# GigaScience

## A draft genome sequence of the elusive giant squid, Architeuthis dux
### --Manuscript Draft--

| Manuscript Number: | GIGA-D-19-00236R1 | |
| --- | --- | --- |
| Full Title: | A draft genome sequence of the elusive giant squid, Architeuthis dux | |
| Article Type: | Data Note | |
| Funding Information: | Villum Fonden (VKR023446) | Dr Rute R. da Fonseca |
| | FP7 People: Marie-Curie Actions (272927) | Dr Rute R. da Fonseca |
| | Fundação para a Ciência e a Tecnologia (PTDC/MAR/115347/2009) | Dr Rute R. da Fonseca |
| | Danmarks Grundforskningsfond (DNRF96) | Dr Rute R. da Fonseca |
| | Programa Operacional Temático Factores de Competitividade (PT) (COMPETE-FCOMP-01-012) | Dr Rute R. da Fonseca |
| | Rede Nacional de Espectrometria de Massa (ROTEIRO/0028/2013) | Dr Hugo Osório |
| | Fundação para a Ciência e a Tecnologia (UID/Multi/04423/2019) | Dr Alexandre Campos |
| | Wellcome Trust (WT108749/Z/15/Z) | Dr Mateus Patricio |
| | Danmarks Grundforskningsfond (DNRF94) | Dr M. Thomas P. Gilbert |
| | Lundbeckfonden (R52-5062) | Dr M. Thomas P. Gilbert |
| | Novo Nordisk Fonden (NNF14CC0001) | Dr Simon Rasmussen |
| | Biotechnology and Biological Sciences Research Council (BB/N020146/1) | Dr Alex Hayward |
| | Biotechnology and Biological Sciences Research Council (BB/M009122/1) | Dr Tobias Baril |
| | Lundbeckfonden (R52-A4895) | Dr Blagoy Blagoev |
| | Nederlandse Organisatie voor Wetenschappelijk Onderzoek (NL) (#825.09.016) | Dr Henk-Jan Hoving |
| | Deutsche Forschungsgemeinschaft (DE) (HO 5569/1-2) | Dr Henk-Jan Hoving |
| | Slovak grant agency VEGA (VEGA 1/0684/16) | Dr Brona Brejova |
| | Slovak grant agency VEGA (VEGA 1/0458/18) | Dr Tomas Vinar |

| Abstract: | Background |
| --- | --- |
| | The giant squid (Architeuthis dux; Steenstrup, 1857) is an enigmatic giant mollusk with a circumglobal distribution in the deep ocean, except in the high Arctic and Antarctic waters. The elusiveness of the species makes it difficult to study. Thus, having a genome assembled for this deep-sea dwelling species will allow unlocking several pending evolutionary questions.Findings |
| | We present a draft genome assembly that includes 200 Gb of Illumina reads, 4 Gb of Moleculo synthetic long-reads and 108 Gb of Chicago libraries, with a final size matching the estimated genome size of 2.7 Gb, and a scaffold N50 of 4.8 Mb. We also present an alternative assembly including 27 Gb raw reads generated using the Pacific Biosciences platform. In addition, we sequenced the proteome of the same individual |

| | |
|---|---|
| | and RNA from three different tissue types from three other species of squid species (Onychoteuthis banksii, Dosidicus gigas, and Sthenoteuthis oualaniensis) to assist genome annotation. We annotated 33,406 protein coding genes supported by evidence and the genome completeness estimated by BUSCO reached 92%. Repetitive regions cover 49.17% of the genome.Conclusions<br><br>This annotated draft genome of A. dux provides a critical resource to investigate the unique traits of this species, including its gigantism and key adaptations to deep-sea environments. |
| Corresponding Author: | Rute R. da Fonseca<br>University of Copenhagen<br>Copenhagen, DENMARK |
| Corresponding Author Secondary Information: | |
| Corresponding Author's Institution: | University of Copenhagen |
| Corresponding Author's Secondary Institution: | |
| First Author: | Rute R. da Fonseca |
| First Author Secondary Information: | |
| Order of Authors: | Rute R. da Fonseca |
| | Alvarina Couto |
| | Andre Machado |
| | Brona Brejova |
| | Caroline B. Albertin |
| | Filipe Silva |
| | Paul Gardner |
| | Tobias Baril |
| | Alex Hayward |
| | Alexandre Campos |
| | Angela Ribeiro |
| | Inigo Barrio Hernandez |
| | Henk-Jan Hoving |
| | Ricardo Tafur-Jimenez |
| | Chong Chu |
| | Barbara Frazão |
| | Bent Petersen |
| | Fernando Peñaloza |
| | Francesco Musacchia |
| | Graham C. Alexander Jr. |
| | Hugo Osório |
| | Inger Winkelmann |
| | Oleg Simakov |
| | Simon Rasmussen |
| | M. Ziaur Rahman |
| | |

| | Davide Pisani |
| --- | --- |
| | Erich Jarvis |
| | Guojie Zhang |
| | Jakob Vinther |
| | Jan Strugnell |
| | L. Filipe C. Castro |
| | Olivier Fedrigo |
| | Mateus Patricio |
| | Qiye Li |
| | Sara Rocha |
| | Agostinho Antunes |
| | Yufeng Wu |
| | Bin Ma |
| | Remo Sanges |
| | Tomas Vinar |
| | Blagoy Blagoev |
| | Thomas Sicheritz-Ponten |
| | Rasmus Nielsen |
| | M. Thomas P. Gilbert |

| Order of Authors Secondary Information: | |
| --- | --- |
| Response to Reviewers: | Dear Editor,<br><br>We herewith submit our revised manuscript 'A draft genome sequence of the elusive giant squid, Architeuthis dux'.<br><br>Regarding the points that you have highlighted, please find the answers below:<br><br>1) Please clarify the rationale for the unconventional assembly strategy in the revised manuscript. If this has "historic" rather than scientific reasons, the reviewer feels this may be fine, but I agree that the reasons should be discussed in the manuscript, for the benefit of readers who are looking for best practice examples.<br><br>The reviewer is correct that there is some degree of history involved. We initially did the assembly without PacBio, and did the presented analyses on this. Later we were offered the chance to try and improve it with PacBio, which we did, but as you can see there was minimal improvement in the assembly statistics (Table 1 and Table S2), but i) an increase of the total genome size to 3.155 Gb, beyond the expected 2.7 Gb estimated in kmergenie, and ii) a slight decrease in the BUSCO completeness assessment. As such, we elected to retain the results based on the original assembly (based on Dovetail), but given that we assume others may wish to use the alternative assembly and explore the differences, we provide both.<br>In the beginning of the "Data analyses" section, we now clearly state which assembly was used in the comparative genomics analyses (from Line 297) and provide an explanation for that choice.<br><br><br>2) Please expand on the methods for protein-coding gene modelling and have another look at your data whether 50K genes may be an overestimate. I also agree with the reviewer's recommendation to analyze gene models in BUSCO to give readers a better idea of their completeness.<br><br>We now expanded the section detailing the filtering of the protein-coding gene set and |

present a total of 33,406 gene annotations in the final set, as these have validation by matching to cephalopod transcripts and/or SwissProt/UniRef90 proteins. We also provide the results from BUSCO when using the gene models as input for comparison (added to Table 1).

Answers to the reviewer's comments:

Reviewer #1: In this study, de Fonseca et al. report the genome of the giant squid as a resource to investigate the unique traits of this fascinating organism. Two assemblies, which are of comparable contiguity to most other recently published molluscan genomes, as well as a set of over 51,000 gene models are reported. Analysis of the genome focuses on repetitive elements (e.g., TEs), non-coding RNAs, and gene families of interest to the authors (WNT genes, Protocadherins, Hox genes, and reflectins). Overall this is a straightforward study that provides a resource that will be broadly useful and I feel it should be published. However, I have a number of suggestions for improvement including a few important issues that need to be addressed.

Major points:

1.1. It is unclear why two different genome assemblies are presented instead of just one most optimal assembly. This is not the way I would have gone about assembling this combination of data but presumably Dovetail scaffolding and gene modelling were performed before PacBio sequencing and scaffolding? Re-doing the assembly would a more logical way would probably have relatively little improvement but a little more explanation of the rationale or 'historical' reasons for two different assemblies and/or this assembly strategy would be a helpful addition to readers looking in the literature for examples on best practices for genome assembly.

Thank you for this comment. The reviewer is correct that there is some degree of history involved. We initially did the assembly without PacBio, and did the presented analyses on this. Later we were offered the chance to try and improve it with PacBio, which we did, but as you can see there was minimal improvement in the assembly statistics (Table 1 and Table S2), but i) an increase of the total genome size to 3.155 Gb, beyond the expected 2.7 Gb estimated in kmergenie, and ii) a slight decrease in the BUSCO completeness assessment. As such, we elected to retain the results based on the Dovetail assembly, but given that we assume others may wish to use the alternative assembly and explore the differences, we provide both.

1.2. Related to this issue, there is little comparison of the two genome assemblies and it is unclear which assembly was used for what analyses and even Table 1 and Table S2's titles are a bit ambiguous with respect to which assembly statistics are presented. Please explicitly state which assembly was used for which analyses.

In the beginning of the "Data analyses" section, we now clearly state which assembly was used in the comparative genomics analyses (from Line 297) and provide an explanation for that choice. Additionally, we also mention the choice in the Methods section (Lines 183 to 185) before describing the strategies for annotation and comparative analyses.

1.3. The approach used for gene annotation is unconventional and the inferred number of protein-coding gene models is very high. This does not mean the gene model set is bad, but I feel that data needed for the reader to assess the quality of the gene models are lacking. Please run BUSCO on the gene models and report these data as well.

We now also provide the results from BUSCO when using the gene models as input for comparison.

1.4. Specimen collection data are not reported in the manuscript.

This information has now been added to Table S1.

Minor points:

1.5. Scientific names of species need to be italicized throughout.

Done.

1.6. Did all the giant squid DNA come from the same individual?

Yes, this is now clear in Line 172.

1.7. Lines 140-141: "currently increasing locally" is a bit awkward and vague.

Replaced by "in some regions".

1.8. Line 176: Which reads? All Illumina reads? PE reads only?

This has now been clarified on Line 176.

1.9. Line 185: Again, this seems to me to be a strange assembly strategy and I think that it should be clearly stated that PacBio data became available 'late in the game' if that is the case. Otherwise, the logic behind this assembly strategy needs to be explained.

In the beginning of the "Data analyses" section, we now clearly state which assembly was used in the comparative genomics analyses (from Line 297) and provide an explanation for that choice.


1.10. Line 199: High-throughput is misspelled.

Done.

1.11. Line 203: Clarify what is meant by reference transcriptome. All reads from all tissues were pooled and assembled together?

This has now been clarified in Lines 203-204.

1.12. Line 205: "EvidencialGene" is a tyo.

Corrected.

1.13. Lines 261-262: Please provide details on exactly what was done in this study in the supplementary material. Description of how the final gene models were selected is vague.

We now further discuss the filters applied in lines 272-275. The total number of protein-coding genes passing all the filters is 33,406.

1.14. Line 277: What is meant by a "bespoke pipeline"? Custom scripts should be made available.

No custom analysis scripts were developed. We simply use 'bespoke' to mean 'tailored to our particular purpose'. Here this refers an analysis pipeline combining: a preliminary analysis using RepeatMasker, followed by a de novo analysis using RepeatModeler and a referenced and publicly available script by Platt et al, followed by a full annotation using RepeatMasker. These steps are fully outlined and referenced in the methods section. We have simplified the sentence which now reads:
"Repetitive elements were first identified using RepeatMasker v.4.0.8"

1.15. Line 450: Correct "Sampling was following"

Done.

1.16. BUSCO results are presented in the methods section (should be in the results by the way) for the pre-PacBio scaffolding genome but not the post-PacBio scaffolding

genome.

The results of BUSCO for post-PacBio step are presented in Table S2 (as indicated in Line 186). We moved the description of the BUSCO results to the "Data analyses" section and added a clarification regarding the choice of the assembly for the overall comparative genomics analyses (from Line 297).

1.17. Table 1: BUSCO should be in all capital letters.

Done.

1.18. Figure 3: What does the note "Gene size only" mean?

This gene was reported to be fully isolated from other Hox genes in a different scaffold but was not alone in the scaffold. There were other non-Hox genes. Figure 3 aims to show both the organisation and the range occupied by Hox genes. Considering the organisation, the gene is isolated such as in O. bimaculoides. Regarding the size, the schematic representation indicates only the Hox "cluster" area. In O. bimaculoides, the scaffolds contain only the Hox genes. This means it could be possible for the cluster to be there but only when considering a very vast distance. In this scenario for C. teleta, the gene is found in the middle of the scaffold, surrounded by other genes. It is not part of the cluster. Indicating the full scaffold size could lead to a wrong interpretation of the gene size and of the Hox gene range. As such, only the gene size is indicated.

1.19. Table S1: Please provide total number of reads and somewhere it should be clarified how many different instrument runs were conducted and if different libraries were multiplexed on the Illumina platform.

This information has now been added to Table S1.

