

1 **RESEARCH ARTICLE**

**Open Access**

2 **A comparison across non-model animals**  
3 **suggests an optimal sequencing depth**  
4 **for *de novo* transcriptome assembly**

5 Warren R Francis<sup>1,2\*</sup>, Lynne M Christianson<sup>1</sup>, Rainer Kiko<sup>3</sup>, Meghan L Powers<sup>1,2</sup>, Nathan C Shaner<sup>4</sup>  
6 and Steven H D Haddock<sup>1\*</sup>

7 **Abstract**

8 **Background:** The lack of genomic resources can present challenges for studies of non-model organisms.  
9 Transcriptome sequencing offers an attractive method to gather information about genes and gene expression  
10 without the need for a reference genome. However, it is unclear what sequencing depth is adequate to assemble the  
11 transcriptome *de novo* for these purposes.

12 **Results:** We assembled transcriptomes of animals from six different phyla (Annelids, Arthropods, Chordates,  
13 Cnidarians, Ctenophores, and Molluscs) at regular increments of reads using Velvet/Oases and Trinity to determine  
14 how read count affects the assembly. This included an assembly of mouse heart reads because we could compare  
15 those against the reference genome that is available. We found qualitative differences in the assemblies of  
16 whole-animals versus tissues. With increasing reads, whole-animal assemblies show rapid increase of transcripts and  
17 discovery of conserved genes, while single-tissue assemblies show a slower discovery of conserved genes though the  
18 assembled transcripts were often longer. A deeper examination of the mouse assemblies shows that with more reads,  
19 assembly errors become more frequent but such errors can be mitigated with more stringent assembly parameters.

20 **Conclusions:** These assembly trends suggest that representative assemblies are generated with as few as 20 million  
21 reads for tissue samples and 30 million reads for whole-animals for RNA-level coverage. These depths provide a good  
22 balance between coverage and noise. Beyond 60 million reads, the discovery of new genes is low and sequencing  
23 errors of highly-expressed genes are likely to accumulate. Finally, siphonophores (polymorphic Cnidarians) are an  
24 exception and possibly require alternate assembly strategies.

25 **Background**

26 RNA-seq has provided a powerful tool for analysis of  
27 transcriptomes. For non-model organisms with limited  
28 genomic information, transcriptome sequencing provides  
29 a cost-saving tool by only sequencing functional and  
30 protein coding RNAs, thus providing direct information  
31 about the genes [1]. There are many benefits of sequencing  
32 a genome, but for relatively large genomes such as human  
33 and mouse, protein coding regions account for under 5%,

34 thus most of the sequencing effort would go to sequenc- 34  
35 ing either regulatory regions or repetitive elements [2]. 35  
36 Smaller genomes could be sequenced and assembled to 36  
37 complement the transcriptomes, though this is not a 37  
38 tractable approach if a genome is quite large. Even still, *de* 38  
39 *novo* genome assembly can produce errors by itself [3]. 39

40 Despite its advantage, transcriptome assembly does 40  
41 present additional challenges when compared to genome 41  
42 assembly. Unlike genomes where most sequences should 42  
43 be approximately equally represented, coverage of any 43  
44 given sequence in a transcriptome can vary over sev- 44  
45 eral orders of magnitude due to expression differences 45  
46 [4]. Because coverage can vary, there is also a question 46  
47 of sequencing depth. Theoretically, there is a sequenc- 47  
48 ing depth beyond which addition of more reads does not 48  
49 provide new information, known as the saturation depth. 49

\*Correspondence: wfrancis@mbari.org; haddock@mbari.org

<sup>1</sup> Monterey Bay Aquarium Research Institute, 7700 Sandholdt Rd, Moss  
Landing, CA 95039, USA

<sup>2</sup> Department of Ocean Sciences, University of California Santa Cruz, Santa  
Cruz, CA, USA

Full list of author information is available at the end of the article

50 Several studies have used approaches which map reads  
51 onto reference genomes and these have suggested saturation  
52 depths at 95% gene coverage ranging from 1.2 million  
53 reads to 50 million for mRNA level coverage, and up to  
54 700 million for splice variants [5-7]. However, these studies  
55 all made use of short reads around 36bp and were not  
56 assembling the transcriptomes *de novo*.

57 Several recent studies have already made use of next-  
58 generation sequencing reads for *de novo* transcriptome  
59 assembly [8-15]. The number of reads used for assembly in  
60 these studies varies widely, ranging from 2.6 million reads  
61 up to 106 million reads [10,11]. The assembly strategies  
62 are equally varied, but share the initial step of removing  
63 low-quality reads and adapters whereupon all remaining  
64 reads are assembled. The assembly quality estimates vary  
65 as well with the most common measure of quality based  
66 on BLAST hits to public databases like Uniprot, though  
67 it was noted that under-representation of many taxa in  
68 public databases limits this approach [8].

69 While many parameters must be optimized for the specific  
70 assembly, it is both inconvenient and costly to acquire  
71 more reads by resequencing. Presently, there is no clear  
72 consensus of what sequencing depth is optimal or what  
73 factors would contribute to the adequate depth. The problems  
74 of omitted genes or variants are obvious with too few  
75 reads. On the other hand, it was suggested that greater  
76 depth may create errors in differential expression analyses,  
77 cost more, and take longer to assemble [16]. Thus, here  
78 we use the same assembly strategy across a diverse set  
79 of organisms to isolate the effects of read count on  
80 assembly quality to attain a general estimate of optimal  
81 read count. We compare trends from *de novo* assemblies  
82 across six phyla. These animals include the mouse (used as  
83 a control for the non-model samples), the Humboldt squid  
84 *Dosidicus gigas*, the scaleworm *Harmothoe imbricata*,  
85 the decapod *Sergestes similis*, the copepod *Pleuromamma*  
86 *robusta*, the ctenophore *Hormiphora californensis*, and  
87 the siphonophore *Chuniphyes multidentata*. To our  
88 knowledge, this is the first study to suggest an optimal  
89 number of reads for *de novo* assembly for the purposes  
90 of mRNA level analysis. These results are applicable to  
91 studies of organisms with limited genomic resources.

## 92 Results and discussion

### 93 *De novo* assembly of transcriptomes

#### 94 *Assembly of mouse heart transcriptome*

95 Raw mouse-transcriptome reads from the ENCODE  
96 project were downloaded from NCBI short-read archive.  
97 Sample SRR453174 (mouse heart RNA-seq) consisted of  
98 82,886,668 x76bp reads as paired-ends. Filtration (see  
99 Methods) removed 11.7% of the reads, almost 95% of  
100 which were due to low quality scores. In order to examine  
101 the role of number of reads on the assembly, we  
102 computationally sub-sampled randomized sets from the

