

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 Author contributions

33 All authors were involved in Conception and design, Acquisition of data,
34 Analysis and inter-pretation of data, Drafting or revising the article.

35

36 **For correspondence**

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

38 pavlopoulosa@janelia.hhmi.org (AP)

39

40 **Competing interests**

41 The authors declare no competing interests

57
42

43 **Funding**

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

58 ABSTRACT

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

70 INTRODUCTION

71 Very few members of the Animal Kingdom hold the esteemed position of major model system for
72 understanding the world we live in. Inventions in molecular and cellular biology increasingly facilitate
73 the development of new model systems for functional genetic studies. Here we analyse the genome
74 sequence of the amphipod crustacean *Parhyale hawaiiensis* (*Parhyale*), in order to underpin its continued
75 development as a model organism. The morphological and ecological diversity of the phylum Arthropoda
76 makes them an ideal group of animals for comparative studies encompassing embryology, adaptation of
77 adult body plans and life history evolution [1–4]. While the most widely studied group are hexapods,
78 reflected by over a hundred sequencing projects available in the NCBI genome database, genomic data in
79 the other three sub-phyla in Arthropoda are still relatively 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 the
83 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
98 be 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 stage of embryogenesis. Embryogenesis takes 10 days at
103 26°C and has been described in detail with an accurate staging system [35]. Embryos display an invariant
104 cell lineage and blastomere at the 8 cell stage already becomes committed to a single germ layer (Figure
105 1D)[35, 36]. Embryonic and post-embryonic stages are amenable to experimental manipulations and direct
106 observation *in vivo* [36–47]. This can be combined with transgenic approaches [23, 45, 48, 49], RNA
107 interference (RNAi) [22] and morpholino-mediated gene knockdown [50], and transgene based lineage
108 tracing [24]. Most recently the utility of the clustered regularly interspaced short palindromic repeats
109 (CRISPR)/CRISPR-associated (Cas) system for targeted genome editing has been elegantly demonstrated
110 during the systematic study of *Parhyale* Hox genes [16, 17]. This arsenal of experimental tools (Table 1)
111 has already established *Parhyale* as an attractive model system for modern research.

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

123 **RESULTS AND DISCUSSION**

124 **Genome assembly, annotation, and validation**

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

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

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

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

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

159 RepeatModeler and RepeatMasker were used on the final assembly to find repetitive regions, which
160 were subsequently classified into families of transposable elements or short tandem repeats (Supplemental
161 HTML:repeat). We found 1,473 different repeat element sequences representing 57% of the assembly
162 (Supplemental Table:repeatClassification).

163 The *Parhyale* assembly comprises of 133,035 scaffolds (90% of assembly), 259,343 unplaced contigs
164 (4% of assembly), and 584,392 potentially allelic contigs (6% of assembly), with a total length of 4.02
165 Gb (Table 2). The N50 length of the scaffolds is 81,190bp. The final genome assembly was annotated

166 with Augustus trained with high confidence gene models derived from assembled transcriptomes, gene
167 homology, and *ab initio* predictions. This resulted in 28,155 final gene models (Figure 4B; Supplemental
168 HTML :annotation) across 14,805 genic scaffolds and 357 unplaced contigs with an N50 of 161,819, bp
169 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 with high confidence gene predictions that were not found in the transcriptome
180 (Figure 4C). The canonical proteome dataset was annotated with both Pfam, KEGG, and BLAST against
181 Uniprot. Assembly quality was further evaluated by alignment to core eukaryotic genes defined by the
182 Core Eukaryotic Genes Mapping Approach (CEGMA) database [56]. We identified 244/248 CEGMA
183 orthology groups from the assembled genome alone and 247/248 with a combination of genome and
184 mapped transcriptome data (Supplemental Figure:cegma). Additionally, 96% of over 280,000 identified
185 transcripts, most of which are fragments that do not contain a large ORF, also mapped to the assembled
186 genome. Together these data suggest that our assembly is close to complete with respect to protein coding
187 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%) than genes that fall within the lower read coverage group (2.68%) (Figure 6B). This is consistent
194 with genes achieving higher mapped genomic read coverage due to having less divergent alleles.

195 The *Parhyale* transcriptome assembled here includes data from a larger laboratory population, hence
196 we expect to see additional polymorphisms beyond the four founder haplotypes of the Chicago-F strain.
197 For a number of developmental genes, we investigated heterozygosity in the genomic reads in addition to
198 extra variants uncovered in the transcriptome or through direct cloning and sequencing of the laboratory
199 population (Supplemental Figure:selectGeneVariant).

200 Applying this analysis to all genes using the transcriptome we found additional variations not found
201 from the genomic reads. We observed that additional variations are not substantially different between

202 genes from the higher (0.88%) versus lower coverage group genes (0.73%; Figure 6C), suggesting
203 heterozygosity and population variance are independent of each other. We also performed an assessment
204 of polymorphism on previously cloned *Parhyale* developmental genes, and found startling levels of
205 variation. (Supplemental Table:devGeneVariant). For example, we found that the cDNAs of the germ line
206 determinants, *nanos* (78 SNPS, 34 non-synonymous substitutions and one 6bp indel) and *vasa* (37 SNPs,
207 7 non-synonymous substitutions and a one 6bp indel) are more distant between *Parhyale* populations than
208 might be observed for orthologs between closely related species.

209 To further evaluate the extent/level of polymorphism across the genome, we mapped the genomic
210 reads to a set of previously published Sanger-sequenced BAC clones of the *Parhyale* HOX cluster from
211 the same line of Chicago-F isofemale line [16]. We detected SNPs at a rate of 1.3 to 2.5% among the
212 BACS (Table 3) and also additional sequence differences between the BACs and genomic reads, again
213 confirming that additional haplotypes exist in the Chicago-F population.

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

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

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

238 Our analysis of shared orthologous groups is equivocal with regard to alternative hypotheses on
239 the relationships among Pancrustacean subgroups: 44 shared groups of orthologous proteins support
240 a Multicrustacea clade (uniting the Malacostraca, Copepoda and Thecostraca), 37 groups support the
241 Allocarida (Branchiopoda and Hexapoda) and 49 groups support the Vericrustacea (Branchiopoda and
242 Multicrustacea)(Supplemental Zip:cladeOrthoGroups).

243 To further analyse the evolution of the *Parhyale* proteome we examined protein families that appeared
244 to be expanded (z-score >2), compared to other taxa (Supplemental Figure:expansion, Supplemental
245 HTML:orthology, Supplemental Txt:orthoGroups). We conservatively identified 29 gene families that
246 are expanded in *Parhyale*. Gene family expansions include the Sidestep (55 genes) and Lachesin (42)
247 immunoglobulin superfamily proteins as well as nephrins (33 genes) and neurotrimins (44 genes), which
248 are thought to be involved in immunity, neural cell adhesion, permeability barriers and axon guidance
249 [58–60]. Other *Parhyale* gene expansions include APN (aminopeptidase N) (38 genes) and cathepsin-like
250 genes (30 genes), involved in proteolytic digestion [61].

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

252 We identified components of all common metazoan cell-signalling pathways are largely conserved in
253 *Parhyale*. At least 13 *Wnt* subfamilies were present in the cnidarian-bilaterian ancestor. *Wnt3* has
254 since been lost in protostomes and Lophotrochozoans retaining 12 *Wnt* genes [62, 63]. Some sampled
255 Ecdysozoans have undergone significant *Wnt* gene loss, for example *C. elegans* has only 5 *Wnt* genes
256 [64]. At most 9 are present in any individual hexapod species [65], with *wnt2* and *wnt4* potentially lost
257 before hexapod radiation. The *Parhyale* genome encodes 6 of the 13 *Wnt* subfamily genes; *wnt1*, *wnt4*,
258 *wnt5*, *wnt10*, *wnt11* and *wnt16* lacks *wnt2*, *wnt6*, *wnt7*, *wnt8*, and *wntA* (Figure 8). While *Wnt* genes are
259 known to have been ancestrally clustered [66]. We observe that *wnt1* and *wnt10* are clustered together on
260 a single scaffold (phaw_30.0003199), given *Wnt9* loss this may be the remnant of the ancient *wnt9-1-6-10*
261 cluster conserved in some protostomes.

262 We could identify 2 *FGF* genes and only a single FGF receptor (*FGFR*) in the *Parhyale* genome,
263 suggesting one *FGFR* has been lost in the Malacostracan lineage (Supplemental Figure:fgf). Within the
264 *TGF-beta* signaling pathway we found 2 genes from the activin subfamily (an activin receptor and a
265 myostatin), 7 genes from the *BMP* subfamily and 2 genes from the inhibin subfamily. Of the *BMP* genes,
266 *Parhyale* has a single decapentaplegic homologue (Supplemental Table:geneClassification). Other compo-
267 nents of the TGF-beta pathway were identified such as the neuroblastoma suppressor of tumorigenicity
268 (present in *Aedes aegypti* and *Tribolium castaneum* but absent in *D. melanogaster* and *D. pulex*) and TGFβ-
269 induced factor homeobox 1 (*TGIF1*) which is a Smad2-binding protein within the pathway present in
270 arthropods but absent in nematodes (*C. elegans* and *Brugia malayi*;Supplemental Table:geneClassification).
271 We identified homologues of *PITX2*, a downstream target of the TGF-beta pathway involved in endoderm
272 and mesoderm formation present [67] in vertebrates and crustaceans (*Parhyale* and *D. pulex*) but not
273 in insects and nematodes. With the exception of *SMAD7* and *SMAD8/9*, all other *SMADs* (*SMAD1*,

274 *SMAD2/3, SMAD4, SMAD6*) are found in arthropods sampled, including *Parhyale*. Components of other
275 pathways interacting with TGF-beta signaling like the *JNK, Par6, ROCK1/RhoA, p38* and *Akt* pathways
276 were also recovered and annotated in the *Parhyale* genome (Supplemental Table:geneClassification).

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

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

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

296 *Parhyale* contains a complement of 9 canonical Hox genes that exhibit both spatial and temporal col-
297 linearity in their expression along the anterior-posterior body axis [16]. Chromosome walking experiments
298 had shown that the Hox genes labial (*lab*) and proboscipedia (*pb*) are linked and that Deformed (*Dfd*),
299 Sex combs reduced (*Scr*), Antennapedia (*Antp*) and Ultrabithorax (*Ubx*) are also contiguous in a cluster
300 [16]. Previous experiments in *D. melanogaster* had shown that the proximity of nascent transcripts in
301 RNA fluorescent *in situ* hybridizations (FISH) coincide with the position of the corresponding genes in
302 the genomic DNA [75, 76]. Thus, we obtained additional information on Hox gene linkage by examining
303 nascent Hox transcripts in cells where Hox genes are co-expressed. We first validated this methodology
304 in *Parhyale* embryos by confirming with FISH, the known linkage of *Dfd* with *Scr* in the first maxillary
305 segment where they are co-expressed (Figure 10A). As a negative control, we detected no linkage between
306 engrailed1 (*en1*) and *Ubx* or *abd-A* transcripts (Figure 10B and C). We then demonstrated the tightly
307 coupled transcripts of *lab* with *Dfd* (co-expressed in the second antennal segment, Figure (Figure 10D),
308 *Ubx* and *abd-A* (co-expressed in the posterior thoracic segments, (Figure 10E), and *abd-A* with *Abd-B*
309 (co-expressed in the anterior abdominal segments, (Figure 10F). Collectively, all evidence supports the
310 linkage of all analysed Hox genes into a single cluster as shown in (Figure 10G). The relative orientation

311 and distance between certain Hox genes still needs to be worked out. So far, we have not been able to
312 confirm that *Hox3* is also part of the cluster due to the difficulty in visualizing nascent transcripts for *Hox3*
313 together with *pb* or *Dfd*. Despite these caveats, *Parhyale* provides an excellent arthropod model system to
314 understand these still enigmatic phenomena of Hox gene clustering and spatio-temporal colinearity, and
315 compare the underlying mechanisms to other well-studied vertebrate and invertebrate models [77].

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

325 **The *Parhyale* genome encodes glycosyl hydrolase enzymes consistent with lignocellu-** 326 **lose digestion (“wood eating”)**

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

342 *Parhyale* is a detritivore that can be sustained on a diet of carrots (Figure 11C), suggesting that they
343 too may be able to depolymerize lignocellulose for energy (Figure 11A and B). We searched for GH
344 family genes in *Parhyale* using the classification system of the CAZy (Carbohydrate-Active enZYmes)
345 database [90] and the annotation of protein domains in predicted genes with PFAM [91]. We identified
346 73 GH genes with complete GH catalytic domains that were classified into 17 families (Supplemental

347 Table:geneClassification) including 3 members of the GH7 family. Phylogenetic analysis of *Parhyale*
348 GH7s show high sequence similarity to the known GH7 genes in *L. quadripunctata* and the amphipod
349 *C. terebrans* [31] (Figure 12A; Supplemental Figure:ghAlignment). GH7 family genes were also iden-
350 tified in the transcriptomes of three more species spanning the Multicrustacea clade: *Echinogammarus*
351 *veneris* (amphipod), *Eucyclops serrulatus* (copepod) and *Calanus finmarchicus* (copepod) (Supplemental
352 Table:geneClassification). As previously reported [92], we also discovered a closely related GH7 gene
353 in the Branchiopod *Daphnia* (Figure 12A). This finding supports the grouping of Branchiopoda with
354 Multicrustacea (rather than with Hexapoda) and the acquisition of a GH7 gene by a Verticrustacean
355 ancestor. Alternatively, this suggests an even earlier acquisition of a GH7 gene by a crustacean ancestor
356 with subsequent loss of the GH7 family gene in the lineage leading to insects.

357 GH families 5,9,10, and 45 encode beta-1,4-endoglucanases which are also required for lignocellulose
358 digestion and are commonly found across Metazoa. We found 3 GH9 family genes with complete catalytic
359 domains in the *Parhyale* genome as well as in the other three Multicrustacean species (Figure 12B).
360 These GH9 enzymes exhibited a high sequence similarity to their homologues in the isopod *Limnoria*
361 and in a number of termites. Beta-glucosidases are the third class of enzyme required for digestion of
362 lignocellulose. They have been classified into a number of GH families: 1, 3, 5, 9 and 30, with GH1
363 representing the largest group [90]. In *Parhyale*, we found 7 beta-glucosidases from the GH30 family and
364 3 from the GH9 family, but none from the GH1 family.

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

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

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

383 The functions of *PGRPs* have been described in detail in insects like *D. melanogaster* [96] and the
384 *PGRP* family has also been ubiquitously reported in vertebrates, molluscs and echinoderms [97, 98].
385 Surprisingly, we found no *PGRP* genes in the *Parhyale* genome. *PGRPs* were also not found in other
386 sequence datasets from Branchiopoda, Copepoda and Malacostraca (Figure 13A), further supporting their
387 close phylogenetic relationship like the *GH7* genes.

388 In the absence of *PGRPs*, the freshwater crayfish *P. leniusculus* relies on a Lysine-type peptidoglycan
389 and serine proteinases, *SPH1* and *SPH2* that forms a complex with *LGBP* during immune response
390 [99]. In an independent analysis In *Parhyale*, we found one *LGBP* gene and two serine proteinases with
391 high sequence identity to *SPH1/2* in *Pacifastacus*. The *D. pulex* genome has also an expanded set of
392 Gram-negative binding proteins (proteins similar to *LGBP*) suggesting a compensatory mechanism for
393 the lost *PGRPs* [100]. Interestingly, we found a putative *PGRP* in the Remipede *Speleonectes tulumensis*
394 (Figure 13A) providing further support for sister group relationship of Remipedia and Hexapoda [14].

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

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

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

420 transmembrane domains) resulting in a total of 38,016 possible combinations.