Reviewer #2: The authors present the genome of the giant squid Architeuthis dux. Several cephalopod genomes have been sequenced, but our genomic understanding of cephalopods living in the deep-sea environment is still poor. The authors sequenced a giant squid species A. dux together with several transcriptomes from the gonad, liver and brain tissues derived from three other squid species including Onychoteuthis banksii, Dosidicus gigas, and Sthenoteuthis oualaniensis.

Having a giant squid genome is an important contribution to the field of cephalopod genomics, especially for further meaningful comparative genomics. The authors provide a decent genome assembly. And the observation of a non-tightly physically linked Hox cluster is interesting. The manuscript is well written in general, however, there are a lot of editing errors throughout the whole manuscript, which distracts the reading. The authors need to carefully fix all these typos and errors during the revision. Further comments are provided below.

Major comments:

2.1.     In the Abstract/Findings, there is a lot of information about "Methods" (e.g. how many raw reads, sequencing of proteome and RNA) instead of what the authors found from the genome itself. Also, the statement "RNA from three different tissue types from three other species of squid to assist genome annotation." is very vague. What tissue types from what species should be clearly described. The authors need to rewrite this section.

In the abstract we followed the format that is usual in a data note, providing detailed information on the data provided by this work. We have now added the names of the three species of squid to the abstract.

2.2.     Line 153: Body patterning system? Usage of body patterning is confusing here since body patterning often refers to the developmental process during embryogenesis but not the skin color pattern.

We have rephrased the sentence to: "Cephalopods can rapidly alter the texture,

pattern, colour and brightness of their skin, and this both enables a complex communication system, as wells as provides exceptional camouflage and mimicry."

2.3.    The authors cited that there is a global proliferation of cephalopods (Lines 140 and 141) but later cited other studies saying that there is a regional extinction. It is a bit confusing whether cephalopods are undergoing proliferation or extinction. Given that the earlier citation is more recent (Doubleday et al., 2016) than others, it is wondering which condition is closer to the current situation.

We have removed the second statement to avoid confusion.

2.4.    Although it is agreeable in general to have genome resources from unexplored species, the authors' argument in the last paragraph of Data description/Context is not convincing. The link between having a genome and aiding conservation efforts as well as ensuring continued existence is not clear.

Without a genome, population genomic studies that provide information regarding the genetic diversity and structure of populations becomes very challenging, with genome-wide data having to be produced from reduced-representation methods that have many biases. In this last paragraph, we state this specifically: "A genome is an important resource for future population genomics studies[…]".

2.5.    Do the authors have any idea why the genome contains so many protein-coding genes (51,225 genes predicted) in comparing to other cephalopod species usually having only 20,000-30,000 genes? For example, is it due to that A. dux has more lineage-specific genes or expansions of certain gene families?

We have revised our gene models and now further discuss the filters applied in lines 272-275. The total number of protein-coding genes passing all the filters is 33,406.

2.6.    Given that genome size and polyploidy of the organisms are often correlated to increased body size (Session et al., 2016), have the authors checked if there is whole-genome duplication or polyploidy in the A. dux genome? Session et al. (2016) Genome evolution in the allotetraploid frog Xenopus laevis. Nature 538, 336-343.

We did confirm that the genome was not polyploid by testing for Hardy–Weinberg equilibrium using re-sequencing data from 32 giant squid individuals (Winkelman et al, unpublished results) and there is no evidence for an ancient duplication since we only found one intact Hox complement.

2.7.    Figure 3: The authors should provide scaffold numbers for the Hox clusters from each species. Also, in most cases, Hox genes in the Hox cluster are adjacent to each other without the insertion of other non-Hox genes. If there is a special case in A. dux and E. scolopes, the authors should show the real gene arrangement on that scaffold, especially for the non-Hox genes (with brief annotation) that are in between Hox genes. This can be achieved by having an additional panel in the same figure. The authors are encouraged to show an illustration on the types of Hox gene organization in order to give the readers a better understanding of this context.

Figure 3 has received new panels. Scaffold information for A. dux was added in panel C (Figure 3-C). As the assemblies of the other species were retrieved from other studies, the readers are directed to the appropriate references for further detail. An extra panel depicting the Hox cluster organisation in more detail has been added. E. scolopes data is shown as reported in its published study. No non-Hox genes were indicated for the area covered in this representation. An additional panel with a simplified version of the various Hox "cluster" types was inserted in panel A (Figure 3-A).

Minor comments:

2.1.1.    Line 149: ~2cm -> "~2 cm"
Done.

2.1.2.     Line 150: 3 orders -> "three orders"
Done.

2.1.3.     Line 150: Architeuthis dux -> "A. dux"
Done.

2.1.4.     Lines 150 and 151: 10-12cm… 20m -> "10-12 cm… 20 m"
Done.

2.1.5.     Line 152: 500kg -> "500 kg"
Done.

2.1.6.     Line 171: a Architeuthis dux sample -> "an A. dux sample"
Done.

2.1.7.     Line 172: What is CTAB?
CTAB = "cetyl trimethylammonium bromide"; this description has been included in the text (Line 172)

2.1.8.     Line 184: For Eukaryota and Metazoa we identified… -> "For Eukaryota and Metazoa, we identified…"
Done.

2.1.9.     Line 184: … 90.4 % and 92.1 %... -> "… 90.4% and 92.1%..."
Done.

2.1.10.    Line 185: 23.38Gb -> "23.38 Gb"
Done.

2.1.11.    Line 186: 14.79kb -> "14.79 kb"

Done.

2.1.12.    "k-mer" (Line 204) or "kmer" (Line 176) to be consistent.
Chose to use "kmer".

2.1.13.    Line 216: 100,000 g -> "100,000×g"
Done.

2.1.14.    Lines 219 and 222: SDS-PAGE -> "SDS-PAGE" (hyphen but not en dash)
Done.

2.1.15.    Line 221: Tris - HCl -> Tris-HCl (single hyphen but not en dash with spaces)
Done.

2.1.16.    Line 226: LC-MS/MS analyses -> "LC-MS/MS analyses" (hyphen but not en dash)
Done.

2.1.17.    Line 254: Using italic for scientific names (i.e. Octopus bimaculatus, Crassostrea gigas, and Lottia gigantea)
Done.

2.1.18.    Line 260: … 200kb (total length 199Mb)… -> "… 200 kb (total length 199 Mb)…"
Done.

2.1.19.    Line 290: Transposable elements -> "transposable elements"
Done.

2.1.20.    Line 300: Architeuthis dux -> "A. dux"
Done.

2.1.21.    Line 323: ~5-8% -> "~5¬-8%" (en dash but not hyphen for a range)

Done.

2.1.22.     Line 381: Octopus bimaculoides -> "O. bimaculoides"
Done.

2.1.23.     Line 383: Euprymna scolopes -> "E. scolopes" (in italic)
Done.

2.1.24.     Line 395: Euprymna scolopes -> "E. scolopes"
Done.

2.1.25.     Lines 397 & 398: 500 - 10,000 kb -> "500-10,000 kb" (en dash but not hyphen for a range)
Done.

2.1.26.     Line 406: ... observed in Hox 1, Hox 4, ANTP, Lox 2, Lox 5, Post 1 and Post 2. Hox 1 did,... -> "... observed in Hox1, Hox4, ANTP, Lox2, Lox5, Post1 and Post2. Hox1 did,..."
Done.

2.1.27.     Line 407: Hox 1 -> "Hox1"
Done.

2.1.28.     Line 408: Hox 3 -> "Hox3"
Done.

2.1.29.     Line 409: Lox 4 -> "Lox4"
Done.

2.1.30.     Lines 410, 412 & 413: Hox 2 -> "Hox2"
Done.

2.1.31.     Line 421: ... contains 7 reflectin genes and 3 reflectin-like genes… -> "... contains seven reflectin genes and three reflectin-like genes…"
Done.

2.1.32.     Line 422: … exception of 1 reflectin gene, … -> "… exception of one reflectin gene, …"
Done.

2.1.33.     Line 436: … (tsa)… -> "… (TSA)…"
Done.

2.1.34.     Lines 647 & 657: Architeuthis dux -> "A. dux"
Done.

2.1.35.     Line 659: Hox 2 -> "Hox2"
Done.

| Additional Information: | |
|---|---|
| Question | Response |
| Are you submitting this manuscript to a special series or article collection? | No |
| **Experimental design and statistics**<br><br><br>Full details of the experimental design and statistical methods used should be given in the Methods section, as detailed in our [Minimum Standards Reporting Checklist](). Information essential to interpreting the | Yes |

| | |
|---|---|
| data presented should be made available in the figure legends.<br><br>Have you included all the information requested in your manuscript? | |
| **Resources**<br><br>A description of all resources used, including antibodies, cell lines, animals and software tools, with enough information to allow them to be uniquely identified, should be included in the Methods section. Authors are strongly encouraged to cite Research Resource Identifiers (RRIDs) for antibodies, model organisms and tools, where possible.<br><br>Have you included the information requested as detailed in our Minimum Standards Reporting Checklist? | Yes |
| **Availability of data and materials**<br><br>All datasets and code on which the conclusions of the paper rely must be either included in your submission or deposited in publicly available repositories (where available and ethically appropriate), referencing such data using a unique identifier in the references and in the "Availability of Data and Materials" section of your manuscript.<br><br>Have you have met the above requirement as detailed in our Minimum Standards Reporting Checklist? | Yes |

1    # A draft genome sequence of the elusive giant squid, *Architeuthis dux*

2

3    Rute R. da Fonseca*[1,2], Alvarina Couto[3], Andre Machado[4], Brona Brejova[5], Carolin B. Albertin[6], Filipe

4    Silva[4], Paul Gardner[7], Toby Baril[8], Alex Hayward[8], Alexandre Campos[4], Ângela Ribeiro[4], Inigo Barrio

5    Hernandez[9], Henk-Jan Hoving[10], Ricardo Tafur-Jimenez[11], Chong Chu[12], Barbara Frazão[4,13], Bent

6    Petersen[14,15], Fernando Peñaloza[16], Francesco Musacchia[17], Graham C. Alexander Jr. [18], Hugo

7    Osório[19,20,21], Inger Winkelmann[22], Oleg Simakov[23], Simon Rasmussen[24], M. Ziaur Rahman[25], Davide

8    Pisani[26], Jakob Vinther[27], Erich Jarvis[28], Guojie Zhang[30,31,32,33], Jan Strugnell[34], L. Filipe C. Castro[4,36], Olivier

9    Fedrigo[28], Mateus Patricio[29], Qiye Li[37], Sara Rocha[3], Agostinho Antunes[4,36], Yufeng Wu[38], Bin Ma[39], Remo

10    Sanges[40,41], Tomas Vinar[5], Blagoy Blagoev[9], Thomas Sicheritz-Ponten[14,15], Rasmus Nielsen[22,42], M. Thomas

11    P. Gilbert[22,43]

12

13    [1]Center for Macroecology, Evolution and Climate, Natural History Museum of Denmark, University of

14    Copenhagen, Copenhagen, Denmark.

15    [2]The Bioinformatics Centre, Department of Biology, University of Copenhagen, Copenhagen, Denmark.

16    [3]Department of Biochemistry, Genetics and Immunology, University of Vigo, Vigo, Spain.

17    [4]CIIMAR, Interdisciplinary Centre of Marine and Environmental Research, University of Porto, Porto,

18    Portugal.

19    [5]Faculty of Mathematics, Physics and Informatics, Comenius University in Bratislava, Bratislava, Slovak

20    Republic.

21    [6]Department of Organismal Biology and Anatomy, University of Chicago, Chicago, USA.

22    [7]Department of Biochemistry, University of Otago, New Zealand.

23    [8]Centre for Ecology and Conservation, University of Exeter, Penryn Campus, Cornwall, UK.

24    [9]Department of Biochemistry and Molecular Biology, University of Southern Denmark, Odense,

25    Denmark.

26    [10]GEOMAR Helmholtz Centre for Ocean Research Kiel, Kiel, Germany.

27    [11]Instituto del Mar del Perú.

28    [12]Department of Biomedical Informatics, Harvard Medical School, Boston, USA.

29    [13]IPMA, Fitoplâncton Lab, Lisboa, Portugal.

30    [14]Centre of Excellence for Omics-Driven Computational Biodiscovery (COMBio), Faculty of Applied

31    Sciences, AIMST University, Kedah, Malaysia.

32    [17]Genomic Medicine, Telethon Institute of Genetics and Medicine, Pozzuoli, Naples, Italy

33    [18]GCB Sequencing and Genomic Technologies Shared Resource, Duke University, Durham, NC, USA.

34    [19]i3S-Instituto de Investigação e Inovação em Saúde, Universidade do Porto, Porto, Portugal.

35    [20]IPATIMUP -Institute of Molecular Pathology and Immunology, University of Porto, Porto, Portugal.

36    [21]Faculty of Medicine of the University of Porto, Porto, Portugal.

37    [22]Centre for GeoGenetics, Natural History Museum of Denmark, University of Copenhagen, Copenhagen,

38    Denmark.

39    [23]Department of Molecular Evolution and Development, University of Vienna, Vienna, Austria.

40    [24]Novo Nordisk Foundation Center for Protein Research, Faculty of Health and Medical Sciences,

41    University of Copenhagen, Copenhagen, Denmark

42    [25]Bioinformatics Solutions Inc, Waterloo, Ontario, Canada.