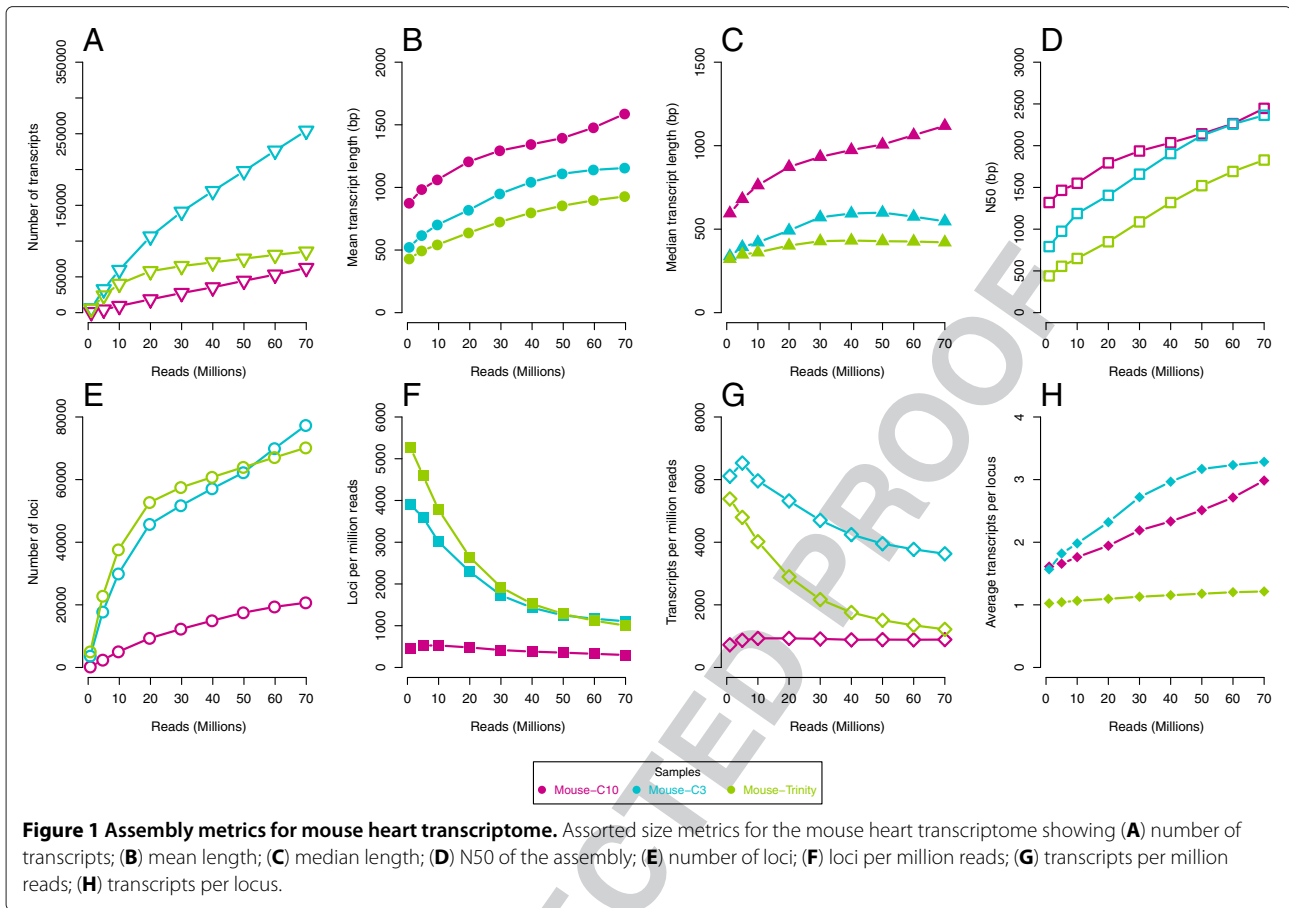
original library. It is suggested that sequencing of very  
small numbers of reads can be most subject to biases  
and that cDNA normalization can improve the uniformity  
of the library at low numbers of reads [17]. Such an  
approach might be quite costly, and the computational  
sub-sampling approach has the advantage of drawing from  
the largest pool of reads and avoid biases which could  
occur at low numbers of reads. Subsets of the filtered  
library were generated containing 1,5,10,20,30,40,50,60,  
and 70 million reads. Reads from each set were included  
in the next largest set, thus all of the reads in the 1 million  
set are included in the 5 million read set, and so forth.  
These sets were assembled with Velvet/Oases [18,19] and  
Trinity [20] (For a detailed comparison of assemblers,  
see [21]).

Schulz *et al.* reported reliable parameters for Oases  
which produced high-quality assemblies of mouse and  
human cell cultures, using 64 million and 30 million reads,  
respectively [19]. This included use of a broad k-mer range  
with a low starting k-mer of 19 or 21 up to a k-mer of 33  
or 35. Accordingly we used k-mers from 21 to 33. Also, a  
minimum k-mer coverage is required by Oases to retain  
any given node during the assembly process; by default  
this is 3 in Oases, that is, any node must have at least  
three-fold coverage for that node to be used. Some differences  
were observed in the output when this parameter was  
changed, and so the same data were assembled with  
coverage cutoff of 3 (referred to hereafter as C3) and a  
stricter cutoff of 10 (C10).

The number of transcripts (Oases terminology for contigs)  
increases steadily for all assemblies (Figure 1A). C10  
also had substantially fewer transcripts and accordingly  
much higher mean and median lengths (Figure 1B-D).  
The pattern of increase for median and N50 (length for  
which half of the total bases are in contigs of this length  
or longer) tracked the mean for the C10 assembly, but not  
the C3 assembly which did not have a clear qualitative  
pattern. The mean, median and N50 were all lower for the  
Trinity assembly than the C3 despite having far fewer  
contigs.

Oases generates transcript "loci", which is Oases terminology  
for the de-Bruijn graph clusters meant to represent genes  
and their splice variants or highly-similar paralogs. Both  
curves approach to a plateau for locus counts (Figure 1E-F).  
The greatest increase in loci was between using 10 million  
to 20 million reads for both C3 and C10. Similarly, the  
C3 assembly shows a decrease in the number of transcripts  
per read (Figure 1G), while the C10 assembly shows an  
almost constant number of transcripts per read. The number  
of transcripts increases while the number of loci tend to  
level off and this means the number of transcripts per  
locus always increases with more reads (Figure 1H). That  
is, on average, more variants will be generated with more  
reads even though some of these are likely due to noise.  
While the Trinity assembly

103  
104  
105  
106  
107  
108  
109  
110  
111  
112  
113  
114  
115  
116  
117  
118  
119  
120  
121  
122  
123  
124  
125  
126  
127  
128  
129  
130  
131  
132  
133 **F1**  
134  
135  
136  
137  
138  
139  
140  
141  
142  
143  
144  
145  
146  
147  
148  
149  
150  
151  
152  
153  
154  
155  
156



157 more closely matches the trends for transcripts per read  
 158 of the C3, the “components” (closest obvious parallel of  
 159 loci) remain close to a unit ratio, suggesting that most  
 160 components have only one associated sequence.

161 **Assembly of invertebrate transcriptomes**

162 Transcriptomes across a broad range of taxa were assem-  
 163 bled as with the mouse and statistics of the largest assem-  
 164 blies are presented in Table 1. The stated GC content of  
 165 the mouse genome is 42% while a subset of conserved  
 166 genes showed a much higher value of 51.24% [22,23].  
 167 Interestingly, for all assemblies except for mouse, the aver-  
 168 age GC content of the assembled contigs was lower than  
 169 that of the raw reads (Figure 2), suggesting either that  
 170 certain genes contribute much more to the overall GC  
 171 content of the library or that biases can be introduced  
 172 from the assembly.

173 For three of six samples (*D.gigas*, *H.imbricata* and  
 174 *S.similis*), only select tissues were used for RNA extrac-  
 175 tion while the rest were whole body (*C.multidentata*,  
 176 *H.californensis* and *Probusta*). It should be noted that  
 177 the *C.multidentata* sample combined sequences from  
 178 the two major tissues, siphosome and nectophore and  
 179 that the *Probusta* sample was a combination of multiple