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

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

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

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

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

455 We also assessed the miRNA and putative long non-coding RNAs (lncRNA) content of *Parhyale* using

456 both MiRPara and Rfam [125, 126]. We annotated 1405 homologues of known non-coding RNAs using
457 Rfam. This includes 980 predicted tRNAs, 45 rRNA of the large ribosomal subunit, 10 rRNA of the small
458 ribosomal subunit, 175 snRNA components of the major spliceosome (U1, U2, U4, U5 and U6), 5 snRNA
459 components of the minor spliceosome (U11, U12, U4atac and U6atac), 43 ribozymes, 38 snoRNAs, 71
460 conserved cis-regulatory element derived RNAs and 42 highly conserved miRNA genes (Supplemental
461 Table:RFAM; Supplemental HTML:rna). *Parhyale* long non-coding RNAs (lncRNAs) were identified
462 from the transcriptome using a series of filters to remove coding transcripts producing a list of 220,284
463 putative lncRNAs (32,223 are multi-exonic). Only one *Parhyale* lncRNA has clear homology to another
464 annotated lncRNA, the sphinx lncRNA from *D. melanogaster* [127].

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

482 Conserved linkages have also been observed between miRNAs and Hox genes in Bilateria [133–137].
483 For example, the phylogenetically conserved *mir-10* is present within both vertebrate and invertebrate Hox
484 clusters between *Hoxb4/Deformed (Dfd)* and *Hoxb5/Scr* [138, 139]. In the *Parhyale* genome and Hox
485 BAC sequences, we find that *mir-10* is also located between *Dfd* and *Src* on BAC clone PA179-K23 and
486 scaffold phaw_30.0001203 (Supplemental Figure:mirnaClusterC,D). However, we could not detect *mir-*
487 *iab-4* near the *Ubx* and *AbdA* genes in *Parhyale*, the location where it is found in other arthropods/insects
488 [140].

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

493 their contribution to morphological diversity.

494 **Methylome analysis of the *Parhyale* genome**

495 Methylation of cytosine residues (m5C) in CpG dinucleotides in animal genomes is regulated by a
496 conserved multi-family group of DNA methyltransferases (DNMTs) with diverse roles in the epigenetic
497 control of gene expression, genome stability and chromosome dynamics [141–143]. The phylogenetic
498 distribution of DNMTs in Metazoa suggests that the bilaterian ancestor had at least one member of the
499 Dnmt1 and Dnmt3 families (involved in *de novo* methylation and maintenance of DNA methylation)
500 and the Dnmt2 family (involved in tRNA methylation), as well as additional RNA methyltransferases
501 [144, 145]. Many animal groups have lost some of these DNA methyltransferases, for example *DNMT1*
502 and 3 are absent from *D. melanogaster* and flatworms [146, 147], while *DNMT2* is absent from nematodes
503 *C. elegans* and *C. briggsae* (Gutierrez and Sommer, 2004). The *Parhyale* genome encodes members of
504 all 3 families *DNMT1*, *DNMT3* and *DNMT2*, as well as 2 orthologs of conserved methyl-CpG-binding
505 proteins and a single orthologue of *Tet2*, an enzyme involved in DNA demethylation [148] (Figure 15A).

506 We used genome wide bisulfite sequencing to confirm the presence and also assess the distribution of
507 CpG dinucleotide methylation. Our results indicated that 20-30% of *Parhyale* DNA is methylated at CpG
508 dinucleotides (Figure 15B). The *Parhyale* methylation pattern is similar to that observed in vertebrates,
509 with high levels of methylation detected in transposable elements and other repetitive elements, in
510 promoters and gene bodies (Figure 15C). A particular class of rolling-circle transposons are very highly
511 methylated in the genome, potentially implicating methylation in silencing these elements. For comparison,
512 about 1% or less of CpG-associated cytosines are methylated in insects like *Drosophila*, *Apis*, *Bombyx*
513 and *Tribolium*. [141, 149, 150]. These data represent the first documentation of a crustacean methylome.
514 Considering the utility of *Parhyale* for genetic and genomic research, we anticipate future investigations to
515 shed light on the functional importance and spatiotemporal dynamics of epigenetic modifications during
516 normal development and regeneration, as well as their relevance to equivalent processes in vertebrate
517 systems.

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

519 *Parhyale* has already emerged as a powerful model for developmental genetic research where the ex-
520 pression and function of genes can be studied in the context of stereotyped cellular processes and with a
521 single-cell resolution. Several experimental approaches and standardized resources have been established
522 to study coding and non-coding sequences (Table 1). These functional studies will be enhanced by
523 the availability of the assembled and annotated genome presented here. As a first application of these
524 resources, we tested the efficiency of CRISPR/Cas system for targeted genome editing in *Parhyale*
525 [15–20, 151, 152]. In these studies, we targeted the Distal-less patterning gene (called *PhDIl-e*) [22] that
526 has a widely-conserved and highly-specific role in animal limb development [153].

527 We first genotyped our wild-type laboratory culture and found two *PhDIl-e* alleles with 23 SNPs
528 and 1 indel in their coding sequences and untranslated regions. For *PhDIl-e* knock-out, two sgRNAs

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

547 The non-homologous end joining (NHEJ) repair mechanism operating in the injected cells can be
548 exploited not only for gene knock-out experiments described above, but also for CRISPR knock-in
549 approaches where an exogenous DNA molecule is inserted into the targeted locus in a homology-
550 independent manner. This homology-independent approach could be particularly useful for *Parhyale* that
551 exhibits high levels of heterozygosity and and lab population polymorphisms, particularly polynucleotide
552 indels in introns and intergenic regions . To this end, we co-injected into 1-cell-stage embryos the Cas9
553 protein together with the strongest sgRNA and a tagging plasmid. The plasmid was designed in such a
554 way that upon its linearization by the same sgRNA and Cas9 and its integration into the *PhDII-e* locus
555 in the appropriate orientation and open reading frame, it would restore the endogenous *PhDII-e* coding
556 sequence in a bicistronic mRNA also expressing a nuclear fluorescent reporter. Among injected F0s, about
557 7% exhibited a nuclear fluorescence signal in the telopodite and exopodite (distal) parts of developing
558 appendages (Figure 16C and Supplemental Figure:funcConstruct E), i.e. in those limb segments that were
559 missing in the knock-out experiments (Figure 16B). Genotyping of one of these embryos demonstrated that
560 the tagged *PhDII-e* locus was indeed encoding a functional *PhDII-e* protein with a small in-frame deletion
561 around the targeted region (Supplemental Figure:funcConstruct F). These results, together with the other
562 recent applications of the CRISPR/Cas system to study Hox genes in *Parhyale* [16, 17], demonstrate that
563 the ability to manipulate the fertilized eggs together with the slow tempo of early cleavages can result in
564 very high targeting frequencies and low levels of mosaicism for both knock-out and knock-in approaches.
565 Considering the availability of the genome-wide resources provided here, we anticipate that the *Parhyale*

566 embryo will prove an extremely powerful system for fast and reliable F0 screens of gene expression and
567 function.

568 CONCLUSION

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

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

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

602 sustainability of this industry by preventing or controlling infectious diseases [93, 154–157].

603 An immune-related problem that will be also interesting to explore in *Parhyale* concerns the possibility
604 of a sterile digestive tract similar to that proposed for limnoriid isopods (REF King et al. PNAS 2012).
605 *Parhyale*, like limnoriid isopods, encodes and expresses all enzymes required for lignocellulose digestion
606 (King et al., 2010), suggesting that it is able to “digest wood” by itself without symbiotic microbial
607 partners. Of course, a lot of work will required to be invested in the characterization of the cellulolytic
608 system in *Parhyale* before any comparisons can be made with other well-established symbiotic digestion
609 systems of lignocellulose. Nevertheless, the possibility of an experimentally tractable animal model
610 that serves as a living bioreactor to convert lignocellulose into simpler metabolites, suggests that future
611 research in *Parhyale* may also have a strong biotechnological potential, especially for the production of
612 biofuels from the most abundant and cheapest raw material, plant biomass.

613 Several of our observations from analysing the *Parhyale* genome, and subsequently other available
614 data sets, suggest that Branchiopoda may not be a more closely related to insects than the Multicrustacea.
615 Parsimonious interpretations of our analyses on immune-related genes and GH enzymes provide support
616 for Branchiopoda as a sister group to Multicrustacea. We observed the absence of *PGRPs* in *D. pulex*,
617 a common feature within Multicrustacea, an observation that is supported by other independent reports
618 [100, 158] (Supplementary table 10). Either *PGRPs* have been lost independently in Multicrustacea
619 and Branchiopoda during arthropod evolution or Branchiopoda are not a sister taxa of insects but are
620 more closely related to the main body of Crustacean taxa. We also identified one glycosyl hydrolase
621 (GH) family 7 gene from the genome of *D. pulex* and this is also supported by other reports [85] and
622 investigation of *D. magna* [30].

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

631 **ACKNOWLEDGMENTS**

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

634 **MATERIALS AND METHODS**

635 A list of software and external datasets used are provided in Supplemental Table:externalDataSoftware.
636 Detailed methodology and codes for each section are provided as supplementary IPython notebooks in

637 HTML format viewable with a web browser. All supplemental data including IPython notebook can be
638 downloaded from this figshare link:
639 [https://figshare.com/articles/supplemental_data_for_Parhyale_hawaniensis_
640 genome/3498104](https://figshare.com/articles/supplemental_data_for_Parhyale_hawaniensis_genome/3498104)

641 **Genome library preparation and sequencing**

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

659 **Karyotyping**

660 For chromosome spreads, tissue was obtained from embryos at stage 14-18 [35]. Eggs were taken from the
661 mother and incubated for 1–2 h in isotonic colchicine solution (0.05% colchicine, ROTH in ASW). After
662 colchicine incubation, embryonic tissue was removed from egg shells and yolk and placed in hypotonic
663 solution (0.075 M KCl) for 25 min.

664 Fixation took place by replacing the hypotonic solution with freshly prepared and ice chilled Carnoy's
665 fixative (six parts ethanol, three parts methanol and one part anhydrous acetic acid) for 25 min. The
666 fixed tissue in Carnoy's fixative was minced with a pair of fine tungsten needles and the resulting cell
667 suspension was dropped with a siliconized Pasteur pipette from a height of about 5 cm onto a carefully
668 degreased and ice chilled microscopic slide. After partial evaporation of the Carnoy's fixative the slides
669 were held few times briefly into a steam of a water bath to rehydrate the tissue. The slides were then
670 dried on a 75°C metal block in a water bath. Finally, the slides with prepared chromosomes were aged
671 overnight at 60°C. After DNA staining either with Hoechst (H33342, Molecular Probes) or with DAPI
672 (Invitrogen), chromosomes were counted on a Zeiss Axioplan II Imaging equipped with C-Apochromat

673 63x/1.2 NA objective and a PCO pixelfly camera. FIJI was used to improve image quality (contrast and
674 brightness) and FIJI plugin 'Cell Counter' was used to determine the number of chromosomes.

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

676 The *Parhyale* raw data and assembled data are available on the NCBI website (project accession
677 SRP066767). Genome assembly was done with Abyss [159] at two different k-mer settings (70, 120) and
678 merged with GAM-NGS. Scaffolding was performed with SSPACE [160]. We chose a cut offs of $\geq 95\%$
679 overlap/95% when removing shorter allelic contigs before scaffolding as these gave better scaffolding
680 results as assessed by assembly metrics. Transcriptome assembly was performed with Trinity [55]. The
681 completeness of the genome and transcriptome was assessed by blasting against CEGMA genes [56] and
682 visualized by plotting the orthologue hit ratio versus e-value. K-mer analysis of variant and repetitive
683 branching was performed with String Graph Assmsembler's preqc module [53]. K-mer intersection analysis
684 was performed using jellyfish2 [161]. An in-depth description of the assembly process is detailed in
685 Supplemental HTML:assembly.

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

687 *Parhyale* transcriptome assembly was generated from Illumina reads collected from diverse embryonic
688 stages (Stages 19, 20, 22, 23, 25, and 28), and adult thoracic limbs and regenerating thoracic limbs (3 and
689 6 days post amputation). For the embryonic samples, RNA was extracted using Trizol; PolyA+ libraries
690 were prepared with the Truseq V1 kit (Illumina), starting with 0.6 - 3.5ug of total mRNA, and sequenced
691 on the Illumina Hiseq 2000 as paired-end 100 base reads, at the QB3 Vincent J. Coates Genomics Sequenc-
692 ing Laboratory. For the limb samples, RNA was extracted using Trizol; PolyA+ libraries were prepared
693 with the Truseq V2 kit (Illumina), starting with 1ug of total mRNA, and sequenced on the Illumina Hiseq
694 2500 as paired-end 100 base reads, at the IGBMC Microarray and Sequencing platform. 260 million
695 reads from embryos and 180 million reads from limbs were used for the transcriptome assembly. Prior to
696 the assembly we trimmed adapter and index sequences using cutadapt [162]. We also removed spliced
697 leader sequences: GAATTTTCACTGTTCCCTTTACCACGTTTTACTG, TTACCAATCACCCCTTTAC-
698 CAAGCGTTTACTG, CCCTTTACCAACTCTTAACTG, CCCTTTACCAACTTTACTG using cutadapt
699 with 0.2 error allowance to remove all potential variants. To assemble the transcriptome we used Trinity
700 (version trinityrnaseq_r20140413) [55] with settings: -min_kmer_cov 2, -path_reinforcement_distance 50.

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

702 Gene prediction was done with a combination of Evidence Modeler [163] and Augustus [164]. The
703 transcriptome was first mapped to the genome using GMAP [165]. A secondary transcriptome reference
704 assembly was performed with STAR/Cufflinks [166, 167]. The transcriptome mapping and Cufflinks
705 assembly was processed through the PASA pipeline [163] to consolidate the annotations. The PASA
706 dataset, a set of Exonerate [168] mapped Uniprot proteins, and Ab initio GeneMark [169] predictions
707 were consolidated with Evidence Modeler to produce a set of gene annotations. A high confidence set

708 of gene models from Evidence Modeler containing evidence from all three sources was used to train
709 Augustus. Evidence from RepeatMasker [170], PASA and Exonerate was then used to generate Augustus
710 gene predictions. A final list of genes for down-stream analysis was generated using both transcriptome
711 and gene predictions (canonical proteome dataset). Detailed methods are described in Supplemental
712 HTML:annotations.

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

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

722 **Polymorphisms in *Parhyale* developmental genes**

723 *Parhyale* genes (nucleotide sequences) were downloaded from GenBank. Each gene was used as a query
724 for blastn against the *Parhyale* genome using the Geneious software [171]. In each case two reference con-
725 tigs hits were observed where both had E values of close to zero. A new sequence called geneX_snp was cre-
726 ated and this sequence was annotated with the snps and/or indels present in the alternative genomic contigs.
727 To determine the occurrence of synonymous and non-synonymous substitution, the original query and the
728 newly created sequence (with polymorphisms annotated) were in silico translated into protein sequences
729 followed by pairwise alignment. Regions showing amino acid changes were annotated as non-synonymous
730 substitutions. Five random genes from the catalogue were selected for PCR, cloning and Sanger sequenc-
731 ing to confirm genomic polymorphisms and assess further polymorphism in the lab population. Primers
732 for genomic PCR designed to capture exon regions are listed as the following: dachshund (PH1F = 5'-
733 GGTGCGCTAAATTGAAGAAATTACG-3' and PH1R = 5'- ACTCAGAGGGTAATAGTAACAGAA-3'),
734 distalless exon 2 (PH2F = 5'-CACGGCCCCGGCACTA ACTATCTC-3' and PH2R = 5'-GTAATATATCTTACAACAACGA
735 3'), distalless exon 3 (PH3F = 5'-GGTGAACGGGCCGGAGTCTC-3' and PH3R = 5'-GCTGTGGGTGCTGTGGGT-
736 3'), homothorax (PH4F = 5'-TCGGGGTGTA AAAAGGACTCTG-3' and PH4R = 5'-AACATAGGAACTCACCTGGTG
737 3'), orthodenticle (PH5F = 5'-TTTGCCACTAACACATATTTGAAA-3' and PH5R = 5'-TCCCAAGTAGATGATCCCT
738 3') and prospero (PH6F = 5'-TACACTGCAACATCCGATGACTTA-3' and PH6R = 5'-CGTGTTATGTTCTCTCGTGGC
739 3').