43    [26]Departments of Biological sciences and Earth Sciences, University of Bristol, Bristol, UK.

44    [27]Departments of Biological sciences and Earth Sciences, University of Bristol, Bristol, UK.

45    [28]The Rockefeller University, New York, USA.

46    [29]European Molecular Biology Laboratory, European Bioinformatics Institute (EMBL-EBI), Wellcome

47    Genome Campus, Hinxton, UK.

48    [30]Section for Ecology and Evolution, Department of Biology, University of Copenhagen, Copenhagen,

49    Denmark.

50    [31]China National Genebank, BGI-Shenzhen, Shenzhen, China.

51    [32]State Key Laboratory of Genetic Resources and Evolution, Kunming Institute of Zoology, Chinese

52    Academy of Sciences, Kunming, China.

53    [33]CAS Center for Excellence in Animal Evolution and Genetics, Chinese Academy of Sciences, Kunming,

54    China.

55    [34]Centre for Sustainable Tropical Fisheries & Aquaculture, James Cook University, Townsville,

56    Queensland, Australia

57    [35]Department of Ecology, Environment and Evolution, School of Life Sciences, La Trobe University,

58    Melbourne, Victoria, Australia

59    [36]Department of Biology, Faculty of Sciences, University of Porto, Portugal.

60    [37]BGI-Shenzhen, Shenzhen, China

61    [38]Department of Computer Science and Engineering, University of Connecticut, Storrs, USA.

62    [39]School of Computer Science, University of Waterloo, Canada.

63    [40]Area of Neuroscience, Scuola Internazionale Superiore di Studi Avanzati (SISSA), Trieste, Italy.

64    [41]Biology and Evolution of Marine Organisms, Stazione Zoologica Anton Dohrn, Napoli, Italy.

65    [42]Departments of Integrative Biology and Statistics, University of California, Berkeley, U.S.A.

66    [43]Norwegian University of Science and Technology, University Museum, Trondheim, Norway

67

68    Email addresses:

69    Rute R da Fonseca: rfonseca@bio.ku.dk (corresponding author)
70    Alvarina Couto: alvarinacouto@gmail.com
71    Andre M.Machado: andre.machado@ciimar.up.pt
72    Brona Brejova: brejova@fmph.uniba.sk
73    Caroline B.Albertin: calbertin@mbl.edu

74   Filipe Silva: filipecgsilva@gmail.com
75   Paul Gardner: paul.gardner@otago.ac.nz
76   Tobias Baril: tb529@exeter.ac.uk
77   Alex Hayward Hayward: Alex.Hayward@exeter.ac.uk
78   Alexandre Campos: acampos@ciimar.up.pt
79   Ângela Ribeiro: ribeiro.angela@gmail.com
80   Inigo Barrio Hernandez: ibarrioh@ebi.ac.uk
81   Henk-Jan Hoving: hhoving@geomar.de
82   Ricardo Tafur-Jiménez: rtafur@imarpe.gob.pe
83   Chong Chu: Chong_Chu@hms.harvard.edu
84   Barbara Frazão: bmfrazao@gmail.com
85   Bent Petersen: bent.petersen@bio.ku.dk
86   Fernando Peñaloza: fpenaloz@lcg.unam.mx
87   Francesco Musacchia: f.musacchia@tigem.it
88   Graham C. Alexander Jr.:gca2@duke.edu
89   Hugo Osório: hosorio@ipatimup.pt
90   Inger E. Winkelmann: inger.winkelmann@gmail.com
91   Oleg Simakov: oleg.simakov@univie.ac.at
92   Simon Rasmussen: simon.rasmussen@cpr.ku.dk
93   M. Ziaur Rahman: zrahman@bioinfor.com
94   Davide Pisani: Davide.Pisani@bristol.ac.uk
95   Erich D. Jarvis: ejarvis@rockefeller.edu
96   Guojie Zhang: zhanggjconi@gmail.com
97   Jakob Vinther: vinther.jakob@gmail.com
98   Jan M. Strugnell: jan.strugnell@jcu.edu.au
99   L. Filipe C. Castro: filipe.castro@ciimar.up.pt
100  Olivier Fedrigo: ofedrigo@rockefeller.edu
101  Mateus Patricio: mateus@ebi.ac.uk
102  Qiye Li: liqiye@genomics.cn
103  Sara Rocha: sprocha@gmail.com
104  Agostinho Antunes: aantunes@ciimar.up.pt
105  Yufeng Wu: ywu@engr.uconn.edu
106  Bin Ma: binma@uwaterloo.ca
107  Remo Sanges: remo.sanges@gmail.com
108  Tomas Vinar: tomas.vinar@fmph.uniba.sk
109  Blagoy Blagoev: bab@bmb.sdu.dk
110  Thomas Sicheritz-Ponten: thomassp@bio.ku.dk
111  Rasmus Nielsen: rasmus_nielsen@berkeley.edu
112  M. Thomas P. Gilbert: tgilbert@snm.ku.dk

## Abstract

### Background

113 

114 

115 The giant squid (*Architeuthis dux*; Steenstrup, 1857) is an enigmatic giant mollusk with a circumglobal

116 distribution in the deep ocean, except in the high Arctic and Antarctic waters. The elusiveness of the

117 species makes it difficult to study. Thus, having a genome assembled for this deep-sea dwelling species

118 will allow unlocking several pending evolutionary questions.

### Findings

119 

120 We present a draft genome assembly that includes 200 Gb of Illumina reads, 4 Gb of Moleculo synthetic

121 long-reads and 108 Gb of Chicago libraries, with a final size matching the estimated genome size of 2.7

122 Gb, and a scaffold N50 of 4.8 Mb. We also present an alternative assembly including 27 Gb raw reads

123 generated using the Pacific Biosciences platform. In addition, we sequenced the proteome of the same

124 individual and RNA from three different tissue types from three other species of squid species

125 (*Onychoteuthis banksii*, *Dosidicus gigas*, and *Sthenoteuthis oualaniensis*) to assist genome annotation.

126 We annotated 33,406 protein coding genes supported by evidence and the genome completeness

127 estimated by BUSCO reached 92%. Repetitive regions cover 49.17% of the genome.

### Conclusions

128 

129 This annotated draft genome of *A. dux* provides a critical resource to investigate the unique traits of this

130 species, including its gigantism and key adaptations to deep-sea environments.

### Keywords

131 

132 Cephalopod, invertebrate, genome assembly.

133

## Data description

### Context

Cephalopods are the most behaviourally complex of the invertebrate protostomes [1] . Their large, highly differentiated brains are comparable in relative size and complexity to those of vertebrates [2], as are their cognitive capabilities [1]. Cephalopods are distributed worldwide from tropical to polar marine habitats, from benthic to pelagic zones and from intertidal areas down to the abyssal parts of the deep sea, with the only exception being the Black Sea. Cephalopod populations are thought to be currently increasing in some regions for a variety of reasons [3], including potential predator release as a consequence of the depletion of fish stocks [4]. The class Cephalopoda contains approximately 800 species, with the vast majority belonging to the soft-bodied subclass Coleoidea (cuttlefishes, octopuses and squids), and a small handful belonging to the Nautiloidea (nautiluses) [5]. Cephalopods are ecologically important as a primary food source for marine mammals, birds and for many fish species. They are also increasingly important as a high-protein food source for humans and are a growing target for commercial fisheries and farming [6].

Cephalopods show a wide variety of morphologies, lifestyles and behaviours [7], but with the exception of the nautiluses they are characterised by having rapid growth and short lifespans, despite a considerable investment in costly sensory adaptations [2]. They range in size from the tiny pygmy squids (~2 cm) to animals that are nearly three orders of magnitude larger, such as the giant squid, *A. du*x (average length 10–12 m, and reported up to 20 m total length) [6,8,9], to the colossal squid, *Mesonychoteuthis hamiltoni* (maximum length remains unclear, but a recorded weight of 500 kg makes it the largest known invertebrate [10]). Cephalopods can rapidly alter the texture, pattern, colour and brightness of their skin, and this both enables a complex communication system, as wells as provides exceptional camouflage and mimicry [11]. Together these allow cephalopods to both avoid predators, and hunt prey highly efficiently, making them some of the top predators in the ocean. The remarkable adaptations of cephalopods also

158    extend to their genome, with recent work demonstrating increased levels of RNA editing to diversify

159    proteins involved in neural functions [12].

160    Over recent years, oceanic warming and acidification, pollution, expanding hypoxia and fishing [13–15]

161    have been shown to affect cephalopod populations. Mercury has been found in high concentrations in

162    the tissue of giant squid specimens [16], and accumulation of flame retardant chemicals has also been

163    detected in the tissue of deep-sea cephalopods [17]. Consequently, there is an urgent need for greater

164    biological understanding of these important, but rarely encountered animals, in order to aid conservation

165    efforts and ensure their continued existence. A genome is an important resource for future population

166    genomics studies aiming at characterizing the diversity of the legendary giant squid, the species which has

167    inspired generations to tell tales of the fabled Kraken.

168

## Methods

169

*DNA extraction, library building, and de novo genome assembly*

170

171    High-molecular-weight genomic DNA was extracted from a single *A. du*x individual (NCBI taxon id:

172    256136) using a cetyl trimethylammonium bromide (CTAB) based buffer followed by organic solvent

173    purification, following Winkelmann et al [18] (details in the Supplementary Information). We generated

174    116 Gb of raw reads from Illumina short-insert libraries, 76 Gb of paired-end reads from libraries ranging

175    from 500 bp to 800 bp in insert size, and 5.4 Gb of mate-pair with a 5 kb insert (Table S1). Furthermore,

176    we generated 3.7 Gb of paired-end reads using Moleculo libraries (3 High-Throughput libraries and 4

177    High-Fidelity libraries). The kmer distribution of the reads under a diploid model in kmergenie [19]

178    predicted the genome size to be 2.7 Gb.

179    An initial assembly generated with Meraculous [20] using Illumina and Moleculo data (N50 of 32 Kb,

180    assembly statistics in Table S2) was used as input for Dovetail Genomic's HiRise scaffolding software

181    together with the Hi-C data generated from two Chicago libraries corresponding to a physical coverage

182    of the genome of 52.1X. This "Meraculous + Dovetail" assembly (statistics in Table 1) was the one used

183    for the genome annotation (non-coding RNAs, protein-coding genes and repeats) and comparative

184    genomics analyses presented in this paper. Further scaffolding was done using 23.38 Gb of PacBio reads

185    (19 SMRT cells, average read length is 14.79 kb) using the default parameters in PBJelly [21] (see

186    assembly statistics in Table S2). The genome gene content completeness was evaluated through the

187    Benchmarking Universal Single-Copy Orthologs (BUSCO v.3.0.2, datasets: Eukaryota, Metazoan) [22].

188    *Transcriptome sequencing and de novo assembly*
189    Given the extreme rarity of live giant squid sightings, we were unable to collect fresh organ samples

190    (following the recommendations in [23]) containing intact RNA from the species to assist with the

191    genome annotation. As an alternative, we extracted total RNA from gonad, liver and brain tissue from

192    live caught specimens of three other oegopsid squid species (*Onychoteuthis banksii*, *Dosidicus gigas*, and

193    *Sthenoteuthis oualaniensis*; NCBI taxon ids 392296, 346249 and 34553, respectively; Supplementary

194    Figure S1), using the Qiagen RNeasy extraction kit (Qiagen,CA, USA). The RNA integrity and quantity was

195    measured on a Qubit fluorometer (Invitrogen, OR, USA) and on the Agilent Bioanalyzer 2100 (Agilent,

196    CA, USA). The Illumina TruSeq Kit v.2.0 was used to isolate the mRNA and prepare cDNA libraries for

197    sequencing, following the recommended protocol. Compatible index sequences were assigned to

198    individual libraries to allow for multiplexing on four lanes of 100bp paired-end technology on an Illumina

199    HiSeq 2000 flow cell. Sequencing of the cDNA libraries was done at the National High-Throughput

200    Sequencing Center at the University of Copenhagen in Denmark. We assessed the quality of the raw

201    reads using FastQC v0.10.0 [24]. After removing indexes and adaptors with CutAdapt [25], we trimmed

202    the reads with the FASTX-toolkit (http://hannonlab.cshl.edu/fastx_toolkit) removing bases with a Phred-

203    scale quality score lower than 25. Reference transcriptomes for each species were built after pooling the

204    reads from all tissues and using these as input in Trinity [26]. This software was used with the default

205    settings including a fixed kmer size of 25 as suggested by the authors. Annotation of coding regions was

206  done with the EvidentialGene pipeline [27].

207  *Protein extraction, separation by 1D SDS–PAGE, MALDI-TOF/TOF and Protein Identification*

208  Given the practical impossibility of obtaining RNA from a giant squid specimen, we produced a library of

209  giant squid peptide sequences to guide the gene annotation process.

