

Author's Response To Reviewer Comments

GIGA-D-17-00041

Draft genome of the Antarctic dragonfish, *Parachaenichthys charcoti* Hyun Park; Do-Hwan Ahn; Seung Chul Shin; Bo-Mi Kim; Seunghyun Kang; Jin-Hyoung Kim; Inhye Ahn; Joonho Park

Reviewer reports:

Reviewer #1: The manuscript "Draft genome of the Antarctic dragonfish, *Parachaenichthys charcoti*" is another important contribution to the scientific community working on comparative teleost genomics. As with similar studies, this manuscript appears to be submitted with the purpose of releasing this valuable dataset to the public, and holds no claim to solve any specific scientific question, but rather put up possibilities for the future use of this dataset.

Major comments

The authors have made a good attempt to conduct thorough sequencing of the Antarctic dragonfish genome, using several paired-end and mate-pair libraries. However, the results, and especially the N50 contig statistic is far below what this reviewer would expected using Celera Assembler (CA) with the sequencing data presented. This reviewer is curious to why this specific sequencing method was applied (i.e three very similar libraries for PE sequencing and 2x300bp).

Author's response:

We planned to assemble sequencing reads into contigs using various assembler programs from the same sequencing libraries: Abyss, ALLPATHS-LG, SOAPdenovo, and Celera assembler were used as assemblers in this study. In case of ALLPATHS-LG, paired-end reads should be merged into single read to assemble using higher k-mer. So, we designed the length of libraries to be shorter than 500 bases. The longer reads were known to be favorable in assemblies using de Bruijn graph methods and overlap-layout-consensus methods. So sequencing libraries were sequenced with 2x300bp mode using Illumine MiSeq.

For all the paired-end libraries the inserts are shorter than the sequencing output, which appears to be quite wasted as the trimmed reads are only 173-212bp on average for these libraries. Would it not have been better to have libraries with an insert size around 700-800bp? This would surely span many more of the repetitive sequences now causing gaps and low continuity. Also, as trimming is part of the CA pipeline, why trim the reads prior to running CA? Additionally, FLASH should have been applied to merge overlapping reads from the paired-end sequencing libraries prior to assembly.

Author's response:

As reviewer's comments, if libraries with an insert size around 700-800bp were used, the assembly statistics would be better from Celera assembler. But we were greedy to create libraries that meets all conditions in assemblers and construct libraries with the an insert size up to 500bp to be merged.

In Celera assembler, windows below average quality value of 12 are trimmed as default. We wanted to use only sequencing reads with high quality in assemblies with Celera assembler and other assemblers using de Bruijn graph, and trimmed the bases with a low quality score < 20 from 3'-end of reads. After that, the reads shorter than 70 in length were also discarded, and the resulting high quality reads were used in all assemblies. The use of FLASH is a good suggestion and we will apply it later to improve the genome assembly.

The authors have also made a fair attempt to annotate this *P. charcoti* draft genome using the MAKER pipeline, and I'm happy to see that effort has been put into RNA sequencing to improve this analysis. However, some shortcuts have been taken in regard to how the annotation was performed. For instance, it is now standard procedure to produce a species specific repeat library, using RepeatModler to aid in the annotation. This was not done. The authors also fail to inform which library that was used for identifying repetitive elements with RepeatMasker. It is also customary to include SNAP, AUGUSTUS and GENEFINDER runs as part of the MAKER pipeline to improve gene prediction. This reviewer cannot see that this has been included in the annotation pipeline, which might explain why the number of predicted genes is so high. I'm also missing information regarding which AED cut-off that was used for the final gene predictions.

