1    Deep, Staged Transcriptomic Resources for the Novel Coleopteran Models *Atrachya menetriesi*

2    and *Callosobruchus maculatus*

3

4    Matthew A. Benton[1+], Nathan J. Kenny[2+], Kai H. Conrads[1], Siegfried Roth[1]* and Jeremy A.

5    Lynch[3]*

6

7    [1] Institute for Developmental Biology, University of Cologne, Zülpicherstrasse 47b, Cologne

8    50674, Germany

9

10   [2] Simon F.S. Li Marine Science Laboratory of School of Life Sciences and Center for Soybean

11   Research of the State Key Laboratory of Agrobiotechnology, The Chinese University of Hong

12   Kong, Shatin, Hong Kong

13

14   [3] Department of Biological Sciences, University of Illinois at Chicago, Chicago, Illinois, United

15   States of America

16

17   [+]Authors contributed equally

18

19   * Co-corresponding authors

20

21   Matthew A. Benton: matthewabenton@gmail.com

22   Nathan J. Kenny: nathanjameskenny@gmail.com

23   Kai Conrads: kai.conrads@gmx.de

24   Siegfried Roth: siegfried.roth@uni-koeln.de

25   Jeremy A. Lynch: jlynch42@uic.edu

26

27

## **Abstract:**

29  *Background:* Despite recent efforts to sample broadly across metazoan and insect diversity,

30  current sequence resources in the Coleoptera do not adequately describe the diversity of the

31  clade. Here we present deep transcriptomic data generated with Illumina platform sequencing of

32  a number of key stages in the development of two coleopteran species, the false melon beetle

33  *Atrachya menetriesi* (Faldermann 1835) and the cowpea weevil *Callosobruchus maculatus*

34  (Fabricius 1775). These species are both agricultural pests and are of great interest to

35  developmental biology, as well as providing additional coleopteran datapoints for comparative

36  genomic analysis.

37

38  *Results:* By sampling at a range of timepoints covering ovary development and several crucial

39  phases of embryonic development we assembled five and three timed transcriptomes for *A.*

40  *menetriesi* and *C. maculatus,* respectively. We utilized this sequence data to build a

41  transcriptomic resource combining all reads into a single, combined assembly for each species,

42  and we analysed each of these resources in detail. The combined *A. menetriesi* assembly

43  consists of 228,096 contigs with an N50 of 1,598 bp, while the combined *C. maculatus*

44  assembly consists of 128,837 contigs with an N50 of 2,263 bp. For each of these assemblies,

45  we identified 78,879 (34.6%) and 41,744 (32.4%) of contigs by BLASTx against the nr database

46  using Blast2GO, and we found that 97% and 98.3% of the BUSCO set of metazoan orthologs

47  were present in these assemblies. We also carried out manual analysis of 38 key

48  developmental genes, and found that nearly all expected genes were present in our

49  transcriptomes. Each of these analyses showed that these transcriptomes are of very high

50  quality. Lastly, we performed read mapping using the individual timed RNA samples, allowing us

51  to identify up- and down-regulated contigs at key stages in the development of these beetles for

52  further analysis.

53

54  *Conclusions:* The data presented here represent a significant increase in the publicly available

55  transcriptomic resources in the Coleoptera. They have also allowed us to note significant

56  changes in expression at important embryonic stages, and will provide a firm basis for a variety

57  of experimentation, both in developmental biology and in wider exploration of the genomic basis

58  of growth and differentiation.

59

60  Key Words: Beetle, Coleopteran, Transcriptome, Ovary, Embryonic, Staged, *Atrachya*

61  *menetriesi*, *Callosobruchus maculatus,* Chrysomelidae

62

## 63 **Background**:

64  The order Coleoptera is the most speciose clade of animals currently known. Despite the best

65  efforts of generations of biologists, its species are only sparsely sampled and are yet to be

66  comprehensively described, with approximately 90% of coleopteran diversity as yet

67  uncategorized (e.g. [1, 2]). A similar discrepancy exists at the molecular level; while several

68  genomic resources are available in this clade, their number and phylogenetic distribution is only

69  just beginning to accurately sample that of the Coleoptera as a whole. A wide range of

70  transcriptomic data is available in whole organisms (for example [3–5], among others),  specific

71  body parts (such as [6, 7]), and in several cases for staged embryos following RNA interference-

72  mediated gene knock-down [8, 9]. The i5K project [10] will also greatly advance our knowledge

73  of the genomic complement of Coleopterans, with 69 species of this Order listed on that

74  database as nominated for genomic sequencing as of the 14/09/15 (url:

75  http://arthropodgenomes.org/wiki/i5K_nominations) and at least one genome publically available

76   [11]. However, the majority of this information is largely still outside of the public domain, the

77   Coleoptera are still relatively undersampled compared to the Diptera and Hymenoptera, and, in

78   particular, timed embryonic resources (e.g. [12]) are rare in the literature.

79   The true phylogeny of the Coleoptera is still under investigation but, in general, four

80   suborders, 17 superfamilies, and 168 families are recognised [13]. The structure of the

81   coleopteran clade can be seen summarised in Fig 1A. Coleopterans have long been used for

82   research into embryology, and in the pre-molecular era, the Chrysomelidae (summarised

83   phylogeny shown in Fig 1B) was one of the best studied superfamilies. For example, the first

84   functional embryonic experiments in any insect were carried out in the Colorado potato beetle,

85   *Leptinotarsa decemlineata* (sub-family Chrysomelinae, see Fig 1B), leading to the discovery of

86   the function of pole cells and the existence of germ plasm in insects [14, 15]. Further, the larval

87   cuticle preparation technique, so vital for arthropod developmental biology, was first perfected in

88   the bean beetle, *Callosobruchus maculatus* [16]. Chrysomelid beetles are also interesting from

89   an ecological and economic viewpoint, as members of the group are usually pest species,

90   perhaps most famously the aforementioned Colorado potato beetle, which is a major pest of

91   potato crops in America, Asia and Europe.

92   The Chrysomelidae are represented in public databases by a number of ongoing

93   transcriptomic and genomic projects. In particular, both *L. decemlineata* (Bioproject

94   PRJNA171749) and *C. maculatus* [17] are the subject of ongoing genomic sequencing.

95   However, while a number of transcriptomes are planned or published in this clade (e.g. [5, 6,

96   18–22]) none yet sample across embryonic time points in a fashion allowing insight into the

97   genetic mechanisms behind key developmental stages and deep sampling of developmentally

98   important genes. Here we present transcriptomic sequences for two species of chrysomelid

99   beetle, the false melon beetle, *Atrachya menetriesi,* and the aforementioned *C. maculatus* (each

100   of which are pictured along with their relative phylogenetic positions in Fig 1).

101       Compared to *C. maculatus,* the false melon beetle *A. menetriesi* (Faldermann), is a

102    relatively unknown and understudied species. It is native to Japan, where it is an agricultural

103    pest, and feeds on a variety of plants such as clover and lettuce. Although there has only been

104    a small amount of research carried out on this beetle, the work that has been done has

105    highlighted several interesting developmental traits, including the possibility of generating twin

106    embryos after egg bisection [23], or up to four, seemingly complete, embryos following

107    treatment with low temperatures [23, 24]. Another interesting trait exhibited by this beetle is that

108    almost all eggs enter diapause at a certain stage, and this is only broken in the wild by winter

109    conditions. However, a small proportion of eggs skip this diapause and continue to adulthood.

110    The ratio of diapause to non-diapause eggs varies in different parts of Japan [25], and is a

111    heritable trait [26]. Further, *A. menetriesi* embryos undergo a very short germ band mode of

112    development [27], contrasting strongly with beetles such as *C. maculatus* (see below). As the

113    last common ancestor of these two species is estimated to have existed only 80 million years

114    ago [28], how their developmental mechanisms have diverged so greatly is a potentially

115    fascinating area for future study. Research on these topics would greatly benefit from modern

116    molecular studies.

117       *C. maculatus* (Fabricius) is native to West Africa [29] but is now found worldwide, and is

118    a common pest of stored legumes. It is also known as the southern cowpea weevil, however, it

119    is not a true weevil. As noted above, this beetle was the focus of active developmental research

120    in the pre-molecular era, with special focus on segmentation [30, 31]. Segmentation in *C.*

121    *maculatus* has also been studied more recently via immunohistochemistry for the even-skipped

122    protein [32]. This recent work confirmed previous reports that the embryos undergo the long

123    germ mode of development, similar to dipterans like *Drosophila melanogaster* and

124    hymenopterans like *Apis mellifera* [33]. It is commonly believed that the long germ mode of

125    development evolved independently in dipterans and hymenopterans, and given the

126    phylogenetic distribution of short and intermediate germ development within the Coleoptera [34,

127    35], it seems likely that the long germ mode of development also evolved independently in the

128    clade to which *C. maculatus* belongs. A comparison of the molecular basis of development in *C.*

129    *maculatus* and other long germ insects, plus with more closely related species that feature short

130    germ development, such as well studied flour beetle, *Tribolium castaneum* (super-family

131    Tenebrionoidea, see Fig.1 A) and *A. menetriesi*, would yield crucial information on how

132    developmental pathways have evolved to generate the long germ mode of development. *C.*

133    *maculatus* has also been studied in other fields and is a useful system for undergraduate lab

134    teaching [17].

135         In order to facilitate research on *A. menetriesi* and *C. maculatus*, as well as wider

136    investigations in the Coleoptera and beyond, we present here deep, multi-stage transcriptomic

137    resources from a range of key time points in the development of these two beetle species.