210  Proteins were solubilised from a giant squid mantle tissue sample according to the procedure described

211  by Kleffmann et al. [28] and employing the following buffers: (1) 40 mM Tris–HCl, 5 mM MgCl2 and 1

212  mM DTT, pH 8.5; (2) 8 M urea, 20 mM Tris, 5 mM $MgCl_2$ and 20 mM DTT; (3) 7 M urea, 2 M thiourea, 20

213  mM Tris, 40 mM DTT, 2% CHAPS (w ∕ v) and 1% Triton X-100 (v ∕ v) and (4) 40 mM Tris, 4% SDS (w ∕ v)

214  and 40 mm DTT. All buffers were augmented with protease inhibitors (Halt™ Protease Inhibitor Cocktail,

215  EDTA-Free, Thermo Scientific). Tissue samples were ground in liquid nitrogen before homogenization, or

216  homogenized directly with ultrasound (probe sonication at 60 Hz, for 3 min) in buffer 1. Solubilised

217  proteins were collected by ultracentrifugation at 100,000xg and 4 ºC. Each extraction was performed in

218  duplicate for each specific buffer and extracts were pooled. Protein extracts were subsequently stored

219  at -20 ºC. Total protein content was estimated according to the Bradford (1976) method [29].

220  Protein separation by 1D SDS-PAGE electrophoresis was carried out as described in Santos et al. [30]. 53

221  μL of sample (39 μg protein) was diluted in 72 μL of Loading Buffer (0.01% bromophenol blue, 2% SDS

222  (Sodium-DodecylSulfate), 20% glycerol, 5% β-mercaptoethanol (w/v/v) in 62.5 mM Tris-HCl, pH 6.8). The

223  resulting solution was heated for 3 min at 99°C. Proteins were separated by SDS–PAGE with 12% (w/v)

224  polyacrylamide gels. Electrophoresis was carried out using the mini Protean Cell (BioRad) at a constant

225  voltage of 150 V. The separated proteins were visualized by staining with Colloidal Coomassie Brilliant

226  Blue (CCB) [31], and lanes were cut into 15 gel sections for subsequent LC-MS/MS analysis.

227  *LC-MS/MS analyses*

228  All samples were analysed with the Easy-nLC system (Thermo Fisher Scientific), connected online to a Q

229  Exactive mass spectrometer (Thermo Fisher Scientific) equipped with a nanoelectrospray ion source

230    (Thermo Fisher Scientific). Tryptic peptides were loaded in a fused silica column (75 µm inner diameter)

231    packed with C18 resin (3-µm beads, Reprosil, Dr. Maisch), with solvent A (0.5% acetic acid). They were

232    then eluted with a 120 minute gradient of solvent B (80% ACN, 0.5% acetic acid) with a constant flow of

233    250 nL/min. The Q exactive was operated in positive mode with a capillary temperature of 250 °C, using

234    the data dependent acquisition method, which switches from full MS scans to MS/MS scans for the 12

235    most intense ions. Fragmentation was achieved by higher-energy collisional dissociation (HCD) with a

236    normalized collisional energy (NCE) of 25. Full MS ranged from 300 to 1750 m/z at a resolution of

237    70,000, an Automatic Gain Control (AGC) of 1e6 and a maximum injection time of 120 ms, whereas

238    MS/MS events were scanned at a resolution of 35,000, an AGC of 1e5, maximum injection time of 124

239    ms, isolation windows of 2 m/z and an exclusion window of 45 seconds.

240    *de novo peptide prediction*
241    Raw LC-MS/MS data were read using Thermo Fisher MSRawFileReader 2.2 library and imported into

242    PEAKS Studio 7.0 and subsequently preprocessed for precursor mass and charge correction, MS/MS de-

243    isotoping, and deconvolution. PEAKS de novo sequencing [31] was performed on each refined MS/MS

244    spectrum with a precursor and fragment ion error tolerance of 7 ppm and 0.02 da respectively.

245    Carbamidomethylation (Cys) was set as a fixed modification and oxidation (Met) and N-terminal

246    Acetylation as variable modifications. At most, five variable modifications per peptide were allowed. For

247    each tandem spectrum, five *de novo* candidates were reported along with their Local Confidence Scores

248    (the likelihood of each amino acid assignment in a de novo candidate peptide). This score was used to

249    determine the accuracy of the de novo peptide sequences. The top de novo peptide for each spectrum

250    was determined by the highest Average Local Confidence score (ALC) among the candidates for that

251    spectrum.

252    *Genome annotation*
253    Protein-coding genes were predicted by ExonHunter [32] , which combines probabilistic models of

254    sequence features with external evidence from alignments. As external evidence, we have used known

255     proteins from *Octopus bimaculatus*, *Crassostrea gigas* (Pacific oyster) and *Lottia gigantea* (Giant owl

256     limpet) and the predicted proteins encoded by the transcriptomes of the three other oegopsid species

257     analysed in this paper (*O. banksii*, *D. gigas*, and *S. oualaniensis)*. These proteins were aligned to the

258     genome by BLASTX. De-novo identified MS/MS-based peptides were initially also considered as external

259     evidence, but were later omitted due to low coverage. Evidence from predicted repeat locations was

260     used to discourage the model to predict genes overlapping repeats. Since no sufficiently close annotated

261     genome was available for training gene finding parameters, ExonHunter was first run using Drosophila

262     melanogaster parameters on a randomly chosen subset of 118 scaffolds longer than 200kb (total length

263     199 Mb). Out of 12,912 exons predicted in this run, 5,716 were supported by protein alignment data

264     and selected to train the parameters of the gene finding model for *A. dux*, using the methods described

265     in [32]. Rerunning ExonHunter with the resulting A. dux  model parameters on the entire genome

266     yielded 51,225 gene predictions genes. Gene prediction in A. dux is challenging due to the fragmentary

267     nature of the genome assembly (60% of predictions span a sequencing gap). This results in a significant

268     number of artifacts, for example short genes with long introns spanning gaps in the assembly. 18,054

269     predictions yield protein product shorter than 100 amino acids, yet the median span of these

270     predictions is more than 4kb and only 32% of them are supported by transcript or protein alignments. In

271     contrast, 83% of genes with product longer than 100aa are supported. In most of the analyses below,

272     we consider only 33,406 genes that were found to have transcript evidence (blastp match to a sequence

273     from a cephalopod transcriptome, with at least 50% of the giant squid coding region covered) and/or

274     matches in Swissprot or UniRef90 databases (Table 1). This supported set contains much fewer

275     extremely short genes (Figure S4).

276     The function of the protein-coding genes was inferred with Annocript 0.2 [33], which is based on the

277     results from blastp [34] runs against the SwissProt (SP) and UniRef90 (Uf). In addition, we performed a

278  rpsblast search using matrices from the conserved domain database (CDD) to annotate specific domains

279  present on the protein queries.

280  Non-coding RNAs were annotated using the cmsearch program from INFERNAL 1.1 and the covariance

281  models (CMs) from the Rfam database v12.0 [35,36]. All matches above the curated GA threshold were

282  included. INFERNAL was selected because it implements the CMs that provide the most accurate

283  bioinformatic annotation tool for ncRNAs available [37]. tRNA-scan v.1.3.1 was subsequently used to

284  refine the annotation of tRNA genes (Table S3). The method uses a number of heuristics to increase the

285  search-speed, annotates the Isoacceptor Type of each prediction, infers if predictions are likely to be

286  functional or tRNA-derived pseudogenes [38,39]. This method uses CMs to identify tRNAs. Rfam

287  matches and the tRNA-scan results for families belonging to the same clan were then "competed", so

288  that only the best match was retained for any genomic region [36].

289  *Transposable element annotation*
290  Repetitive elements were first identified using RepeatMasker v.4.0.8 [40] with the eukaryota RepBase

291  [41] repeat library. Low-complexity repeats were ignored (-nolow) and a sensitive (-s) search was

292  performed. Following this, a de novo repeat library was constructed using RepeatModeler v.1.0.11 [42] ,

293  including RECON v.1.08 [43] and RepeatScout v.1.0.5 [44]. Novel repeats identified by RepeatModeler

294  were analyzed with a 'BLAST, Extract, Extend' process to characterise elements along their entire length

295  [45]. Consensus sequences and classification information for each repeat family were generated. The

296  resulting de novo repeat library was utilized to identify repetitive elements using RepeatMasker.

297   Data analyses

298   We present a main draft genome assembly produced using 200 Gb of Illumina reads, 4 Gb of Moleculo

299   synthetic long-reads and 108 Gb of Chicago libraries, with a final size matching the estimated genome

300   size of 2.7 Gb, and a scaffold N50 of 4.8 Mb (assembly and annotation statistics in Table 1). Genome

301   completeness estimated by BUSCO reached 90.4% (Eukaryota) and 92.1% (Metazoa), and the

302   completeness for the 33,406 protein-coding genes was 91.2% (Eukaryota) and 84.0 (Metazoa).

303   We also produced an alternative assembly including 27 Gb raw reads generated using the Pacific

304   Biosciences platform, but this showed minimal improvement in assembly statistics, genome size larger

305   than the predicted and lower BUSCO completeness (Table S2).

306   *Comparative analyses of transposable elements*
307   We estimated the total repeat content of the giant squid genome to be approximately half its total size

308   (~49.1%) (Figure 1, Supplementary Table S4). Out of all the repeats present in the giant squid genome,

309   only a few were predicted to be small RNAs, satellites, simple or low complexity repeats (~0.89% of the

310   total genome), with the vast majority (~48.21%) instead consisting of transposable elements (TEs; i.e.

311   SINEs, LINEs, LTR retrotransposons, and DNA transposons; Figure 1, Supplementary Table S4). Of the TE

312   portion of the giant squid genome, the main contribution from annotated TEs is from DNA elements

313   (11.06%) and LINEs (6.96%), with only a small contribution from SINEs (1.99%) and LTR elements

314   (0.72%). TEs are a nearly universal feature of eukaryotic genomes, often comprising a large proportion

315   of the total genomic DNA (e.g. the maize genome is ~85% TEs [46], stick insect genome is ~52% TEs [47],

316   and the human genome is >45% TEs [48]), consequently these account for the majority of observed

317   genome size variation among animals.

318   In Figure 1, we summarise the recently reported TE analyses performed on assembled cephalopod

319   genomes, as follows: California two-spot octopus (*Octopus bimaculatus*) [11] and long-arm octopus (*O.*

320   *minor*) [49], Hawaiian bobtail squid (*Euprymna scolopes*) [50], and giant squid (*A. dux*). The varying

321   sequencing strategies employed to generate currently available cephalopod genomes (and

322    accompanying variation in assembly quality) complicates the comparative analysis of TE content for this

323    group. However, notwithstanding this caveat, it does seem clear that TEs make up a large fraction of the

324    total genomic content across all cephalopod genomes published to date (Figure 1). DNA transposons

325    and LINEs dominate in available cephalopod genomes, while LTR elements and SINEs generally

326    represent a minor portion of cephalopod TEs (Figure 1). Within decapod cephalopods (i.e. squid and

327    cuttlefish), patterns in TE content are generally similar, however, the giant squid has a notably larger

328    proportion of DNA transposons (1,626,482 elements, 11.06% of the total genome) than the Hawaiian

329    bobtail squid (855,308 elements, 4.05% of the total genome), with the bobtail squid in turn having a

330    similar proportion of LINEs (752,629 elements, 6.83% of the total genome) than the giant squid (766,382

331    elements, 6.96% of the total genome; Figure 1).

332    The defining ability of TEs to mobilise, in other words, to transfer copies of themselves into other parts

333    of the genome, can result in harmful mutations. However, TEs can also facilitate the generation of

334    genomic novelty, and there is increasing evidence of their importance for the evolution of host-adaptive

335    processes [51]. In the giant squid genome, all classes of TEs were more frequent (~38.23) in intergenic

336    regions (here defined as regions >2kb upstream or downstream of an annotated gene), than in genic

337    regions versus % of the genome in intergenic regions (~16.6%; Figure 2A). These findings are broadly

338    similar to those reported for other cephalopods, although a larger proportion of the giant squid genome

339    is composed of repeats located within genic regions (percentage of the genome represented by TEs for

340    *O. bimaculoides*: ~6% genic versus ~30% intergenic, and for *O. minor* ~6% genic versus ~40% intergenic

341    [49]).

342    A Kimura distance-based copy divergence analysis revealed that the most frequent TE sequence

343    divergence relative to the TE consensus sequence in the giant squid genome was ~5–8% across all

344    repeat classes, suggesting a relatively recent transposition burst across all major TE types (Figure 2B).

345    Divergence peaks were most pronounced in LINE RTE elements, Tc/Mar and hAT DNA transposons, and

14

346 unclassified TEs, with smaller divergence peaks in SINE tRNA elements and Penelope LINE elements

347 (Figure 2B). Divergence peaks were most pronounced in LINE RTE elements, Tc/Mar and hAT DNA

348 transposons, and unclassified TEs, with smaller divergence peaks in SINE tRNA elements and Penelope

349 LINE elements (Figure 2B). In comparison to observations from other cephalopods, these results suggest

350 a shorter and more intense burst of recent TE activity in the giant squid genome. Overall, further

351 genomic sampling within each of the cephalopod clades will be needed to understand TE evolution, as

352 closely related species can show significant differences (*e.g.*, *O. bimaculoides* to *O. vulgaris*) [52].