Author's response:

We used de novo repeat library to identify repetitive elements using RepeatMasker, and the de novo repeat library were produced using repeatModeler with the Repbase library (Ver. 20140131). We also selected the SNAP in MAKER annotation pipeline. Because predicted genes with AED score less than 0.75 were about 3%, we used 1 as AED cut-off values for the final gene predictions. The number of genes with AED value below 0.75 and below 0.25 was 31,642 and 19,708, respectively. We remained the gene with high AED value for manual review, added AED value into the file called "Blast2Go_annotation_with_AED.tab" in GigaDB, and we changed manuscript as follows:

"MAKER2 annotation pipeline was used for genome annotation with default parameters [12]. It first identified repetitive elements using RepeatMasker (ver. 3.3.0) with a de novo repeat library [13], which was constructed using RepeatModeler (Ver. 1.0.3) [14] with the Repbase library (Ver. 20140131). The SNAP gene finder [15] was selected to perform ab initio gene prediction from this masked genome sequence. Alignment of transcriptome assembly results using BLASTn and homologous protein information from tBLASTx were considered for gene annotation as RNA and protein evidence, respectively. Transcriptome assembly was performed by using the program CLC Genomics Workbench 8.0 with default parameters, and sequencing reads from PE500 (Table 1) were used. Proteins from six species were used in the analysis: *Notothenia coriiceps* (NCBI reference sequence NC_015653.1) and *Danio rerio*, *Gasterosteus aculeatus*, *Takifugu rubripes*, *Tetraodon nigroviridis*, and *Gadus morhua* (all from Ensembl release 69). MAKER2 include integration of the Annotation Edit Distance (AED) metric for controlling the quality of annotation [16]. AED values are bounded between 0 and 1, an AED value of 0 indicated that its aligned evidence and annotated gene showed an exact match, Conversely, a value of 1 indicated no evidence support. But the AED cut-off was not applied for this gene predictions. Instead, AED values were denoted in gene annotation and were considered for orthologous gene analysis and gene gain and loss."

The authors have further investigated the gene space completeness using BUSCO, which is

good. However, there is reasons to believe that the gene sets reported are not up to date, especially since there is now a Actinopterygii specific gene set available (http://busco.ezlab.org/frame_meta.html). This should be quick to run and the results can easily be implemented in Table 3.

Author's response:

We did re-run BUSCO analysis to the Actinopterygii DB, and changed Table 3 with new results as reviewer's comments.

In an attempt to conduct comparative genomics, the authors have grouped orthologous genes from several species into orthologous groups using OrthoMCL. This is an OK starting point for a comparative analysis, however, their analysis is based on unfiltered data for the ENSEMBLE (which is know to include thousands of duplicates and Gene ID's without any sequence data available). For instance, would 24,460 genes be a much more adequate dataset to use for the zebrafish. It also included all of the 32,712 *P. charcoti* gene predictions, which leads me to believe that most of the 333 othologous groups (according to Figure 3a, yet referred to as 333 genes" in the text) contain false positives and/or repeats. Based on these results, the authors also produce a "gain-and-loss" figure for the investigated species, yet there is no mentioning on how this analysis was performed.

Author's response:

We filtered the data from the ENSEMBLE, we selected one gene among transcript variants, and discarded the gene without any sequence data available. Then, we have grouped orthologous genes with filtered data into orthologous groups. 25,637 zebra fish gene were used in this analysis, and the number of filtered genes for the ENSEMBLE were indicated in Figure 3A. We did not filtered the 32,712 *P. charcoti* gene predictions completely. Instead, we performed the analysis with the genes corresponding to the three conditions (AED cut-off;1, 0.75, 0.25) and the results corresponding to each case were shown in Figure 3A and 3B. The method producing a "gain-and-loss" was added to manuscript. We also changed the manuscripts as follow:

"Likelihood analysis of gene gain and loss

We estimated differences in the size of orthologs to identify gene families that have undergone significant size changes through evolution [20, 21]. We used the program CAFE3.0 [22] and performed analyses against three groups including the coding sequence of *P. charcoti* with different AED threshold separately. We performed phylogenetic analyses among seven representative fishes with the protein-coding gene in the orthologous groups to obtain the Newick description of a rooted and bifurcating phylogenetic tree. 8,951 orthologous gene sets were selected using the criterion of reciprocal best BLASTP hit and were aligned using PRANK (Ver. 130820) under a codon model with the "-dna -codon" option [23], poor alignment sites were eliminated using Gblock (Ver. 0.91) under a codon model with the "-t = c" option [24]. The remaining alignment regions were concatenated, and used in the construction of the phylogenetic tree by using the neighbor-joining method {Saitou, 1987 #51} in the MEGA (Ver. 6) program [25]. The ultrametric tree of the species with branch lengths in units of time were prepared by referring TimeTree [26] for CAFE3.0 (Figure 3B). The program was performed using $p < 0.05$, and estimated rates of birth (λ) and death (μ) were calculated using the program LambdaMu with the "-s" option. The number of gene gains and losses were calculated on each branch of the tree with the "-t" option. *P. charcoti* gained 937 and lost 1916 gene families (Figure 3B)."