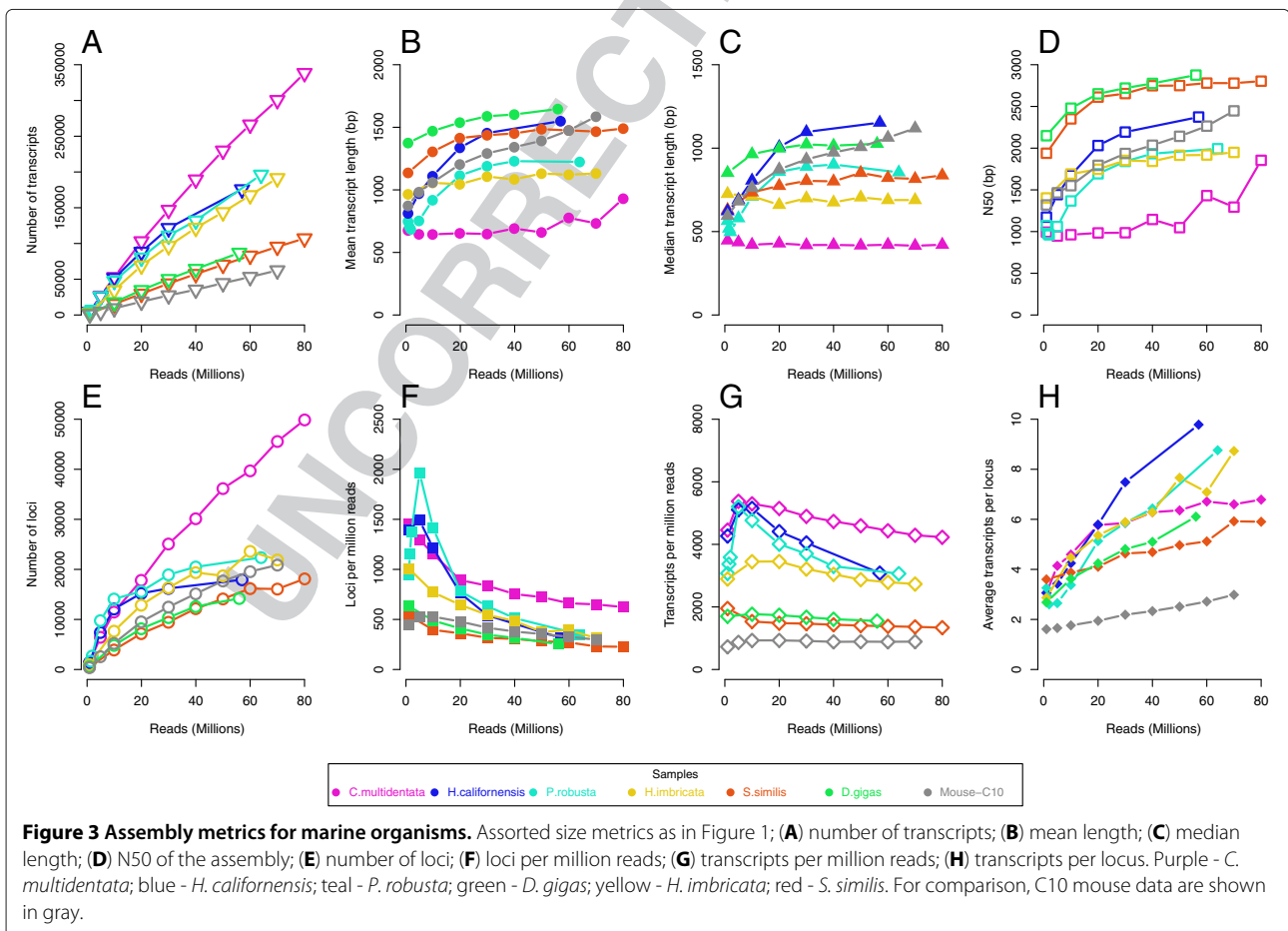
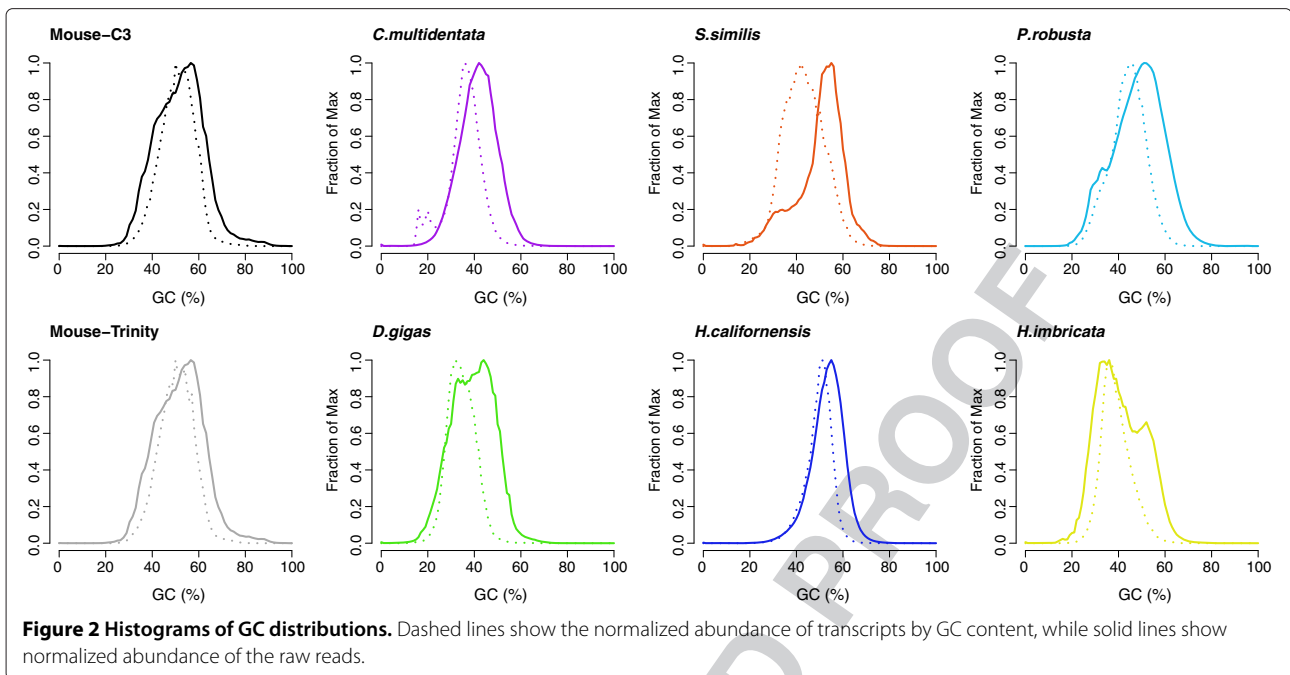
180 individuals. This decision was based on size of the ani- 180  
 181 mals since very small organisms are difficult to dissect. 181  
 182 Assembly trends analogous to Figure 1 for the six animals 182  
 183 are shown in Figure 3. Mouse C10 data from Figure 1 are 183 F3  
 184 shown in gray as reference. Three main trends emerged. 184  
 185 Whole-body samples were characterized by a rapid gain 185  
 186 of transcripts and increases in transcript size through 40 186  
 187 million reads, while all other parameters level off after 187  
 188 40 million reads. Single tissue samples showed a slow 188  
 189 gain of relatively long transcripts across fewer loci. Lastly, 189  
 190 the whole-body siphonophore showed continuous gain 190  
 191 of both short transcripts and loci without reaching an 191  
 192 asymptote at the maximum number of reads assembled. 192

193 Four of the animals showed modest gains in mean, 193  
 194 median and N50 with more reads (average 20% from 194  
 195 fewest to most reads), while *P.robusta* and *H.californensis* 195  
 196 nearly doubled from the fewest to the most reads 196  
 197 (Figure 3B-D). Most of the transcript-length increase 197  
 198 occurred before 30 million reads, suggesting that adding 198  
 199 more reads did not produce longer sequences beyond 199  
 200 that threshold, or that they became longer at the same 200  
 201 rate that new, short transcripts were generated. As with 201  
 202 the mouse samples, transcripts were added continually 202  
 203 with more reads (Figure 3A). Compared to the mouse, 203

t1.1 **Table 1 Assembly Statistics**

| t1.2  | Organism           | Mouse cov-cutoff-3 | Mouse cov-cutoff-10 | Mouse-Trinity | Chuniphyes multidentata | Sergestes similis | Pleuromamma robusta | Dosidicus gigas | Hormiphora californensis | Harmothoe imbricata |
|-------|--------------------|--------------------|---------------------|---------------|-------------------------|-------------------|---------------------|-----------------|--------------------------|---------------------|
| t1.3  | Phylum             | Chordata           | Chordata            | Chordata      | Cnidaria                | Arthropoda        | Arthropoda          | Mollusca        | Ctenophora               | Annelida            |
| t1.4  | Tissue             | Heart              | Heart               | Heart         | Whole body              | Legs              | Whole body          | Mantle          | Whole body               | Scale               |
| t1.5  | Raw Reads          | 82,886,668         | 82,886,668          | 82,886,668    | 103,415,276             | 93,597,558        | 64,116,306          | 60,661,588      | 64,675,964               | 75,608,018          |
| t1.6  | Raw GC (%)         | 51.90              | 51.90               | 51.90         | 42.29                   | 50.74             | 48.86               | 39.89           | 53.71                    | 41.52               |
| t1.7  | Filtered Reads     | 73,187,048         | 73,187,048          | 73,187,048    | 102,366,438             | 92,423,904        | 63,867,922          | 56,264,099      | 57,583,204               | 70,340,105          |
| t1.8  | Assembled Reads    | 70,000,000         | 70,000,000          | 70,000,000    | 80,000,000              | 80,000,000        | 63,867,922          | 56,264,099      | 57,583,204               | 70,340,105          |
| t1.9  | Transcripts        | 254,215            | 62,353              | 85,294        | 338,254                 | 107,082           | 196,104             | 86,897          | 175,701                  | 191,290             |
| t1.10 | Total Length (Mbp) | 293.55             | 98.84               | 79.12         | 314.99                  | 159.59            | 240.05              | 143.09          | 272.23                   | 216.66              |
| t1.11 | Mean (bp)          | 1,154              | 1,585               | 927           | 931                     | 1,490             | 1,224               | 1,646           | 1,549                    | 1,132               |
| t1.12 | Median (bp)        | 547                | 1,119               | 421           | 421                     | 837               | 855                 | 1,026           | 1,153                    | 689                 |
| t1.13 | N50 (bp)           | 2,364              | 2,447               | 1,828         | 1,854                   | 2,803             | 1,993               | 2,876           | 2,373                    | 1,949               |
| t1.14 | Oases Loci         | 77,411             | 20,889              | 70272         | 49,831                  | 18,139            | 22,385              | 14,227          | 17,960                   | 21,914              |
| t1.15 | GC (%)             | 54.08              | 53.95               | 53.46         | 31.24                   | 44.66             | 45.78               | 36.55           | 51.66                    | 40.53               |

t1.16 Summary statistics of the largest transcriptome assembly for each organism.