353 *Non-coding RNAs*
354 We identified 50,598 ncRNA associated loci in the squid sequencing data, using curated homology-based

355 probabilistic models from the Rfam database[53] and the specialized tRNAscan-SE transfer RNA

356 annotation tool [38]. The essential and well conserved Metazoan ncRNAs: tRNAs, rRNAs (5S, 5.8S, SSU

357 and LSU), RNase P, RNase MRP, SRP and the major spliceosomal snRNAs (U1, U2, U4, U5, U6), as well as

358 the minor spliceosomal snRNAs (U11, U12, U4atac & U6atac), are all found in the *A. dux* genome. Some

359 of the copy numbers associated with the core ncRNAs are extreme. For example, we identified: i)

360 approximately 24,000 loci that appear to derive from 5S rRNA; ii) approximately 17,000 loci that are

361 predicted to be tRNA derived; iii) approximately 3,200 Valine tRNAs isotypes and approximately 1,300

362 U2 spliceosomal RNAs. The microRNA mir-598 also exhibits high copy-numbers at 172. Many of these

363 are likely to be SINEs derived by transposition. All 20 tRNA isotypes were identified in *A. dux* genome.

364 Again, many of these had relatively large copy numbers (summarised in Table 1). These ranged from 46

365 (Cys) up to 2,541 (Val). We identified 174 loci that share homology with 34 known snoRNA families,

366 these included 15 scaRNA, 41 H/ACA box and 118 C/D box snoRNA associated loci [10]. The snoRNAs are

367 predominantly involved in rRNA maturation. We identified 7,049 loci that share homology with 283

368 families of microRNA. Some of these may be of limited reliability, as CMs for simple hairpin structures

369 can also match other, non-homologous, hairpin-like structures in the genome e.g. inverted repeats. A

370    number of cis-regulatory elements were also identified. These included 235 hammerhead 1 ribozymes,

371    133 Histone 30 UTR stem-loops, and 14 Potassium channel RNA editing signal sequences. There are very

372    few matches to obvious non-metazoan RNA families in the current assemblies. The only notable

373    exceptions are bablM, IMES-2, PhotoRC-II and rspL. Each of these families are also found in marine

374    metagenomic datasets, possibly explaining their presence as "contamination" from the environment.

375

376    *Analyses of specific gene families*
377    Several gene families involved in development, such as transcription factors or signaling ligands, are

378    highly conserved across metazoans and may therefore reveal signatures of genomic events, such as a

379    whole genome duplication.

380    WNT is a family of secreted lipid-modified signaling glycoproteins with a key role during development

381    [54]. Comparative analysis of molluscan genomes indicates that the ancestral state was 12 *WNT* genes,

382    as *Wnt3* is absent in all protostomes examined thus far [55]. The giant squid has the typical 12

383    lophotrochozoan WNTs (1, 2, 4, A, 5, 6, 7, 8, 9, 10, 11 and 16; Supplementary Figure S2), and therefore

384    has retained the ancestral molluscan complement, including *Wnt8*, which is absent, for instance, in the

385    genome of the slipper snail *Lottia gigantea* [56].

386    Protocadherins are a family of cell adhesion molecules that appear to play an important role in

387    vertebrate brain development [57]. It is thought that they act as multimers at the cell surface in a

388    manner akin to DSCAM in flies, which lack protocadherins [58]. Cephalopods have massively expanded

389    this family, with 168 identified in the *O. bimaculoides* genome, whereas only 17-25 protocadherins have

390    been identified in the genomes of annelids and non-cephalopod molluscs [11]. We identified

391    approximately 135 protocadherin genes in *A. dux*, many of which are located in clusters in the genome.

392    The possibility that this gene family plays a developmental role parallel to that of protocadherins in

393    vertebrate neurodevelopment thus remains a compelling hypothesis.

394    Development organisation of the highly diverse body plans found in the Metazoa is controlled by a

395    conserved cluster of homeotic genes, which includes, among others, the Hox genes. These are

396    characterized by a DNA sequence referred to as the homeobox, comprising 180 nucleotides that encode

397    the homeodomain [59]. Hox genes are usually found in tight physical clusters in the genome and are

398    sequentially expressed in the same chronological order as they are physically located in the DNA

399    (temporal and spatial collinearity) [60]. Different combinations of Hox gene expression in the same

400    tissue type can lead to a wide variety of different structures [61]. This makes the Hox genes a key subject

401    for understanding the origins of the multitude of forms found in the cephalopods. In *O. bimaculoides*

402    genome assembly no scaffold contained more than a single Hox gene, meaning that they are fully

403    atomised [11]. However, in *E. scolopes*, the Hox cluster was found spanning two scaffolds [50]. In the

404    giant squid, we recovered a full Hox gene cluster in a single scaffold (Figure 3-B). The Hox gene

405    organization found in the giant squid genome suggests either the presence of a disorganised cluster, so-

406    called type D, or atomised clusters, type A [61], or possibly a combination of the two (the genes are still

407    organized, but physically distant from each other). The existence of a "true" cluster seems unlikely, given

408    the presence of other unrelated genes in between and the relatively large distances (Figure 3-C). The

409    classification as type A (atomised) might seem most obvious, despite the co-presence of the genes in a

410    single scaffold, due to these large distances. However, the definition of type D (disorganised) does allow

411    for the presence of non-Hox genes in between members of the cluster (Figure 3-A). Thus, it is difficult to

412    clearly categorise the recovered "cluster", but it does remain clear that these genes are not as tightly

413    bundled as they are in other Bilateria lineages. The *A. dux* Hox "cluster" is spread across 11 Mb of a 38

414    Mb scaffold, and this suggests a far larger size range in the cephalopods than in other described animals,

415    as recently suggested based on the genome of *E. scolopes* [50]. It is possible that this is the reason for

416    the apparent atomisation of Hox genes in the more fragmented *O. bimaculoides* assembly. Hox clusters

417    are usually found in contigs of around 100 kb length in vertebrates [6, 7] and between 500 – 10,000 kb

418    in invertebrates [8] An assembled contig easily containing the complete cluster for these smaller cluster

419    sizes, would manage to cover only one member of the Hox gene cluster in the studied coleoids. As such,

420    our results suggest that the Hox cluster may not be fully atomised in *O. bimaculoides* as previously

421    hypothesised. Further improvements of genome assemblies in cephalopods will be required to address

422    this question. The biological reason for this dramatic increase in the distance between the genes in the

423    Hox cluster presents an intriguing avenue of future research. The homeodomain of all the obtained Hox

424    genes in cephalopods were compared with those of other mollusks. Few differences were found relative

425    to a previous study [62], as no significant modifications were observed in Hox1, Hox4, ANTP, Lox2, Lox5,

426    Post1 and Post2. Hox1 did, however, show reduced conservation in residues 22 to 25 in the *A. dux*

427    sequence. This observation for Hox1 in *A. dux* is visible only in the Pacbio assembly. Additionally, the

428    Hox3 homeodomain analysis supports a basal placement of the nautiloids within cephalopods. The Lox4

429    gene was the most variable among all groups. As of to date, Hox2 still remains undetected in the coleoid

430    cephalopods [63]. Assembly errors notwithstanding, gain and loss of Hox genes has been attributed to

431    fundamental changes in animal body plans, and the apparent loss of Hox2 may therefore be significant.

432    For example, Hox gene loss has been associated with the reduced body-plan segmentation of spider

433    mites [42]. The circumstance that Hox2 has been readily found in *Nautilus*, but remains undetected in all

434    coleoids sequenced thus far, might signify an important developmental split within the Cephalopoda.

435    Alternatively, and equally intriguing, this Hox gene may have undergone such drastic evolutionary

436    modifications that it is presently undetectable by conventional means.

437    On a final note, we analyzed genes encoding reflectins, a class of cephalopod-specific proteins first

438    described in *E. scolopes* [64]. Reflectins form flat structures that reflect ambient light (other marine

439    animals use purine-based platelets), thus modulating iridescence for communication or camouflage

440    purposes [65]. The giant squid genome contains seven reflectin genes and three reflectin-like genes

441    (Supplementary Figure S3). All of these genes, with the exception of one reflectin gene, appear on the

442    same scaffold, which corresponds very well with the distribution pattern of octopus reflectin genes

443    [11]).

## Conclusions

445    Not only because of its astonishing proportions, but also for the lack of knowledge of the key facets of

446    its deep-sea lifestyle, the giant squid has long captured the imagination of scientists and the general

447    public alike. With the release of this annotated giant squid genome, we set the stage for future research

448    into the enigmas that enshroud this truly awe-inspiring creature. Further, given the paucity of available

449    cephalopod genomes, we provide a valuable contribution to the genomic description of cephalopods,

450    and more widely to the growing number of fields that are recognizing the potential, which this group of

451    behaviourally advanced invertebrates holds for improving our understanding of the diversity of life on

452    Earth in general.

## Availability of supporting data

454    The data sets supporting the results of this article are available in the NCBI database a Bioproject

455    PRJNA534469. The three transcriptome data sets (TSA) have ids GHKK01000000, GHKL01000000 and

456    GHKH01000000 and the sequence data used for the genome assemblies has id VCCN01000000.

## Additional files

458    Supplement.txt. Supplementary methods, tables and figures.

## Declarations

## Abbreviations

461    Gb: gigabase pairs; Mb: megabase pairs; BUSCO: Benchmarking Universal Single-copy Orthologs; bp:

462    base pair; NCBI: National Center for Biotechnology Information; LC-MS/MS: liquid chromatography (LC)

463    tandem mass spectrometry (MS); CCB: Colloidal Coomassie Brilliant Blue; HCD: higher-energy collisional

464    dissociation; NCE: normalized collisional energy; AGC: Automatic Gain Control; ALC: Average Local

465    Confidence; SP: SwissProt; Uf: UniRef90; CDD: conserved domain database; CM: covariance model; TE:

466    transposable element; LINE: Long interspersed nuclear element; SINE: Short interspersed nuclear

467    element; LRT: long terminal repeat.

## Ethics statement

469    Sampling followed the recommendations from Moltschaniwskyj et al., 2007 [23].

## Consent for publication

471    Not applicable.

## Competing interests

473    The authors declare that they have no competing interests.

## Funding

## Authors contributions

R.D.F. and M.T.P.G. designed the study. J.S., H-J.H. AND R.T. carried out the sampling. Alex.C., A.R., B.F.,

G.C.A.Jr, H.O. and I.W. performed the laboratory work. R.D.F., Alv.C., A.M., C.B.A., F.S., P.G., T.B., A.H.,

I.B.H., C.C., B.P., F.P., M.P., F.M., O.S., S.R., M.Z.R. and D.P. analyzed the data. E.J., G.Z., J.V., O.F. and Q.L.

contributed with genomic resources. R.D.F., L.F.C.C., A.A., Y.W., B.M., R.S., T.V., B.B., T.S-P., M.T.P.G.

contributed with supervision and computational resources. R.R.F., T.S-P., R.N., M.T.P.G paid for

sequencing. R.D.F. wrote the manuscript with contributions from all authors. All authors have read and

approved the manuscript.

## Acknowledgments

## References

510    References

511    1. Zullo L, Hochner B. A new perspective on the organization of an invertebrate brain. Commun Integr

512    Biol [Internet]. Taylor & Francis; 2011 [cited 2019 Feb 19];4:26–9. Available from:

513    http://www.ncbi.nlm.nih.gov/pubmed/21509172

514    2. Nixon M, Young JZ. The brains and lives of cephalopods. Oxford: Oxford University Press, Oxford;

515    2003.

516    3. Doubleday ZA, Prowse TAA, Arkhipkin A, Pierce GJ, Semmens J, Steer M, et al. Global proliferation of

517    cephalopods. Curr Biol [Internet]. 2016 [cited 2019 May 2];26:R406–7. Available from:

518    http://www.ncbi.nlm.nih.gov/pubmed/27218844

519    4. Vecchione M, Allcock L, Piatkowski U, Jorgensen E, Barratt I. Persistent Elevated Abundance of

520    Octopods in an Overfished Antarctic Area. Smithson Poles  Contrib to Int Polar Year Sci [Internet].

521    Smithsonian Institution Scholarly Press; 2009 [cited 2019 May 2]. p. 197–204. Available from:

522    https://repository.si.edu/handle/10088/6827

523    5. Young RE, Vecchione M, Mangold KM. Cephalopoda, Cuvier 1797 [Internet]. Tree Life. 2018. Available

524    from: http://tolweb.org/Cephalopoda/19386

525    6. Roper CF, Sweeney MJ, Nauen CE. FAO Species Catalogue Vol. 3. Cephalopods of the world. An

526    annotated and illustrated catalogue of species of interest to fisheries. FAO Fish Synopsis [Internet].

