May 2005.
- 1467 [164] Cole Trapnell, Brian A Williams, Geo Pertea, Ali Mortazavi, Gordon Kwan, Marijke J van Baren,
1468 Steven L Salzberg, Barbara J Wold, and Lior Pachter. Transcript assembly and quantification by
1469 RNA-Seq reveals unannotated transcripts and isoform switching during cell differentiation. *Nature*
1470 *Biotechnology*, 28(5):516–520, May 2010.
- 1471 [165] Alexander Dobin, Carrie A Davis, Felix Schlesinger, Jorg Drenkow, Chris Zaleski, Sonali Jha,
1472 Philippe Batut, Mark Chaisson, and Thomas R Gingeras. STAR: ultrafast universal RNA-seq aligner.
1473 *Bioinformatics*, 29(1):15–21, January 2013.
- 1474 [166] Guy St C Slater and Ewan Birney. Automated generation of heuristics for biological sequence
1475 comparison. *BMC bioinformatics*, 6:31, 2005.
- 1476 [167] A V Lukashin and M Borodovsky. GeneMark.hmm: new solutions for gene finding. *Nucleic Acids*
1477 *Research*, 26(4):1107–1115, 1998.
- 1478 [168] A F A Smit, R Hubley, and P Green. *RepeatMasker Open-4.0.*, 2013.
- 1479 [169] Matthew Kears, Richard Moir, Amy Wilson, Steven Stones-Havas, Matthew Cheung, Shane
1480 Sturrock, Simon Buxton, Alex Cooper, Sidney Markowitz, Chris Duran, Tobias Thierer, Bruce
1481 Ashton, Peter Meintjes, and Alexei Drummond. Geneious Basic: an integrated and extendable
1482 desktop software platform for the organization and analysis of sequence data. *Bioinformatics*,
1483 28(12):1647–1649, June 2012.

- 1484 [170] E J Rehm, R L Hannibal, R C Chaw, M A Vargas-Vila, and N H Patel. In Situ Hybridization
1485 of Labeled RNA Probes to Fixed Parhyale hawaiiensis Embryos. *Cold Spring Harbor Protocols*,
1486 2009(1):pdb.prot5130–pdb.prot5130, January 2009.
- 1487 [171] Robert C Edgar. MUSCLE: multiple sequence alignment with high accuracy and high throughput.
1488 *Nucleic Acids Research*, 32(5):1792–1797, 2004.
- 1489 [172] A Stamatakis. RAxML version 8: a tool for phylogenetic analysis and post-analysis of large
1490 phylogenies. *Bioinformatics*, 2014.
- 1491 [173] Felix Krueger and Simon R Andrews. Bismark: a flexible aligner and methylation caller for
1492 Bisulfite-Seq applications. *Bioinformatics*, 27(11):1571–1572, June 2011.
- 1493 [174] Pin-Hsiang Chou, Hao-Shuo Chang, I Tung Chen, Han-You Lin, Yi-Min Chen, Huey-Lang Yang,
1494 and K C Han-Ching Wang. The putative invertebrate adaptive immune protein *Litopenaeus vannamei*
1495 Dscam (LvDscam) is the first reported Dscam to lack a transmembrane domain and cytoplasmic tail.
1496 *Developmental & Comparative Immunology*, 33(12):1258–1267, December 2009.
- 1497 [175] E P Nawrocki and S R Eddy. Infernal 1.1: 100-fold faster RNA homology searches. *Bioinformatics*,
1498 29(22):2933–2935, October 2013.
- 1499 [176] Sam Griffiths-Jones, Harpreet Kaur Saini, Stijn van Dongen, and Anton J Enright. miRBase: tools
1500 for microRNA genomics. *Nucleic Acids Research*, 36(Database issue):D154–8, January 2008.