740 **Evolutionary analyses of orthologous groups**

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

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

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

751 Embryo fixation and in-situ hybridization was performed according to [38, 172]. To enhance the nascent
752 nuclear signal over mature cytoplasmic transcript, we used either early germband embryos (Stages 11
753 – 15) in which expression of *lab*, *Dfd*, and *Scr* are just starting [16], or probes that contain almost
754 exclusively intron sequence (*Ubx*, *abd-A*, *Abd-B*, and *en1*). *Lab*, *Dfd*, and *Scr* probes are described
755 in [16]. Template for the intron-spanning probes were amplified using the following primers: *en1*-
756 Intron1, AAGACACGACGAGCATCCTG and CTGTGTATGGCTACCCGTCC; *Ubx*-Intron1, GGTAT-
757 GACAGCCGTCCAACA and AGAGTGCCAAGGATACCCGA; *abd-A*, CGATATACCCAGTCCGGTGC
758 and TCATCAGCGAGGGCACAATT; *Abd-B*, GCTGCAGGATATCCACACGA and TGCAGTTGC-
759 CGCCATAGTAA. A T7-adapter was appended to the 5' end of each reverse primer to enable direct
760 transcription from PCR product. Probes were labeled with either Digoxigenin (DIG) or Dinitrophenol
761 (DNP) conjugated UTPs, and visualized using sheep α -DIG (Roche) and donkey α -Sheep AlexaFluor
762 555 (Thermo Fischer Scientific), or Rabbit α -DNP (Thermo Fischer Scientific) and Donkey α -Rabbit
763 AlexaFluor 488 (Jackson ImmunoResearch), respectively following the procedure of Ronshaugen and
764 Levine (2004). Preparations were imaged on an LSM 780 scanning laser confocal (Zeiss), and processed
765 using Volocity software (Perkin-Elmer).

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

770 **Phylogenetic tree construction**

771 Multiple sequence alignments of protein sequences for gene families of *FGF*, *FGFR*, *CERS*, *GH7*,
772 *GH9*, *PGRP*, Toll-like receptors, *DICER*, Piwi and Argonaute were performed using MUSCLE [173].
773 Phylogenetic tree construction was performed with RAxML [174] using the WAG+G model from
774 MUSCLE multiple alignments.

775 **Bisulfite sequencing**

776 Libraries for DNA methylation analysis by bisulfite sequencing were constructed from 100ng of genomic
777 DNA extracted from one *Parhyale* male individual, using the Illumina Truseq DNA methylation kit

778 according to manufacturers instructions. Alignments to the *Parhyale* genome were generated using the
779 core Bismark module from the program Bismark [175], having first artificially joined the *Parhyale* contigs
780 to generate 10 pseudo-contigs as the program is limited as to the number of separate contigs it can analyse.
781 We then generated genome-wide cytosine coverage maps using the bismark_methylation_extraction
782 module with the parameter `-CX` specified to generate annotations of CG, CHH and CHG sites. In order
783 to analyse genome-wide methylation patterns, cytosines with more than 10 read depth coverage were
784 selected. Overall methylation levels at CG, CHH and CHG sites were generated using a custom Perl
785 script. To analyse which regions were methylated we mapped back from the joined contigs to the original
786 contigs and assigned these to functional regions based on RepeatMasker [170] and transcript annotations
787 of repeats and genes respectively. To generate overall plots of methylation levels in different features we
788 averaged over all sites mapping to particular features, focusing on CG methylation and measuring the
789 %methylation at each site as the number of reads showing methylation divided by the total number of
790 reads covering the site. Meta gene plots over particular features were generated similarly except that sites
791 mapping within a series of 100bp wide bins from 1000bp upstream of the feature start site onwards were
792 collated.

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

794 For the identification of *Dscam* in the *Parhyale*, we used the *Dscam* protein sequence from crustaceans *D.*
795 *pulex* [110] and *L. vannamei* [176] as queries to probe the assembled genome using tBlastN. A 300kb
796 region on scaffold phaw_30.0003392 was found corresponding to the *Parhyale Dscam* extending from
797 IG1 to FN6 exons. This sequence was annotated using transcriptome data together with manual searches
798 for open reading frames to identify IG, FN exons and exon-intron boundaries (Figure 10). Hypervariable
799 regions of IG2, IG3 and IG7 were also annotated accordingly on the scaffold (Figure 8). This region
800 represents a bona fide *Dscam* paralog as it matches the canonical extracellular *Dscam* domain structure
801 of nine IGs – four FNs – one IG and two FNs. *Parhyale* mRNA extractions were performed using
802 the Zymo Research Direct-zol RNA MiniPrep kit according to manufacturer's instructions. Total RNA
803 extract was used for cDNA synthesis using the Qiagen QuantiTect Reverse Transcription Kit according to
804 manufacturer's instructions. To identify and confirm potential hypervariable regions from the *Parhyale*
805 *Dscam* (PhDscam) transcript, three regions of PhDscam was corresponding to IG2, IG3 and IG7 exons
806 respectively were amplified using the following primer pairs. IG2 region:

807 DF1 = 5'-CCCTCGTGTTCCCGCCCTTCAAC-3'

808 DR1 = 5'-GCGATGTGCAGCTCTCCAGAGGG-3'

809 IG3 region:

810 DF2 = 5'-TCTGGAGAGCTGCACATCGCTAAT-3'

811 DR2 = 5'-GTGGTCATTGCGTACGAAGCACTG-3'

812 IG7 region:

813 DF3 = 5'-CGGATACCCCATCGACTCCATCG-3'

814 DR3 = 5'-GAAGCCGTCAGCCTTGCATTCAA-3'

815 PCR of each region was performed using Phusion High-fidelity polymerase from Thermo Fisher Scientific
816 and thermal cycling was done as the following: 98°C 30s, followed by 30 cycles of 98°C 10s, 67°C 30s,
817 72°C 1m30s, and then 72°C 5m. PCR products were cloned into pGEMT-Easy vector and a total of 81
818 clones were selected and Sanger sequenced and in silico translated in the correct reading frame using
819 Geneious (R7; [171] for multiple sequence alignment.

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

821 *Parhyale* non-protein-coding RNAs were identified using two independent approaches. Infernal 1.1.1
822 [177] was used with the RFAM 12.0 database [126] to scan the genome to identified potential non-protein-
823 coding RNAs according. Additionally, MiRPara [125] was used to scan the genome for potential miRNA
824 precursors. These potential precursors were further filtered using small RNA read mapping and miRBase
825 mapping [178]. Putative lncRNAs were identified from the transcriptome by applying filtering criteria
826 including removal of known coding proteins and removal of predicted proteins. Detailed methods are
827 available in Supp_rna.

828 **CRISPR/Cas genome editing**

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

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

837 5'-GAAATTAATACGACTCACTATA

838 AGAGTTGTTACCAAAGAAGTTTTAGAGCTAGAAATAGC-3'

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

840 5'-GAAATTAATACGACTCACTAT

841 AGGCTTCCCCGCCCATGTAGTTTTAGAGCTAGAAATAGC-3'

842 together with the universal reverse primer:

843 5'-AAAAGCACCGACTCGGTGCCACTTTTTCAAGTTGATAA

844 CGGACTAGCCTTATTTAACTTGCTATTTCTAGCTCTAAAAC-3'

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

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

850 Before microinjection, a small aliquot of the sgRNA was centrifuged, the pellet was washed with 70%
851 ethanol, resuspended in nuclease-free water and quantified on a Nanodrop spectrophotometer (Thermo
852 Scientific). The Cas9 was provided either as in vitro synthesized capped mRNA or as recombinant protein.
853 Cas9 mRNA synthesis was carried out as previously described [45] using plasmid T7-Cas9 (a gift from
854 David Stern and Justin Crocker) linearized with EcoRI digestion. The lyophilized Cas9 protein (PNA
855 Bio Inc) was resuspended in nuclease-free water at a concentration of 1.25 µg/µl and small aliquots were
856 stored at -80°C. For microinjections, we mixed 400 ng/µl of Cas9 protein with 40-200 ng/µl sgRNA,
857 incubated at 37°C for 5 min, transferred on ice, added the inert dye phenol red (5x from Sigma-Aldrich)
858 and, for knock-in experiments, the tagging plasmid at a concentration of 10 ng/µl. The injection mix was
859 centrifuged for 20 min at 4°C and the cleared solution was microinjected into 1-cell-stage embryos as
860 previously described [45].

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

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

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

887 Ruby2 tagging plasmid were screened for nuclear fluorescence in the developing appendages under an
888 Olympus MVX10 epi-fluorescence stereomicroscope. To image expression, live embryos at stage S22
889 were mounted in 0.5% SeaPlaque low-melting agarose (Lonza) in glass bottom microwell dishes (MatTek
890 Corporation) and scanned as described above acquiring both the fluorescence and transmitted light on an
891 inverted Zeiss 880 confocal microscope. To recover the chromosome-plasmid junctions, genomic DNA
892 was extracted from transgenic siblings with fluorescent limbs and used as template in PCR reaction as
893 described above with primer pair 313For and H2BRev (5'-TTACTTAGAAGAAGTGTACTTTG-3') for
894 the left junction and primer pair M13 forward and 557Rev for the right junction. Amplified products were
895 purified and cloned as described above and sequenced with the M13 forward and M13 reverse primers.

896 FIGURES AND TABLES

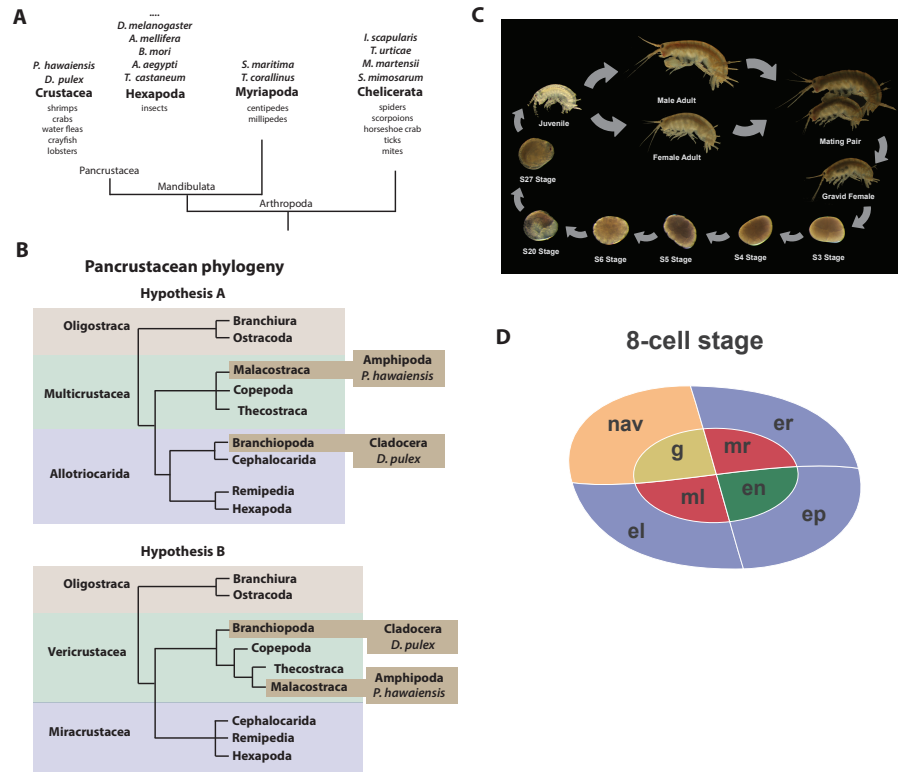


Figure 1. Introduction. (A) Phylogenetic relationship of Arthropods showing the traditional view of Chelicerata as an outgroup to Mandibulata and the Pancrustacea clade which includes crustaceans and insects. Species listed for each clade have ongoing or complete genomes. Species for Crustacea include: *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) Alternative hypotheses of Pancrustacean phylogeny. Hypothesis A depicts Branchiopoda as part of the Allotricarida clade that includes remipedes and insects. Hypothesis B depicts two of the four Pancrustacea clades (Vericrustacea and Miracrustacea). According to hypothesis B, Branchiopoda is a sister group to Multicrustacea (Copepoda, Thecostraca and Malacostraca). (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.

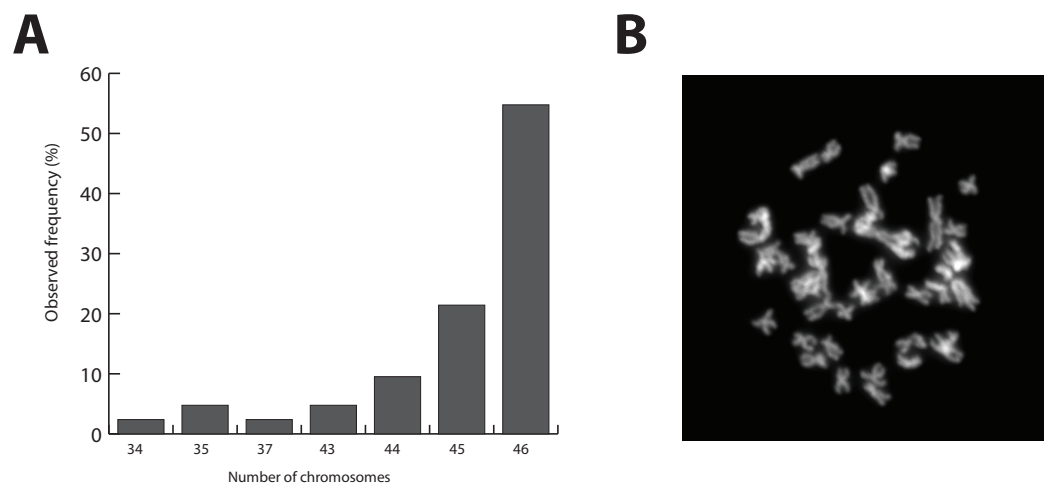


Figure 2. *Parhyale* Karyotype. (A) Frequency of the number of chromosomes observed in 42 mitotic spreads. Forty-six chromosomes were observed in more than half preparations. (B) Representative image of Hoechst-stained chromosomes.

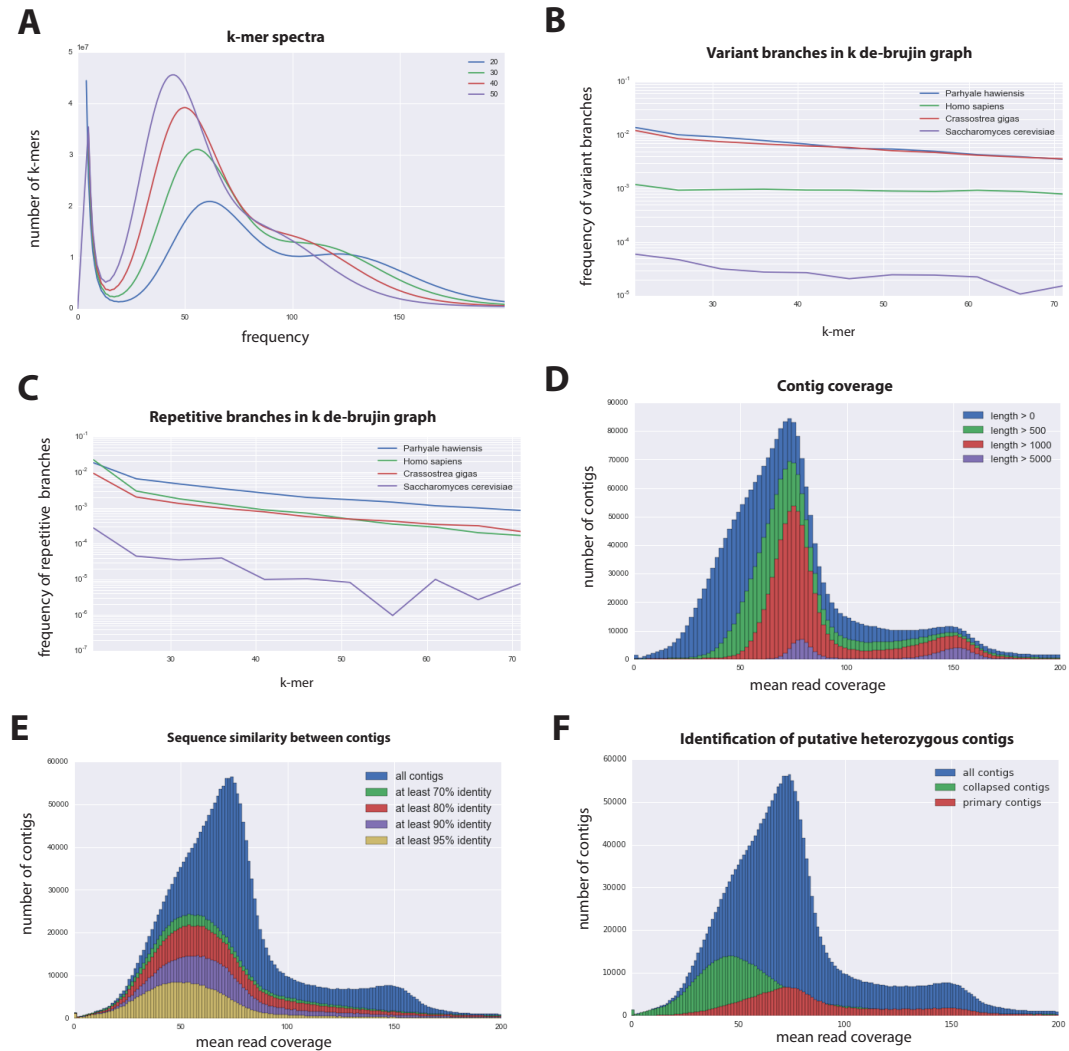


Figure 3. Assembly metrics. (A) K-mer frequency spectra of all reads for k-length from 20 to 50. (B) K-mer branching analysis performed with String Graph Assembler’s pre-qc module showing the frequency of k-mer branches classified as variants compared to *Homo sapiens*, *Crassostrea gigas*, and *Saccharomyces cerevisiae*. (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.