204 on average these six animals all had more transcripts  
205 per locus (Figure 3H). It is unclear why this would  
206 be the case, though the C10 assembly had the fewest  
207 number of transcripts overall for all numbers of reads.  
208 The most pronounced gains in loci happened within  
209 the first 10 million reads, particularly for *Probusta* and  
210 *H.californensis* (Figure 3E-F). Gains in loci tended to level  
211 out between 40 and 60 million reads, suggesting most  
212 genes (or parts of genes) were assembled by 60 million  
213 reads.

214 A very high number of transcripts for *C.multidentata*  
215 (Figure 3, purple) led to the lowest mean, median, and  
216 N50. The number of removed, low-quality reads is com-  
217 parable in this sample to others, so low quality is unlikely  
218 to be the cause. As two sets of reads were combined  
219 into a whole animal, this may have created artifacts.  
220 However, another *C.multidentata* siphosome sample pro-  
221 duced assemblies with large numbers of relatively short  
222 sequences (data unpublished). One possible explanation  
223 is that siphonophores have continuously developing dif-  
224 ferentiated zooids [24]. These zooids have specialized  
225 functions which are in some ways analogous to organs,  
226 and a whole organism can contain multiple developmental  
227 stages and express a large part of the genome, possi-  
228 bly confounding the assembly process. Assemblies of a  
229 number other siphonophores (data unpublished) similarly  
230 had many short transcripts. We speculate that alternate  
231 assembly strategies or very careful dissections might be  
232 required for animals in this lineage.

### 233 Discovery of conserved genes

#### 234 Conserved mouse genes

235 One approach used to assess genome completeness is to  
236 search only for conserved eukaryotic orthologous genes  
237 (KOGs). The current NCBI KOG database has 860 gene  
238 clusters across 7 eukaryotes with over 16000 proteins  
239 [25]. The KOG reference genes did not include mouse  
240 sequences, and this provided an opportunity to test pre-  
241 dictions about *de novo* transcriptome quality while still  
242 having a reference in the end to confirm the reliability  
243 of the sequences. For each KOG, the transcripts were  
244 aligned against the reference KOGs with *tblastn*, and the  
245 best coding sequence was kept. The putative proteins were  
246 classified by length relative to the range of sizes of the  
247 reference KOGs. The size range allowed some flexibility,  
248 as 12 mouse proteins were larger than the longest refer-  
249 ence protein for that KOG, and 5 were shorter than  
250 the shortest reference protein. Finally the proteins were  
251 aligned with *blastp* against reviewed mouse proteins in  
252 Uniprot to determine accuracy. One protein was unre-  
253 viewed (Q3UWL8, Mouse Prefoldin 4). For this test, Trin-  
254 ity and Oases are comparable at assembling full-length  
255 proteins, though Trinity appears to be slightly better at  
256 reconstructing canonical proteins (Figure 4A).

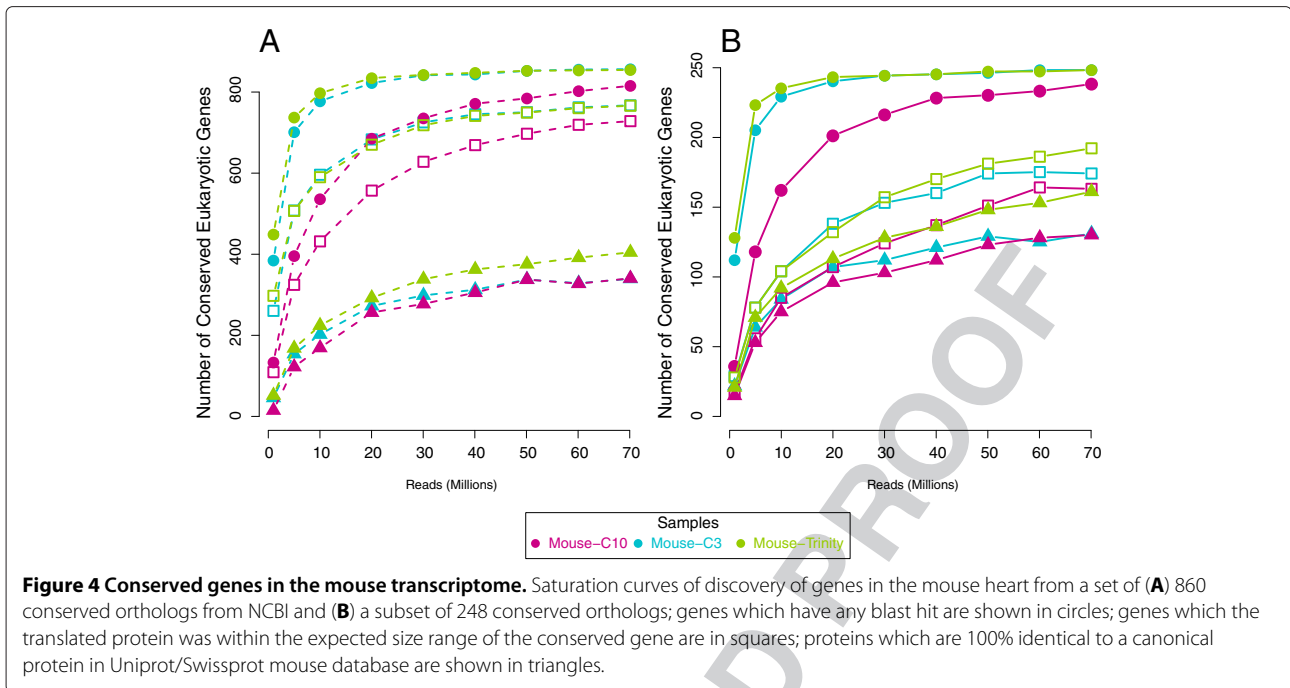
257 However, gene duplications present difficulties for such  
258 assessments unless one had *a priori* knowledge of how  
259 many copies should be present in the genome. For this  
260 study, we also used the subset of eukaryotic KOGs con-  
261 taining 248 genes from the CEGMA pipeline which  
262 were identified as single-copy orthologs in most genomes  
263 [26,27]. Almost one third of these KOGs are involved  
264 in processes like transcription and translation and were  
265 expected to be expressed in many tissues. Trinity and  
266 Oases with a lower coverage cutoff of 3 found simi-  
267 lar numbers of KOGs at much lower numbers of reads  
268 (Figure 4B) than compared to the C10 assembly. Also  
269 more KOGs were found within expected length much  
270 faster with C3 than with the higher cutoff of 10, and  
271 the Trinity assembly outperformed both of these. These  
272 results suggest that it is better to have a lower cutoff and  
273 assemble more sequences. Likewise, the Trinity assem-  
274 bly had more transcripts than C10 and were shorter than  
275 those in C3, yet more KOGs were found with fewer  
276 reads and more coding transcripts were correctly assem-  
277 bled at greater numbers of reads. However, for the Oases  
278 assemblies this had remarkably little effect on the number  
279 of correct canonical proteins that were found (Figure 4,  
280 triangles). Although there is some overestimation, no pro-  
281 tein designated as too short was ever correct. Regarding  
282 the fate of the other full-length proteins, for C3 at 70 mil-  
283 lion reads, 186 KOGs were found within the expected  
284 range, though only 131 were correct. Eight of the 186  
285 KOGs had only 1 mismatch in the amino-acid sequence  
286 compared to the reference protein which could be due  
287 to errors, splice variants, tissue-specific modifications or  
288 alleles. The remaining KOGs had at least two amino-  
289 acid changes but were within the size range. Thus for  
290 the mouse, the size range was a reliable predictor of true  
291 full-length proteins.