Finally, the authors also present analyses based on (crude) Gene Ontology analyses which offer little scientific value. The entire paragraph on GO enrichment testing (including the results) is not very interesting. So, unless there is any biological meaning applied to the genes or pathways identified, this could/should be removed.

Author's response:

We removed the paragraph and tables for the gene ontology analyses according to reviewer's suggestion.

Minor comments

i) Please use an appropriate "thousands separator" for all values across the manuscript

Content were corrected: content of Table 1-3.

ii) Please make sure that the genus name is not spelled out several times.

Content were corrected: *Parachaenichthys charcoti* to *P. charcoti*

iii) Exchange "illumine" for "Illumina" prior to Table 1

Modification of content: "illumine" to (Illumina, San Diego, USA)

Reviewer #2: Review of Manuscript GIGA-D-17-00041

Overview

Hyun Park's group present the first genome sequence for *Parachaenichthys charcoti*, a member of the bathydraconid (Antarctic dragonfish) clade of the notothenioid group of Antarctic teleosts. This is the second notothenioid genome to be made publically available, following the publication of the Antarctic bullhead *Notothenia coriiceps* (Shin SC et al. *Genome Biology*. 2014;15:468). As a fish biologist interested in physiological evolution, the availability of multiple notothenioid genomes presents a great opportunity for deciphering the genomic basis of adaptive/non-adaptive changes made possible by the extreme cold environment and unusual evolutionary history linked to the notothenioid radiation. As a resource, the *P. charcoti* genome will be used for comparative analyses with *N. coriiceps* and other teleost genomes. I am particularly excited about the eventual publication of a genome for an Antarctic icefish species (*Channichthyidae*), for which the most extreme physiological traits linked to cold conditions are observed (e.g. total loss of haemoglobin). The genomes of *N. coriiceps* and *P. charcoti* will be crucial for such comparative analyses. It is important to note that the Antarctic dragonfishes and *Notothenia* lineages are relatively distant, so the availability of both genome sequences allows both shared-ancestral and lineage-specific changes or adaptations to be disentangled. Moreover, these genomes are generally important in the context of understanding the physiological capacity of notothenioids - key to the overall fauna of Antarctica - to respond to contemporary changes in climate. The manuscript is generally well written.

Thus, overall, I support the publication of this Data Note in GigaScience and I think the paper will encourage the uptake of the *P. charcoti* genome for a range of physiological and evolutionary questions. The data provided by the authors is generally comprehensive and relevant. I offer a number of comments/suggestions, aiming to either increase the clarity surrounding the manuscript's organization and the data and its applications, or requesting more details on aspects of the methodology. I split my comments into general suggestions and a larger set of minor points, the latter linked to particular text in the paper.

General suggestions

1. The authors might consider adding an informative heading to the first paragraph of the Data Description section, such as 'Context' or "Background". This would increase the clarity of the manuscript's organization.

Author's response:

We added "Introduction" as an informative heading.

2. I suggest that authors include an additional dedicated section at the end of the manuscript along the lines of the "Re-use potential" subheading suggested in the Journal guidelines. At the moment, the paper does not do a very proficient job in helping the reader envisage specific uses for the Data Set presented. Hence, in current form, the wider interest of the data set is not fully justified. I would like to see elaboration of the author's stance concerning data re-use, which I feel is necessary to meet the Journal's aim to "contextualize exceptional datasets to encourage reuse". This could provide more context in light of the findings of Shin et al. 2014 (e.g. the new *P. charcoti* genome will allow questions such as, which genomic traits are ancestral to all notothenioids? Which are lineage-specific? Which evolved by convergence? etc.), or give more context on interesting physiological traits observed in notothenioids for which researchers are seeking to clarify the underlying genomic basis.