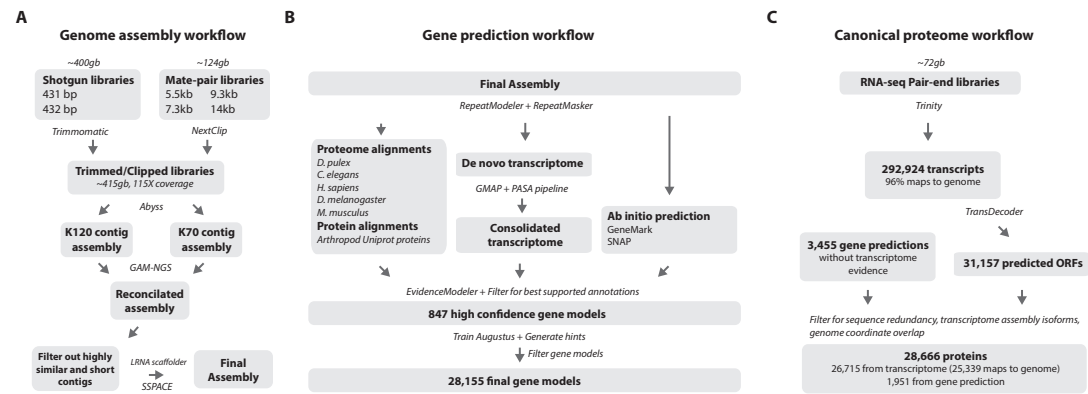


Figure 4. Workflows of assembly, annotation, and proteome generation. (A) Flowchart of the genome assembly. Two shotgun libraries and four mate-pair libraries with the indicated average sizes were prepared from a single male animal and sequenced at a 115x coverage after read filtering. Contigs were assembled 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.

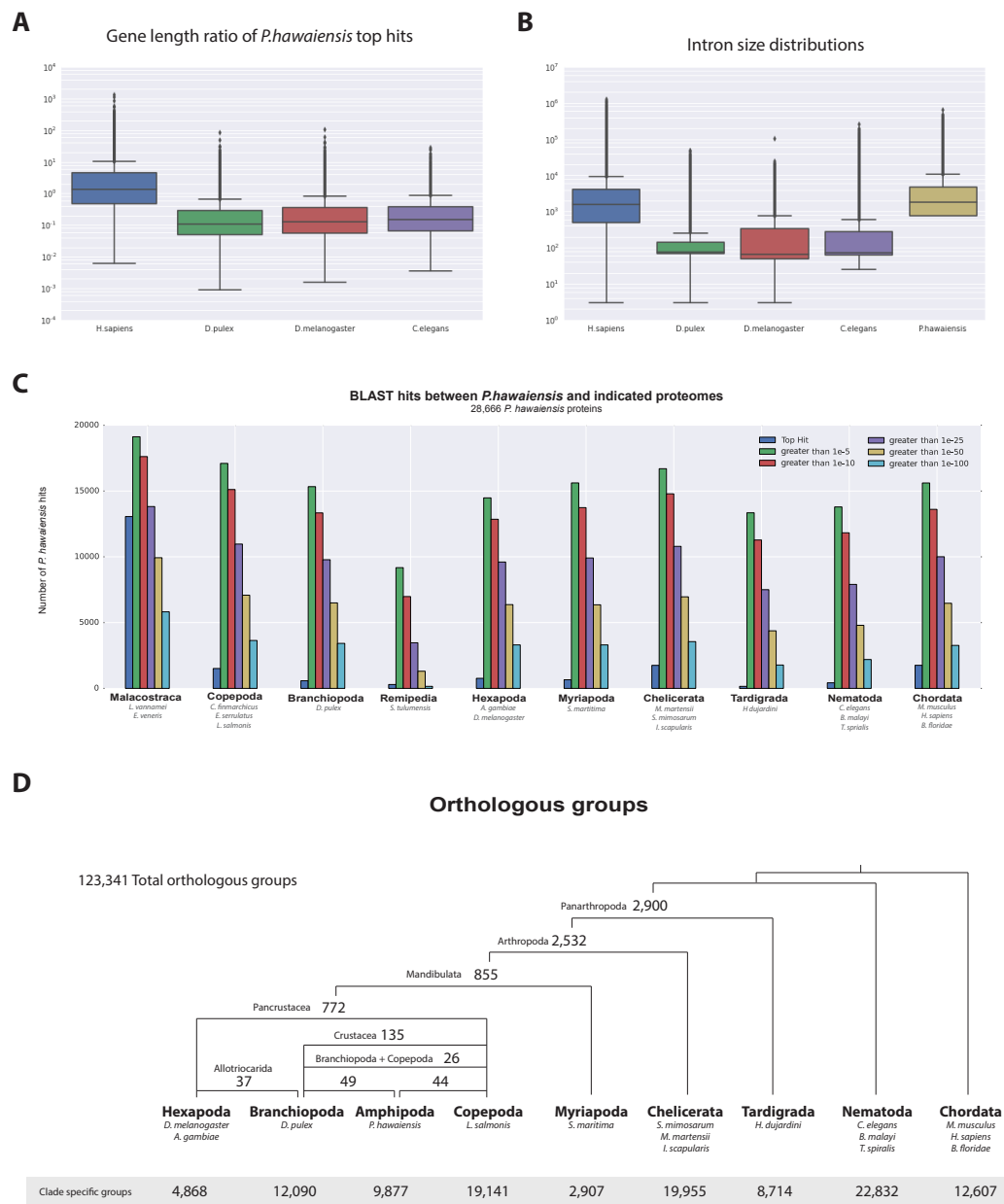


Figure 5. *Parhyale* genome comparisons. (A) Box plots comparing gene size between *Parhyale* and Human (*H. sapiens*), water fleas (*D. pulex*), flies (*D. melanogaster*) and nematodes (*C. elegans*). Ratios were calculated by dividing the size of the top blast hits in each species with the corresponding *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).

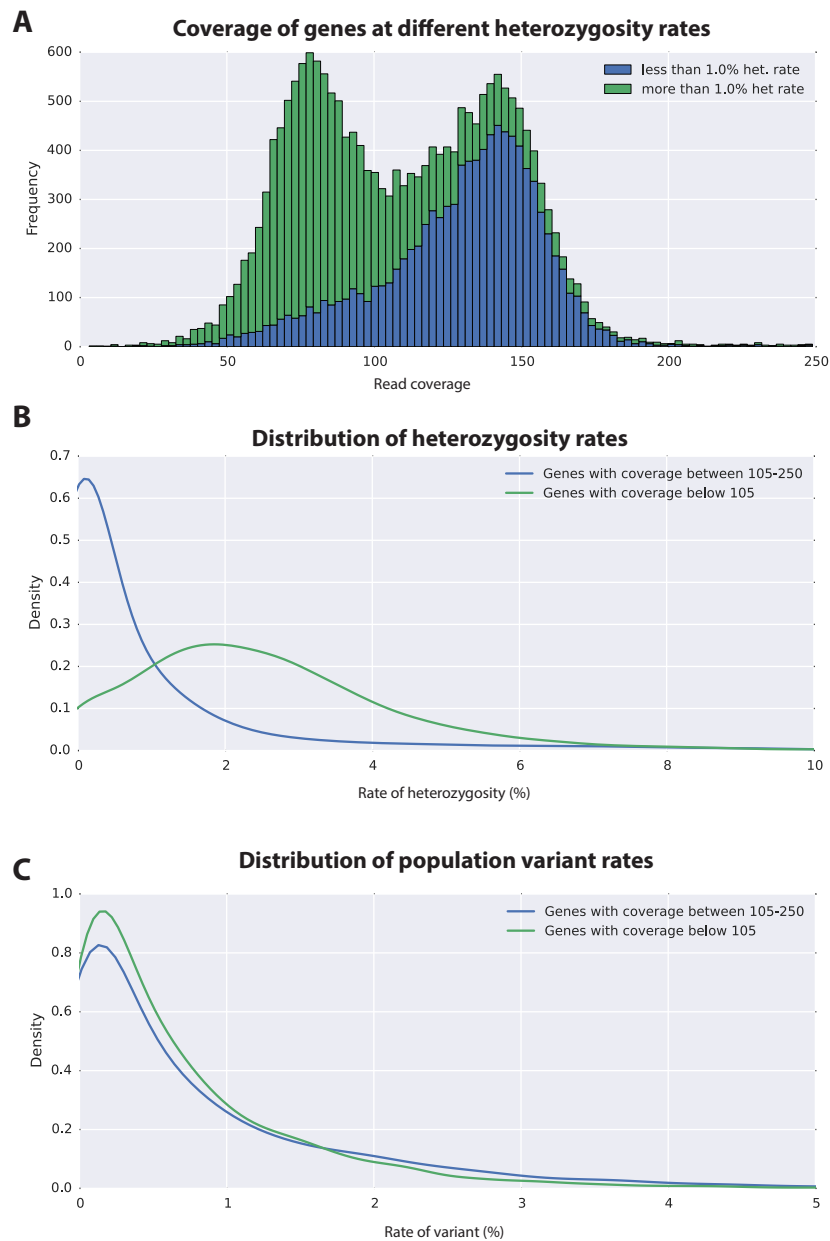


Figure 6. Variation analyses in gene coding regions. (A) A read coverage histogram of the gene predictions. Reads were first mapped to the genome, then coverage of the loci defined by the gene predictions were extracted to calculate mean coverage values. (B) Distribution plot shows that genes in the 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 both 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.	
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 >= 1bp	2	1	5	1	1	0	0	6		

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

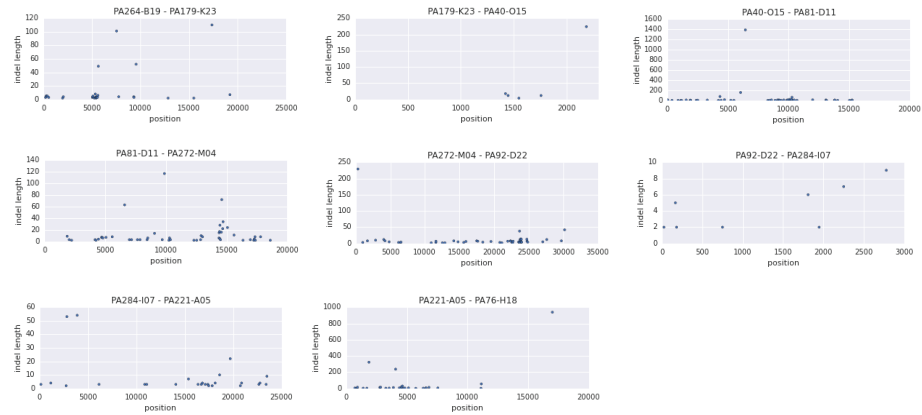


Figure 7. Variation observed in contiguous BAC sequences. (A) Schematic diagram of the contiguous BAC clones and their % sequence identities. “Overlap length” refers to the lengths (bp) of the overlapping regions between two BAC clones. “BAC supported single nucleotide polymorphisms (SNPs)” refer to 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. (B) Position versus indel lengths across each overlapping BAC region.

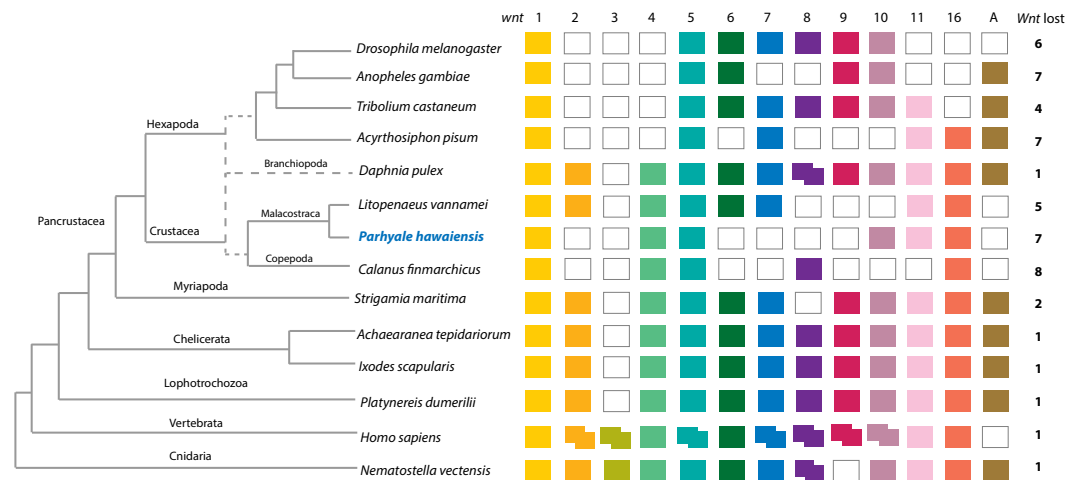


Figure 8. Comparison of Wnt family members across Metazoa. Comparison of Wnt family members across Metazoa. Tree on the left illustrates the phylogenetic relationships of species used. Dotted lines in the phylogenetic tree illustrate the alternative hypothesis of Branchiopoda + Hexapoda versus Branchiopoda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