#### 292 Conserved invertebrate genes

293 We then examined our invertebrate transcriptomes for  
294 completion using the same set of KOGs. There was a  
295 clear, qualitative difference between whole-body organ-  
296 isms (Figure 5A) and dissected tissues (Figure 5B). C10  
297 mouse data are included for reference. For whole-body  
298 transcriptomes, over 90% of the KOGs were detectable at  
299 20 million reads, yet the number of within-length KOGs  
300 went down with higher numbers of reads past 20 mil-  
301 lion. This could be caused if proteins declared to be  
302 within-range were longer than the true protein due to mis-  
303 assembly causing addition of pieces, or if the true protein  
304 became mis-assembled with addition of noisy reads. In  
305 nearly all of our assemblies, it was the latter: mis-assembly  
306 of the putative protein which generated stop codons.  
307 *C.multidentata* (Figure 5A, purple) was again exceptional,  
308 as the number of within-length KOGs increased more  
309 slowly with addition of more reads than the other two

F5

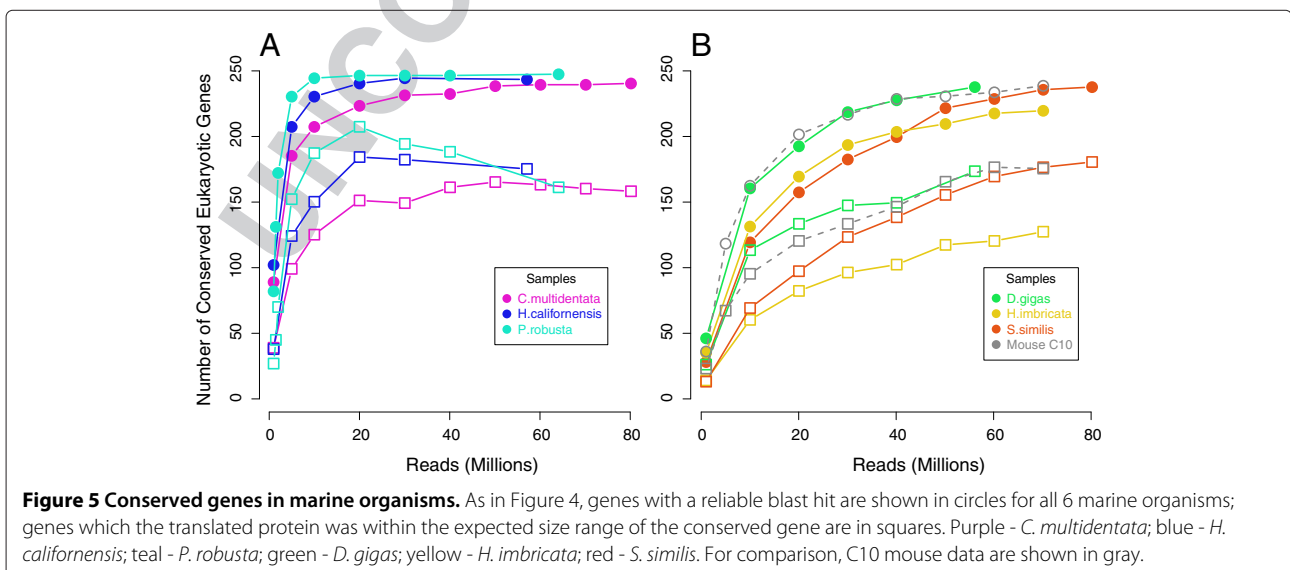
F4



310 whole-body animals (*H.californensis* and *Probusta*) and  
 311 only decreased after 50 million reads rather than 20  
 312 million.

313 For dissected-tissue transcriptomes (*Dosidicus gigas*,  
 314 *Harmothoe imbricata*, and *Sergestes similis*), the rate  
 315 of discovery of KOGs was much slower with between  
 316 63% and 81% of KOGs detectable at 20 million reads  
 317 (Figure 5B). This was not surprising since those genes  
 318 may not be highly-expressed in all tissues and it is likely  
 319 tissue-specific genes account for the bulk of the assembly  
 320 at low numbers of reads. Isolated tissues may express  
 321 fewer universal KOGs that we selected in our test, and

we expected that other abundant transcripts should mis-  
 assemble at high numbers of reads in that tissue. However,  
 the dissected-tissue transcriptomes had longer transcripts  
 and fewer loci, suggesting this was not the case. Since  
 whole-animal transcriptomes include all tissues, a greater  
 proportion of the genome is expressed so coverage of any  
 given transcript or splice-variant is proportionally much  
 lower. The length saturation patterns appear to be different  
 between whole-animal and tissue transcriptomes. However,  
 using conserved genes as a metric, there appears to be  
 limited benefit of sequencing beyond 60 million reads.



### 334 Mis-assembly at high numbers of reads

335 KOGs with single-exon coding sequences in the mouse  
336 were examined for mis-assembly. To increase the number  
337 of genes examined, another set of KOGs from only  
338 metazoans (*C.elegans*, *D.melanogaster* and *H.sapiens*,  
339 CDH) was used. The KOG database at NCBI contained  
340 1147 clusters common to CDH. Again, only genes that  
341 were annotated as single copy in all three animals were  
342 used, leaving a final set of 202 KOGs specific to  
343 metazoans. These combined sets of 450 had 12 genes  
344 in mouse which were presumed single-copy and annotated  
345 in NCBI to have a single-exon coding sequence  
346 (GenBank:NP\_062724.1, NP\_666327.2, NP\_082281.2,  
347 NP\_058612.3, XP\_899832.1, NP\_001153802.1, NP\_001104758.1,  
348 NP\_077152.1, XP\_486217.2, NP\_598737.1, NP\_032025.2,  
349 NP\_075969.1). At 70 million reads, 3 genes in C3  
350 had alternate erroneous coding sequences: NAT6, CHMP1B1/  
351 DID2, FTSJ (N-acetyl transferase 6, Charged multivesicular  
352 body protein 1b-1, Ribosomal RNA methyltransferase,  
353 respectively). The sequence of CHMP1B1 was never  
354 assembled correctly for any number of reads and the best  
355 version was missing 9 amino acids at the N-terminus  
356 including the start codon. Only NAT6 had extraneous  
357 coding sequence in C10, suggesting that such errors can  
358 be controlled by limiting read count as well as increasing  
359 k-mer coverage thresholds.

360 While some mis-assemblies can occur with more reads,  
361 overall this is not a problem, as shown by the curves in  
362 Figures 4 and 5. However, select cases of mis-assembly of  
363 the mouse genes are shown in Figure 6. AlaRS (Alanyl-tRNA  
364 synthetase) presents an example of the optimal scenario,  
365 whereby the protein is not found at all with few reads,  
366 but then pieces come together with the addition of more  
367 reads until the final protein is correctly assembled. The  
368 majority of proteins follow this trend. 2-OGDH shows an  
369 unusual oscillation between the reference protein and  
370 alternate forms. EF2 is assembled correctly with few  
371 reads, then errors accumulate as more reads are added.  
372 From this, it cannot be assumed that the largest set of  
373 reads will produce the best contigs. Schulz *et al.* indicated  
374 that between 10 and 20% of Oases transcripts had some  
375 degree of misassembly [19]. This value was found to  
376 correlate with the smallest k-mer used in assembly and  
377 the authors suggest using larger k-mers if problems arise  
378 due to chimeric transcripts. Thus if using more reads,  
379 it may be advisable to use larger k-mers or a higher static  
380 coverage cutoff.

### 381 Conclusions

382 In this study, number of whole animals and tissues from  
383 non-model organisms and one mouse organ were assembled  
384 and the completeness was assessed using a set of conserved  
385 genes. Additionally, a comparison was made between two  
386 high-performing assemblers with respect to

the mouse data. Oases required much greater memory  
usage while Trinity had much longer run times (approximately  
2-fold longer). Both Trinity and Oases perform comparably  
at assembling conserved genes across a large set, indicating  
that the saturation depth is not greatly affected by assembler  
choice.