527    Rome; 1984;125:277. Available from: http://www.fao.org/3/ac479e/ac479e00.htm

528    7. Jereb P, Roper CFE. Cephalopods of the world. An annotated and illustrated catalogue of cephalopod

529    species known to date. Myopsid and Oegopsid Squids. FAO Species Cat Fish Purp [Internet]. Food and

530    Agriculture Organization of the United Nations; 2010 [cited 2019 Feb 19];2:605. Available from:

531    http://www.fao.org/3/i1920e/i1920e00.htm

532     8. McClain CR, Balk MA, Benfield MC, Branch TA, Chen C, Cosgrove J, et al. Sizing ocean giants: patterns

533     of intraspecific size variation in marine megafauna. PeerJ [Internet]. PeerJ Inc.; 2015 [cited 2019 May

534     15];3:e715. Available from: https://peerj.com/articles/715

535     9. Paxton CGM. Unleashing the Kraken: on the maximum length in giant squid (Architeuthis sp.). J Zool

536     [Internet]. 2016 [cited 2019 May 15];300:82–8. Available from:

537     https://zslpublications.onlinelibrary.wiley.com/doi/pdf/10.1111/jzo.12347

538     10. Rosa R, Seibel BA. Slow pace of life of the Antarctic colossal squid. J Mar Biol Assoc United Kingdom

539     [Internet]. Cambridge University Press; 2010 [cited 2019 May 2];90:1375–8. Available from:

540     https://www.cambridge.org/core/product/identifier/S0025315409991494/type/journal_article

541     11. Albertin CB, Simakov O, Mitros T, Wang ZY, Pungor JR, Edsinger-Gonzales E, et al. The octopus

542     genome and the evolution of cephalopod neural and morphological novelties. Nature [Internet]. Nature

543     Publishing Group; 2015 [cited 2018 May 2];524:220–4. Available from:

544     http://www.nature.com/articles/nature14668

545     12. Liscovitch-Brauer N, Alon S, Porath HT, Elstein B, Unger R, Ziv T, et al. Trade-off between

546     Transcriptome Plasticity and Genome Evolution in Cephalopods. Cell [Internet]. 2017 [cited 2019 May

547     2];169:191–202.e11. Available from: http://www.ncbi.nlm.nih.gov/pubmed/28388405

548     13. Gilly WF, Beman JM, Litvin SY, Robison BH. Oceanographic and Biological Effects of Shoaling of the

549     Oxygen Minimum Zone. Ann Rev Mar Sci [Internet]. 2013 [cited 2019 Feb 19];5:393–420. Available from:

550     http://www.ncbi.nlm.nih.gov/pubmed/22809177

551     14. Golikov AV, Sabirov RM, Lubin PA, Jørgensen LL. Changes in distribution and range structure of Arctic

552     cephalopods due to climatic changes of the last decades. Biodiversity [Internet].  Taylor & Francis Group

553     ; 2013 [cited 2019 Feb 19];14:28–35. Available from:

554    http://www.tandfonline.com/doi/abs/10.1080/14888386.2012.702301

555    15. Balmaseda MA, Trenberth KE, Källén E. Distinctive climate signals in reanalysis of global ocean heat

556    content. Geophys Res Lett [Internet]. John Wiley & Sons, Ltd; 2013 [cited 2019 Feb 19];40:1754–9.

557    Available from: http://doi.wiley.com/10.1002/grl.50382

558    16. Bustamante P, González AF, Rocha F, Miramand P, Guerra A. Metal and metalloid concentrations in

559    the giant squid Architeuthis dux from Iberian waters. Mar Environ Res [Internet]. 2008 [cited 2019 Feb

560    19];66:278–87. Available from: http://www.ncbi.nlm.nih.gov/pubmed/18514304

561    17. Unger MA, Harvey E, Vadas GG, Vecchione M. Persistent pollutants in nine species of deep-sea

562    cephalopods. Mar Pollut Bull [Internet]. 2008 [cited 2019 Feb 19];56:1498–500. Available from:

563    https://linkinghub.elsevier.com/retrieve/pii/S0025326X0800218X

564    18. Winkelmann I, Campos PF, Strugnell J, Cherel Y, Smith PJ, Kubodera T, et al. Mitochondrial genome

565    diversity and population structure of the giant squid Architeuthis: genetics sheds new light on one of the

566    most enigmatic marine species. Proceedings Biol Sci [Internet]. 2013 [cited 2019 Apr 4];280:20130273.

567    Available from: http://rspb.royalsocietypublishing.org/cgi/doi/10.1098/rspb.2013.0273

568    19. Chikhi R, Medvedev P. Informed and automated k-mer size selection for genome assembly.

569    Bioinformatics [Internet]. Oxford University Press; 2014 [cited 2016 Aug 11];30:31–7. Available from:

570    http://www.ncbi.nlm.nih.gov/pubmed/23732276

571    20. Chapman J a, Ho I, Sunkara S, Luo S, Schroth GP, Rokhsar DS. Meraculous: de novo genome assembly

572    with short paired-end reads. PLoS One [Internet]. 2011 [cited 2013 Feb 28];6:e23501. Available from:

573    http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=3158087&tool=pmcentrez&rendertype=ab

574    stract

575    21. English AC, Richards S, Han Y, Wang M, Vee V, Qu J, et al. Mind the Gap: Upgrading Genomes with

576     Pacific Biosciences RS Long-Read Sequencing Technology. Liu Z, editor. PLoS One [Internet]. 2012 [cited

577     2019 Apr 4];7:e47768. Available from: http://www.ncbi.nlm.nih.gov/pubmed/23185243

578     22. Simão FA, Waterhouse RM, Ioannidis P, Kriventseva E V., Zdobnov EM. BUSCO: assessing genome

579     assembly and annotation completeness with single-copy orthologs. Bioinformatics [Internet]. 2015

580     [cited 2019 Apr 8];31:3210–2. Available from: https://academic.oup.com/bioinformatics/article-

581     lookup/doi/10.1093/bioinformatics/btv351

582     23. Moltschaniwskyj NA, Hall K, Lipinski MR, Marian JEAR, Nishiguchi M, Sakai M, et al. Ethical and

583     welfare considerations when using cephalopods as experimental animals. Rev Fish Biol Fish [Internet].

584     Kluwer Academic Publishers; 2007 [cited 2019 Jun 22];17:455–76. Available from:

585     http://link.springer.com/10.1007/s11160-007-9056-8

586     24. Patel RK, Jain M. NGS QC Toolkit: A Toolkit for Quality Control of Next Generation Sequencing Data.

587     Liu Z, editor. PLoS One [Internet]. Public Library of Science; 2012 [cited 2018 Jul 31];7:e30619. Available

588     from: http://dx.plos.org/10.1371/journal.pone.0030619

589     25. Martin M. Cutadapt removes adapter sequences from high-throughput sequencing reads.

590     EMBnet.journal [Internet]. 2011;17:10–2. Available from:

591     http://journal.embnet.org/index.php/embnetjournal/article/view/200

592     26. Haas BJ, Papanicolaou A, Yassour M, Grabherr M, Blood PD, Bowden J, et al. De novo transcript

593     sequence reconstruction from RNA-seq using the Trinity platform for reference generation and analysis.

594     Nat Protoc [Internet]. NIH Public Access; 2013 [cited 2018 Jun 24];8:1494–512. Available from:

595     http://www.ncbi.nlm.nih.gov/pubmed/23845962

596     27. Gilbert D. Gene-omes built from mRNA seq not genome DNA [Internet]. Notre Dame: 7th annual

597     arthropod genomics symposium; 2013. Available from: http://globalhealth.nd.edu/7th-annual-

598    arthropod-genomics-symposium/

599    28. Kleffmann T, Russenberger D, von Zychlinski A, Christopher W, Sjölander K, Gruissem W, et al. The

600    Arabidopsis thaliana chloroplast proteome reveals pathway abundance and novel protein functions.

601    Curr Biol [Internet]. 2004 [cited 2014 Aug 23];14:354–62. Available from:

602    http://www.ncbi.nlm.nih.gov/pubmed/15028209

603    29. Bradford MM. A rapid and sensitive method for the quantitation of microgram quantities of protein

604    utilizing the principle of protein-dye binding. Anal Biochem [Internet]. 1976 [cited 2019 Mar 28];72:248–

605    54. Available from: http://www.ncbi.nlm.nih.gov/pubmed/942051

606    30. Santos R, da Costa G, Franco C, Gomes-Alves P, Flammang P, Coelho A V. First Insights into the

607    Biochemistry of Tube Foot Adhesive from the Sea Urchin Paracentrotus lividus (Echinoidea,

608    Echinodermata). Mar Biotechnol [Internet]. 2009 [cited 2019 Mar 28];11:686–98. Available from:

609    http://www.ncbi.nlm.nih.gov/pubmed/19221839

610    31. Neuhoff V, Arold N, Taube D, Ehrhardt W. Improved staining of proteins in polyacrylamide gels

611    including isoelectric focusing gels with clear background at nanogram sensitivity using Coomassie

612    Brilliant Blue G-250 and R-250. Electrophoresis [Internet]. John Wiley & Sons, Ltd; 1988 [cited 2019 Mar

613    28];9:255–62. Available from: http://doi.wiley.com/10.1002/elps.1150090603

614    32. Brejová B, Vinar T, Chen Y, Wang S, Zhao G, Brown DG, et al. Finding genes in Schistosoma

615    japonicum: annotating novel genomes with help of extrinsic evidence. Nucleic Acids Res [Internet]. 2009

616    [cited 2016 Mar 10];37:e52. Available from:

617    http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=2673418&tool=pmcentrez&rendertype=ab

618    stract

619    33. Musacchia F, Basu S, Petrosino G, Salvemini M, Sanges R. Annocript: a flexible pipeline for the

620  annotation of transcriptomes able to identify putative long noncoding RNAs. Bioinformatics [Internet].

621  2015 [cited 2016 Mar 14];31:2199–201. Available from:

622  http://bioinformatics.oxfordjournals.org/content/early/2015/02/19/bioinformatics.btv106

623  34. Camacho C, Coulouris G, Avagyan V, Ma N, Papadopoulos J, Bealer K, et al. BLAST+: architecture and

624  applications. BMC Bioinformatics [Internet]. 2009 [cited 2014 Jul 9];10:421. Available from:

625  http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=2803857&tool=pmcentrez&rendertype=ab

626  stract

627  35. Burge SW, Daub J, Eberhardt R, Tate J, Barquist L, Nawrocki EP, et al. Rfam 11.0: 10 years of RNA

628  families. Nucleic Acids Res [Internet]. Narnia; 2013 [cited 2019 Apr 4];41:D226–32. Available from:

629  http://academic.oup.com/nar/article/41/D1/D226/1050811/Rfam-110-10-years-of-RNA-families

630  36. Gardner PP, Daub J, Tate J, Moore BL, Osuch IH, Griffiths-Jones S, et al. Rfam: Wikipedia, clans and

631  the &quot;decimal&quot; release. Nucleic Acids Res [Internet]. Narnia; 2011 [cited 2019 Apr

632  4];39:D141–5. Available from: https://academic.oup.com/nar/article-lookup/doi/10.1093/nar/gkq1129

633  37. Freyhult EK, Bollback JP, Gardner PP. Exploring genomic dark matter: a critical assessment of the

634  performance of homology search methods on noncoding RNA. Genome Res [Internet]. Cold Spring

635  Harbor Laboratory Press; 2007 [cited 2019 Apr 4];17:117–25. Available from:

636  http://www.ncbi.nlm.nih.gov/pubmed/17151342

637  38. Chan PP, Lowe TM. GtRNAdb: a database of transfer RNA genes detected in genomic sequence.

638  Nucleic Acids Res [Internet]. Narnia; 2009 [cited 2019 Apr 4];37:D93–7. Available from:

639  https://academic.oup.com/nar/article-lookup/doi/10.1093/nar/gkn787

640  39. Lowe TM, Eddy SR. tRNAscan-SE: a program for improved detection of transfer RNA genes in

641  genomic sequence. Nucleic Acids Res [Internet]. 1997 [cited 2019 Apr 4];25:955–64. Available from:

642     http://www.ncbi.nlm.nih.gov/pubmed/9023104

643     40. Smit AFA, Hubley RR, Green PR. RepeatMasker Open-4.0. 2013.

644     41. Bao W, Kojima KK, Kohany O. Repbase Update, a database of repetitive elements in eukaryotic

645     genomes. Mob DNA [Internet]. 2015 [cited 2019 Apr 17];6:11. Available from:

646     http://www.ncbi.nlm.nih.gov/pubmed/26045719

647     42. Smit A, Hubley R. RepeatModeler Open-1.0. 2015.

648     43. Bao Z, Eddy SR. Automated de novo identification of repeat sequence families in sequenced

649     genomes. Genome Res [Internet]. 2002 [cited 2019 Apr 17];12:1269–76. Available from:

650     http://www.genome.org/cgi/doi/10.1101/gr.88502

651     44. Price AL, Jones NC, Pevzner PA. De novo identification of repeat families in large genomes.