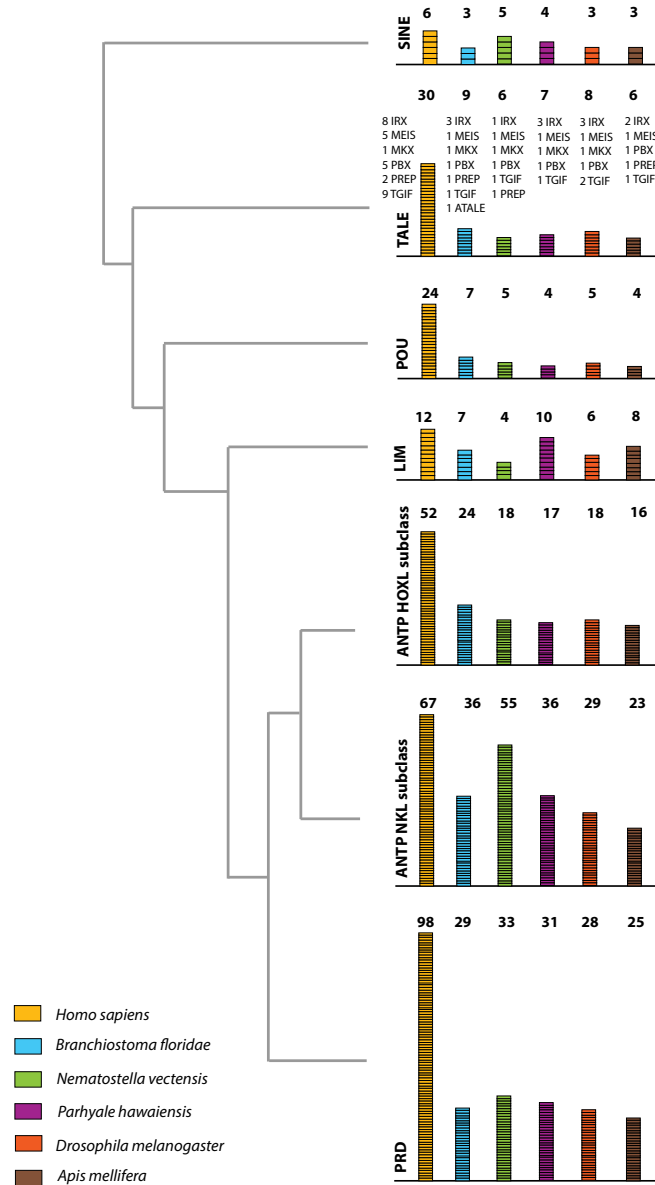


Figure 9. The overview of homeodomain radiation and phylogenetic relationships among homeodomain proteins from Arthropoda (*P. hawaiiensis*, *D. melanogaster* and *A. mellifera*) Chordata (*H. sapiens* and *B. floridae*) Cnidaria (*N. vectensis*). Six major homeodomain classes are illustrated (SINE,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

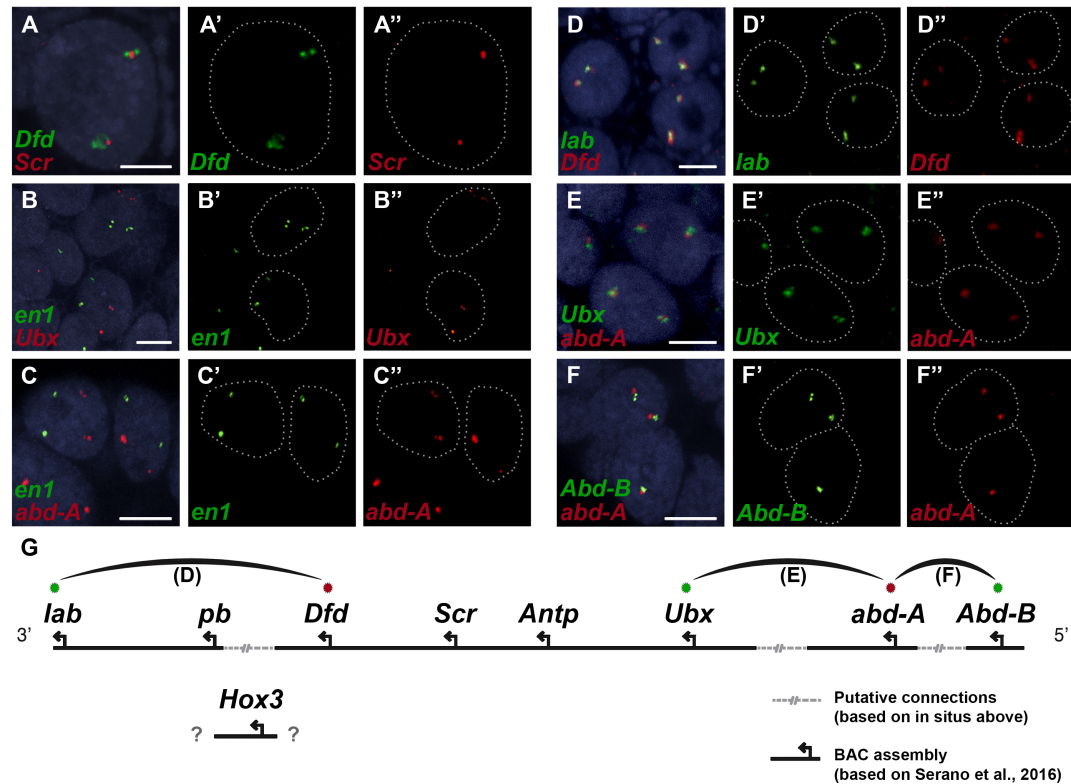


Figure 10. Variation observed in contiguous BAC sequences. Shown in panels A-F are nascent transcripts for pairs of genes expressed in *Parhyale* embryos as detected by two-color fluorescent in situ hybridization. Cell nuclei are stained with DAPI (blue) or outlined with a white dotted line. Panel A shows the co-localization of Hox sequences Dfd (A, A') and Scr (A, A''), which were previously known to be adjacent to one another from BAC assembly and sequencing. Thus, the adjacent positioning of nascent Dfd and Scr transcripts shown here serves as a positive control. Panel B and C show negative controls in which nascent transcript of Ubx are not located near those of engrailed1 (en1) (B, B', B''), and nascent transcript of abd-A are also not located near those of en1 (C, C', C''). Panels D – F show the co-localization of Hox genes not previously connected together by BAC data – lab (D, D') with Dfd (D, D''), Ubx (E, E') with abd-A (E, E''), and Abd-B (F, F') with abd-A (F, F''), establishing their proximate positioning on the same chromosome. Panel G shows a schematic of the predicted configuration of the Hox complex in *Parhyale*. Previously known genomic assembly is represented by the solid black lines, whereas linkages established by in situ results (described above) are shown as arcs. The relative orientation and order of all the genes is not known with certainty, and remains to be confirmed, but our data is consistent with the collinear orientation depicted here. Scale bars, 5 μ m (A-F).

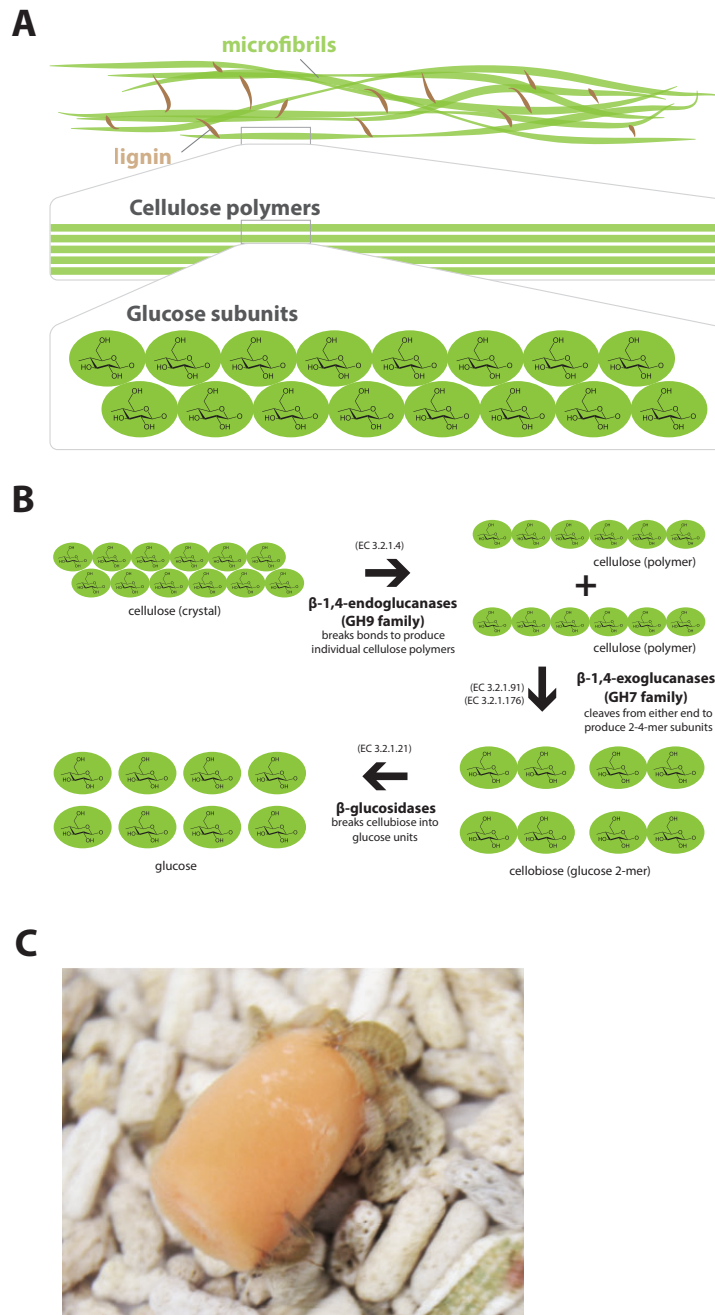


Figure 11. Phylogenetic analysis of GH7 and GH9 family proteins. (A) Structure of lignocellulosic biomass showing carbohydrate polymers and sugar monomers. **(B)** Schematic drawing illustrating mechanisms of lignocellulose degradation involving glycosyl hydrolases: β -1,4-endoglucanases, β -1,4-exoglucanases and β -glucosidases. **(C)** *Parhyale* feeding on carrots.

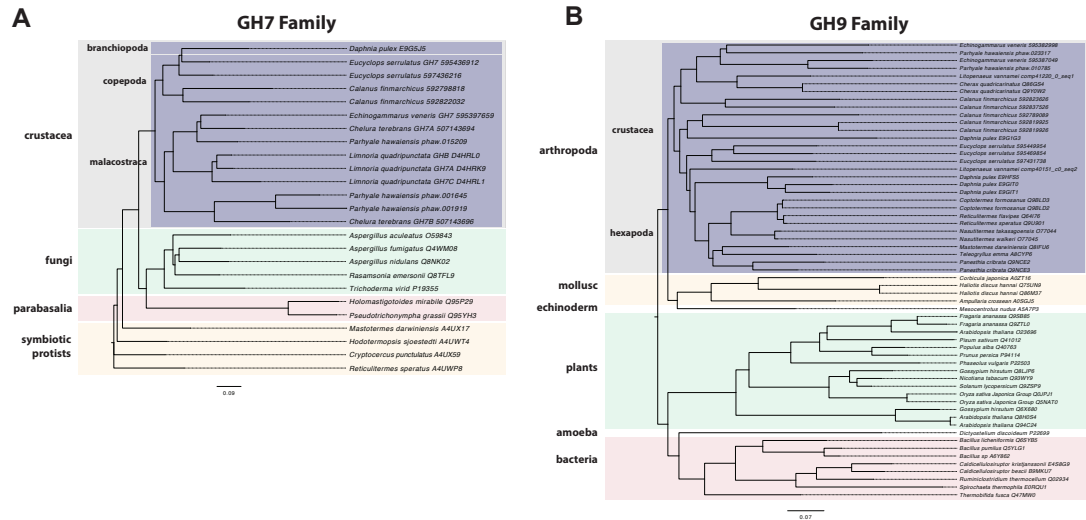


Figure 12. Phylogenetic analysis of GH7 and GH9 family proteins. (A) Phylogenetic tree showing the relationship between GH7 family proteins of *Parhyale*, other crustaceans from Vericrustacea (Malacostraca, Branchiopoda, Copepoda), fungi and symbiotic protists (root). UniProt and GenBank accessions are 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.

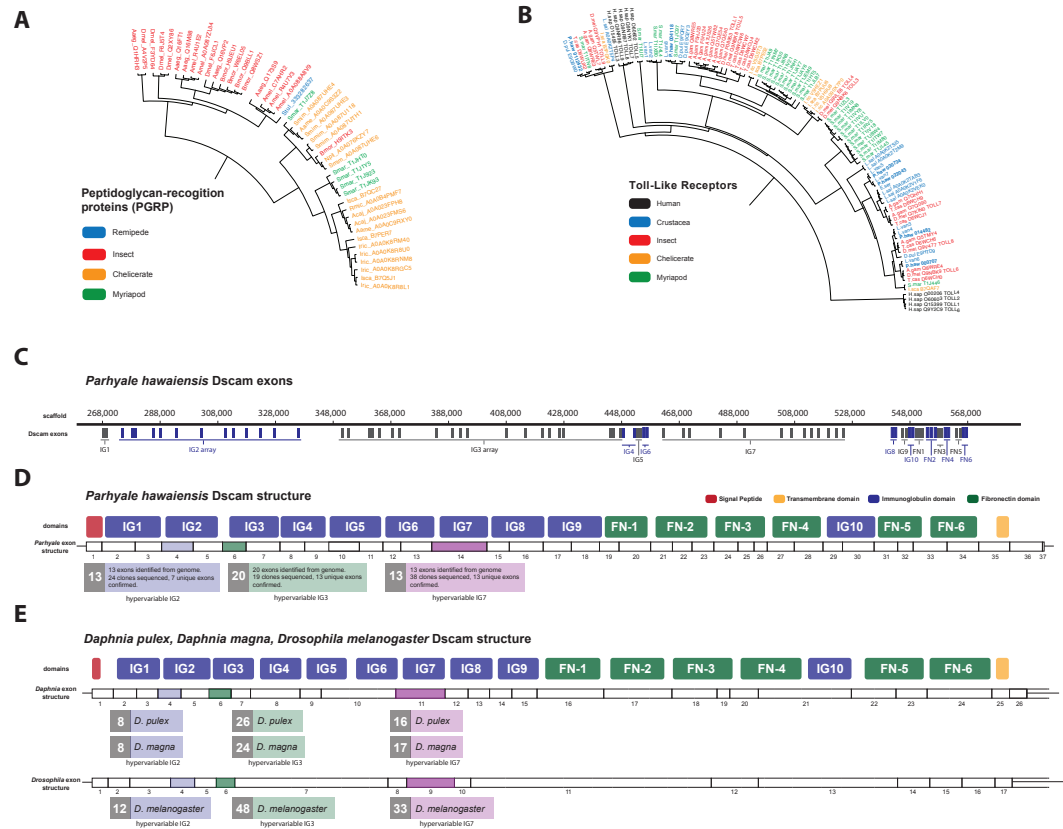


Figure 13. Peptidoglycan recognition proteins (PGRPs) and Toll-like receptors (TLRs) phylogeny. (A) Phylogenetic tree of peptidoglycan recognition proteins (PGRPs). With the exception of remipedes, PGRPs were not found in crustaceans. PGRPs have been found in the rest arthropods, including insects, myriapods 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.