138    Using RSEM-based methods we have compared transcript abundance across these life stages,

139    which will allow the investigation of genes that play key roles in developmental changes in these

140    species, particularly at the maternal/zygotic transition. We have carried out extensive searches

141    for key genes in developmental patterning and cell signalling pathways, and from our analyses

142    we conclude that the transcriptomes for both species are of very high coverage, with almost all

143    expected genes being present with full open reading frames. We have already made extensive

144    use of these resources for our own studies on the embryonic development of these two beetles,

145    and are confident that they will be of broad utility to a range of fields in genomics and

146    developmental biology.

147

## Results and Discussion

149    *RNA Extraction and Sample Selection*

150    Ovaries and embryos were collected as described in Materials and Methods, and as seen in Fig

151    2. The chosen time-windows cover a variety of important stages in the development of these

6

152    species, and can be expected to contain the majority of embryonically active transcripts in their

153    expressed complement. Briefly, RNA was collected from dissected ovaries from each species

154    and from four embryonic time windows for *A. menetriesi* and two embryonic time windows for *C.*

155    *maculatus*. Samples were sequenced on the Illumina HiSeq platform (one lane per species),

156    adaptor trimming was performed, along with preliminary assessment of read quality, and data

157    was made available for download from an external server.

158

159    *Read Quality and Assembly Metrics:*

160    Table 1 summarises initial read information, and these reads are available from the NCBI SRA,

161    Bioproject Accession numbers: PRJNA293391 and PRJNA293393. To confirm quality of read

162    data, FastQC [36] was run on all read data. This showed Phred quality scores were high, with

163    median scores always exceeding 28 through to the 101st base, and generally in the mid-30s. A

164    slight bias was found in initial nucleotide sequence. To ensure that Illumina adaptors were

165    removed in their entirety, Trimmomatic v.0.32 [37] was used to check for residual adaptor

166    sequence, but none were observed. We therefore posit that this bias is due to known biases in

167    Illumina hexamer binding [38] rather than sequencing-based artifacts.

168          Using Trinity [39] as an assembler, Trimmomatic-corrected read assemblies were then

169    compared with assemblies based on uncorrected reads. Initial assays of Trimmomatic-treated

170    read assemblies empirically found them to be less well assembled than those using un-trimmed

171    reads, with a shorter N50 and less overall sequence recovery. Trimmed assemblies were

172    therefore discarded in favour of untrimmed assemblies, which were used for all further analyses.

173          Assays of initial assemblies noted a small amount of fungal contamination in the data for

174    both species. DeconSeq [40] was therefore run on the Trinity output for both species, with high

175    stringency as noted in the methods. A total of 336 and 206 contigs with some homology to

176    fungal sequence was removed from the *A. menetriesi* and *C. maculatus* total assemblies as a

177    result, before read mapping was performed. In the *A. menetriesi* assembly, we observed

7

178    particularly high contamination with the protists, notably *Dictyostelium discoideum* and

179    *Naegleria gruberi,* and some of this almost certainly remains in the transcriptome. Metrics for

180    final assemblies after the removal of contamination can be seen in Table 2, alongside the

181    results of assemblies for individual time points.

182        The final, combined timepoint, contamination-removed assemblies, comprising 228,096

183    and 128,837 contigs for *A. menetriesi* and *C. maculatus* respectively, contain a large number of

184    well-assembled transcripts, with the number of contigs greater than 1 kb in length (52,217 in *A.*

185    *menetriesi,* 33,250 in *C. maculatus*) and a high N50 (1,598 bp *A. menetriesi,* 2,263 bp *C.*

186    *maculatus*) indicating a very well assembled dataset. This size is sufficient to span most protein

187    coding domains, allowing easy inference of homology, and will be the full length of many

188    transcripts. It is important to note that many of the contigs in our assembly will represent splice

189    variants of single genes, and some genes will have multiple splice variants, which will affect

190    these statistics. However, excellent recovery of splicing variation itself will be useful for a range

191    of later analysis. GC content of the final assembled transcriptomes closely mirrors that of reads

192    (32.82% in *A. menetriesi,* 38.7% in *C. maculatus*; reads 36-38%/40-43% respectively depending

193    on library).

194        To gain an understanding of whether full-length transcripts were present in our data, we

195    ran TransDecoder v 2.01 [41] to identify open reading frames (ORFs) and filtered for the results

196    that were at least 100 amino acids long. For the *A. menetriesi* assembly, this analysis yielded

197    71,961 raw and 51,912 filtered contigs, while for the *C. maculatus* assembly, the analysis

198    yielded 65,433 raw and 36,535 filtered contigs. The mean average length of the predicted

199    polypeptides, after the filtering step, was 351 (*A. menetriesi*) and 448 (*C. maculatus*) amino

200    acids. This is long enough for us to be confident that our assembly adequately spans coding

201    regions, as the average eukaryotic protein length is 361 amino acids [42]. Together, these

202    analyses suggest that we have recovered the vast majority of coding sequence in our combined

203    assemblies, with sufficient length to adequately span ORFs, a conclusion further supported by

204     gene annotation data as discussed further below. The numbers of ORFs presented here are

205     considerably more than the 16,404 gene models observed in the *T. castaneum* genome

206     (*Tribolium* Genome Sequencing Consortium, 2008), and this is likely due to both spurious ORFs

207     in our dataset and to multiple splice variants. All assemblies are available from Figshare online

208     (*A. menetriesi* DOI: 10.6084/m9.figshare.2056464, *C. maculatus* DOI:

209     10.6084/m9.figshare.2056467).

210

211     *Timed RNA Expression - Differential Expression Analysis*

212     Key developmental stages for both *A. menetriesi* and *C. maculatus* can be easily observed

213     following fixation and nuclei staining (data not shown). Briefly, egg lay to uniform blastoderm

214     stage takes approximately 24 hours in *A. menetriesi*, and 7 hours in *C. maculatus*. Germband

215     formation, gastrulation and elongation occur from 24-100 hours in *A. menetriesi*, and from 7-24

216     hours in *C. maculatus*. The latter period was subdivided in *A. menetriesi* according to

217     characteristic stages of short germ development pertinent to our research interests. These

218     stages, along with assayed sample periods, are shown diagramatically in Fig 2. As well as being

219     used for combined assemblies as described earlier, RNA extracted from mature ovaries and

220     from embryos collected during the aforementioned time periods was also individually assembled

221     using Trinity.

222         The timed transcriptome assemblies for each species often possess better assembly

223     quality when compared to the combined assemblies by metrics such as N50 and mean contig

224     length (Fig 2, Table 2). As a result of being made up of a subset of the total reads, the individual

225     assemblies do not possess the breadth of the combined assemblies, with fewer contigs,

226     especially at long contig length (e.g. 1kb <). As such, the staged transcriptomes were used

227     solely for comparison of expression levels across time, while the combined assemblies were

228     used for gene family analyses. We emphasise that no technical replicate was performed for

229     these comparisons, and any conclusions drawn from them should hold this consideration in

9

230    regard. With this limitation in mind, we carried out differential expression analyses and observed

231    broad trends in expression, as can be seen in Fig 3.

232        First, we generated matrices from general comparison between time points in order to

233    find the most similar samples (Fig 3 A,B). Generally, these results are congruent with steady

234    changes in gene expression across the course of development, with most time points being

235    most similar to those immediately preceding and following them. However, a split can be seen in

236    *A. menetriesi* (Fig 3 A,B) between the ovary and 0-24 hour samples and all others, with the

237    three later samples resembling each other more closely than the 0-24 hour dataset resembles

238    the 24-48 hour sample. This could be due to the maternal:zygotic transition, which will occur at

239    some point in this time frame. This cannot be seen for *C. maculatus*, and could be due to more

240    admixture of RNA within the last 7-24 hour sample. Focused analysis on when the

241    maternal:zygotic transition occurs in each species is required to resolve this question.

242        Next, we clustered the results of our differential expression calculations (Fig 3 C,D). The

243    results shown are those with RSEM considering each isoform separately (rather than taking into

244    account clustering into genes performed by Trinity). Numerous up- and down-regulated contigs

245    can be seen at each time point, with some time points more obviously possessing or lacking a

246    subset of genes found in the combined transcriptomes.

247        While replicates have not been performed and we have not analysed up and down

248    regulated transcripts in detail, these data are available to download from Additional Files 6 and

249    7 attached to this document online. These data will act as a good initial guide for those

250    interested in tracking differential expression of specific genes across development and will be

251    useful for hypothesis building.

252

253    *Basic Gene Annotation*

254    To gain an understanding of the depth of coverage of our datasets we used the BUSCO library

255    of well-annotated genes [43], which are known to be highly conserved in single copy across the

10

256   Metazoa, as a basis for comparison with our combined, assembled transcriptomes. Of the

257   BUSCO set of 843 metazoan orthologs, the *A. menetriesi* assembly possesses 801 (95%)

258   complete (of which 225, 26%, are duplicated), 16 fragmented (1.8%) and 26 missing genes

259   (3.0%), for a total recovery of 97% of the BUSCO dataset. The *C. maculatus* assembly contains

260   815 (96%) complete (with 202, 23%, duplicated), 13 (1.5%) fragmented and 15 (1.7%) missing

261   genes (98.3% recovery). This extremely high level of recovery gives us confidence that at least

262   all housekeeping genes expected to be present in these species are found in our datasets,

263   which strongly suggests that these transcriptomes contain the vast majority of the gene cassette

264   of these species.

265   The number of duplicate genes in our BUSCO analyses likely reflects the construction of

266   our transcriptomes from mixed RNA samples, with the allelic variation that this implies.

267   Discerning true duplicates from allelic and splice variant data is largely contingent on the

268   availability of well-assembled genomes. However, the recovery of these putative duplicates in

269   our assemblies underlines that our RNA sequencing and assembly was of good depth and

270   quality (respectively). With this information in-hand, a range of investigations will be made

271   possible, particularly into the regulation and expression of developmentally important genes.