Author's response:

We added an additional section at the end of the manuscript to satisfy for "Re-use potential" as follow:

"The Antarctic dragonfish *P. charcoti* is a species in the sister lineage of icefishes [27-29] which is the only hemoglobinless vertebrates. The dragonfish (Bathypagrus) and the icefish (Channichthyidae) were generally considered to be evolved from common notothenioid ancestor, which was characterized by decreased hematocrit and blood hemoglobin concentrations [30-34]. The dragonfish showed most similar patterns in these trends among red-blooded notothenioid taxa [34]. The globin complex of the dragonfish *P. charcoti* was hypothesized to be similar in length and organization to that of ancestral icefish prior to loss of functionality [35]. Along with the recently published *N. coriiceps* genome [36], the genome of *P. charcoti* will broaden our understanding of how Antarctic fish have evolved to survive in sub-zero temperatures, and might provide an important clue to understand the process of evolution to the hemoglobinless Antarctic fish and their distinct phenotypes (an increase of blood volume, low blood viscosity, large bore capillaries, increased vascularity with great capacitance, cardiomegaly, and high blood flow)."

3. I find the GO analyses to have tangential relevance as a dataset of meaningful future use, unless it is dissected considerably more than presented within this Data Note, where it appears much as a 'bolt-on'. The biological meaning of data presented in Table 4 (overrepresented GO terms in *P. charcoti*) does not add much insight to fuel on-going research. The data in Table 5 may be misleading in terms of its potential meaning for notothenioid-specific evolution, since the gene list was defined on the basis of comparing two notothenioids with stickleback as the next nearest phylogenetic lineage. As tens of millions of years separate notothenioids and stickleback, it is impossible to say the genes are restricted to notothenioids. This is a minor point, but for me, the paper would be clearer without the GO analyses.

Author's response:

We removed the paragraph and tables for the gene ontology analyses according to reviewer's suggestion.

4. The authors should use species abbreviations consistently throughout the manuscript, which is not the case currently.

Author's response:

We corrected species abbreviation throughout the manuscript.

5. The authors used Celera to assemble the paired end MiSeq reads. As this is an OLC assembler, I would not have naturally considered this to be an optimal approach using relatively short read lengths provided by MiSeq. However, the authors provide evidence that a reasonable draft genome and annotation was nonetheless generated. I am intrigued, in a comparative sense, to know how the same data would have performed using the best-performing assemblers built on the de Bruijn Graph approach. Did the authors attempt any such assemblies, and if so, why did they eventually choose to go with the Celera assembly? To clarify, I am not requesting this as a necessary revision, though if the authors had some available data, I feel it would be of wider interest to contrast the performance of different assemblers.

Author's response:

We assembled sequencing reads into contigs using various assembler from the same sequencing libraries: Abyss, ALLPATHS-LG, SOAPdenovo, and Celera assembler. The assembly statistics from Celera assembler were best among assemblers.

	P. charcoti	CA	8.3	Abyss	2.0.2	SOAPdenovo2	Allpath-LG		
Scaffold Total scaffold length (bases)	794	596	176	1	460	857	469	1,130,003,516	685,815,544
Gap size (bases)	86	840	902	385	080	136	529,475,795	172,038,706	
Scaffolds (n)	12	602	5	921	399	785,432	29,613		
N50 scaffold length (bases)	178	362	10	786	50,086	74,560			
Max scaffold length (bases)	1	318	127	993	314	691,673	716,090		
Contig Total contig length (bases)	709,540,340	1,076,189,796	607,268,662	529,876,330					
Contigs (n)	153,398	6,198,487	2,431,352	139,649					
N50 contig length (bases)	6,145	279	313	6,067					
Max contig length (bases)	65,864	32,177	3,493	67,562					
Gene Number									
(n)	32,712								