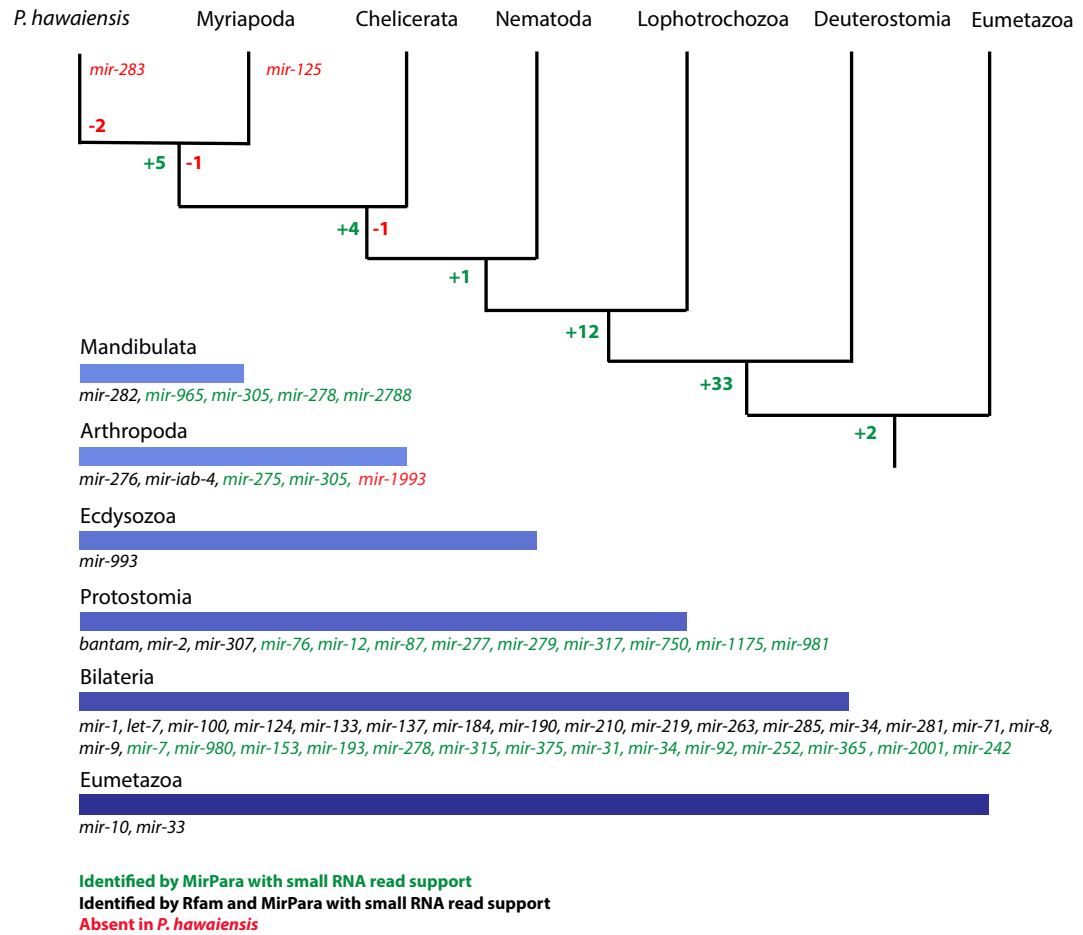


Figure 14. Conserved miRNA families found in *Parhyale* consistent with miRNAs identified in Eumetazoa, Bilateria, Protostomia, Ecdysozoa, Arthropoda and Mandibulata. miRNAs marked in red were not found in *Parhyale*. miRNAs marked in green were identified by MirPara with small RNA sequencing read support. miRNAs marked in black were identified by Rfam and MirPara with small RNA sequencing read support.

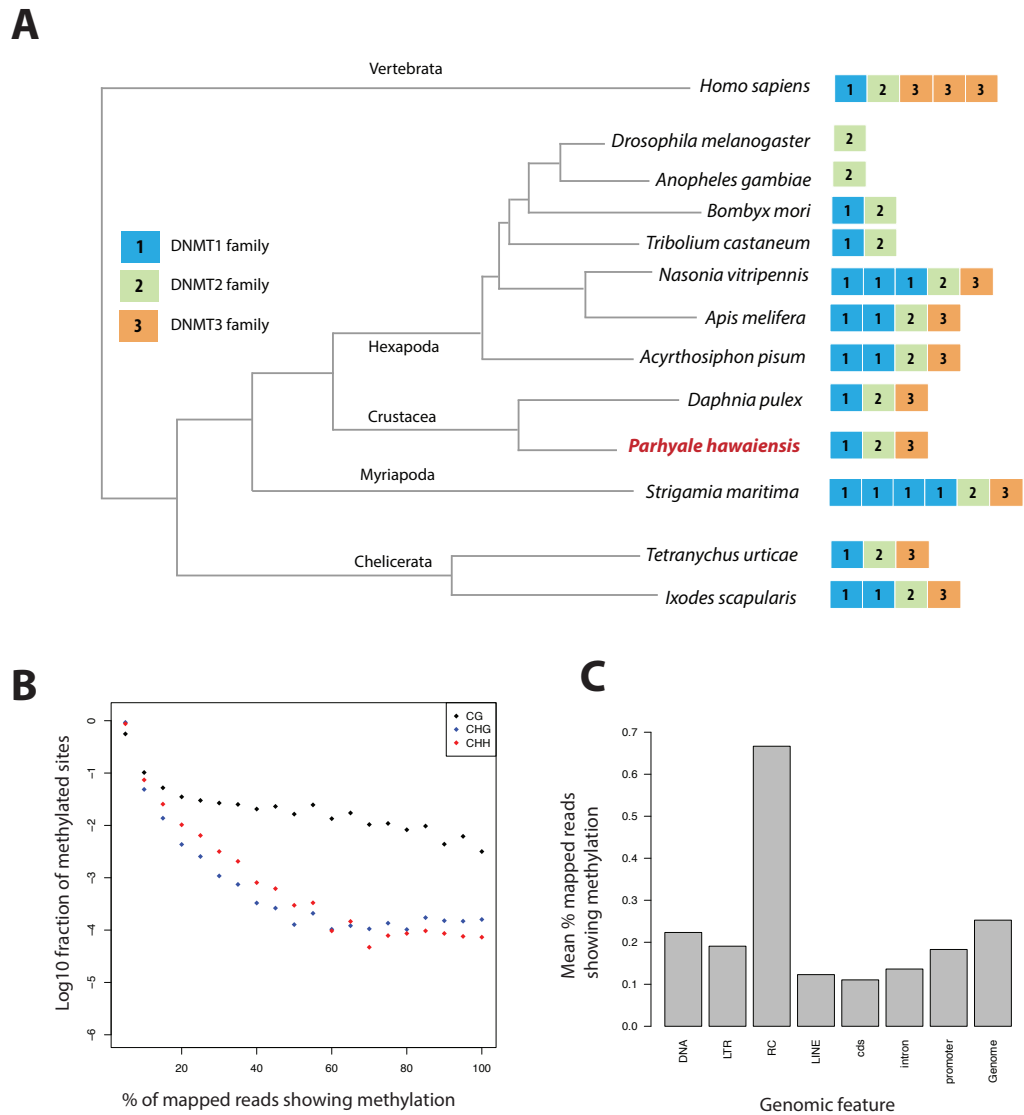


Figure 15. Comparison of *Parhyale* DNMT proteins with other arthropods and *H. sapiens*. (A) Comparison of *Parhyale* DNMT proteins with other arthropods and *H. sapiens*. Tree on the left illustrates the phylogenetic relationships of species used. Colour boxes indicate the presence of a particular DNMT subfamily for a given species. Paralogs of DNMT are indicated accordingly. (B) Amount of methylation is presented as percentage of reads showing methylation at a site in which CpG sites showed preferential methylation. (C) Methylation of various genomic features: DNA transposons (DNA), long terminal repeats (LTR), rolling circle (RC) repeats, long interspersed element (LINE) transposons, coding sequence, introns, promoters and the rest of the genome.

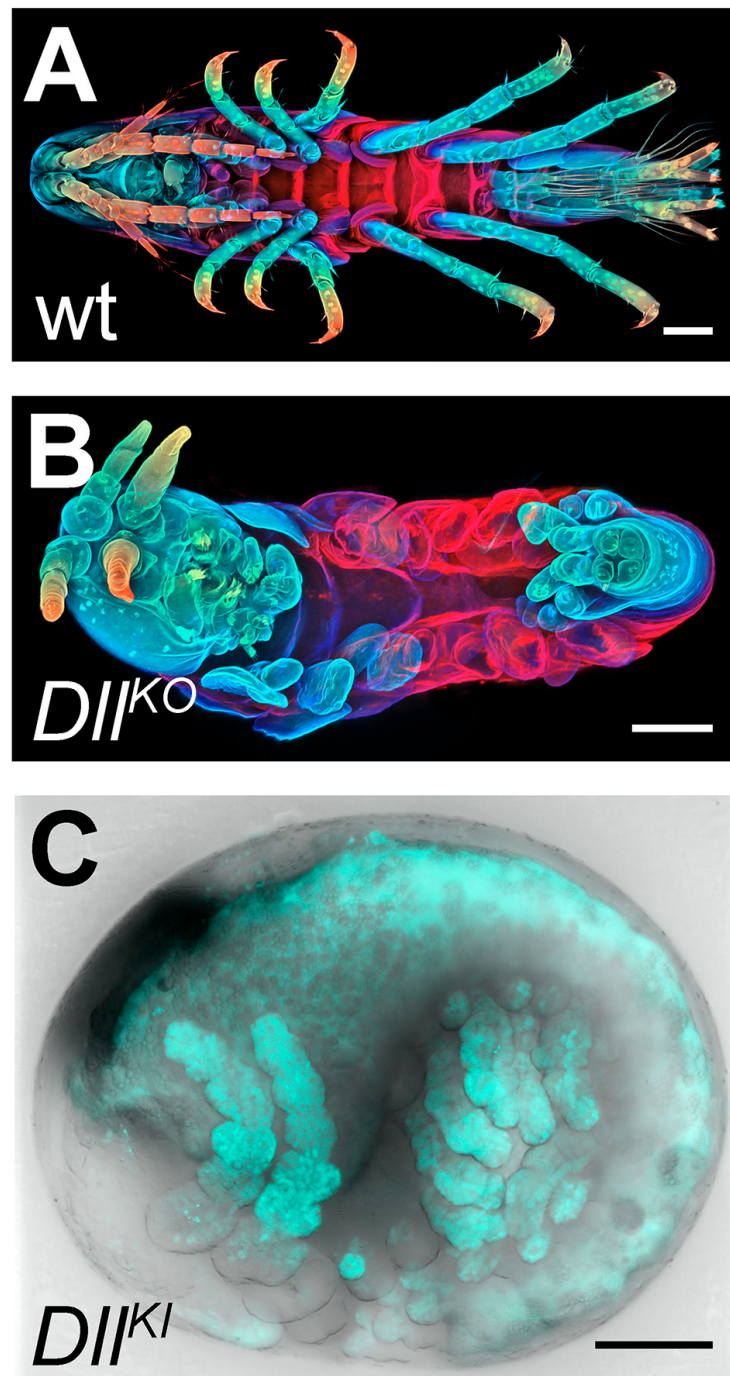


Figure 16. *Parhyale hawaiiensis*, an emerging model system in developmental genetics research. (A) Wild-type morphology and (B) mutant *Parhyale* with truncated limbs after CRISPR-mediated knock-out of the limb patterning gene Distal-less (DII). Panels show ventral views of juveniles stained for cuticle and color-coded by depth with anterior to the left. (C) Fluorescent tagging of DII expressed in most limbs (shown in cyan) by CRISPR-mediated knock-in 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 μm.

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
		phaw_30_tra_m.014329
		phaw_30_tra_m.012297
rm62	7	phaw_30_tra_m.004444
		phaw_30_tra_m.012605
		phaw_30_tra_m.001849
		phaw_30_tra_m.006468
Piwi/aubergine	2	phaw_30_tra_m.023485
		phaw_30_tra_m.011247
Dicer 1	1	phaw_30_tra_m.016012
Dicer 2	1	phaw_30_tra_m.001257
argonaute 1	1	phaw_30_tra_m.021619
		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

897 REFERENCES

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

- 930 Justen Andrews, Teresa J Crease, Haixu Tang, Susan M Lucas, Hugh M Robertson, Peer Bork,
931 Eugene V Koonin, Evgeny M Zdobnov, Igor V Grigoriev, Michael Lynch, and Jeffrey L Boore. The
932 Ecoresponsive Genome of *Daphnia pulex*. *Science*, 331(6017):555–561, 2011.
- 933 [11] K Meusemann, B M von Reumont, S Simon, F Roeding, S Strauss, P Kuck, I Ebersberger, M Walz,
934 G Pass, S Breuers, V Achter, A von Haeseler, T Burmester, H Hadrys, J W Wagele, and B Misof. A
935 Phylogenomic Approach to Resolve the Arthropod Tree of Life. *Molecular biology and evolution*,
936 27(11):2451–2464, October 2010.
- 937 [12] Jerome C Regier, Jeffrey W Shultz, Andreas Zwick, April Hussey, Bernard Ball, Regina Wetzer,
938 Joel W Martin, and Clifford W Cunningham. Arthropod relationships revealed by phylogenomic
939 analysis of nuclear protein-coding sequences. *Nature*, 463(7284):1079–1083, February 2010.
- 940 [13] T H Oakley, J M Wolfe, A R Lindgren, and A K Zaharoff. Phylotranscriptomics to Bring the Under-
941 studied into the Fold: Monophyletic Ostracoda, Fossil Placement, and Pancrustacean Phylogeny.
942 *Molecular biology and evolution*, 30(1):215–233, December 2012.
- 943 [14] Bjoern M von Reumont, Ronald A Jenner, Matthew A Wills, Emiliano Dell’ampio, Günther Pass,
944 Ingo Ebersberger, Benjamin Meyer, Stefan Koenemann, Thomas M Iliffe, Alexandros Stamatakis,
945 Oliver Niehuis, Karen Meusemann, and Bernhard Misof. Pancrustacean phylogeny in the light of
946 new phylogenomic data: support for Remipedia as the possible sister group of Hexapoda. *Molecular
947 biology and evolution*, 29(3):1031–1045, March 2012.
- 948 [15] Le Cong, F Ann Ran, David Cox, Shuailiang Lin, Robert Barretto, Naomi Habib, Patrick D Hsu,
949 Xuebing Wu, Wenyan Jiang, Luciano A Marraffini, and Feng Zhang. Multiplex genome engineering
950 using CRISPR/Cas systems. *Science*, 339(6121):819–823, February 2013.
- 951 [16] Julia M Serano, Arnaud Martin, Danielle M Liubicich, Erin Jarvis, Heather S Bruce, Konnor La,
952 William E Browne, Jane Grimwood, and Nipam H Patel. Comprehensive analysis of Hox gene
953 expression in the amphipod crustacean *Parhyale hawaiiensis*. *Developmental Biology*, pages 1–13,
954 November 2015.
- 955 [17] Arnaud Martin, Julia M Serano, Erin Jarvis, Heather S Bruce, Jennifer Wang, Shagnik Ray, Carryn A
956 Barker, Liam C O’Connell, and Nipam H Patel. CRISPR/Cas9 Mutagenesis Reveals Versatile Roles
957 of Hox Genes in Crustacean Limb Specification and Evolution. *Current biology : CB*, December
958 2015.
- 959 [18] Prashant Mali, Luhan Yang, Kevin M Esvelt, John Aach, Marc Guell, James E DiCarlo, Julie E
960 Norville, and George M Church. RNA-guided human genome engineering via Cas9. *Science*,
961 339(6121):823–826, February 2013.