272   A further understanding of the content of our assemblies was gained from Blast2GO

273   analysis [44, 45]. Genes were annotated using b2gpipe, on the basis of BLASTx (*E* value cutoff,

274   $10^{-3}$) comparisons made against the nr database as downloaded on the 26 January 2015. Of

275   228,096 (*A. menetriesi)* and 128,837 (*C. maculatus)* contigs in each assembly, 78,879 (34.6%)

276   and 41,744 (32.4%) possessed a hit in the nr databases above the threshold. After further

277   annotation with ANNEX and Interproscan, a total of 36,315 (15.9%) and 13,096 (10.2%) contigs

278   were assigned to one or more GO categories. These numbers, while only a fraction of the total

279   number of contigs within our transcriptomes, more closely reflect the expected eukaryotic

280   protein complement in number. Blast2GO annotations are given in Additional Files 2 and 3

281    Fig 4A shows the distribution of species best hit by BLASTx comparison of contigs from

282    *A. mentriesi and C. maculatus* with the nr database. For both species, *T. castaneum* is the best

283    represented species - a reflection of both the phylogenetic position of this species and its well

284    annotated genome. The fact that contigs in our transcriptomes match *T. castaneum* and other

285    coleopteran and insect species more closely than those of other species suggests that gene

286    orthology will be easy to assign in many cases, and that a higher than normal rate of molecular

287    evolution is not observed in our species.

288    The distribution of GO terms within our datasets are shown in Fig 4B, alongside those of

289    the well-annotated *D. melanogaster* and *Mus musculus* proteomes. In general, our datasets

290    resemble one another more than they mirror that of the two sequenced genomes noted. Our

291    transcriptomic resources empirically seem under-represented relative to *D. melanogaster* and

292    *M. musculus* in developmentally interesting categories such as 'Protein Binding' (Molecular

293    Function), 'Multicellular Organism Development' and 'Cell Differentiation' (Biological Process).

294    Given our results for more targeted investigations, found below, we feel this is likely a result of

295    poor annotation of these by Blast2GO, rather than true absence from the transcriptome.

296    At a gross level, in the intracellular 'Cellular Component' GO categories, our data

297    appears more similar to that seen in the two 'model' species than in Molecular Function or

298    Biological Process categories, although this has not been tested statistically. This would

299    suggest these well-conserved structural components were more readily assigned GO categories

300    than the developmentally interesting categories listed above. While not all GO categories are as

301    well-annotated by this process as may be desired, the broad classification of our data into a

302    wide range of GO categories of all levels of GO distribution demonstrates that Blast2GO

303    annotations of our data are a useful starting point for more focussed investigations and

304    identification of specific genes and pathways of interest.

305

306    *Gene Family Recovery*

307     To extend the semi-automated analyses presented above, we performed more targeted

308     analysis of individual, developmentally important gene families. Both the Hox family of

309     transcription factors and the TGF-β cassette were exceptionally well recovered in our dataset.

310     The Hox genes, and in particular the ANTP HOXL class, which pattern the anterior-posterior

311     body axis, are recovered almost in their entirety in both species when compared to well-

312     catalogued databases (e.g. [46]). This can be seen in Fig 5, which shows the phylogenetic

313     distribution of recovered ANTP HOXL class sequences from our transcriptomes alongside

314     previously annotated members of this class.

315         In *A. menetriesi,* sequences for all the ANTP HOXL class genes are recovered in our

316     transcriptome, as can be seen in Fig 5, with sequences and alignments in Additional File 1. We

317     note, however, that the *A. menetriesi Hox 2 / maxillopedia (mxp)/ proboscipedia (pb)* sequence

318     bears some BLAST similarity with *Hox 4 / Dfd* sequences, while the putative *Hox 4 / Deformed*

319     *(Dfd)* recovered does not include the homeodomain sequence - whether it has lost this crucial

320     domain, or if this portion of the sequence is simply not recovered in our assembly is at present

321     unknown. The putative *A. menetriesi Hox 4 / Dfd* sequence is given in Additional File 1,

322     although it is not shown in the phylogeny in Fig 5 due to its truncated length.

323         While 2 (*A. menetriesi)* and 4 (*C. maculatus*) *Hox 3 / zerknüllt (zen)* variants are seen in

324     our species, these are identical at the coding level, and therefore seem to be splice or allelic

325     variants, rather than the two paralogous genes seen in *T. castaneum* [47]. Similarly, no

326     evidence of the *caudal* (*cad*) duplication seen in *T. castaneum* [48] can be found in our

327     transcriptomic assemblies, suggesting this is perhaps specific to the flour beetle lineage. Both

328     *zen* and *caudal* are important embryonic patterning genes, and comparison of these genes in

329     our two species and *Tribolium* would be an excellent situation in which to study how sub- and

330     neo-functionalisation occurs. Likely allelic or splice variants are also observed for other HOXL

331     genes in both *A. menetriesi* and *C. maculatus*. It should be noted that these could represent

332     very recent duplications, or the effect of gene conversion, although this can only be tested fully

13

333   with the advent of a complete genomic resource. No *Pdx/Xlox* gene was seen, adding further

334   circumstantial evidence for the broad scale loss of this gene across the Arthropoda.

335       Our *C. maculatus* transcriptome also contains the full complement of ANTP HOXL

336   genes, although, similar to the case of *Hox 4 / Dfd* in *A. menetriesi*, the 4 recovered *abdominal-*

337   *A (abd-A)* homologs lack the whole homeodomain sequence, with several residues missing.

338   These truncations excluded them from the phylogeny shown in Fig 5, but the sequences for

339   these putative homologs are given in Additional File 1.

340       Even more so than in *A. menetriesi,* a remarkable diversity of potential splice/allelic

341   variants are noted in *C. maculatus,* particularly for the *Hox 6/8* superfamily and *Hox 9-13* /

342   *Abdominal-B (Abd-B)* gene family. Of the *Hox 6/8* gene superfamily, normally represented by

343   four genes in the beetle (*prothoraxless (ptl), fushi tarazu (ftz), Ultrathorax (Utx)* and *abd-A*), 12

344   *ptl,* 1 *ftz,* 4 *Utx* and 4 *abd-A* representatives were found in our analysis. Furthermore, up to 26

345   different potential allelic or splice variants of abdB are recorded. As our transcriptome is made

346   of mixed embryonic samples, it is perhaps not surprising that a diversity of splice/allelic variants

347   are observed, but the excellent recovery of this data confirms the deep coverage provided by

348   our sequencing and assembly.

349       The TGF-β cassettes of the insects have been very well described previously (e.g. [49,

350   50]). Our datasets recover almost the full expected complement of the Coleoptera. A slight

351   exception to this is *Activin (Act)*, a partial sequence of which is recovered for both species: a

352   portion of the propeptide which does not span the mature signal peptide sequence. Whether this

353   is a consequence of loss of the mature domain in these species or low levels of expression at

354   the sampled timepoints remains to be established. A *BMPx* ortholog of clear homology to genes

355   of that family can be found in *C. maculatus,* but has been excluded from the tree seen in Fig 6

356   as it is incomplete in length. Its sequence can be found in Additional File 1, and we have no

357   doubt as to its identity due to high levels of sequence conservation between it and the *T.*

358   *castaneum* and *A. menetriesi* orthologs of this gene.

359     The *glass bottom boat (gbb)* duplication observed in *T. castaneum* cannot be found in

360     our data, but we can recover a range of splice or allelic variants for other genes, especially in *A.*

361     *menetriesi*. These do not differ in the protein coding regions, which leads us to suspect that

362     these are not from gene duplications (unless the duplication(s) occurred very recently). The

363     phylogeny shown in Fig 6 confirms the homology of all genes, and splice/allelic variant numbers

364     observed are given there in brackets, with all sequences available in Additional File 1.

365     We also note the discovery of an additional TGF-β ligand in *C. maculatus*. This gene has

366     been previously automatically annotated as *derriere* in *T. castaneum*, (XP_008191586.1) and if

367     it is truly of this family, which is also known as *GDF1/3/Univin/Vg1*, this would be a surprise, as

368     its presence outside the Deuterostomia is controversial [51]. If proof could be found for this

369     being a *bona fide GDF1/3/Univin/Vg1*, the presence of this gene in more than one coleopteran

370     could suggest that this might in fact be ancestrally present in all bilaterian species, but further

371     investigation is warranted before strong conclusions can be drawn in this regard. We could gain

372     no phylogenetic support for placing either of these beetle sequences in the GDF1/3 clade, and it

373     may well be that these sequences instead represent a coleopteran novelty.

374     Our recovery of not only the full expected complements of these vital developmental

375     genes, but also a remarkable diversity of alternative variants, demonstrates the depth of our

376     assemblies as a resource. Whether used as the basis for simple cloning or more sophisticated

377     analysis of patterns of gene variation and diversification, these transcriptomes will be of wide

378     utility to the field of coleopteran and insect developmental biology.

379

380     *Pathway Recovery*

381     As well as examining specific gene families, we investigated a number of broader pathways

382     commonly studied in insects [52]. This allows us to both note how well-recovered such

383     pathways are in our species as a measure of transcriptome utility, as well as note interesting

384     differences between these pathways in our species when compared to others. We did this using

385    both automated (KEGG KAAS mapping) and manual (BLAST based) methods. Some

386    representative results of KEGG KAAS mapping can be seen in Fig 7, and all KEGG annotations

387    can be downloaded from Additional Files 4 and 5.