Repeat content

(% of genome) 19.4

BUSCO Complete BUSCOs (%) 88.6 75.9 78.9 65

Complete and single-copy BUSCOs (%) 86.3 74 76.9 64

Complete and duplicated BUSCOs (%) 2.3 1.9 2 2

Fragmented BUSCOs (%) 5.7 13.3 9.4 18

Missing BUSCOs (%) 5.7 10.8 11.7 17

Total BUSCO groups searched (n) 4 584*

* Total number of Actinopterygii database

Specific minor points

1. Abstract: "... and *P. charcoti* has undergone molecular and ecological diversification to survive in this cold environment". The wording here can be misconstrued, as the same statement is true for the wider notothenioid lineage. Better to write "... and all notothenioids have undergone molecular and ecological diversification to survive in this cold environment.

Modification of content:

Notothenioids dominate Antarctic fish, making up 90% biomass and all notothenioids have undergone molecular and ecological diversification to survive in this cold environment.

2. "However, little is known about the biology of this species, except that globin intergenic regulatory regions play a role in its low levels of alpha-globin expression". I found this sentence a little disappointing as an upfront motivation for the Data. I feel the abstract could more strongly communicate the importance of the target species for our comparative understanding of evolution in Antarctic fish. Which genomic traits are ancestral to notothenioids, which are lineage-specific, which evolved by convergence, etc.? I suspect these are the motivating questions and in my opinion, the paper would be stronger if this came through more strongly generally, including the abstract.

Modification of content:

Therefore, they (notothenioids) are considered an attractive Antarctic fish model for evolutionary and ancestral genomic studies. Bathydraconidae is a speciose family of the Notothenioidei, the dominant taxonomic component of Antarctic teleosts. To understand the process of evolution of Antarctic fish, we select a typical Antarctic bathydraconid dragonfish, *P. charcoti*.

3. Keywords: the authors might consider elaborating this list, for example to include mention of a genome assembly. Currently the keyword list could be linked to almost any field where Antarctic fish are studied, so it should better represent a genome biology paper.

Modification of content:

Keywords: *Parachaenichthys charcoti*, Antarctic dragonfish, Notothenioid, De novo genome assembly, Genome annotation.

4. Data description paragraph 1: "Antarctic notothenioid teleosts have evolved to adapt to the

extreme Antarctic marine environment. The fish fauna of the Southern Ocean is dominated by a single lineage belonging to the perciform suborder Notothenioidei, consisting of 132 species and 8 families. They survive in the extreme Antarctic marine environment, which includes large seasonal changes in food availability and cold ocean water.

These first few sentences have an issue with the flow of information, which jumps about abruptly, as if thrown together. Consider a reformulation: "The fish fauna of the Southern Ocean is dominated by a single lineage belonging to the perciform suborder Notothenioidei, consisting of 132 species and 8 families. All Antarctic notothenioids have evolved to adapt to the extreme Antarctic marine environment, which includes large seasonal changes in food availability and stably cold water temperature."

Modification of content:

The fish fauna of the Southern Ocean is dominated by a single lineage belonging to the perciform suborder Notothenioidei, consisting of 132 species and 8 families. All Antarctic notothenioids have evolved to adapt to the extreme Antarctic marine environment, which includes large seasonal changes in food availability and stably cold water temperature.

5. Data description paragraph 1: "Nototheniidae is the most speciose family of the Notothenioidei, the dominant taxonomic component of Antarctic teleosts, making up 90% of the fish biomass of the continental shelf and upper slope [1-4]. *Parachaenichthys charcoti*, the Antarctic bathydraconid dragonfish, was first described by Vaillant in 1906".

I find the construction of these sentences to be unusual - when first reading the information, the implication I got was that *P. charcoti* is a member of Nototheniidae, which is not the case. Can the authors please address the construction of the text to improve the clarity of the information?