- 962 [19] Martin Jinek, Krzysztof Chylinski, Ines Fonfara, Michael Hauer, Jennifer A Doudna, and Em-
963 manuelle Charpentier. A programmable dual-RNA-guided DNA endonuclease in adaptive bacterial
964 immunity. *Science*, 337(6096):816–821, August 2012.
- 965 [20] Anna F Gilles and Michalis Averof. Functional genetics for all: engineered nucleases, CRISPR and
966 the gene editing revolution. *EvoDevo*, 5(1):43–13, 2014.
- 967 [21] M Averof and N H Patel. Crustacean appendage evolution associated with changes in Hox gene
968 expression. *Nature*, 388(6643):682–686, 1997.
- 969 [22] Danielle M Liubicich, Julia M Serano, Anastasios Pavlopoulos, Zacharias Kontarakis, Meredith E
970 Protas, Elaine Kwan, Sandip Chatterjee, Khoa D Tran, Michalis Averof, and Nipam H Patel.
971 Knockdown of Parhyale Ultrabithorax recapitulates evolutionary changes in crustacean appendage
972 morphology. *Proceedings of the National Academy of Sciences of the United States of America*,
973 106(33):13892–13896, August 2009.
- 974 [23] Anastasios Pavlopoulos, Zacharias Kontarakis, Danielle M Liubicich, Julia M Serano, Michael
975 Akam, Nipam H Patel, and Michalis Averof. Probing the evolution of appendage specialization by
976 Hox gene misexpression in an emerging model crustacean. *Proceedings of the National Academy of
977 Sciences of the United States of America*, 106(33):13897–13902, August 2009.
- 978 [24] Nikolaos Konstantinides and Michalis Averof. A common cellular basis for muscle regeneration in
979 arthropods and vertebrates. *Science*, 343(6172):788–791, February 2014.
- 980 [25] Jeanne L Benton, Rachel Kery, Jingjing Li, Chadanat Noonin, Irene Söderhäll, and Barbara S Beltz.
981 Cells from the Immune System Generate Adult-Born Neurons in Crayfish. 30(3):322–333, August
982 2014.
- 983 [26] L Vazquez, J Alpuche, G Maldonado, C Agundis, A Pereyra-Morales, and E Zenteno. Review:
984 Immunity mechanisms in crustaceans. *Innate Immunity*, 15(3):179–188, May 2009.
- 985 [27] Chris Hauton. The scope of the crustacean immune system for disease control. *Journal of Inverte-
986 brate Pathology*, 110(2):251–260, June 2012.
- 987 [28] T L Maginnis. The costs of autotomy and regeneration in animals: a review and framework for
988 future research. *Behavioral Ecology*, 17(5):857–872, June 2006.
- 989 [29] Sunetra Das and David S Durica. Ecdysteroid receptor signaling disruption obstructs blastemal cell
990 proliferation during limb regeneration in the fiddler crab, *Uca pugilator*. *Molecular and cellular
991 endocrinology*, 365(2):249–259, January 2013.
- 992 [30] Andrew J King, Simon M Cragg, Yi Li, Jo Dymond, Matthew J Guille, Dianna J Bowles, Neil C
993 Bruce, Ian A Graham, and Simon J McQueen-Mason. Molecular insight into lignocellulose digestion

- 994 by a marine isopod in the absence of gut microbes. *Proceedings of the National Academy of Sciences*,
995 107(12):5345–5350, March 2010.
- 996 [31] M Kern, J E McGeehan, and S D Streeter. Structural characterization of a unique marine animal
997 family 7 cellobiohydrolase suggests a mechanism of cellulase salt tolerance. In *Proceedings of the*
998 *...*, 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, Xuanting 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 Karagogeos. 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 Gorojankina, 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] R D Finn, J Mistry, and B Schuster-Böckler. Pfam: clans, web tools and services. *Nucleic acids . . .*,
1208 2006.
- 1209 [92] Simon M Cragg, Gregg T Beckham, Neil C Bruce, Timothy D H Bugg, Daniel L Distel, Paul Dupree,
1210 Amaia Green Etxabe, Barry S Goodell, Jody Jellison, John E McGeehan, Simon J McQueen-Mason,
1211 Kirk Schnorr, Paul H Walton, Joy E M Watts, and Martin Zimmer. Lignocellulose degradation
1212 mechanisms across the Tree of Life. *Current Opinion in Chemical Biology*, 29:108–119, December
1213 2015.
- 1214 [93] G D Stentiford, D M Neil, E J Peeler, J D Shields, H J Small, T W Flegel, J M Vlak, B Jones,
1215 F Morado, S Moss, J Lotz, L Bartholomay, D C Behringer, C Hauton, and D V Lightner. Disease
1216 will limit future food supply from the global crustacean fishery and aquaculture sectors. *Journal of*
1217 *Invertebrate Pathology*, 110(2):141–157, June 2012.
- 1218 [94] Robert M Waterhouse, Evgenia V Kriventseva, Stephan Meister, Zhiyong Xi, Kanwal S Alvarez,
1219 Lyric C Bartholomay, Carolina Barillas-Mury, Guowu Bian, Stephanie Blandin, Bruce M Chris-
1220 tensen, Yuemei Dong, Haobo Jiang, Michael R Kanost, Anastasios C Koutsos, Elena A Levashina,
1221 Jianyong Li, Petros Ligoxygakis, Robert M Maccallum, George F Mayhew, Antonio Mendes, Kristin
1222 Michel, Mike A Osta, Susan Paskewitz, Sang Woon Shin, Dina Vlachou, Lihui Wang, Weiqi Wei,
1223 Liangbiao Zheng, Zhen Zou, David W Severson, Alexander S Raikhel, Fotis C Kafatos, George
1224 Dimopoulos, Evgeny M Zdobnov, and George K Christophides. Evolutionary dynamics of immune-
1225 related genes and pathways in disease-vector mosquitoes. *Science*, 316(5832):1738–1743, June
1226 2007.

- 1227 [95] Charles A Janeway and Ruslan Medzhitov. Innate immune recognition. *Annual review of immunol-*
1228 *ogy*, 20:197–216, 2002.
- 1229 [96] T Werner, K Borge-Renberg, P Mellroth, H Steiner, and D Hultmark. Functional Diversity of
1230 the *Drosophila* PGRP-LC Gene Cluster in the Response to Lipopolysaccharide and Peptidoglycan.
1231 *Journal of Biological Chemistry*, 278(29):26319–26322, July 2003.
- 1232 [97] C Liu, Z Xu, D Gupta, and R Dziarski. Peptidoglycan Recognition Proteins: A novel family
1233 of four human innate immunity pattern recognition molecules. *Journal of Biological Chemistry*,
1234 276(37):34686–34694, September 2001.
- 1235 [98] Abdur Rehman, Ping Taishi, Jidong Fang, Jeannine A Majde, and James M Krueger. The cloning
1236 of a rat peptidoglycan recognition protein (PGRP) and its induction in brain by sleep deprivation.
1237 *Cytokine*, 13(1):8–17, January 2001.
- 1238 [99] Haipeng Liu, Chenglin Wu, Yasuyuki Matsuda, Shun-ichiro Kawabata, Bok Luel Lee, Kenneth
1239 Söderhäll, and Irene Söderhäll. Peptidoglycan activation of the proPO-system without a peptidogly-
1240 can receptor protein (PGRP)? *Developmental & Comparative Immunology*, 35(1):51–61, January
1241 2011.
- 1242 [100] Seanna J McTaggart, Claire Conlon, John K Colbourne, Mark L Blaxter, and Tom J Little. The
1243 components of the *Daphnia pulex* immune system as revealed by complete genome sequencing.
1244 *BMC Genomics*, 10(1):175–119, 2009.
- 1245 [101] Catherine Dostert, Emmanuelle Jouanguy, Phil Irving, Laurent Troxler, Delphine Galiana-Arnoux,
1246 Charles Hetru, Jules A Hoffmann, and Jean-Luc Imler. The Jak-STAT signaling pathway is required
1247 but not sufficient for the antiviral response of *drosophila*. *Nature Immunology*, 6(9):946–953, August
1248 2005.
- 1249 [102] T Tanji, X Hu, A N R Weber, and Y T Ip. Toll and IMD Pathways Synergistically Activate an Innate
1250 Immune Response in *Drosophila melanogaster*. *Molecular and cellular biology*, 27(12):4578–4588,
1251 May 2007.
- 1252 [103] Matthew A. Benton, Matthias Pechmann, Nadine Frey, Dominik Stappert, Kai H. Conrads, Yen-
1253 Ta Chen, Evangelia Stamataki, Anastasios Pavlopoulos, and Siegfried Roth. Toll genes have an
1254 ancestral role in axis elongation. *Current Biology*, 26(12):1609 – 1615, 2016.
- 1255 [104] Natalia I Arbouzova and Martin P Zeidler. JAK/STAT signalling in *Drosophila*: insights into
1256 conserved regulatory and cellular functions. *Development*, 133(14):2605–2616, July 2006.
- 1257 [105] E A Levashina, L F Moita, S Blandin, G Vriend, M Lagueux, and F C Kafatos. Conserved role of
1258 a complement-like protein in phagocytosis revealed by dsRNA knockout in cultured cells of the
1259 mosquito, *Anopheles gambiae*. *Cell*, 104(5):709–718, 2001.

- 1260 [106] H Decker. Recent findings on phenoloxidase activity and antimicrobial activity of hemocyanins.
1261 *Developmental & Comparative Immunology*, 28(7-8):673–687, June 2004.
- 1262 [107] So Young Lee, Bok Luel Lee, and Kenneth Söderhäll. Processing of crayfish hemocyanin subunits
1263 into phenoloxidase. *Biochemical and Biophysical Research Communications*, 322(2):490–496,
1264 September 2004.
- 1265 [108] D Schmucker, J C Clemens, H Shu, C A Worby, J Xiao, M Muda, J E Dixon, and S L Zipursky.
1266 *Drosophila Dscam* is an axon guidance receptor exhibiting extraordinary molecular diversity. *Cell*,
1267 101(6):671–684, June 2000.
- 1268 [109] Fiona L Watson, Roland Püttmann-Holgado, Franziska Thomas, David L Lamar, Michael Hughes,
1269 Masahiro Kondo, Vivienne I Rebel, and Dietmar Schmucker. Extensive diversity of Ig-superfamily
1270 proteins in the immune system of insects. *Science*, 309(5742):1874–1878, September 2005.
- 1271 [110] Daniela Brites, Seanna McTaggart, Krystalynne Morris, Jobriah Anderson, Kelley Thomas, Isabelle
1272 Colson, Thomas Fabbro, Tom J Little, Dieter Ebert, and Louis Du Pasquier. The *Dscam* homologue
1273 of the crustacean *Daphnia* is diversified by alternative splicing like in insects. *Molecular biology
1274 and evolution*, 25(7):1429–1439, July 2008.
- 1275 [111] Stephane E Castel and Robert A Martienssen. RNA interference in the nucleus: roles for small
1276 RNAs in transcription, epigenetics and beyond. *Nature reviews. Genetics*, 14(2):100–112, February
1277 2013.
- 1278 [112] Alexei A Aravin, Natalia M Naumova, Alexei V Tulin, Vasilii V Vagin, Yakov M Rozovsky,
1279 and Vladimir A Gvozdev. Double-stranded RNA-mediated silencing of genomic tandem repeats
1280 and transposable elements in the *D. melanogaster* germline Alexei A. Aravin*. *Current Biology*,
1281 11(13):1–11, July 2001.
- 1282 [113] N J Caplen, S Parrish, F Imani, A Fire, and R A Morgan. Specific inhibition of gene expression by
1283 small double-stranded RNAs in invertebrate and vertebrate systems. *Proceedings of the National
1284 Academy of Sciences of the United States of America*, 98(17):1–7, August 2001.
- 1285 [114] Julius Brennecke, Alexei A Aravin, Alexander Stark, Monica Dus, Manolis Kellis, Ravi Sachidanandam,
1286 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 [115] Weifeng Gu, Masaki Shirayama, Darryl Conte Jr, Jessica Vasale, Pedro J Batista, Julie M Claycomb,
1289 James J Moresco, Elaine M Youngman, Jennifer Keys, Matthew J Stoltz, Chun-Chieh G Chen,
1290 Daniel A Chaves, Shenghua Duan, Kristin D Kasschau, Noah Fahlgren, John R Yates III, Shohei
1291 Mitani, James C Carrington, and Craig C Mello. Distinct Argonaute-Mediated 22G-RNA Pathways
1292 Direct Genome Surveillance in the *C. elegans* Germline. *Molecular cell*, 36(2):231–244, October
1293 2009.

- 1294 [116] Heng-Chi Lee, Weifeng Gu, Masaki Shirayama, Elaine Youngman, Darryl Conte, and Craig C
1295 Mello. *C. elegans* piRNAs mediate the genome-wide surveillance of germline transcripts. *Cell*,
1296 150(1):78–87, July 2012.
- 1297 [117] Lin He and Gregory J Hannon. MicroRNAs: small RNAs with a big role in gene regulation. *Nature*
1298 *reviews. Genetics*, 5(7):522–531, July 2004.
- 1299 [118] J Michael Thomson, Martin Newman, Joel S Parker, Elizabeth M Morin-Kensicki, Tricia Wright,
1300 and Scott M Hammond. Extensive post-transcriptional regulation of microRNAs and its implications
1301 for cancer. *Genes & development*, 20(16):2202–2207, August 2006.
- 1302 [119] Witold Filipowicz, Suvendra N Bhattacharyya, and Nahum Sonenberg. Mechanisms of post-
1303 transcriptional regulation by microRNAs: are the answers in sight? *Nature reviews. Genetics*,
1304 2008(2):102–114, February 2008.
- 1305 [120] Peter Sarkies, Murray E Selkirk, John T Jones, Vivian Blok, Thomas Boothby, Bob Goldstein,
1306 Ben Hanelt, Alex Ardila-Garcia, Naomi M Fast, Phillip M Schiffer, Christopher Kraus, Mark J
1307 Taylor, Georgios Koutsovoulos, Mark L Blaxter, and Eric A Miska. Ancient and Novel Small RNA
1308 Pathways Compensate for the Loss of piRNAs in Multiple Independent Nematode Lineages. *PLoS*
1309 *biology*, 13(2):e1002061–20, February 2015.
- 1310 [121] Ying Dong and Markus Friedrich. Nymphal RNAi: systemic RNAi mediated gene knockdown in
1311 juvenile grasshopper. *BMC Biotechnology*, 5:25, 2005.
- 1312 [122] George M Weinstock, Gene E Robinson, Richard A Gibbs, George M Weinstock, George M
1313 Weinstock, Gene E Robinson, Kim C Worley, Hugh M Robertson, Daniel B Weaver, Martin Beye,
1314 Peer Bork, Jay D Evans, Klaus Hartfelder, Greg J Hunt, Gene E Robinson, Ryszard Maleszka,
1315 George M Weinstock, Klaus Hartfelder, Gro V Amdam, Mrcia M G Bitondi, Anita M Collins,
1316 Alexandre S Cristino, H Michael, G Lattorff, Carlos H Lobo, Robin F A Moritz, Francis M F Nunes,
1317 Robert E Page, Zil L P Simões, Diana Wheeler, Piero Carninci, Shiro Fukuda, Yoshihide Hayashizaki,
1318 Chikatoshi Kai, Jun Kawai, Naoko Sakazume, Daisuke Sasaki, Michihira Tagami, Gro V Amdam,
1319 Stefan Albert, Geert Baggerman, Kyle T Beggs, Guy Bloch, Giuseppe Cazzamali, Mira Cohen,
1320 Mark David Drapeau, Dorothea Eisenhardt, Christine Emore, Michael A Ewing, Susan E Fahrbach,
1321 Sylvain Foret, Cornelis J P Grimmelikhuijzen, Frank Hauser, Amanda B Hummon, Greg J Hunt,
1322 Jurgen Huybrechts, Andrew K Jones, Noam Kaplan, Gérard Lebouille, Michal Linial, J Troy
1323 Littleton, Alison R Mercer, Robert E Page, Gene E Robinson, Timothy A Richmond, Sandra L
1324 RodriguezZas, Elad B Rubin, David B Sattelle, David Schlipalius, Liliane Schoofs, Yair Shemesh,
1325 Jonathan V Sweedler, Rodrigo Velarde, Peter Verleyen, Evy Vierstraete, Michael R Williamson,
1326 Martin Beye, Seth A Ament, Susan J Brown, Miguel Corona, Peter K Dearden, W Augustine
1327 Dunn, Michelle M Elekonich, Christine G Elsik, Tomoko Fujiyuki, Irene Gattermeier, Tanja Gempe,
1328 Martin Hasselmann, Tatsuhiko Kadowaki, Eriko Kage, Azusa Kamikouchi, Takeo Kubo, Robert