388        KEGG KAAS mapping uses BLAST results to annotate known pathways, and gives a

389    rapid overview of the recovery of these. Here we have shown the well-known Wnt, Notch and

390    Hedgehog pathways to indicate the depth of our transcriptomes, and show how they may be

391    useful for future research. However, these maps often use terminology based on vertebrate

392    nomenclature, and contain genes known to be absent from particular clades. We have therefore

393    indicated in Fig 7 (using unshaded boxes as shown in the Key) genes that may be absent

394    ancestrally in the Coleoptera, based on their absence from the *T. castaneum* pathway. Of genes

395    expected to be present in the Coleoptera we find almost total recovery in our transcriptomes. In

396    the three pathways examined, only three genes noted to be present in *T. castaneum* were noted

397    as absent from both of our transcriptome datasets, all in the Wnt cascades (Fig 7A). The

398    expected Hedgehog cassette was recovered *in toto* (Fig 7B) and in the Notch signalling

399    cascade (Fig 7C), only APH-1 was noted to be absent, and only from the *C. maculatus*

400    transcriptome. We must note that these may not be true absences - KEGG KAAS mapping is

401    based on automatic BLAST assignation, and if these sequences are divergent in our

402    transcriptomes they may have been missed by this analysis.

403        We also examined pathways manually, using reciprocal BLAST hits and closer manual

404    investigation to confirm identity of individual genes, the results of which can be seen in Table 3.

405    The anterior-posterior patterning genes *cad* (mentioned earlier) and *hunchback (hb)* are present

406    in both species. Of the germline establishment and localization genes examined, *nanos* was

407    surprisingly absent in both species, while *bruno (bru;* also known as *arrest)*, *exuperantia (exu)*,

408    *tudor* (*tud;* 2 copies in *A. menetriesi)*, *oskar (osk)*, *vasa (vas)* and *valois (vls)* were present.

409    Interestingly, *pumilio (pum)* is present in *C. maculatus* in single copy (although it is divided

410    across two contigs), while *A. menetriesi* possesses a total of seven copies. The different *A.*

16

411    *menetriesi pum* copies varied both at the nucleotide level and in their amino acid sequences,

412    strongly suggesting that they are in fact paralogs. In depth analysis of these genes is required to

413    uncover why they have undergone several rounds of duplication. Orthologs of the *Drosophila*

414    gene *swallow (swa)* could not be found in either of our transcriptome resources, nor is it present

415    in several other insects (data not shown) and we suggest it may therefore be a schizophoran

416    novelty.

417           Canonical gap genes *Krüppel (Kr)*, *knirps (kni)*, *giant (gt)*, *huckebein (hkb)*, *tailless* (*tll;* 2

418    paralogs in *C. maculatus*), *buttonhead (btd)*, *empty spiracles (ems)* and both *orthodenticle*

419    orthologs (*Otd* and *Otd2*) were recovered in both species examined here. The *C. maculatus*

420    specific duplication of *tll* would be another excellent opportunity to study gene duplication and

421    evolution. The pair rule genes e*ven skipped (eve), hairy (h), fushi tarazu (ftz), odd paired (opa),*

422    *odd skipped (odd), paired (prd), runt (run)* and *sloppy paired 1 (slp1)* were present in single

423    copy. The segment polarity gene cassette was also present in both species, with the notable

424    absence of *gooseberry-neuro (gsb-n)* from our datasets. The genes *armadillo (arm), cubitis*

425    *interruptus (ci), Engrailed (En), fused (fu), gooseberry (gsb), hedgehog (hh), pangolin (pan),*

426    *patched (ptc)* and *wingless (wg)* were all present, and their sequences can be found in

427    Additional File 1.

428           All of these pathways are commonly studied in insects, and the annotations provided

429    here, along with preliminary timed expression data, will provide a basis for a wide range of

430    targeted investigations into the embryonic development of these two species, and how these

431    pathways have changed over the course of evolution. Furthermore, the excellent recovery of

432    these pathways by both automated (KEGG-KAAS) and manual annotation gives us high

433    confidence in the completeness of our transcriptomic resources. This confirms the results of our

434    BUSCO analysis, and our datasets are therefore likely to contain the vast majority of transcribed

435    genes in these two species, with only lowly expressed and temporally restricted genes absent

436    from these transcriptome resources.

17

437

## Conclusions:

439    Our production of deep transcriptomic sequence data for *A. menetriesi* and *C. maculatus* will

440    assist in the inference of character gain and loss across the Coleoptera, aid in future

441    phylogenetic efforts, and allow a range of investigations into the embryonic development of

442    these species at the molecular level. The status of these organisms as common agricultural

443    pests also suggests that such resources may allow targeted control mechanisms to be

444    developed for these species. This data will be another key building block in our understanding of

445    the transcriptomic basis to embryological development, and provide a window into the basic

446    biology of the most successful clade of animals.

447

## Methods:

449    *Animal Husbandry*

450    *A. menetriesi* eggs were kindly provided by Dr Yoshikazu Ando, and were reared at 25°C on wet

451    sand or soil and fed fresh lettuce. *C. maculatus* beetles were kindly provided by Dr Joel Savard,

452    and were reared at 30°C on dry black-eyed peas.

453

454    *RNA Extraction and Sequencing*

455    RNA was extracted from ovary and timed embryonic samples using a TRIzol RNA extraction kit

456    according to the manufacturer's protocols. RNA quantity and quality was checked using a

457    Thermo Scientific Nanodrop 2000C Spectrophotometer and 1 µg was sent for sequencing by

458    the Cologne Centre for Genomics. Adaptor trimming and initial quality control was performed by

459    the commercial provider according to their internal standards. This cleaned data was then made

460    available to us for download from an external server. Paired end read quality after sequencing

461    was assessed using the FastQC program [36].

462

463 *Transcriptome Assembly and Comparative Expression Analyses*

464 Assemblies used in our final analysis were made using Trinity version 2013_08_14 [39], with the

465 default settings. Full assemblies were made using reads from all time points, and individual

466 assemblies were then constructed using each sampled time point individually. All assemblies

467 are available from Figshare online (*A. menetriesi* DOI: 10.6084/m9.figshare.2056464, *C.*

468 *maculatus* DOI: 10.6084/m9.figshare.2056467). DeconSeq standalone version 0.4.3 [40] was

469 run on full assemblies with settings -i 95 -c 95, using the bact, fungi, hsref, and prot databases.

470 Comparative expression analyses were performed using RSEM [53] as packaged in the Trinity

471 module, and results shown here are the 'as-isoform' data, although 'as-gene' data is provided in

472 Additional Files. BUSCO v1.1b1 [43] was used to assess gene complement completeness.

473

474 *Functional Annotation*

475 Our combined assemblies were automatically assigned homologs and annotated according to

476 gene ontology (GO) terms using Blast2GO [44, 45]. Initially, BLASTx was run using BLAST

477 2.2.29+ against the NCBI nr database as downloaded to a local server on the 17/01/2015, with

478 settings -evalue 0.001 -max_target_seqs 5 -outfmt 5. GO term distribution within the *D.*

479 *melanogaster* and *H. sapiens* genomes were downloaded from B2GO-FAR [54] and calculated

480 using the Combined Graph function of Blast2GO. KEGG KAAS mapping was automatically

481 performed using the KEGG KAAS tool (http:// www.genome.jp/tools/kaas/), single-directional

482 best hit with default BLAST settings, and with the eukaryote dataset as a basis for annotation.

483

484 *Gene Identification*

485 Gene sequences were manually identified and their homology confirmed by independently using

486 tBLASTn [55] searches using gene sequences of known homology downloaded from the NCBI

487 nr database as queries against standalone databases created on a local server using BLAST

19

488    2.2.29+ or the CLC Main workbench. Genes putatively identified using this method were

489    reciprocally BLASTed against the online NCBI nr database using BLASTx to confirm their

490    identity. Where identity was uncertain, phylogenetic analysis was used to confirm identity.

491

492    *Phylogenetic Tree Construction*

493    Sequences were aligned using MAFFT 7 [56] unless otherwise stated under the L-INS-i

494    strategy. Alignments were then saved and exported to MEGA 6, where regions of poor

495    alignment were manually excluded and maximum likelihood phylogenetic trees were

496    constructed using the LG model, 1000 bootstrap replicates as indicated, 4 gamma categories

497    and invariant sites, and all other default prior settings [57].

498

499

## List of Abbreviations Used:

501    BLAST: Basic Local Alignment Search Tool, BUSCO: Benchmarking Universal Single-Copy

502    Orthologs, GDF: Growth Differentiation Factor, KEGG-KAAS: Kyoto Encyclopaedia of Genes

503    and Genomes - Automatic Annotation Server, MAFFT: Multiple Alignment using Fast Fourier

504    Transform, NCBI: National Centre for Biotechnology Information,ORF: Open Reading Frame,

505    RNA: Ribonucleic Acid, RSEM: RNA-Seq by Expectation-Maximization, SRA: Short Read

506    Archive, TGF: Transforming Growth Factor

507

## Competing Interests:

509    The authors declare that they have no competing interests.

510

## Authors' contributions:

512     MAB, KC and JAL maintained animals, extracted RNA and arranged sequencing. MAB and NJK

513     assembled transcriptomes, analysed raw data and wrote the manuscript. MAB, NJK, KC, JAL

514     and SR designed experiments and contributed to gene family analyses presented in this

515     manuscript. All authors read and approved the final manuscript.

516

## Acknowledgements:

525

## Availability of data and materials

527     The datasets supporting the conclusions of this article are available in the NCBI SRA repository

528     [Bioproject Accession numbers: PRJNA293391 and PRJNA293393] and in the Figshare

529     repository [DOIs: 10.6084/m9.figshare.2056464 and 10.6084/m9.figshare.2056467]. NOTE:

530     DATA EMBARGOED UNTIL PUBLICATION

531

## Bibliography:

533     1. Stork NE: **Insect diversity - facts, fiction and speculation**. *Biol J Linn Soc* 1988, **35**:321–

534     337.

535     2. Grimaldi D, Engel M: *Evolution of the Insects*. Cambridge University Press; 2005.