Overall, these results suggest that for whole-body transcriptomes  
and individual organs or cells, 30 and 20 million reads are  
sufficient for mRNA level coverage, respectively. For the read  
length used in this study, that would produce 2-3 gigabases  
of sequence. It should be noted that the mouse data consisted  
of shorter reads than used for the invertebrates, but this did  
not appear to have substantial effect as this difference was  
only between 75bp reads and 100bp reads. Assembly errors  
are evident in whole-body transcriptomes after 30 million  
reads, and the average length appeared to level off at the  
same depth. Presumably this depth would apply for studies  
of differential expression as well, as the highly expressed  
transcripts should be present and distinguishable at that  
sequencing depth. In our experience, we find it is optimal  
to acquire between 50 and 60 million reads, and then sub-  
sample up around 20 or 30 million. This approach reliably  
assembles nearly all proteins of interest. There are still  
observable differences between assemblies, although some of  
these differences may ultimately be due to variations in  
RNA quality or properties of the animal.

### Methods

#### Samples and sequencing

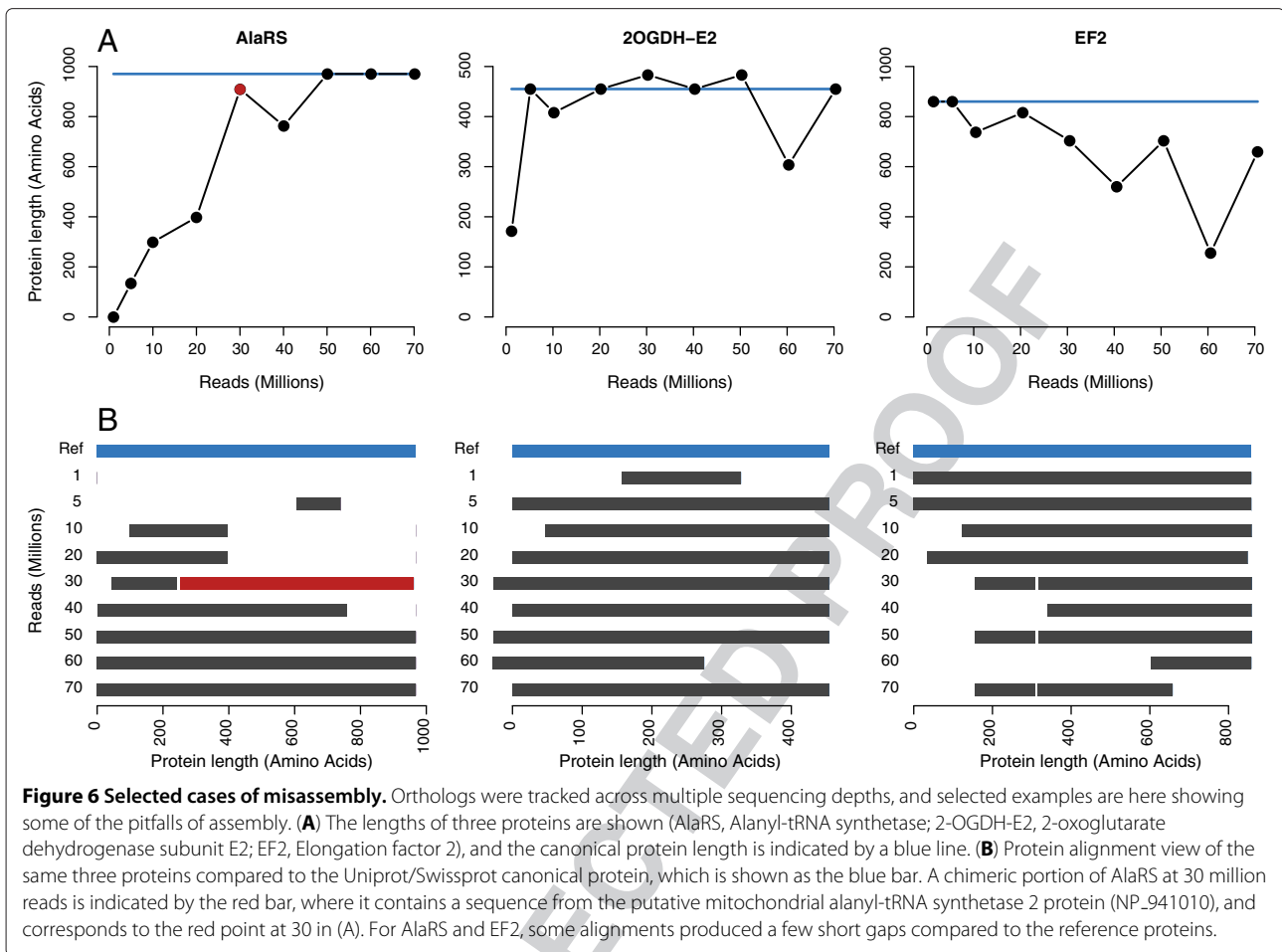
*D.gigas* and *H.californensis* were collected in the Gulf  
of California by jig and trawl net, respectively. *C.multidentata*  
and *S.similis* were collected in the Monterey Bay using  
remotely-operated-underwater vehicles. *H.imbricata* samples  
were given courtesy of T. Rivers. All samples were flash  
frozen in liquid nitrogen immediately following collection.  
Total RNA was extracted using RNeasy kit (Qiagen) as per  
instructions. *C.multidentata* RNA was extracted with Trizol  
and purified with the RNeasy kit. Preparation of RNA-seq  
libraries was done using Illumina TruSeq kit for paired end  
reads. Total RNA was sent for sequencing at University of  
Utah. Multiple individuals of *P.robusta* were sampled off the  
coast of Namibia and sequenced at the Institute for Clinical  
Molecular Biology, (IKMB, Kiel University). Sequencing  
was done using the Illumina HiSeq2000 platform on a  
paired-end protocol with 100 cycles. Mouse heart data were  
downloaded from NCBI accession GSE36025, sample SRR453174.

#### Transcriptome assembly

All computations were done on a computer with two quad-  
core processors and 96GB RAM. For each sample, the orders  
of all raw reads were randomized with the

F6





439 randomize.cpp program and processed with a modified  
 440 version of the filter\_illumina.cpp program in the Agalma  
 441 transcriptome package (<https://github.com/caseywdunn/agalma>). This removed low-quality reads (with mean  
 442 Phred score < 28), as well as reads containing adapters  
 443 and reads that were mostly repeated bases, such as polyT  
 444 tracts. Reads from pairs with one good read and one bad  
 445 read retained the good read for the largest assembly. Oth-  
 446 erwise, only good pairs were used in other assemblies. The  
 447 transcriptome for each set was assembled *de novo* using  
 448 Velvet v1.2.06 /Oases v0.2.06. Identical assembly param-  
 449 eters were used unless otherwise noted. Multiple k-mer  
 450 assemblies were generated (21,25,29,33) and merged with  
 451 Oases-M (k-mer of 27). A static coverage cutoff of 10 was  
 452 used and insert size of the paired ends was estimated with  
 453 the “-exp\_cov auto” parameter, typically around 180bp,  
 454 as expected. The minimum contig length was set to 100,  
 455 which is the read length. The Trinity assembler was also  
 456 used for comparison of mouse assemblies using the same  
 457 filtered subsets of reads. Other than insert length being  
 458 specified as the upper limit rather than the mean, default  
 459 assembly parameters were used including a minimum  
 460

transcript length of 200bp. Transcript lengths and GC 461  
 content were measured with an in-house python script, 462  
 sizecutter.py, available at the MBARI public repository 463  
 ([bitbucket.org/beroe/mbari-public/src](http://bitbucket.org/beroe/mbari-public/src)). 464

### Conserved gene analyses 465

All blast searches were done using the NCBI blast 2.2.25+ 466  
 package [28]. We generated a script to blast and ana- 467  
 lyze the matches, kogblaster.py (on the public repository, 468  
 as above). Briefly, the reference KOGs (860 orthologous 469  
 groups from NCBI, or 248 orthologous groups, from 470  
<http://korflab.ucdavis.edu/Datasets/cegma/>) were aligned 471  
 to each assembly with tblastn with an e-value cutoff of 472  
 $10^{-6}$ . For each alignment, the subject hit was translated 473  
 and coding sequences were only kept if they contained 474  
 both start and stop codons. From this subset, the best 475  
 alignment was declared to be the correct sequence. Next, 476  
 the length of the correct sequence was used to estimate 477  
 whether that sequence was full-length relative to the 478  
 conserved orthologs. For each KOG in the CEGMA dataset, 479  
 there were 6 proteins from 6 species and there was some 480  
 variability in protein length (average 11.8% from longest 481

Q2

482 to shortest). The variability from the the reference set was  
483 used to establish boundaries for size classifications which  
484 were made to watch the progression of assembly of indi-  
485 vidual genes: (1) within the size range of the KOG; (2)  
486 within the range but where the alignment was less than  
487 90% of the length of the protein; (3) longer than those in  
488 the size range; (4) shorter than the size range; (5) shorter  
489 than the size range and shorter than the alignment, often  
490 indicative of a stop codon bridged by the alignment. The  
491 full-length size range was defined by ratios of the short-  
492 est protein to the second shortest, and analogously for  
493 the longest protein and second longest. For example, if  
494 the shortest protein within a KOG was 80AAs, and the  
495 second shortest was 100AAs, the lower bound would be  
496  $(80 * (80/100))$ , and thus 64AAs. This was calculated for  
497 each KOG, and was to account for proteins which could  
498 potentially become the 'new' shortest or longest. Ulti-  
499 mately, only those within the size range (1) were declared  
500 as full-length sequences.