- 1329 Kucharski, Takekazu Kunieda, Marcé Lorenzen, Natalia V Milshina, Mizue Morioka, Kazuaki
1330 Ohashi, Ross Overbeek, Robert E Page, Gene E Robinson, Christian A Ross, Morten Schioett, Teresa
1331 Shippy, Hideaki Takeuchi, Amy L Toth, Judith H Willis, Megan J Wilson, Evgeny M Zdobnov,
1332 Karl H J Gordon, Ivica Letunic, Kevin Hackett, Jane Peterson, Adam Felsenfeld, Mark Guyer,
1333 Michel Solignac, Richa Agarwala, Jean Marie Cornuet, Christine Emore, Greg J Hunt, Monique
1334 Monnerot, Florence Mougél, Justin T Reese, David Schlipalius, Dominique Vautrin, Daniel B
1335 Weaver, Joseph J Gillespie, Jamie J Cannone, Robin R Gutell, J Spencer Johnston, Michael B
1336 Eisen, Amanda B Hummon, Venky N Iyer, Vivek Iyer, Peter Kosarev, Aaron J Mackey, Timothy A
1337 Richmond, Victor Solovyev, Alexandre Souvorov, George M Weinstock, Michael R Williamson,
1338 Katherine A Aronstein, Katarina Bilikova, Yan Ping Chen, Andrew G Clark, Laura I Decanini,
1339 William M Gelbart, Charles Hetru, Dan Hultmark, Jean-Luc Imler, Haobo Jiang, Michael Kanost,
1340 Kiyoshi Kimura, Brian P Lazzaro, Dawn L Lopez, Jozef Simuth, Graham J Thompson, Zhen Zou,
1341 Pieter De Jong, Erica Sodergren, Miklós Csűrös, Aleksandar Milosavljevic, J Spencer Johnston,
1342 Kazutoyo Osoegawa, Stephen Richards, Chung-Li Shu, George M Weinstock, Laurent Duret, Eran
1343 Elhaik, Dan Graur, Daniel B Weaver, Gro V Amdam, Juan M Anzola, Kathryn S Campbell, Kevin L
1344 Childs, Derek Collinge, Madeline A Crosby, C Michael Dickens, Karl H J Gordon, L Sian Gramates,
1345 Christina M Grozinger, Peter L Jones, Mireia Jorda, Xu Ling, Beverly B Matthews, Jonathan Miller,
1346 Natalia V Milshina, Craig Mizzen, Miguel A Peinado, Jeffrey G Reid, Gene E Robinson, Susan M
1347 Russo, Andrew J Schroeder, Susan E St Pierre, Ying Wang, Pinglei Zhou, Richa Agarwala, Natalia V
1348 Milshina, Daniel B Weaver, Kevin L Childs, C Michael Dickens, William M Gelbart, Huaiyang Jiang,
1349 Paul Kitts, Natalia V Milshina, Barbara Ruef, Susan M Russo, Anand Venkatraman, George M
1350 Weinstock, Lan Zhang, Pinglei Zhou, J Spencer Johnston, Gildardo Aquino-Perez, Jean Marie
1351 Cornuet, Monique Monnerot, Michel Solignac, Dominique Vautrin, Charles W Whitfield, Susanta K
1352 Behura, Stewart H Berlocher, Andrew G Clark, J Spencer Johnston, Walter S Sheppard, Deborah R
1353 Smith, Andrew V Suarez, Neil D Tsutsui, and Daniel B and... Weaver. Insights into social insects
1354 from the genome of the honeybee *Apis mellifera*. *Nature*, 443(7114):931–949, October 2006.
- 1355 [123] Weina Xu and Zhaojun Han. Cloning and phylogenetic analysis of sid-1-like genes from aphids.
1356 *Journal of insect science (Online)*, 8(30):1–6, 2008.
- 1357 [124] J Y Roignant, C Carre, R Mugat, D Szymczak, J A Lepesant, and C Antoniewski. Absence of
1358 transitive and systemic pathways allows cell-specific and isoform-specific RNAi in *Drosophila*. *RNA*,
1359 9(3):299–308, March 2003.
- 1360 [125] Yonggan Wu, Bo Wei, Haizhou Liu, Tianxian Li, and Simon Rayner. MiRPara: a SVM-based
1361 software tool for prediction of most probable microRNA coding regions in genome scale sequences.
1362 *BMC bioinformatics*, 12(1):107, 2011.
- 1363 [126] Eric P Nawrocki, Sarah W Burge, Alex Bateman, Jennifer Daub, Ruth Y Eberhardt, Sean R Eddy,
1364 Evan W Floden, Paul P Gardner, Thomas A Jones, John Tate, and Robert D Finn. Rfam 12.0:

- 1365 updates to the RNA families database. *Nucleic Acids Research*, 43(Database issue):D130–7, January
1366 2015.
- 1367 [127] W Wang, F G Brunet, E Nevo, and M Long. Origin of sphinx, a young chimeric RNA gene in
1368 *Drosophila melanogaster*. *Proceedings of the National Academy of Sciences of the United States of*
1369 *America*, 99(7):4448–4453, 2002.
- 1370 [128] Martin J Blythe, Damian Kao, Sunir Malla, Joanna Rowsell, Ray Wilson, Deborah Evans, Jamie
1371 Jowett, Amy Hall, Virginie Lemay, Sabrina Lam, and A Aziz Aboobaker. A Dual Platform Approach
1372 to Transcript Discovery for the Planarian Schmidtea Mediterranea to Establish RNAseq for Stem
1373 Cell and Regeneration Biology. *PLoS ONE*, 5(12):e15617, December 2010.
- 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 Degnan, 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 within Hox 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] Derek Lemons and William McGinnis. Gene Regulatory Networks in the Evolution and Development
1401 of the Heart. *Science*, 313(5795):1918–1922, September 2006.
- 1402 [140] S Cumberledge, A Zaratzian, and S Sakonju. Characterization of two RNAs transcribed from the
1403 cis-regulatory region of the abd-A domain within the Drosophila bithorax complex. *Proceedings of
1404 the National Academy of Sciences of the United States of America*, 87(9):3259–3263, May 1990.
- 1405 [141] Assaf Zemach, Ivy E McDaniel, Pedro Silva, and Daniel Zilberman. Genome-Wide Evolutionary
1406 Analysis of Eukaryotic DNA Methylation. *Science*, 328(5980):916–919, 2010.
- 1407 [142] Julie A Law and Steven E Jacobsen. Establishing, maintaining and modifying DNA methylation
1408 patterns in plants and animals. *Nature reviews. Genetics*, 11(3):204–220, February 2010.
- 1409 [143] Peter A Jones. Functions of DNA methylation: islands, start sites, gene bodies and beyond. *Nature
1410 reviews. Genetics*, 13(7):484–492, May 2012.
- 1411 [144] Peter A. Jones and Gangning Liang. Rethinking how DNA methylation patterns are maintained.
1412 *Nature Reviews Genetics*, 10(11):805–811, September 2009.
- 1413 [145] Albert Jeltsch, Ann Ehrenhofer-Murray, Tomasz P. Jurkowski, Frank Lyko, Gunter Reuter, Serge
1414 Ankri, Wolfgang Nellen, Matthias Schaefer, and Mark Helm. Mechanism and biological role of
1415 dnmt2 in nucleic acid methylation. *RNA Biology*, 0(0):1–16, 0. PMID: 27232191.
- 1416 [146] Mary Grace Goll, Finn Kirpekar, Keith A Maggert, Jeffrey A Yoder, Chih-Lin Hsieh, Xiaoyu Zhang,
1417 Kent G Golic, Steven E Jacobsen, and Timothy H Bestor. Methylation of tRNA^{Asp} by the DNA
1418 methyltransferase homolog Dnmt2. *Science*, 311(5759):395–398, January 2006.
- 1419 [147] Farah Jaber-Hijazi, Priscilla J K P Lo, Yuliana Mihaylova, Jeremy M Foster, Jack S Benner,
1420 Belen Tejada Romero, Chen Chen, Sunir Malla, Jordi Solana, Alexey Ruzov, and A Aziz Aboobaker.
1421 Planarian MBD2/3 is required for adult stem cell pluripotency independently of DNA methylation.
1422 *Developmental Biology*, 384(1):141–153, December 2013.
- 1423 [148] Jamie A Hackett, Roopsha Sengupta, Jan J Zyllicz, Kazuhiro Murakami, Caroline Lee, Thomas A
1424 Down, and M Azim Surani. Germline DNA Demethylation Dynamics and Imprint Erasure Through
1425 5-Hydroxymethylcytosine. *Science*, 339(6118):448–452, 2013.
- 1426 [149] Suhua Feng, Shawn J. Cokus, Xiaoyu Zhang, Pao-Yang Chen, Magnolia Bostick, Mary G. Goll,
1427 Jonathan Hetzel, Jayati Jain, Steven H. Strauss, Marnie E. Halpern, Chinweike Ukomadu, Kirsten C.
1428 Sadler, Sriharsa Pradhan, Matteo Pellegrini, and Steven E. Jacobsen. Conservation and divergence
1429 of methylation patterning in plants and animals. *Proceedings of the National Academy of Sciences*,
1430 107(19):8689–8694, 2010.

- 1431 [150] Albert Jeltsch. Phylogeny of methylomes. *Science*, 328(5980):837–838, 2010.
- 1432 [151] Haoyi Wang, Hui Yang, Chikdu S Shivalila, Meelad M Dawlaty, Albert W Cheng, Feng Zhang,
1433 and Rudolf Jaenisch. One-Step Generation of Mice Carrying Mutations in Multiple Genes by
1434 CRISPR/Cas-Mediated Genome Engineering. *Cell*, 153(4):910–918, May 2013.
- 1435 [152] Hui Yang, Haoyi Wang, Chikdu S Shivalila, Albert W Cheng, Linyu Shi, and Rudolf Jaenisch.
1436 One-Step Generation of Mice Carrying Reporter and Conditional Alleles by CRISPR/Cas-Mediated
1437 Genome Engineering. *Cell*, 154(6):1370–1379, September 2013.
- 1438 [153] G Panganiban, S M Irvine, C Lowe, H Roehl, L S Corley, B Sherbon, J K Grenier, J F Fallon,
1439 J Kimble, M Walker, G A Wray, B J Swalla, M Q Martindale, and S B Carroll. The origin and
1440 evolution of animal appendages. *Proceedings of the National Academy of Sciences of the United*
1441 *States of America*, 94(10):5162–5166, 1997.
- 1442 [154] Karyn N Johnson, Marielle C W van Hulten, and Andrew C Barnes. “Vaccination” of shrimp
1443 against viral pathogens: Phenomenology and underlying mechanisms. *Vaccine*, 26(38):4885–4892,
1444 September 2008.
- 1445 [155] Yanan Lu, Junjun Liu, Liji Jin, Xiaoyu Li, YuHong Zhen, Hongyu Xue, Jiansong You, and Yongping
1446 Xu. Passive protection of shrimp against white spot syndrome virus (WSSV) using specific antibody
1447 from egg yolk of chickens immunized with inactivated virus or a WSSV-DNA vaccine. *Fish and*
1448 *Shellfish Immunology*, 25(5):604–610, November 2008.
- 1449 [156] S Rajesh Kumar, V P Ishaq Ahamed, M Sarathi, A Nazeer Basha, and A S Sahul Hameed. Immuno-
1450 logical responses of *Penaeus monodon* to DNA vaccine and its efficacy to protect shrimp against
1451 white spot syndrome virus (WSSV). *Fish and Shellfish Immunology*, 24(4):467–478, April 2008.
- 1452 [157] Andrew F Rowley and Edward C Pope. Vaccines and crustacean aquaculture—A mechanistic
1453 exploration. *Aquaculture*, 334-337(C):1–11, March 2012.
- 1454 [158] William J Palmer and Francis M Jiggins. Comparative Genomics Reveals the Origins and Diversity
1455 of Arthropod Immune Systems. *Molecular biology and evolution*, 32(8):2111–2129, August 2015.
- 1456 [159] J T Simpson, K Wong, S D Jackman, J E Schein, S J M Jones, and I Birol. ABySS: A parallel
1457 assembler for short read sequence data. *Genome Research*, 19(6):1117–1123, June 2009.
- 1458 [160] M Boetzer, C V Henkel, H J Jansen, D Butler, and W Pirovano. Scaffolding pre-assembled contigs
1459 using SSPACE. *Bioinformatics*, 27(4):578–579, February 2011.
- 1460 [161] Guillaume Marçais and Carl Kingsford. A fast, lock-free approach for efficient parallel counting of
1461 occurrences of k-mers. *Bioinformatics*, 27(6):764–770, March 2011.
- 1462 [162] Marcel Martin. Cutadapt removes adapter sequences from high-throughput sequencing reads.
1463 *EMBnet*, 17(1):10–12, August 2011.

- 1464 [163] Brian J Haas, Steven L Salzberg, Wei Zhu, Mihaela Pertea, Jonathan E Allen, Joshua Orvis, Owen
1465 White, C Robin Buell, and Jennifer R Wortman. Automated eukaryotic gene structure annotation
1466 using EvidenceModeler and the Program to Assemble Spliced Alignments. *Genome biology*,
1467 9(1):R7, 2008.
- 1468 [164] M Stanke and S Waack. Gene prediction with a hidden Markov model and a new intron submodel.
1469 *Bioinformatics*, 19(Suppl 2):ii215–ii225, October 2003.
- 1470 [165] Thomas D Wu and Colin K Watanabe. GMAP: a genomic mapping and alignment program for
1471 mRNA and EST sequences. *Bioinformatics*, 21(9):1859–1875, May 2005.
- 1472 [166] Cole Trapnell, Brian A Williams, Geo Pertea, Ali Mortazavi, Gordon Kwan, Marijke J van Baren,
1473 Steven L Salzberg, Barbara J Wold, and Lior Pachter. Transcript assembly and quantification by
1474 RNA-Seq reveals unannotated transcripts and isoform switching during cell differentiation. *Nature*
1475 *Biotechnology*, 28(5):516–520, May 2010.
- 1476 [167] Alexander Dobin, Carrie A Davis, Felix Schlesinger, Jorg Drenkow, Chris Zaleski, Sonali Jha,
1477 Philippe Batut, Mark Chaisson, and Thomas R Gingeras. STAR: ultrafast universal RNA-seq aligner.
1478 *Bioinformatics*, 29(1):15–21, January 2013.
- 1479 [168] Guy St C Slater and Ewan Birney. Automated generation of heuristics for biological sequence
1480 comparison. *BMC bioinformatics*, 6:31, 2005.
- 1481 [169] A V Lukashin and M Borodovsky. GeneMark.hmm: new solutions for gene finding. *Nucleic Acids*
1482 *Research*, 26(4):1107–1115, 1998.
- 1483 [170] A F A Smit, R Hubley, and P Green. *RepeatMasker Open-4.0.*, 2013.
- 1484 [171] Matthew Kears, Richard Moir, Amy Wilson, Steven Stones-Havas, Matthew Cheung, Shane
1485 Sturrock, Simon Buxton, Alex Cooper, Sidney Markowitz, Chris Duran, Tobias Thierer, Bruce
1486 Ashton, Peter Meintjes, and Alexei Drummond. Geneious Basic: an integrated and extendable
1487 desktop software platform for the organization and analysis of sequence data. *Bioinformatics*,
1488 28(12):1647–1649, June 2012.
- 1489 [172] E J Rehm, R L Hannibal, R C Chaw, M A Vargas-Vila, and N H Patel. In Situ Hybridization
1490 of Labeled RNA Probes to Fixed Parhyale hawaiiensis Embryos. *Cold Spring Harbor Protocols*,
1491 2009(1):pdb.prot5130–pdb.prot5130, January 2009.
- 1492 [173] Robert C Edgar. MUSCLE: multiple sequence alignment with high accuracy and high throughput.
1493 *Nucleic Acids Research*, 32(5):1792–1797, 2004.
- 1494 [174] A Stamatakis. RAxML version 8: a tool for phylogenetic analysis and post-analysis of large
1495 phylogenies. *Bioinformatics*, 2014.

- 1496 [175] Felix Krueger and Simon R Andrews. Bismark: a flexible aligner and methylation caller for
1497 Bisulfite-Seq applications. *Bioinformatics*, 27(11):1571–1572, June 2011.
- 1498 [176] Pin-Hsiang Chou, Hao-Shuo Chang, I Tung Chen, Han-You Lin, Yi-Min Chen, Huey-Lang Yang,
1499 and K C Han-Ching Wang. The putative invertebrate adaptive immune protein *Litopenaeus vannamei*
1500 Dscam (LvDscam) is the first reported Dscam to lack a transmembrane domain and cytoplasmic tail.
1501 *Developmental & Comparative Immunology*, 33(12):1258–1267, December 2009.
- 1502 [177] E P Nawrocki and S R Eddy. Infernal 1.1: 100-fold faster RNA homology searches. *Bioinformatics*,
1503 29(22):2933–2935, October 2013.
- 1504 [178] Sam Griffiths-Jones, Harpreet Kaur Saini, Stijn van Dongen, and Anton J Enright. miRBase: tools
1505 for microRNA genomics. *Nucleic Acids Research*, 36(Database issue):D154–8, January 2008.