536     3. **1KITE: 1000 Insect Transcriptome Evolution** [http://www.1kite.org/]

537    4. Keeling CI, Henderson H, Li M, Yuen M, Clark EL, Fraser JD, Huber DPW, Liao NY, Docking

538    TR, Birol I, Chan SK, Taylor GA, Palmquist D, Jones SJM, Bohlmann J: **Transcriptome and**

539    **full-length cDNA resources for the mountain pine beetle, Dendroctonus ponderosae**

540    **Hopkins, a major insect pest of pine forests.** *Insect Biochem Mol Biol* 2012, **42**:525–536.

541    5. Kumar A, Congiu L, Lindstrom L, Piiroinen S, Vidotto M, Grapputo A: **Sequencing, De Novo**

542    **assembly and annotation of the Colorado Potato Beetle, Leptinotarsa decemlineata,**

543    **Transcriptome.** *PLoS One* 2014, **9**:e86012.

544    6. Pauchet Y, Wilkinson P, van Munster M, Augustin S, Pauron D, ffrench-Constant RH:

545    **Pyrosequencing of the midgut transcriptome of the poplar leaf beetle Chrysomela**

546    **tremulae reveals new gene families in Coleoptera.** *Insect Biochem Mol Biol* 2009, **39**:403–

547    413.

548    7. Chen H, Lin L, Xie M, Zhang G, Su W: **De novo sequencing, assembly and**

549    **characterization of antennal transcriptome of Anomala corpulenta Motschulsky**

550    **(Coleoptera: Rutelidae).** *PLoS One* 2014, **9**:e114238.

551    8. Oberhofer G, Grossmann D, Siemanowski JL, Beissbarth T, Bucher G: **Wnt/ -catenin**

552    **signaling integrates patterning and metabolism of the insect growth zone**. *Development*

553    2014, **141**:4740–4750.

554    9. Jacobs CGC, Braak N, Lamers GEM, van der Zee M: **Elucidation of the serosal cuticle**

555    **machinery in the beetle Tribolium by RNA sequencing and functional analysis of**

556    **Knickkopf1, Retroactive and Laccase2**. *Insect Biochem Mol Biol* 2015, **60**:7–12.

557    10. i5k-Consortium: **The i5K Initiative: advancing arthropod genomics for knowledge,**

558    **human health, agriculture, and the environment.** *J Hered* 2013, **104**:595–600.

559    11. Richards S, Gibbs R a, Weinstock GM, Brown SJ, Denell R, Beeman RW, Gibbs R, Bucher

560    G, Friedrich M, Grimmelikhuijzen CJP, Klingler M, Lorenzen M, Roth S, Schröder R, Tautz D,

561    Zdobnov EM, Muzny D, Attaway T, Bell S, Buhay CJ, Chandrabose MN, Chavez D, Clerk-

562    Blankenburg KP, Cree A, Dao M, Davis C, Chacko J, Dinh H, Dugan-Rocha S, Fowler G, et al.:

563   **The genome of the model beetle and pest Tribolium castaneum.** *Nature* 2008, **452**:949–55.

564   12. Yin A, Pan L, Zhang X, Wang L, Yin Y, Jia S, Liu W, Xin C, Liu K, Yu X, Sun G, Al-hudaib K,

565   Hu S, Al-Mssallem IS, Yu J: **Transcriptomic study of the red palm weevil Rhynchophorus**

566   **ferrugineus embryogenesis.** *Insect Sci* 2015, **22**:65–82.

567   13. Hunt T, Bergsten J, Levkanicova Z: **A comprehensive phylogeny of beetles reveals the**

568   **evolutionary origins of a superradiation**. *Science (80- )* 2007, **438**:1–4.

569   14. Hegner RW: **Effects of Removing the Germ-Cell Determinants from the Eggs of Some**

570   **Chrysomelid Beetles. Preliminary Report**. 1908, **16**:19–26.

571   15. Hegner RW: **The origin and early history of the germ-cells in some chrysomelid**

572   **beetles**. *J Morphol* 1909, **20**:231–296.

573   16. Meer JM Van Der: **Optical clean and permanent whole mount preparation for phase-**

574   **contrast microscopy of cuticular structures of insect larvae**. *Dros Inf Serv* 1977, **52**.

575   17. **Bean Beetles: A Model Organism for Inquiry-based Undergraduate Laboratories**

576   [http://www.beanbeetles.org/]

577   18. Pauchet Y, Wilkinson P, Chauhan R, Ffrench-Constant RH: **Diversity of beetle genes**

578   **encoding novel plant cell wall degrading enzymes.** *PLoS One* 2010, **5**:e15635.

579   19. Kirsch R, Wielsch N, Vogel H, Svatos A, Heckel DG, Pauchet Y: **Combining proteomics**

580   **and transcriptome sequencing to identify active plant-cell-wall-degrading enzymes in a**

581   **leaf beetle.** *BMC Genomics* 2012, **13**:587.

582   20. Flagel LE, Bansal R, Kerstetter RA, Chen M, Carroll M, Flannagan R, Clark T, Goldman BS,

583   Michel AP: **Western corn rootworm (Diabrotica virgifera virgifera) transcriptome assembly**

584   **and genomic analysis of population structure.** *BMC Genomics* 2014, **15**:195.

585   21. Strauss AS, Wang D, Stock M, Gretscher RR, Groth M, Boland W, Burse A: **Tissue-**

586   **specific transcript profiling for ABC transporters in the sequestering larvae of the**

587   **phytophagous leaf beetle Chrysomela populi.** *PLoS One* 2014, **9**:e98637.

588   22. Chi YH, Salzman R a, Balfe S, Ahn J-E, Sun W, Moon J, Yun D-J, Lee SY, Higgins TJ V,

23

589     Pittendrigh B, Murdock LL, Zhu-Salzman K: **Cowpea bruchid midgut transcriptome**

590     **response to a soybean cystatin--costs and benefits of counter-defence.** *Insect Mol Biol*

591     2009, **18**:97–110.

592     23. Miya K, Kobayashi K: **The embryonic development of Atrachya menetriesi. Faldermann**

593     **(Coleoptera, Chrysomelidae). II. Analyses of early development by ligation and low**

594     **temperature treatment.** *J Fac Agric Iwate Univ* 1974, **12**:39–55.

595     24. Miya K, Ando Y, Kurihara M: **Formation of duplicated embryos by treatment of low**

596     **temperature in Atrachya menetriesi Faldermann (Chrysomelidae, Coleoptera)**. *Proc 26th*

597     *Ann Meet Ent Soc Japan* 1966, **9**.

598     25. Ando Y: **Geographic-Variation In The Incidence Of Non-Diapause Eggs Of The False**

599     **Melon Beetle, Atrachya-Menetriesi Faldermann (Coleoptera, Chrysomelidae)**. *Appl*

600     *Entomol Zool* 1979, **14**:193–202.

601     26. Ando Y, Miya K: **Diapause character in the false melon beetle, Atrachya menetriesi**

602     **Faldermann, produced by crossing between diapause and non diapause strains.** *Bull Fac*

603     *Agri Iwate Univ* 1968, **9**:87–96.

604     27. Miya K: **The embryonic development of a Chrysomelid Beetle, Atrachya menetriesi.**

605     **Faldermann (Coleoptera) I. The stages of development and changes of external form**. *J*

606     *Fac Agric Iwate Univ* 1965, **7**:155–166.

607     28. Gómez-Zurita J, Hunt T, Kopliku F, Vogler AP: **Recalibrated tree of leaf beetles**

608     **(Chrysomelidae) indicates independent diversification of angiosperms and their insect**

609     **herbivores.** *PLoS One* 2007, **2**:e360.

610     29. Tran BMD, Credland PF: **Consequences of inbreeding for the cowpea seed beetle,**

611     **Callosobruchus maculatus (F)(Coleoptera: Bruchidae)**. *Biol J Linn Soc* 1995, **56**:483–503.

612     30. Meer J van der: **of metameric order in the insect Callosobruchus maculatus**

613     **Fabr.(Coleoptera) I. Incomplete segment patterns can result from constriction-induced**

614     **cytological damage**. *J Embryol Exp …* 1979, **51**:1–26.

24

615    31. Meer JM Van Der: **Parameters influencing reversal of segment sequence in posterior**

616    **egg fragments of Callosobruchus (Coleoptera)**. *Roux's Arch Dev Biol* 1984:339–356.

617    32. Patel NH, Condron BG, Zinn K: **Pair-rule expression patterns of even-skipped are found**

618    **in both short- and long-germ beetles.** *Nature* 1994, **367**:429–434.

619    33. Osborne PW, Dearden PK: **Expression of Pax group III genes in the honeybee (Apis**

620    **mellifera).** *Dev Genes Evol* 2005, **215**:499–508.

621    34. Anderson DT: **The Development of Holometabolous Insects**. In *Developmental systems.*

622    *Insects, Vol 1*. Edited by Counce SJ, Waddinton CH. New York: Academic Press; 1972:165–

623    242.

624    35. Davis GK, Patel NH: **SHORT, LONG, AND BEYOND: Molecular and Embryological**

625    **Approaches to**. *Annu Rev Entomol* 2002:669–99.

626    36. **FastQC: A quality control tool for high throughput sequence data**

627    [http://www.bioinformatics.babraham.ac.uk/projects/fastqc/]

628    37. Bolger AM, Lohse M, Usadel B: **Trimmomatic: a flexible trimmer for Illumina sequence**

629    **data.** *Bioinformatics* 2014, **30**:2114–2120.

630    38. Hansen KD, Brenner SE, Dudoit S: **Biases in Illumina transcriptome sequencing caused**

631    **by random hexamer priming.** *Nucleic Acids Res* 2010, **38**:e131.