501 The animals in this study were treated ethically and  
502 responsibly. Because no vertebrates or octopus were  
503 involved, no formal certification is required per the  
504 Helsinki Declaration. The mouse data presented in the  
505 paper were not obtained from our experiments, but were  
506 downloaded from a database.

#### 507 Competing interests

508 The authors declare no competing financial interests.

#### 509 Authors' contributions

510 WF, RK and SH designed experiments. LC, RK, MP and SH caught animals. LC,  
511 RK, MP and NS processed animals and extracted RNA. WF assembled  
512 transcriptomes. WF, RK and SH analyzed data. WF wrote the paper.  
513 All authors read and approved the final manuscript.

#### 514 Acknowledgements

515 WRF would like to thank J. Maitin-Shepard for help optimizing the Python  
516 scripts, D. Zerbino and M. Schulz for numerous tips and bug corrections to  
517 Velvet and Oases during the process of assembly, C. Dunn and M. Howison for  
518 preliminary versions of the Agalma package. The NIH National Institute of  
519 General Medical Sciences (ROI-GMO87198) to S. H. D. H. supported our work.  
520 This research was also supported by the David and Lucile Packard Foundation  
521 through the Monterey Bay Aquarium Research Institute. RK was supported by  
522 a Postdoc-Stipend of the German Academic Exchange Service and through  
523 the DFG funded Cluster of Excellence "Future Ocean"; project CP0923.

#### 524 Author details

525 <sup>1</sup>Monterey Bay Aquarium Research Institute, 7700 Sandholdt Rd, Moss  
526 Landing, CA 95039, USA. <sup>2</sup>Department of Ocean Sciences, University of  
527 California Santa Cruz, Santa Cruz, CA, USA. <sup>3</sup>Helmholtz Center for Ocean  
528 Research Kiel, GEOMAR, Hohenbergstr. 2, 24105 Kiel, Germany. <sup>4</sup>The Scintillon  
529 Institute, 9924 Mesa Rim Rd., San Diego, CA 92121, USA.

530 Received: 15 August 2012 Accepted: 23 January 2013

531 Published: 12 March 2013

#### 532 References

533 1. Wang Z, Gerstein M, Snyder M: **RNA-Seq: a revolutionary tool for**  
534 **transcriptomics.** *Nature Rev Genet* 2009, **10**:57–63. [http://www.  
535 pubmedcentral.nih.gov/articlerender.fcgi?artid=2949280&tool=  
536 pmcentrez&rendertype=abstract]

2. Sakharkar MK, Perumal BS, Sakharkar KR, Kanguane P: **An analysis on**  
537 **gene architecture in human and mouse genomes.** *In Silico Biol* 2005,  
538 **5**(4):347–365. [http://www.ncbi.nlm.nih.gov/pubmed/16268780] 539
3. Salzberg SL, Yorke Ja: **Beware of mis-assembled genomes.**  
540 *Bioinformatics (Oxford, England)* 2005, **21**(24):4320–4321. [http://www.  
541 ncbi.nlm.nih.gov/pubmed/16332717] 542
4. Bullard JH, Purdom E, Hansen KD, Dudoit S: **Evaluation of statistical**  
543 **methods for normalization and differential expression in mRNA-Seq**  
544 **experiments.** *BMC Bioinformatics* 2010, **11**:94. [http://www.  
545 pubmedcentral.nih.gov/articlerender.fcgi?artid=2838869&tool=  
546 pmcentrez&rendertype=abstract] 547
5. Li H, Lovci MT, Kwon YS, Rosenfeld MG, Fu XD, Yeo GW: **Determination of**  
548 **tag density required for digital transcriptome analysis: application**  
549 **to an androgen-sensitive prostate cancer model.** *Proc Natl Acad Sci*  
550 *USA* 2008, **105**(51):20179–20184. [http://www.pubmedcentral.nih.gov/  
551 articlerender.fcgi?artid=2603435&tool=pmcentrez&rendertype=abstract] 552
6. Blencowe BJ, Ahmad S, Lee LJ: **Current-generation high-throughput**  
553 **sequencing: deepening insights into mammalian transcriptomes.**  
554 *Genes Dev* 2009, **23**(12):1379–1386. [http://www.ncbi.nlm.nih.gov/  
555 pubmed/19528315] 556
7. Cloonan N, Forrest ARR, Kollé G, Gardiner BBA, Faulkner GJ, Brown MK,  
557 Taylor DF, Steptoe AL, Wani S, Bethel G, Robertson AJ, Perkins AC,  
558 Bruce SJ, Lee CC, Ranade SS, Peckham HE, Manning JM, McKernan KJ,  
559 Grimmond SM, Mellissa K, Andrew C, Kevin J: **Stem cell transcriptome**  
560 **profiling via massive-scale mRNA sequencing.** *Nat Methods* 2008,  
561 **5**(7):613–619. [http://www.ncbi.nlm.nih.gov/pubmed/18516046] 562
8. Feldmeyer B, Wheat CW, Krezdorn N, Rotter B, Pfenninger M: **Short read**  
563 **illumina data for the de novo assembly of a non-model snail species**  
564 **transcriptome (Radix balthica, Basommatophora, Pulmonata), and a**  
565 **comparison of assembler performance.** *BMC Genomics* 2011, **12**:317.  
566 [http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=3128070&  
567 tool=pmcentrez&rendertype=abstract] 568
9. Mattila TM, Bechsgaard JS, Hansen TT, Schierup MH, Bilde T: **Orthologous**  
569 **genes identified by transcriptome sequencing in the spider genus**  
570 **Stegodyphus.** *BMC Genomics* 2012, **13**:70. [http://www.pubmedcentral.  
571 nih.gov/articlerender.fcgi?artid=3350440&tool=pmcentrez&rendertype=  
572 abstract] 573
10. Garg R, Patel RK, Tyagi AK, Jain M: **De novo assembly of chickpea**  
574 **transcriptome using short reads for gene discovery and marker**  
575 **identification.** *DNA Res : Int J Rapid Publ Reports Genes Genomes* 2011,  
576 **18**:53–63. [http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=  
577 3041503&tool=pmcentrez&rendertype=abstract] 578
11. Yang D, Fu Y, Wu X, Xie Y, Nie H, Chen L, Nong X, Gu X, Wang S, Peng X,  
579 Yan N, Zhang R, Zheng W, Yang G: **Annotation of the transcriptome**  
580 **from Taenia pisiformis and its comparative analysis with three**  
581 **Taeniidae species.** *PLoS One* 2012, **7**(4):e32283. [http://www.  
582 pubmedcentral.nih.gov/articlerender.fcgi?artid=3326008&tool=  
583 pmcentrez&rendertype=abstract] 584
12. Shi CY, Yang H, Wei CL, Yu O, Zhang ZZ, Jiang CJ, Sun J, Li YY, Chen Q,  
585 Xia T, Wan XC: **Deep sequencing of the Camellia sinensis**  
586 **transcriptome revealed candidate genes for major metabolic**  
587 **pathways of tea-specific compounds.** *BMC Genomics* 2011, **12**:131.  
588 [http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=3056800&  
589 tool=pmcentrez&rendertype=abstract] 590
13. Barrero Ra, Chapman B, Yang Y, Moolhuijzen P, Keeble-Gagnère G,  
591 Zhang N, Tang Q, Bellgard MI, Qiu D: **De novo assembly of Euphorbia**  
592 **fischeriana root transcriptome identifies prostratin pathway related**  
593 **genes.** *BMC Genomics* 2011, **12**:600. [http://www.pubmedcentral.nih.  
594 gov/articlerender.fcgi?artid=3273484&tool=pmcentrez&rendertype=  
595 abstract] 596
14. Wang XW, Luan JB, Li JM, Bao YY, Zhang CX, Liu SS: **De novo**  
597 **characterization of a whitefly transcriptome and analysis of its gene**  
598 **expression during development.** *BMC Genomics* 2010, **11**:400. [http://  
599 www.pubmedcentral.nih.gov/articlerender.fcgi?artid=2898760&tool=  
600 pmcentrez&rendertype=abstract] 601
15. Crawford JE, Guelbeogo WM, Sanou A, Traoré A, Vernick KD, Sagnon N,  
602 Lazzaro BP: **De novo transcriptome sequencing in Anopheles**  
603 **funestus using Illumina RNA-seq technology.** *PLoS One* 2010,  
604 **5**(12):e14202. [http://www.pubmedcentral.nih.gov/articlerender.fcgi?  
605 artid=2996306&tool=pmcentrez&rendertype=abstract] 606