652     Bioinformatics [Internet]. 2005 [cited 2019 Apr 17];21:i351–8. Available from:

653     http://www.ncbi.nlm.nih.gov/pubmed/15961478

654     45. Platt RN, Blanco-Berdugo L, Ray DA. Accurate Transposable Element Annotation Is Vital When

655     Analyzing New Genome Assemblies. Genome Biol Evol [Internet]. 2016 [cited 2019 Apr 17];8:403–10.

656     Available from: https://academic.oup.com/gbe/article-lookup/doi/10.1093/gbe/evw009

657     46. Schnable PS, Ware D, Fulton RS, Stein JC, Wei F, Pasternak S, et al. The B73 Maize Genome:

658     Complexity, Diversity, and Dynamics. Science (80- ) [Internet]. 2009 [cited 2019 Apr 17];326:1112–5.

659     Available from: http://www.ncbi.nlm.nih.gov/pubmed/19965430

660     47. Wu C, Twort VG, Crowhurst RN, Newcomb RD, Buckley TR. Assembling large genomes: analysis of

661     the stick insect (Clitarchus hookeri) genome reveals a high repeat content and sex-biased genes

662     associated with reproduction. BMC Genomics [Internet]. 2017 [cited 2019 Apr 17];18:884. Available

663     from: http://www.ncbi.nlm.nih.gov/pubmed/29145825

664    48. Initial sequencing and analysis of the human genome. Nature [Internet]. Nature Publishing Group;

665    2001 [cited 2019 Apr 17];409:860–921. Available from: http://www.nature.com/articles/35057062

666    49. Kim B-M, Kang S, Ahn D-H, Jung S-H, Rhee H, Yoo JS, et al. The genome of common long-arm octopus

667    Octopus minor. Gigascience [Internet]. 2018 [cited 2019 Apr 17];7. Available from:

668    http://www.ncbi.nlm.nih.gov/pubmed/30256935

669    50. Belcaid M, Casaburi G, McAnulty SJ, Schmidbaur H, Suria AM, Moriano-Gutierrez S, et al. Symbiotic

670    organs shaped by distinct modes of genome evolution in cephalopods. Proc Natl Acad Sci U S A

671    [Internet]. National Academy of Sciences; 2019 [cited 2019 Apr 17];116:3030–5. Available from:

672    http://www.ncbi.nlm.nih.gov/pubmed/30635418

673    51. Schrader L, Schmitz J. The impact of transposable elements in adaptive evolution. Mol Ecol

674    [Internet]. 2018 [cited 2019 Apr 17];28:1537–49. Available from:

675    http://www.ncbi.nlm.nih.gov/pubmed/30003608

676    52. Zarrella I, Herten K, Maes GE, Tai S, Yang M, Seuntjens E, et al. The survey and reference assisted

677    assembly of the Octopus vulgaris genome. Sci Data [Internet]. Nature Publishing Group; 2019 [cited

678    2019 Jun 9];6:13. Available from: http://www.nature.com/articles/s41597-019-0017-6

679    53. Nawrocki EP, Burge SW, Bateman A, Daub J, Eberhardt RY, Eddy SR, et al. Rfam 12.0: updates to the

680    RNA families database. Nucleic Acids Res [Internet]. 2015 [cited 2019 Apr 4];43:D130–7. Available from:

681    http://www.ncbi.nlm.nih.gov/pubmed/25392425

682    54. Cadigan KM, Nusse R. Wnt signaling: a common theme in animal development. Genes Dev [Internet].

683    Cold Spring Harbor Laboratory Press; 1997 [cited 2019 May 8];11:3286–305. Available from:

684    http://www.ncbi.nlm.nih.gov/pubmed/9407023

685    55. Cho S-J, Valles Y, Giani VC, Seaver EC, Weisblat DA. Evolutionary Dynamics of the wnt Gene Family: A

686    Lophotrochozoan Perspective. Mol Biol Evol [Internet]. 2010 [cited 2019 May 16];27:1645–58. Available

687    from: http://www.ncbi.nlm.nih.gov/pubmed/20176615

688    56. Simakov O, Marletaz F, Cho S-J, Edsinger-Gonzales E, Havlak P, Hellsten U, et al. Insights into

689    bilaterian evolution from three spiralian genomes. Nature [Internet]. 2012 [cited 2019 May

690    16];493:526–31. Available from: http://www.ncbi.nlm.nih.gov/pubmed/23254933

691    57. Chen W V, Maniatis T. Clustered protocadherins. Development [Internet]. Company of Biologists;

692    2013 [cited 2018 Oct 21];140:3297–302. Available from:

693    http://www.ncbi.nlm.nih.gov/pubmed/23900538

694    58. Zipursky SL, Sanes JR. Chemoaffinity revisited: dscams, protocadherins, and neural circuit assembly.

695    Cell [Internet]. 2010 [cited 2019 May 16];143:343–53. Available from:

696    https://linkinghub.elsevier.com/retrieve/pii/S0092867410011451

697    59. Pratihar S, Prasad Nath R, Kumar Kundu J. Hox genes and its role in animal development. Int J Sci Nat

698    [Internet]. 2010 [cited 2019 Apr 4];1:101–3. Available from:

699    http://www.scienceandnature.org/IJSN_V1(2)_D2010/IJSN_V1(2)_2.pdf

700    60. Fröbius AC, Matus DQ, Seaver EC. Genomic Organization and Expression Demonstrate Spatial and

701    Temporal Hox Gene Colinearity in the Lophotrochozoan Capitella sp. I. Butler G, editor. PLoS One

702    [Internet]. 2008 [cited 2019 Apr 4];3:e4004. Available from:

703    http://www.ncbi.nlm.nih.gov/pubmed/19104667

704    61. Mallo M, Wellik DM, Deschamps J. Hox genes and regional patterning of the vertebrate body plan.

705    Dev Biol [Internet]. 2010 [cited 2019 Apr 4];344:7–15. Available from:

706    http://www.ncbi.nlm.nih.gov/pubmed/20435029

707    62. Pernice M, Deutsch JS, Andouche A, Boucher-Rodoni R, Bonnaud L. Unexpected variation of Hox

708    genes' homeodomains in cephalopods. Mol Phylogenet Evol [Internet]. Academic Press; 2006 [cited

709    2019 Apr 4];40:872–9. Available from:

710    https://www.sciencedirect.com/science/article/pii/S1055790306001369?via%3Dihub

711    63. Barucca M, Canapa A, Biscotti MA, Zappavigna V. An Overview of Hox Genes in Lophotrochozoa:

712    Evolution and Functionality. J Dev Biol [Internet]. Multidisciplinary Digital Publishing Institute (MDPI);

713    2016 [cited 2019 Apr 4];4:1–15. Available from: http://www.ncbi.nlm.nih.gov/pubmed/29615580

714    64. Crookes WJ, Ding L-L, Huang QL, Kimbell JR, Horwitz J, McFall-Ngai MJ. Reflectins: The Unusual

715    Proteins of Squid Reflective Tissues. Science (80- ) [Internet]. 2004 [cited 2019 Feb 19];303:235–8.

716    Available from: http://www.ncbi.nlm.nih.gov/pubmed/14716016

717    65. Wardill TJ, Gonzalez-Bellido PT, Crook RJ, Hanlon RT. Neural control of tuneable skin iridescence in

718    squid. Proc R Soc B Biol Sci [Internet]. The Royal Society; 2012 [cited 2019 May 16];279:4243–52.

719    Available from: http://rspb.royalsocietypublishing.org/cgi/doi/10.1098/rspb.2012.1374

720    66. Pace RM, Grbić M, Nagy LM. Composition and genomic organization of arthropod Hox clusters.

721    Evodevo [Internet]. BioMed Central; 2016 [cited 2019 Oct 27];7:11. Available from:

722    http://evodevojournal.biomedcentral.com/articles/10.1186/s13227-016-0048-4

723

724

726    **Table 1.** Statistics of the giant squid genome assembly (Meraculous + Dovetail) and corresponding gene
727    prediction and functional annotation. The transcript evidence was confirmed by blastp hits with e-value
728    < 10E[-6] using the transcriptomes of three other species of squid (see the "Transcriptome sequencing"
729    section).

730

| Global Statistics | | |
|---|---|---|
| **Genome assembly*** | **Genome** | **Gene models with evidence** |
| Input assembly | Meraculous | |
| Contig N50 length (Mb) | 0.005 | |
| Longest contig (Mb) | 0.120 | |
| Scaffold N50 length (Mb) | 4.852 | |
| Longest scaffold (Mb) | 32.889 | |
| Total length (Gb) | 2.693 | |
| **BUSCO statistics ([1]Euk / [2]Met)** | | |
| Complete BUSCOs, (%) | 86.1 / 88.5 | 81.6 / 78.3 |
| Complete and single-copy, (%) | 85.1 / 87.6 | 79.9 / 77.7 |
| Complete and duplicated, (%) | 1.0 / 0.9 | 1.7 / 0.6 |
| Partial, (%) | 4.3 / 3.6 | 9.6 / 5.7 |
| Missing, (%) | 9.6 / 7.9 | 8.8 / 16.0 |
| Total Buscos found, (%) | 90.4 / 92.1 | 91.2 / 84.0 |
| **Genome annotation / Gene Prediction** | | |
| Protein-coding gene number | 33,406 | |
| Transcript evidence | 30,472 | |
| Average Protein length, (aa) | 339 | |
| Longest Protein, (aa) | 17,047 | |
| Average CDS length, (bp) | 1,015 | |
| Longest CDS, (bp) | 51,138 | |

| | |
|---|---|
| Average exon length, (bp) | 199 |
| Average exons per gene | 5 |

**Functional annotation (Number of Hits)**

| | |
|---|---|
| Swissprot | 15,749 |
| Uniref90 | 29,553 |
| GO Terms | 4,712 |
| Conserved Domains Database (CDD) | 15,280 |

*The presented statistics are to contigs/scaffolds with length >= 500 bp.

[1]Euk: Database of Eukaryota orthologs genes, containing a total of 303 BUSCO groups.

[2]Met: Database of Metazoa orthologs genes, containing a total of 978 BUSCO groups.

731
732

733

734

**Figure 1.** Comparison of genome repeat content among available cephalopod genomes with assembled

genomes (repeat data for *O. minor* and *O. bimaculoides* from [49] and for *E. scolopes* from [50]). The

tree indicates evolutionary relationships among the two available octopod cephalopods and the two

available decapod cephalopods. Pie charts are scaled according to genome size (*O. bimaculoides*: 2.7Gb,

*O. minor*: 5.09Gb, *E. scolopes*: 5.1Gb, *A. dux*: 2.7Gb), with different repeat types indicated by the colours

presented in the key.

**Figure 2. A)** Stacked bar chart illustrating the proportions (expressed as percentage of the total genome)

of repeats found in genic (≤2kb from an annotated gene) and intergenic regions (>2kb from an

annotated gene) for the giant squid genome. **B)** Transposable element (TE) accumulation history in the

giant squid genome, based on a Kimura distance-based copy divergence analysis of TEs, with Kimura

substitution level (CpG adjusted) illustrated on the x-axis, and percentage of the genome represented by

each repeat type on the y-axis. Repeat type is indicated by the colour chart below the x-axis.

**Figure 3.** Schematic representation of the Hox gene clusters. Different scaffolds are separated by two

slashes. **A)** Simplified classification of the Hox clusters genomic organisation. Type A identifies the lack of

a "typical" Hox cluster configuration, i.e. genes are scattered through the genome (not closely placed);

Type S indicates a Hox cluster that is separated by a chromosomal breakpoint; Type D clusters

comprehend all the genes in the same location but encompassing a larger region than in organised

clusters and may display non-Hox genes and repeats in between; Type O indicates a very compact

cluster embracing a short region with only Hox genes. Non-coding RNA and miRNA can be found. **B)**

Simplified scheme of the chromosomal organisation in various invertebrates. Scaffold length is shown

underneath. Unlike in other coleoids, for *Architeuthis dux* all Hox genes were found in the same scaffold.

However, the distance between the genes was larger than expected for invertebrate organisms, and

758    non-homeobox genes were also present within the cluster. Hox2 remains undetected in coleoids. *A. dux*

759    cluster can be found in scaffold25. *E. scolopes*, *O. bimaculoides*, *L. gigantea*, *C. teleta* and *D.*

760    *melanogaster* assemblies and Hox cluster details can be found in [11,50,56,66]. (*) This gene was

761    reported in a different scaffold, adjacent to non-Hox genes (the length corresponds to the size of the

762    gene). **C)** Complete representation of the Hox cluster found in *A. dux* including the non-Hox genes. PO –

763    Predicted open reading frame; TATDN2 – Putative deoxyribonuclease TATDN2; ZMYM1 – Zinc finger

764    MYM-type protein 1; POGK – Pogo transposable element with KRAB; Zinc finger – Zinc finger protein;

765    MYB-like – Putative Myb-like DNA-binding domain protein; MAPRE1 – Microtubule-associated protein

766    RP/EB family member 1; MGC12965 – Similar to Cytochrome c, somatic.

767

Figure 1

Figure 1

Figure 2

Figure 3

Click here to access/download
**Supplementary Material**
RFonseca_supplement_RF1.docx

Dear Editor,

We herewith submit our revised manuscript 'A draft genome sequence of the elusive giant squid, *Architeuthis dux*'.

Regarding the points that you have highlighted, please find the answers below:

1) Please clarify the rationale for the unconventional assembly strategy in the revised manuscript. If this has "historic" rather than scientific reasons, the reviewer feels this may be fine, but I agree that the reasons should be discussed in the manuscript, for the benefit of readers who are looking for best practice examples.