632    39. Grabherr MG, Haas BJ, Yassour M, Levin JZ, Thompson DA, Amit I, Adiconis X, Fan L,

633    Raychowdhury R, Zeng Q, Chen Z, Mauceli E, Hacohen N, Gnirke A, Rhind N, di Palma F,

634    Birren BW, Nusbaum C, Lindblad-Toh K, Friedman N, Regev A: **Full-length transcriptome**

635    **assembly from RNA-Seq data without a reference genome.** *Nat Biotechnol* 2011, **29**:644–

636    652.

637    40. Schmieder R, Edwards R: **Fast identification and removal of sequence contamination**

638    **from genomic and metagenomic datasets.** *PLoS One* 2011, **6**:e17288.

639    41. Haas BJ, Papanicolaou A, Yassour M, Grabherr M, Blood PD, Bowden J, Couger MB,

640    Eccles D, Li B, Lieber M, Macmanes MD, Ott M, Orvis J, Pochet N, Strozzi F, Weeks N,

641   Westerman R, William T, Dewey CN, Henschel R, Leduc RD, Friedman N, Regev A: **De novo**

642   **transcript sequence reconstruction from RNA-seq using the Trinity platform for reference**

643   **generation and analysis.** *Nat Protoc* 2013, **8**:1494–1512.

644   42. Brocchieri L, Karlin S: **Protein length in eukaryotic and prokaryotic proteomes.** *Nucleic*

645   *Acids Res* 2005, **33**:3390–3400.

646   43. Simao FA, Waterhouse RM, Ioannidis P, Kriventseva E V, Zdobnov EM: **BUSCO:**

647   **assessing genome assembly and annotation completeness with single-copy orthologs.**

648   *Bioinformatics* 2015, **31**:3210–3212.

649   44. Conesa A, Gotz S, Garcia-Gomez JM, Terol J, Talon M, Robles M: **Blast2GO: a universal**

650   **tool for annotation, visualization and analysis in functional genomics research.**

651   *Bioinformatics* 2005, **21**:3674–3676.

652   45. Gotz S, Garcia-Gomez JM, Terol J, Williams TD, Nagaraj SH, Nueda MJ, Robles M, Talon

653   M, Dopazo J, Conesa A: **High-throughput functional annotation and data mining with the**

654   **Blast2GO suite.** *Nucleic Acids Res* 2008, **36**:3420–3435.

655   46. Zhong Y-F, Butts T, Holland PWH: **HomeoDB: a database of homeobox gene diversity.**

656   *Evol Dev* 2008, **10**:516–518.

657   47. van der Zee M, Berns N, Roth S: **Distinct Functions of the Tribolium zerknullt Genes in**

658   **Serosa Specification and Dorsal Closure**. *Curr Biol* 2005, **15**:624–636.

659   48. Schulz C, Schröder R, Hausdorf B, Wolff C, Tautz D: **A caudal homologue in the short**

660   **germ band beetle Tribolium shows similarities to both, the Drosopila and the vertebrate**

661   **caudal expression patterns**. *Dev Genes Evol* 1998, **208**:283–289.

662   49. Van der Zee M, da Fonseca RN, Roth S: **TGFbeta signaling in Tribolium: vertebrate-like**

663   **components in a beetle.** *Dev Genes Evol* 2008, **218**:203–213.

664   50. Ozuak O, Buchta T, Roth S, Lynch JA: **Ancient and diverged TGF-beta signaling**

665   **components in Nasonia vitripennis.** *Dev Genes Evol* 2014, **224**:223–233.

666   51. Kenny NJ, Namigai EKO, Dearden PK, Hui JHL, Grande C, Shimeld SM: **The**

667    **Lophotrochozoan TGF-beta signalling cassette - diversification and conservation in a key**

668    **signalling pathway.** *Int J Dev Biol* 2014, **58**:533–549.

669    52. Gilbert SF (Ed): *Developmental Biology*. 10th edition. Sunderland, MA: Sinauer Associates,

670    Inc.; 2013.

671    53. Li B, Dewey CN: **RSEM: accurate transcript quantification from RNA-Seq data with or**

672    **without a reference genome.** *BMC Bioinformatics* 2011, **12**:323.

673    54. Gotz S, Arnold R, Sebastian-Leon P, Martin-Rodriguez S, Tischler P, Jehl M-A, Dopazo J,

674    Rattei T, Conesa A: **B2G-FAR, a species-centered GO annotation repository.** *Bioinformatics*

675    2011, **27**:919–924.

676    55. Altschul SF, Gish W, Miller W, Myers EW, Lipman DJ: **Basic local alignment search tool.**

677    *J Mol Biol* 1990, **215**:403–410.

678    56. Katoh K, Standley DM: **MAFFT multiple sequence alignment software version 7:**

679    **improvements in performance and usability.** *Mol Biol Evol* 2013, **30**:772–780.

680    57. Tamura K, Stecher G, Peterson D, Filipski A, Kumar S: **MEGA6: Molecular Evolutionary**

681    **Genetics Analysis version 6.0.** *Mol Biol Evol* 2013, **30**:2725–2729.

682

683    **Figure legends**

684    *Figure 1:* A) Cladogram of Coleoptera simplified from that determined by Hunt *et al.* (2007) [13].

685    Red asterisk indicates the superfamily in which the Chrysomelidae are located. B) Cladogram of

686    Chrysomelidae simplified from that determined by Gomez-Zurita *et al.* (2007) [28]. Black and

687    grey asterisks indicate the sub-families in which *A. menetriesi* and *C. maculatus* are located

688    (respectively). C) *A. menetriesi* and D) *C. maculatus* adults.

689

690    *Figure 2*: Summary of RNA sources, life stages sampled and sequencing results, with

691    *A.menetriesi* data presented at left and *C. maculatus* data at right.

692

693   *Figure 3:* Overview of results of differential expression analysis performed by RSEM within the

694   Trinity framework. A) and B) show the sample correlation matrix for *A. menetriesi* and *C.*

695   *maculatus* respectively. C) and D) show relative expression of each differentially expressed

696   contig, considered as isoforms, across time.

697

698   *Figure 4:* Blast2GO results A) Distribution of BLASTx best hits by species, showing metazoans

699   only, for *A. menetriesi* in orange, *C. maculatus* in blue. B) Distribution of GO terms expressed as

700   a percentage of annotated contigs which were assigned a term within each of the three

701   (Molecular Function, Cellular Component, Biological Process) categories of GO ID.

702

703   *Figure 5:* Phylogenetic inter-relationships of ANTP HOXL class genes*,* as reconstructed by

704   MEGA 6 using the LG+Freqs model with 4 gamma categories and invariant sites, based on a 59

705   amino acid alignment spanning the homeodomain. Numbers at base of nodes represent

706   bootstrap percentages of 1000 replicates. Scale bar at base of phylogeny gives substitutions

707   per site at given unit distance. Red underline indicates *A. menetriesi* and *C. maculatus*

708   sequences, coloured boxes are used to delineate known gene families (and in the case of Hox

709   6/7/8, a superfamily).

710

711   *Figure 6:* Maximum Likelihood Phylogeny of TGF-β ligands, as determined using MEGA under

712   the LG+Freqs model with 4 gamma categories and invariant sites, on the basis of a 72 amino

713   acid alignment of mature peptide sequences. The given scale depicts the number of

714   substitutions per site per unit length. Bootstrap percentages (of 1000 replicates) are given at

715   base of nodes. *A. menetriesi* and *C. maculatus* sequences are underlined in red. Coloured

716   boxes represent known gene families with representatives in our transcriptomic resources, while

717   all gene families, including those not found in our datasets, are indicated at right.

718

719    *Figure 7:* KEGG style pathway maps showing recovery in our transcriptome resources of A) the

720    Wnt signalling pathway in canonical and non-canonical contexts, B) the Hedgehog signalling

721    pathway and C) the Notch signalling pathway. Coloration of genes indicates presence, absence

722    or ancestral absence from the Coleoptera as detailed in the key, which also gives other

723    information as noted.

724

725    **Tables**

726    *Table 1:* Basic read data metrics

|  | A. menetriesi | | | | | C. maculatus | | |
|---|---|---|---|---|---|---|---|---|
| *Library: (hpf= hours post fertilization)* | *Ovary* | *0-24 hpf* | *24-48 hpf* | *48-72 hpf* | *72-100 hpf* | *Ovary* | *0-7 hpf* | *7-24 hpf* |
| *Number of Paired-end Reads* | 44,695,199 | 56,139,626 | 39,911,618 | 46,217,614 | 45,213,862 | 67,026,913 | 61,060,357 | 55,273,927 |
| *Read Length (bp)* | 101 | 101 | 101 | 101 | 101 | 101 | 101 | 101 |
| *Amount of Data (Gigabytes)* | 7.3 | 9.5 | 6.8 | 7.8 | 7.7 | 10.8 | 9.8 | 8.8 |
| *GC content (%)* | 38 | 36 | 37 | 37 | 37 | 43 | 41 | 40 |