Modification of content:

Notothenioids dominate Antarctic fish, making up 90% biomass and all notothenioids have undergone molecular and ecological diversification to survive in this cold environment. Therefore, they are considered an attractive Antarctic fish model for evolutionary and ancestral genomic studies. Bathydraconidae is a speciose family of the Notothenioidei, the dominant taxonomic component of Antarctic teleosts [1-4]. *Parachaenichthys charcoti*, the Antarctic bathydraconid dragonfish, was first described by Vaillant in 1906 (Notothenioidei: Bathydraconidae) (AphiaID: 234687; Fishbase ID: 7102).

6. Page 4, "All sequencing (Table 1)", would read more clearly as "All sequencing (data statistics provided in Table 1)". In the current form, the table citation is not clearly linked to the provided text about 'sequencing processes'.

Content were corrected:

For genomic DNA sequencing, three paired-end libraries (PE300, PE400 and PE450) were constructed from sheared genomic DNA (consisting of 300, 400 and 450 bp fragments) and subsequently prepared using standard Illumina sample preparation methods. Mate-pair libraries (MP3K, MP5K, MP8K and MP20K) were prepared for scaffolding, and sequencing was performed according to the manufacturer's instructions (consisting of 3 kb, 5 kb, 8 kb and 20 kb fragments) (Illumina, San Diego, USA).

Because expressed sequence tags are essential for gene annotation in draft genomes, transcriptome library was conducted using TruSeq® Sample Preparation v2 (Illumina) with total RNA. Total RNA were extracted from liver tissue and purified using the RNeasy Mini Kit (Qiagen) with the RNase-Free DNaseI Kit (Qiagen). Extracted sample quality and concentration were determined with 2100 Bioanalyzer (Agilent Technologies, Santa Clara, CA). mRNA was isolated from 2 µg of the total RNA for double-stranded cDNA library construction with poly-A selection. For transcriptome sequencing, paired-end libraries (PE500) were constructed from sheared cDNA consisting of 500 bp fragments and subsequently prepared using standard Illumina sample preparation methods. Final transcriptome libraries length and concentration were determined with 2100 Bioanalyzer. Transcriptome libraries were sequenced using runs of 300×2 paired-end reads (Table 1).

All resulting Illumina reads were trimmed using the FASTX-Toolkit (ver. 0.0.11) (http://hannonlab.cshl.edu/fastx_toolkit) with the parameters -t 20, -l 70 and -Q 33, after which a paired sequence from the trimmed Illumina reads was selected. All sequencing processes for three paired-end libraries (genomic DNA), four mate-pair libraries (genomic DNA) and one paired-end libraries (transcriptome) were performed by Korea Polar Research Institutes (data statistics provided in Table 1).

7. Page 4, "illumine, Carlsbad, USA". Please correct the typo.

Content were corrected:

"illumine, Carlsbad, USA" to "Illumina, San Diego, USA"

8. Page 4, "Finally, paired-end trimmed reads data with 73-fold coverage were obtained (Table 1).

How was the fold-coverage estimated in this case? Also, why present coverage just for the paired-end libraries and not the mate pair libraries?

Author's response:

We divided the sum of paired-end trimmed sequence by the predicted genome size to calculate the fold-coverage. Because the mate-pair libraries were used only in scaffolding, we did not considered it as coverage. But this sentence was not informative. So we deleted this sentence.

9. Page 5: "The assembled contig revealed a contig coverage of approximately 36.57x". By what approach was this assessed?

Author's response:

A contig coverage were calculated by Celera assembler, so we added "in Celera assembler" at the end of the sentence as follow:

"The assembled contig revealed a contig coverage of approximately 36.57x from Celera assembler."

10. Page 5: Why were the selected parameters in Celera selected? Are these simply generally optimized default parameters?

Author's response:

We had tried some optimized Celera assembler parameters, but default option generated best result, although some parameter was optimized for our computer power. Our experience was identical to other genome cases.

11. Page 5: "Contigs from the initial assembly were used for scaffolding using the stand-alone scaffolding tool SSPACE (ver. 2.0) [11]. Trimmed mate-pair reads created using the FASTX-Toolkit were used in the scaffolding process".

Can the authors please provide enough information on the SSPACE parameters employed to allow the reader to repeat the analysis?