The reviewer is correct that there is some degree of history involved. We initially did the assembly without PacBio, and did the presented analyses on this. Later we were offered the chance to try and improve it with PacBio, which we did, but as you can see there was minimal improvement in the assembly statistics (Table 1 and Table S2), but i) an increase of the total genome size to 3.155 Gb, beyond the expected 2.7 Gb estimated in kmergenie, and ii) a slight decrease in the BUSCO completeness assessment. As such, we elected to retain the results based on the original assembly (based on Dovetail), but given that we assume others may wish to use the alternative assembly and explore the differences, we provide both.
In the beginning of the "Data analyses" section, we now clearly state which assembly was used in the comparative genomics analyses (from Line 297) and provide an explanation for that choice.

2) Please expand on the methods for protein-coding gene modelling and have another look at your data whether 50K genes may be an overestimate. I also agree with the reviewer's recommendation to analyze gene models in BUSCO to give readers a better idea of their completeness.

We now expanded the section detailing the filtering of the protein-coding gene set and present a total of 33,406 gene annotations in the final set, as these have validation by matching to cephalopod transcripts and/or SwissProt/UniRef90 proteins. We also provide the results from BUSCO when using the gene models as input for comparison (added to Table 1).

Answers to the reviewer's comments:

Reviewer #1: In this study, de Fonseca et al. report the genome of the giant squid as a resource to investigate the unique traits of this fascinating organism. Two assemblies, which are of comparable contiguity to most other recently published molluscan genomes, as well as a set of over 51,000 gene models are reported. Analysis of the genome focuses on repetitive elements (e.g., TEs), non-coding RNAs, and gene families of interest to the authors (WNT genes, Protocadherins, Hox genes, and reflectins). Overall this is a straightforward study that provides a resource that will be broadly useful and I feel it should be published. However, I have a number of suggestions for improvement including a few

Major points:

1.1. It is unclear why two different genome assemblies are presented instead of just one most optimal assembly. This is not the way I would have gone about assembling this combination of data but presumably Dovetail scaffolding and gene modelling were performed before PacBio sequencing and scaffolding? Re-doing the assembly would a more logical way would probably have relatively little improvement but a little more explanation of the rationale or 'historical' reasons for two different assemblies and/or this assembly strategy would be a helpful addition to readers looking in the literature for examples on best practices for genome assembly.

Thank you for this comment. The reviewer is correct that there is some degree of history involved. We initially did the assembly without PacBio, and did the presented analyses on this. Later we were offered the chance to try and improve it with PacBio, which we did, but as you can see there was minimal improvement in the assembly statistics (Table 1 and Table S2), but i) an increase of the total genome size to 3.155 Gb, beyond the expected 2.7 Gb estimated in kmergenie, and ii) a slight decrease in the BUSCO completeness assessment. As such, we elected to retain the results based on the Dovetail assembly, but given that we assume others may wish to use the alternative assembly and explore the differences, we provide both.

1.2. Related to this issue, there is little comparison of the two genome assemblies and it is unclear which assembly was used for what analyses and even Table 1 and Table S2's titles are a bit ambiguous with respect to which assembly statistics are presented. Please explicitly state which assembly was used for which analyses.

In the beginning of the "Data analyses" section, we now clearly state which assembly was used in the comparative genomics analyses (from Line 297) and provide an explanation for that choice. Additionally, we also mention the choice in the Methods section (Lines 183 to 185) before describing the strategies for annotation and comparative analyses.

1.3. The approach used for gene annotation is unconventional and the inferred number of protein-coding gene models is very high. This does not mean the gene model set is bad, but I feel that data needed for the reader to assess the quality of the gene models are lacking. Please run BUSCO on the gene models and report these data as well.

We now also provide the results from BUSCO when using the gene models as input for comparison.

1.4. Specimen collection data are not reported in the manuscript.

This information has now been added to Table S1.

Minor points:
1.5. Scientific names of species need to be italicized throughout.

Done.

Yes, this is now clear in Line 172.

Replace by "in some regions".

This has now been clarified on Line 176.

In the beginning of the "Data analyses" section, we now clearly state which assembly was used in the comparative genomics analyses (from Line 297) and provide an explanation for that choice.

Done.

This has now been clarified in Lines 203-204.

Corrected.

We now further discuss the filters applied in lines 272-275. The total number of protein-coding genes passing all the filters is 33,406.

No custom analysis scripts were developed. We simply use 'bespoke' to mean 'tailored to our particular purpose'. Here this refers an analysis pipeline combining: a preliminary analysis using RepeatMasker, followed by a *de novo* analysis using RepeatModeler and a referenced and publicly available script by Platt et al, followed by a full annotation using RepeatMasker. These steps are fully outlined and referenced in the methods section. We have simplified the sentence which now reads:

"Repetitive elements were first identified using RepeatMasker v.4.0.8"

Done.

The results of BUSCO for post-PacBio step are presented in Table S2 (as indicated in Line 186). We moved the description of the BUSCO results to the "Data analyses" section and added a clarification regarding the choice of the assembly for the overall comparative genomics analyses (from Line 297).

Done.

This gene was reported to be fully isolated from other Hox genes in a different scaffold but was not alone in the scaffold. There were other non-Hox genes. Figure 3 aims to show both the organisation and the range occupied by Hox genes. Considering the organisation, the gene is isolated such as in *O. bimaculoides*. Regarding the size, the schematic representation indicates only the Hox "cluster" area. In *O. bimaculoides*, the scaffolds contain only the Hox genes. This means it could be possible for the cluster to be there but only when considering a very vast distance. In this scenario for *C. teleta*, the gene is found in the middle of the scaffold, surrounded by other genes. It is not part of the cluster. Indicating the full scaffold size could lead to a wrong interpretation of the gene size and of the Hox gene range. As such, only the gene size is indicated.

This information has now been added to Table S1.

manuscript, which distracts the reading. The authors need to carefully fix all these typos and errors during the revision. Further comments are provided below.

Major comments:

2.1.    In the Abstract/Findings, there is a lot of information about "Methods" (e.g. how many raw reads, sequencing of proteome and RNA) instead of what the authors found from the genome itself. Also, the statement "RNA from three different tissue types from three other species of squid to assist genome annotation." is very vague. What tissue types from what species should be clearly described. The authors need to rewrite this section.

In the abstract we followed the format that is usual in a data note, providing detailed information on the data provided by this work. We have now added the names of the three species of squid to the abstract.

2.2.    Line 153: Body patterning system? Usage of body patterning is confusing here since body patterning often refers to the developmental process during embryogenesis but not the skin color pattern.

We have rephrased the sentence to: "Cephalopods can rapidly alter the texture, pattern, colour and brightness of their skin, and this both enables a complex communication system, as wells as provides exceptional camouflage and mimicry."

2.3.    The authors cited that there is a global proliferation of cephalopods (Lines 140 and 141) but later cited other studies saying that there is a regional extinction. It is a bit confusing whether cephalopods are undergoing proliferation or extinction. Given that the earlier citation is more recent (Doubleday et al., 2016) than others, it is wondering which condition is closer to the current situation.

We have removed the second statement to avoid confusion.

2.4.    Although it is agreeable in general to have genome resources from unexplored species, the authors' argument in the last paragraph of Data description/Context is not convincing. The link between having a genome and aiding conservation efforts as well as ensuring continued existence is not clear.

Without a genome, population genomic studies that provide information regarding the genetic diversity and structure of populations becomes very challenging, with genome-wide data having to be produced from reduced-representation methods that have many biases. In this last paragraph, we state this specifically: "A genome is an important resource for future population genomics studies[…]".

2.5.    Do the authors have any idea why the genome contains so many protein-coding genes (51,225 genes predicted) in comparing to other cephalopod species usually having only 20,000-30,000 genes? For example, is it due to that A. dux has more lineage-specific genes or expansions of certain gene families?

We have revised our gene models and now further discuss the filters applied in lines 272-275. The total number of protein-coding genes passing all the filters is 33,406.

2.6.    Given that genome size and polyploidy of the organisms are often correlated to increased body size (Session et al., 2016), have the authors checked if there is whole-genome duplication or polyploidy in the A. dux genome? Session et al. (2016) Genome evolution in the allotetraploid frog Xenopus laevis. Nature 538, 336-343.

We did confirm that the genome was not polyploid by testing for Hardy–Weinberg equilibrium using re-sequencing data from 32 giant squid individuals (Winkelman et al, unpublished results) and there is no evidence for an ancient duplication since we only found one intact Hox complement.

2.7.    Figure 3: The authors should provide scaffold numbers for the Hox clusters from each species. Also, in most cases, Hox genes in the Hox cluster are adjacent to each other without the insertion of other non-Hox genes. If there is a special case in A. dux and E. scolopes, the authors should show the real gene arrangement on that scaffold, especially for the non-Hox genes (with brief annotation) that are in between Hox genes. This can be achieved by having an additional panel in the same figure. The authors are encouraged to show an illustration on the types of Hox gene organization in order to give the readers a better understanding of this context.

Figure 3 has received new panels. Scaffold information for *A. dux* was added in panel C (Figure 3-C). As the assemblies of the other species were retrieved from other studies, the readers are directed to the appropriate references for further detail. An extra panel depicting the Hox cluster organisation in more detail has been added. *E. scolopes* data is shown as reported in its published study. No non-Hox genes were indicated for the area covered in this representation. An additional panel with a simplified version of the various Hox "cluster" types was inserted in panel A (Figure 3-A).

Minor comments:

2.1.1.    Line 149: ~2cm -> "~2 cm"

Done.

2.1.2.    Line 150: 3 orders -> "three orders"

Done.

2.1.3.    Line 150: Architeuthis dux -> "A. dux"

Done.

2.1.4.    Lines 150 and 151: 10-12cm… 20m -> "10-12 cm… 20 m"

Done.

2.1.5.    Line 152: 500kg -> "500 kg"

Done.

2.1.6.    Line 171: a Architeuthis dux sample -> "an A. dux sample"

Done.

2.1.7.    Line 172: What is CTAB?

CTAB = "cetyl trimethylammonium bromide"; this description has been included in the text (Line 172)

2.1.8.    Line 184: For Eukaryota and Metazoa we identified… -> "For Eukaryota and Metazoa, we identified…"

Done.

2.1.9.    Line 184: … 90.4 % and 92.1 %... -> "… 90.4% and 92.1%..."

Done.

2.1.10.    Line 185: 23.38Gb -> "23.38 Gb"

Done.

2.1.11.    Line 186: 14.79kb -> "14.79 kb"

Done.

2.1.12.    "k-mer" (Line 204) or "kmer" (Line 176) to be consistent.

Chose to use "kmer".

2.1.13.    Line 216: 100,000 g -> "100,000×g"

Done.

2.1.14.    Lines 219 and 222: SDS-PAGE -> "SDS-PAGE" (hyphen but not en dash)

Done.

2.1.15.    Line 221: Tris - HCl -> Tris-HCl (single hyphen but not en dash with spaces)

Done.

2.1.16.    Line 226: LC-MS/MS analyses -> "LC-MS/MS analyses" (hyphen but not en dash)

Done.

2.1.17.    Line 254: Using italic for scientific names (i.e. Octopus bimaculatus, Crassostrea gigas, and Lottia gigantea)

**Done.**

2.1.18.    Line 260: … 200kb (total length 199Mb)… -> "… 200 kb (total length 199 Mb)…"

**Done.**

2.1.19.    Line 290: Transposable elements -> "transposable elements"

**Done.**

2.1.20.    Line 300: Architeuthis dux -> "A. dux"

**Done.**

2.1.21.    Line 323: ~5-8% -> "~5¬-8%" (en dash but not hyphen for a range)

**Done.**

2.1.22.    Line 381: Octopus bimaculoides -> "O. bimaculoides"

**Done.**

2.1.23.    Line 383: Euprymna scolopes -> "E. scolopes" (in italic)

**Done.**

2.1.24.    Line 395: Euprymna scolopes -> "E. scolopes"

**Done.**

2.1.25.    Lines 397 & 398: 500 - 10,000 kb -> "500-10,000 kb" (en dash but not hyphen for a range)

**Done.**

2.1.26.    Line 406: … observed in Hox 1, Hox 4, ANTP, Lox 2, Lox 5, Post 1 and Post 2. Hox 1 did,… -> "… observed in Hox1, Hox4, ANTP, Lox2, Lox5, Post1 and Post2. Hox1 did,…"

**Done.**

2.1.27.    Line 407: Hox 1 -> "Hox1"

**Done.**

2.1.28.    Line 408: Hox 3 -> "Hox3"

Done.

2.1.29.    Line 409: Lox 4 -> "Lox4"

Done.

2.1.30.    Lines 410, 412 & 413: Hox 2 -> "Hox2"

Done.

2.1.31.    Line 421: … contains 7 reflectin genes and 3 reflectin-like genes… -> "… contains seven reflectin genes and three reflectin-like genes…"

Done.

2.1.32.    Line 422: … exception of 1 reflectin gene, … -> "… exception of one reflectin gene, …"

Done.

2.1.33.    Line 436: … (tsa)… -> "… (TSA)…"
Done.

2.1.34.    Lines 647 & 657: Architeuthis dux -> "A. dux"

Done.

2.1.35.    Line 659: Hox 2 -> "Hox2"

Done.