727

728

729

730

731

732    *Table 2:* Assembly metrics for individual time point and combined transcriptomic resources

| | *A. menetriesi* | | | | | | *C. maculatus* | | | |
|---|---|---|---|---|---|---|---|---|---|---|
| | Ovary | 0-24 hpf | 24-48 hpf | 48-72 hpf | 72-100 hpf | Combined | Ovary | 0-7 hpf | 7-24 hpf | Combined |
| Number of contigs | 100,084 | 106,548 | 120,497 | 140,868 | 126,400 | 228,096 | 98,082 | 57,879 | 80,665 | 128,837 |
| Max contig length (bp) | 20,759 | 17,765 | 18,398 | 17,591 | 17,180 | 20,757 | 30,742 | 22,941 | 22,956 | 29,933 |
| Mean contig length (bp) | 983.95 | 1027.36 | 924.27 | 879.03 | 914.2 | 844.66 | 1074.86 | 1216.55 | 1106.56 | 991.04 |
| Median contig length (bp) | 471 | 490 | 453 | 426 | 448 | 401 | 416 | 526 | 474 | 391 |
| N50 contig length (bp) | 1,894 | 1,993 | 1,721 | 1,643 | 1,686 | 1,598 | 2,450 | 2,570 | 2,308 | 2,263 |
| # contigs in N50 | 14,196 | 15,200 | 17,511 | 19,814 | 18,105 | 31,196 | 11,695 | 7,874 | 10,978 | 15,155 |
| # contigs > 1kb | 28,255 | 31,944 | 32,307 | 34,520 | 33,002 | 52,217 | 27,763 | 19,919 | 25,529 | 33,250 |
| # bases, total | 98,477,173 | 109,463,362 | 111,371,655 | 123,827,226 | 115,555,146 | 192,664,213 | 105,424,220 | 70,412,726 | 89,260,499 | 127,682,507 |
| # bases in contigs > 1kb | 68,737,074 | 78,526,348 | 75,208,106 | 80,796,007 | 77,199,414 | 122,986,547 | 78,519,052 | 55,032,546 | 67,296,566 | 91,466,099 |

| GC Content % (2dp) | 33.54 | 33.14 | 33.14 | 33 | 33.17 | 32.82 | 39.26 | 39.21 | 38.59 | 38.7 |
|---|---|---|---|---|---|---|---|---|---|---|

733

734

735

736

737

738 Table 3: Genes identified by manual annotation

| Pathway/Gene | A. menetriesi | C.maculatus | Pathway/Gene | A. menetriesi | C.maculatus |
|---|---|---|---|---|---|
| Maternal Effect: | | | Pair rule: | | |
| caudal | present | present | even skipped | present | present |
| hunchback | present | present | fushi tarazu | present (see Hox) | present (see Hox) |
| | | | hairy | present | present |
| Germline: | | | odd paired | present | present |
| bruno/arrest | present | present | odd skipped | present | present |
| exuperantia | present | present | paired | present | present |
| nanos | absent | absent | runt | present | present |
| oskar | present | present | sloppy paired 1 | present | present |
| pumilio | present - 7 copies | present | | | |
| swallow | absent | absent | Segment Polarity: | | |
| tudor | present - 2 copies | present | armadillo | present | present |
| valois | present | present | cubitis interruptus | present | present |
| vasa | present | present | Engrailed | present | present |

|  |  |  | *fused* | present | present |
|---|---|---|---|---|---|
| Gap genes: |  |  | *gooseberry* | present | present |
| *buttonhead* | present | present | *gooseberry-neuro* | absent | absent |
| *empty spiracles* | present | present | *hedgehog* | present | present (but on 3 contigs) |
| *giant* | present | present | *pangolin* | present | present |
| *huckebein* | present | present | *patched* | present | present |
| *knirps* | present | present | *wingless* | present | present |
| *Krüppel* | present | present |  |  |  |
| *orthodenticle 1* | present | present |  |  |  |
| *orthodenticle 2* | present | present |  |  |  |
| *tailless* | present | present - 2 copies |  |  |  |

739

740

## Additional Files:

742 Additional File 1: Sequences of all genes referred to in text, and alignments used in

743 phylogenetic analyses (.xls)

744 Additional File 2: Annotations of transcriptome (Blast2GO .annot file) *A. menetriesi*

745 Additional File 3: Annotations of transcriptome (Blast2GO .annot file) *C. maculatus*

746 Additional File 4: KEGG data, *A. menetriesi*

747 Additional File 5: KEGG data, *C. maculatus*

748 Additional File 6: Comparative expression data, *A. menetriesi*

749 Additional File 7: Comparative expression data, *C. maculatus*

**A)** Coleoptera

**B)** Chrysomelidae

**C)** *Atrachya menetriesi*

2 cm

**D)** *Callosobruchus maculatus*

1 cm

| | Tenebrionoidea | | Cerylonid Series |
|---|---|---|---|
| | Cleroidea | | Cucujoidea |
| | Chrysomeloidea | | Curculionoidea |
| | Elateriformia | | Bostrichiformia |
| | Scarabaeiformia | | Staphyliniformia |
| | Lymexyloidea | | Hydradephaga |
| | Geadephaga | | |

| | Alticinae | | Galerucinae |
|---|---|---|---|
| | Chrysomelinae | | Caddidinae |
| | Hispinae | | Clytrinae |
| | Cryptocephalinae | | Chlamisinae |
| | Eumolpinae | | Spilopyrinae |
| | Criocerinae | | Donaciinae |
| | Synetinae | | Bruchinae |

| *Atrachya menetriesi* | Developing Gamete | *Callosobruchus maculatus* |
|---|---|---|
| **Sample 1: Mixed Ovarian**<br>44,695,199 paired end reads<br>100,084 contigs | **Ovarian Development** | **Sample 1: Mixed Ovarian**<br>67,026,913 paired end reads<br>98,082 contigs |
| **Sample 2: 0-24 hpf**<br>56,139,626 paired end reads<br>106,548 contigs | **Early Egg to Blastoderm Stage** | **Sample 2: 0-7 hpf**<br>61,060,357 paired end reads<br>57,879 contigs |
| **Sample 3: 24-48 hpf**<br>39,911,618 paired end reads<br>120,497 contigs | **Embryo Formation, Gastrulation, Germband Elongation** | **Sample 3: 7-24 hpf**<br>55,273,927 paired end reads<br>80,665 contigs |
| **Sample 4: 48-72 hpf**<br>46,217,614 paired end reads<br>140,868 contigs | | |
| **Sample 5: 72-100 hpf**<br>45,213,862 paired end reads<br>126,400 contigs | | |
| **Samples 1-5 Combined:**<br>232,177,919 paired end reads<br>228,096 contigs | **Final, Combined Assemblies** | **Samples 1-3 Combined:**<br>183,361,197 paired end reads<br>128,837 contigs |

**A)** *Atrachya menetriesi* Sample Correlation Matrix

Color Key
Relative Similarity
0.2 0.4 0.6 0.8 1

72-100 hour
48-72 hour
24-48 hour
0-24 hour
ovary

ovary
0-24 hour
24-48 hour
48-72 hour
72-100 hour

**B)** *Callosobruchus maculatus* Sample Correlation Matrix

Color Key
Relative Similarity
0.2 0.4 0.6 0.8 1

7-24 hour
0-7 hour
ovary

ovary
0-7 hour
7-24 hour

**C)** *Atrachya menetriesi* Contig Relative Expression

Color Key
Relative expression
-10 -5 0 5 10

ovary
0-24 hour
24-48 hour
48-72 hour
72-100 hour

**D)** *Callosobruchus maculatus* Contig Relative Expression

Color Key
Relative expression
-5 0 5

ovary
0-7 hour
7-24 hour

**A)** Legend:
- *A. menetriesi*
- *C. maculatus*

Y-axis: Number of Top Hits

X-axis categories: Tribolium castaneum, Dendroctonus ponderosae, Acyrthosiphon pisum, Stegodyphus mimosarum, Bombyx mori, Diaphorina citri, Cerapachys biroi, Camponotus floridanus, Hydra vulgaris, Microplitis demolitor, Harpegnathos saltator, Danaus plexippus, Acromyrmex echinatior, Zootermopsis nevadensis, Drosophila melanogaster, Nasonia vitripennis, Strongylocentrotus purpuratus, Aedes aegypti, Capitella teleta, Branchiostoma floridae, Other

**B)** Legend:
- *Atrachya menetriesi*
- *Callosobruchus maculatus*
- *Drosophila melanogaster*
- *Mus musculus*

Y-axis: Percentage of GO Category (BP, CC or MF) with Given Annotation

Molecular Function: Catalytic Activity, Signal Transducer Activity, Receptor Activity, Structural Molecule Activity, Transporter Activity, Binding, Protein Binding, Enzyme Regulator Activity, Oxidoreductase Activity, Transferase Activity, Hydrolase Activity, Ligase Activity, Ion Transmembrane Transporter Activity, Channel Activity

Cellular Component: Extracellular Region, Extracellular Space, Cell, Nucleus, Chromosome, Cytoplasm, Membrane, Macromolecular Complex

Biological Process: Transport, Multicellular Organismal Development, Metabolic Process, Catabolic Process, Biosynthetic Process, Cellular Process, Cell Differentiation, Macromolecule Metabolic Process, Secretion, Regulation Of Biological Process, Response To Stimulus, Biological Regulation, Cellular Metabolic Process, Oxidation-Reduction Process

*Atrachya menetriesi comp103230 c0 seg2 Utx (10 variants)*
*Callosobruchus maculatus comp41872 c0 seg15 Utx (4 variants)*
*Tribolium castaneum Utx*
*Drosophila melanogaster Ubx*
*Drosophila melanogaster abd-A*
*Tribolium castaneum abd-A*
*Atrachya menetriesi comp139998 c0 seg1 abd-A*
*Branchiostoma floridae Hox8*
*Branchiostoma floridae Hox6*
*Branchiostoma floridae Hox7*
*Drosophila melanogaster Antp*
*Tribolium castaneum ptl*
*Atrachya menetriesi comp109961 c0 seq1 ptl (2 variants)*
*Callosobruchus maculatus comp41872 c0 seq2 ptl (12 variants)*
*Drosophila melanogaster ftz*
*Tribolium castaneum ftz*
*Atrachya menetriesi comp92490 c0 seq1 ftz*
*Callosobruchus maculatus comp42069 c2 seq12 ftz*