- 607 16. Tarazona S, García-Alcalde F, Dopazo J, Ferrer A, Conesa A: **Differential**  
608 **expression in RNA-seq: a matter of depth.** *Genome Res* 2011,  
609 **21**(12):2213–2223. [http://www.pubmedcentral.nih.gov/articlerender.  
610 fcgi?artid=3227109&tool=pmcentrez&rendertype=abstract]
- 611 17. Hale MC, McCormick CR, Jackson JR, Dewoody JA: **Next-generation**  
612 **pyrosequencing of gonad transcriptomes in the polyploid lake**  
613 **sturgeon (*Acipenser fulvescens*): the relative merits of**  
614 **normalization and rarefaction in gene discovery.** *BMC Genomics* 2009,  
615 **10**:203. [http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=  
616 2688523&tool=pmcentrez&rendertype=abstract]
- 617 18. Zerbino DR, Birney E: **Velvet: algorithms for de novo short read**  
618 **assembly using de Bruijn graphs.** *Genome Res* 2008, **18**(5):821–829.  
619 [http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=2336801&  
620 tool=pmcentrez&rendertype=abstract]
- 621 19. Schulz MH, Zerbino DR, Vingron M, Birney E: **Oases: Robust de novo**  
622 **RNA-seq assembly across the dynamic range of expression levels.**  
623 *Bioinformatics (Oxford, England)* 2012:1–12. [http://www.ncbi.nlm.nih.gov/  
624 pubmed/22368243]
- 625 20. Grabherr MG, Haas BJ, Yassour M, Levin JZ, Thompson Da, Amit I,  
626 Adiconis X, Fan L, Raychowdhury R, Zeng Q, Chen Z, Mauceli E,  
627 Hacohen N, Gnirke A, Rhind N, di Palma F, Birren BW, Nusbaum C,  
628 Lindblad-Toh K, Friedman N, Regev A: **Full-length transcriptome**  
629 **assembly from RNA-Seq data without a reference genome.**  
630 *Nat Biotechnol* 2011, **29**(7):644–652. [http://www.ncbi.nlm.nih.gov/  
631 pubmed/21572440]
- 632 21. Zhao QY, Wang Y, Kong YM, Luo D, Li X, Hao P: **Optimizing de novo**  
633 **transcriptome assembly from short-read RNA-Seq data: a**  
634 **comparative study.** *BMC Bioinformatics* 2011, **12**(Suppl 14):S2.  
635 [http://www.biomedcentral.com/1471-2105/12/S14/S2]
- 636 22. Waterston RH, Lindblad-Toh K, Birney E, Rogers J, Abril JF, Agarwal P,  
637 Agarwala R, Ainscough R, Alexandersson M, An P, Antonarakis SE,  
638 Attwood J, Baertsch R, Bailey J, Barlow K, Beck S, Berry E, Birren B, Bloom T,  
639 Bork P, Botcherby M, Bray N, Brent MR, Brown DG, Brown SD, Bult C,  
640 Burton J, Butler J, Campbell RD, Carninci P, et al: **Initial sequencing and**  
641 **comparative analysis of the mouse genome.** *Nature* 2002,  
642 **420**(6915):520–562. [http://www.ncbi.nlm.nih.gov/pubmed/12466850]
- 643 23. Romiguier J, Ranwez V, Douzery EJP, Galtier N: **Contrasting GC-content**  
644 **dynamics across 33 mammalian genomes: relationship with**  
645 **life-history traits and chromosome sizes.** *Genome Res* 2010,  
646 **20**(8):1001–1009. [http://www.pubmedcentral.nih.gov/articlerender.fcgi?  
647 artid=2909565&tool=pmcentrez&rendertype=abstract]
- 648 24. Dunn C: **Siphonophores.** *Current Biol: CB* 2009, **19**(6):R233–R234.  
649 [http://www.ncbi.nlm.nih.gov/pubmed/19321136]
- 650 25. Tatusov RL, Fedorova ND, Jackson JD, Jacobs AR, Kiryutin B, Koonin EV,  
651 Krylov DM, Mazumder R, Mekhedov SL, Nikolskaya AN, Rao BS, Smirnov S,  
652 Sverdlov AV, Vasudevan S, Wolf YI, Yin JJ, Natale Da: **The COG database:**  
653 **an updated version includes eukaryotes.** *BMC Bioinformatics* 2003,  
654 **4**:41. [http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=  
655 222959&tool=pmcentrez&rendertype=abstract]
- 656 26. Parra G, Bradnam K, Korf I: **CEGMA: a pipeline to accurately annotate**  
657 **core genes in eukaryotic genomes.** *Bioinformatics (Oxford, England)*  
658 **2007, 23**(9):1061–1067. [http://www.ncbi.nlm.nih.gov/pubmed/  
659 17332020]
- 660 27. Parra G, Bradnam K, Ning Z, Keane T, Korf I: **Assessing the gene space in**  
661 **draft genomes.** *Nucleic Acids Res* 2009, **37**:289–297. [http://www.  
662 pubmedcentral.nih.gov/articlerender.fcgi?artid=2615622&tool=  
663 pmcentrez&rendertype=abstract]
- 664 28. Camacho C, Coulouris G, Avagyan V, Ma N, Papadopoulos J, Bealer K,  
665 Madden TL: **BLAST+: architecture and applications.** *BMC Bioinformatics*  
666 **2009, 10**:421. [http://www.pubmedcentral.nih.gov/articlerender.fcgi?  
667 artid=2803857&tool=pmcentrez&rendertype=abstract]

doi:10.1186/1471-2164-14-167

**Cite this article as:** Francis et al.: A comparison across non-model animals suggests an optimal sequencing depth for *de novo* transcriptome assembly. *BMC Genomics* 2013 **14**:167.

**Submit your next manuscript to BioMed Central  
and take full advantage of:**

- Convenient online submission
- Thorough peer review
- No space constraints or color figure charges
- Immediate publication on acceptance
- Inclusion in PubMed, CAS, Scopus and Google Scholar
- Research which is freely available for redistribution

Submit your manuscript at  
www.biomedcentral.com/submit



# Author Query Form

---

**Journal:** BMC Genomics

**Article:** A comparison across non-model animals suggests an optimal sequencing depth for *de novo* transcriptome assembly

Dear Author,

During the copyediting of your paper, the following queries arose. Please respond to these by annotating your proofs with the necessary changes/additions.

- . If you intend to annotate your proof electronically, please refer to the E-annotation guidelines.
- . If you intend to annotate your proof by means of hard-copy mark-up, please refer to the proof mark-up symbols guidelines. If manually writing corrections on your proof and returning it by fax, do not write too close to the edge of the paper. Please remember that illegible mark-ups may delay publication.

Whether you opt for hard-copy or electronic annotation of your proofs, we recommend that you provide additional clarification of answers to queries by entering your answers on the query sheet, in addition to the text mark-up.

| Query No. | Query  | Remark |
|-----------|--|--------|
| Q1        | Upon checking, it is noticed that there are panels inside image of Figure 5 however it is not located within its corresponding caption. Please also mention panels within the figure caption to correspond with the image. |        |
| Q2        | URLs: Please check all URLs if it is working. If not, please provide alternatives.   |        |