We added the parameters at the end of the sentence as follow:

"Contigs from the initial assembly were used for scaffolding using the stand-alone scaffolding tool SSPACE (ver. 2.0) with the following parameters: -x 0, -k 3, -a 0.8, and -T 60 [11]."

12. Page 5: "After scaffolding, the number of scaffolds decreased from 153,398 to 12,381, and the N50 scaffold length increased from 6,135 to 166,726 bp (Table 2)."

The authors might consider stating the total size of the final scaffolds (~795 Mb), which is approaching the genome size according to the K-mer analysis.

We added this sentence at the end of paragraph:

"The total size of the final scaffolds (~795 Mb) was consistent with the estimated genome size (805 Mb)."

13. Page 6: "We first identified repetitive elements using RepeatMasker (ver. 3.3.0) [13], and this masked genome sequence was used for ab initio gene prediction using the SNAP software [14]"

Can the authors please provide more details on their use of RepeatMasker? Which repeats were used? How were they generated bioinformatically?

We changed manuscript as follow:

"MAKER2 annotation pipeline was used for genome annotation with default parameters [12]. It first identified repetitive elements using RepeatMasker (ver. 3.3.0) with a de novo repeat library [13], which was constructed using RepeatModeler (Ver. 1.0.3) [14] with the Repbase library (Ver. 20140131). The SNAP gene finder [15] was selected to perform ab initio gene prediction from this masked genome sequence."

14. Page 6: "Transcriptome assembly results, which were generated using CLC Genomics Workbench 8.0, were used for expressed sequence tags"

Some more details are needed here. Can the authors please clarify the information in terms of the parameters used in CLC? Also, was there not a step to go from a raw transcriptome to a reference transcriptome assembly used for annotation?

We changed manuscript as follow:

"Transcriptome assembly was performed by using the program CLC Genomics Workbench 8.0 with default parameters, and sequencing reads from PE500 (Table 1) were used."

15. Page 6: "A total of 32,712 genes were predicted in *P. charcoti* using MAKER, and 61,709 ab initio prediction, with insufficient evidence were generated (Table 2)."

Much of the information listed in the text is not linked to Table 2. Can the authors please check they have included all information intended in Table 2?

We deleted the ab initio prediction in manuscript and added more information into Table 2.

16. Page 7: Minor point - consider using the term 'partial' rather than 'Fragmented' in Table 1, to be better aligned to information given in the text (or used 'fragmented' in the text). Would the authors also like to comment on why the number of vertebrate BUSCO genes is substantially lower than the eukaryotic or metazoan set?

We change "fragmented" with "Partial" in Table 3.

We did re-run BUSCO analysis to the Actinopterygii DB, and changed Table 3 with new results as reviewer's comments.

17. Page 8: "We identified 8,951 orthologous groups common to all seven fish; 288 of 32,636 *N. coriiceps* genes and 333 of 32,712 *P. charcoti* genes were not identified in any other species, and 2,519 groups were identified only in the two Antarctic fish (Fig. 3A). Subsequently, gene gain-and loss was analyzed in seven representative fish species, *P. charcoti* gained 937 and lost 1916 gene families (Fig. 3B)."

The authors must provide methods to explain how the phylogenetic tree provided in Figure 3 was produced and how they performed the gene gain/loss approach. I suspect the methods are the same as presented in Shin et al. 2014, but this should be clarified. I also must request that the authors either directly provide (or offer some easy way) for an interested reader to extract the relevant subsets of the 8,951 orthogroups (e.g. 333 genes specific to *P. charcoti*; 258 genes specific to *N. coriiceps*; 2,519 common to the two Antarctic fish) as these will be a useful start point for future investigations. Looking at the current data provided in the GigaDB repository, I can only see the 8,951 orthogroups.

The method producing a "gain-and-loss" was added to manuscript, and we uploaded additional orthogroups data into GigaDB (orthologues_List_specific_Antarctic_fish.txt)

18. Page 10: "Availability of supporting data". The authors should break down the full set of data attached in the GigaDB online repository.

We mended as comment.