***Hox6/7/8***
***Utx/Ubx***
***Abd-A***
***Ptl/Antp***
***Ftz***

*Branchiostoma floridae Hox5*
*Tribolium castaneum Cx*
*Drosophila melanogaster Scr*
*Atrachya menetriesi comp75090 c0 seq2 Cx*
*Callosobruchus maculatus comp41872 c0 seq9 Cx (5 variants)*

***Hox5 / Cx/Scr***

*Branchiostoma floridae Hox4*
*Callosobruchus maculatus comp37257 c0 seq2 comp37257 c1 seq1 Dfd*
*Drosophila melanogaster Dfd*
*Tribolium castaneum Dfd*

***Hox4 / Dfd***

*Branchiostoma floridae Xlox*
*Branchiostoma floridae Hox3*
*Tribolium castaneum zen1*
*Atrachya menetriesi comp102511 c0 seq2 zen (2 variants)*
*Callosobruchus maculatus comp34354 c0 seq6 zen (4 variants)*
*Tribolium castaneum zen2*

***Hox3 / Zen***

*Drosophila melanogaster Abd-B*
*Tribolium castaneum Abd-B*
*Callosobruchus maculatus comp44633 c0 seq26 Abd-B (26 variants)*
*Atrachya menetriesi comp146889 c0 seq1 Abd-B*
*Branchiostoma floridae Hox10*
*Branchiostoma floridae Hox9*

***Hox9-13 / Abd-B***

*Branchiostoma floridae Cdx*
*Drosophila melanogaster cad*
*Tribolium castaneum Cad2*
*Tribolium castaneum Cad1*
*Atrachya menetriesi comp39955 c0 seq1 Cad*
*Callosobruchus maculatus comp45520 c0 seq1 Cad*

***Cdx / Cad***

*Atrachya menetriesi comp42215 c0 seq1 mxp (NB Dfd like)*
*Branchiostoma floridae Hox2*
*Callosobruchus maculatus comp37995 c0 seq3 mxp (3 variants)*
*Tribolium castaneum mxp*
*Drosophila melanogaster pb*

***Hox2 / Mxp/Pb***

*Branchiostoma floridae Hox1*
*Drosophila melanogaster lab*
*Tribolium castaneum lab*
*Atrachya menetriesi comp74830 c0 seq2 lab (2 variants)*
*Callosobruchus maculatus comp42069 c2 seq15 seq1 lab*

***Hox1 / Lab***

*Drosophila melanogaster zen*
*Drosophila melanogaster zen2*
*Drosophila melanogaster ind*
*Branchiostoma floridae Gsx*
*Callosobruchus maculatus comp26552 c0 seq2 Ind*
*Tribolium castaneum Ind*
*Atrachya menetriesi comp108136 c0 seq1 Ind (2 variants)*

***Gsx / Ind***

*Callosobruchus maculatus comp38822 c1 seq3 Btn*
*Drosophila melanogaster btn*
*Atrachya menetriesi comp68233 c0 seq1 Btn*
*Tribolium castaneum Btn*
*Branchiostoma floridae Mox*

***Mox / Btn***

*Branchiostoma floridae Evxb*
*Branchiostoma floridae Evxa*
*Tribolium castaneum Eve*
*Drosophila melanogaster eve*
*Atrachya menetriesi comp98599 c0 seq1 Eve (2 variants)*
*Callosobruchus maculatus comp39876 c0 seq1 Eve.*

***Evx / Eve***

*Tribolium castaneum Unpg*
*Atrachya menetriesi comp127691 c0 seq1 Unpg*
*Callosobruchus maculatus comp2890 c0 seq2 Unpg*
*Drosophila melanogaster unpg*
*Branchiostoma floridae Gbx*

***Gbx / Unpg***

*Branchiostoma floridae Mnxa*
*Branchiostoma floridae Mnxb*
*Drosophila melanogaster exex*
*Callosobruchus maculatus comp12077 c1 seq1 Exex*
*Tribolium castaneum Exex*
*Atrachya menetriesi comp84723 c0 seq1 Exex*

***Mnx / Exex***

*Drosophila melanogaster bcd*
*Branchiostoma floridae Pitx*
*Drosophila melanogaster Ptx1*
*Tribolium castaneum Ptx1*

0.2

Branchiostoma floridae BMP 2/4 AAC97488.1
Mus musculus BMP4 AAC37698.1
Saccoglossus kowalevskii BMP24
Mus musculus BMP2 NP031579.2
Apis mellifera Dpp BMP2/4 XP001122815.2
Drosophila melanogaster Dpp AAN10431.1
Lottia gigantea Dpp/BMP2/4 205842
Nematostella vectensis BMP2/4 AAR13362.1
Tribolium castaneum NP 001034540.1 Dpp
Callosobruchus maculatus comp33819 c0 seq2 Dpp
Atrachya menetriesi comp100246 c0 seq7 Dpp (6 variants)

*Dpp / BMP2/4*

Saccoglossus kowalevskii univin-like
Branchiostoma floridae GDF 5/6/7 XP002602867.1
Nematostella vectensis GDF5 AAR27581.1
Mus musculus GDF5 NP032135.2
Mus musculus GDF7 NP038555.1
Mus musculus GDF6 NP038554.1

*GDF 5/6/7*

Mus musculus BMP8 NP031584.1
Nematostella vectensis BMP5/8 XP 001638358.1
Callosobruchus maculatus comp40430 c0 seq1 Gbb
Atrachya menetriesi comp89425 c0 seq1 Gbb
Tribolium castaneum NP 001107813.1 Gbb1
Apis mellifera Gbb XP394252.1
Tribolium castaneum EFA04646.1 Gbb2
Drosophila melanogaster Gbb AAF47075.1
Saccoglossus kowalevskii BMP5/8
Lottia gigantea BMP5/8 195882
Branchiostoma floridae BMP5/7/8 C3ZTY2
Mus musculus BMP7 NP031583.2
Mus musculus BMP5 NP031581.2
Mus musculus BMP6 NP031582.1

*Gbb /BMP 5/6/7/8*

Drosophila melanogaster Scw AAN11056.2
Lottia gigantea BMP10 111943
Mus musculus BMP10 NP033886.2
Mus musculus BMP9 AAD56961.1
Apis mellifera BMPx (BMP9/10) XP001120039.1
Atrachya menetriesi comp170252 c0 seq1 BMPx (BMP9/10)
Tribolium castaneum XP 973577.1 BMPx (BMP9/10)

*BMPx / BMP 9/10*

Nematostella vectensis ADMP related AFP87424.1
Mus musculus GDF3 NP032134.2
Mus musculus GDF1 NP032133.2
Mus musculus BMP3 NP775580.1

*GDF 1/3*

Saccoglossus kowalevskii BMP3
Branchiostoma floridae ADMP XP002604737.1
Lottia gigantea ADMP 110168
Saccoglossus kowalevskii ADMP
Danio rerio ADMP NP571951.2

*ADMP / BMP3*

Saccoglossus kowalevskii Nodal A
Lottia gigantea Nodal ACB42423.1
Branchiostoma floridae Nodal AAL99367.1
Mus musculus Nodal AAI28019.1

*Nodal*

Callosobruchus maculatus comp48645 c0 seq1 Maverick
Atrachya menetriesi comp106563 c0 seq1 Maverick (3 variants)
Tribolium castaneum XP 001811434.1 Maverick
Apis mellifera BMP2B/Maverick XP001122118.2
Drosophila melanogaster Maverick AAF59328.1

*Maverick*

Saccoglossus kowalevskii TGFB2 ADB22639.1
Mus musculus TGFB 2 NP033393.2
Mus musculus TGFB 1 NP035707.1
Mus musculus TGFB 3 NP033394.2

*TGF β*

Mus musculus GDF11 NP034402.1
Mus musculus Myostatin AAO46885.1
Branchiostoma floridae Myostatin XP002599461.1
Tribolium castaneum XP 966819.1 Myostatin
Callosobruchus maculatus comp39049 c0 seq1 Myo
Atrachya menetriesi comp104142 c0 seq2 Myo (2 variants)
Nematostella vectensis Myostatin XP 001641598.1
Saccoglossus kowalevskii GDF8/11 XP 002734819.1
Lottia gigantea Myostatin 82990
Drosophila melanogaster Myostatin AAF59319.1

*Myostatin / GDF 8/11*

Mus musculus Inhibin e NP032408.2
Apis mellifera Activin XP001123044.2
Tribolium castaneum XP 008194984.1 Act
Drosophila melanogaster Activin NP651942.2
Mus musculus Inhibin a NP032406.1
Mus musculus Inhibin b NP032407.1

*Inhibin / Activin*

Saccoglossus kowalevskii Lefty NP001164679.1
Branchiostoma floridae Inhibin beta XP002606835.1
Saccoglossus kowalevskii Inhibin/Activin NP001161496.1
Nematostella vectensis Activin ABF61781.1

*Lefty*

*Inhibin / Activin*

Lottia gigantea Activin 151945
Apis mellifera ALP XP001122210.1
Drosophila melanogaster Dawdle (ALP) NP722840.1
Tribolium castaneum XP 970355.1 Dawdle (ALP)
Callosobruchus maculatus comp24632 c3 seq1 Dawdle (ALP)
Atrachya menetriesi comp62923 c0 seq1 Dawdle (ALP) (2 variants)

*Dawdle / ALP*

Callosobruchus maculatus comp30559 c0 seq2 putative Derriere (2 variants)
Tribolium castaneum XP 008191586.1 putative Derriere

*Unknown Derriere-Like*

Mus musculus BMP15 NP033887.1
Mus musculus GDF9 NP032136.2

*BMP15 / GDF9*

Mus musculus Neurturin NP 032764.1
Danio rerio GDNF NP 571807.1
Mus musculus GDNF NP 034405.1
Xenopus laevis GDNF NP 001090196.1

0.5

A) Wnt Signalling:

B) Hedgehog Signalling:

C) Notch Signalling: