1    # *De novo* genome and transcriptome analyses provide insights

2    # into the biology of the trematode human parasite *Fasciolopsis*

3    # *buski*

4    Devendra K. Biswal[1]†, Tanmoy Roychowdhury[2]†, Priyatama Pandey[2], Veena

5    Tandon[1,3]§

6    [1]Bioinformatics Centre, North-Eastern Hill University, Shillong 793022,

7    Meghalaya, India

8    [2]School of Computational and Integrative Sciences, Jawaharlal Nehru University,

9    New Delhi 110067, India

10   [3]Department of Zoology, North-Eastern Hill University, Shillong 793022,

11   Meghalaya, India

12   †Equal contributors

13   §Corresponding author

14

15   Email addresses:

16   DKB: devbioinfo@gmail.com

17   TR: tanmoy63@gmail.com

18   PP: priyatamapandey@gmail.com

19   VT: tandonveena@gmail.com

20

21

22   # Abstract

23  Many trematode parasites cause infection in humans and are thought to be a

24  major public health problem. Their ecological diversity in different regions

25  provides challenging questions on evolution of these organisms. In this report,

26  we perform transcriptome analysis of the giant intestinal fluke, *Fasciolopsis*

27  *buski*, using next generation sequencing technology. Short read sequences

28  derived from polyA containing RNA of this organism were assembled into 30677

29  unigenes that led to the annotation of 12380 genes. Annotation of the assembled

30  transcripts enabled insight into processes and pathways in the intestinal fluke,

31  such as RNAi pathway and energy metabolism. The expressed kinome of the

32  organism was characterized by identifying all protein kinases. We have also

33  carried out whole genome sequencing and used the sequences to confirm

34  absence of some of the genes, not observed in transcriptome data, such as

35  genes involved in fatty acid biosynthetic pathway. Transcriptome data also

36  helped us to identify some of the expressed transposable elements. Though

37  many Long Interspersed elements (LINEs) were identified, only two Short

38  Interspersed Elements (SINEs) were visible. Overall transcriptome and draft

39  genome analysis of *F. buski* helped us to characterize some its important

40  biological characteristics and provided enormous resources for development of a

41  suitable diagnostic system and anti-parasitic therapeutic molecules.

42

43

44  **Introduction**

45 *Fasciolopsis buski* (giant intestinal fluke) is one of the largest digenean

46 trematode flatworms infecting humans causing the disease fasciolopsiasis, with

47 epidemiological records from Asia including the Indian subcontinent. Clinical

48 manifestation of the disease is diarrhea, with presence of ulceration, intestinal

49 wall abscess and hemorrhage, with the possibility of death if not treated in time

50 [1]. Majority of the infected people (up to 60% in India and the mainland China)

51 remain asymptomatic without manifestation of clinical invasive disease [1,2].

52 In India, these trematodes have been reported from different regions

53 including the Northeast, associated with animal rearing (mainly pigs), or in

54 underwater vegetables such as water chestnut. Although *F. buski* inhabits warm

55 and moist regions and has been reported as single species in the genus,

56 morphological variations among flukes from different geographical isolates have

57 been observed suggesting genetic polymorphism or host-specific parasite

58 adaptation [3,4].

59 Infective metacercarial cysts usually occur in 15-20 groups on the surface

60 of aquatic vegetation [5,6] and once consumed, the juvenile adult stage of *F.*

61 *buski* emerges and adheres to the small intestine of its host, attached until the

62 host dies or is removed [1,2,7]. Unembryonated eggs are discharged into the

63 intestine in the fecal material, with release of parasitic organisms (miracidia) two-

64 week post infection. These hatch between 27–32°C, and invade snails

65 (intermediate host) where developmental stages (sporocysts, rediae, and

66 cercariae) are reported. The cercariae released and encysted as metacercariae

67    (on aquatic plants), infect the mammalian host and develop into adult flukes

68    completing its life cycle [8].

69         The current strategy for control of fasciolopsiasis, similar to that for other

70    trematodes, such as *Gastrodiscoides hominis* is by blocking transmission among

71    different hosts, i.e., human, animal reservoir and intermediate host (molluscan)

72    [5]. Many drugs have been used to treat fasciolopsiasis. The drug of choice is

73    praziquantel since it has high efficacy even in cases of severe fasciolopsiasis [9,

74    10, 11, 12, 13, 14]. More recently, the efficacy of triclabendazole, oxyclozanide

75    and rafoxanide has been evaluated in pigs with favorable response to treatment

76    [15]. Despite control programs, public health is a concern in endemic areas, and

77    there is a need for new control measures with reports of sporadic re-emergence

78    of infection (in Uttar Pradesh, India) from non-endemic regions [16, 17,18, 19].

79         Development of novel therapeutic agents targeting intestinal flukes is

80    complex. The parasite biology occurs in multi-host life cycles and mechanisms

81    for parasite host immune evasion strategies are not well understood. Genomics

82    approaches are likely to be more successful in identifying new targets for

83    intervention. Unfortunately, little genomic information is available for *F. buski* in

84    the public domain other than a PCR-based molecular characterization using ITS

85    1 & 2 regions of ribosomal DNA (rDNA) genes [20,21]. Developments in next

86    generation sequencing (NGS) technologies and computational analysis tools

87    enable rapid data generation to decipher organismal biology [22, 23, 24, 25, 26].

88    NGS analysis to understand transcriptomes of related organisms, such as

89    *Fasciola gigantica*, *Fasciola hepatica*, *Schistosoma mansoni*, *Schistosoma*

90     *japonicum*, *Clonorchis sinensis, Opisthorchis viverrini, Fascioloides magna* and

91     *Echinostoma caproni* reveal genes involved in adult parasite-host responses,

92     such as antioxidants, heat shock proteins and cysteine proteases [27, 28, 29, 30,

93     31, 32]. Lack of molecular analysis of *F. buski* has hampered understanding of

94     evolutionary placement of this organism among other *Fasciola*, as well as our

95     understanding of mechanisms that regulate host-pathogen relationships.

96     Previously we have reported the mitochondrial DNA (mt DNA) of the intestinal

97     fluke for the first time that would help investigate Fasciolidae taxonomy and

98     systematics with the aid of mtDNA NGS data [33]. Here, we report results of our

99     effort to characterize the transcriptome of the adult stage of *F. buski.* NGS

100     technology was used to get RNA-seq and bioinformatics analysis was carried out

101     on the sequence output. In a few cases we have used the *F. buski* draft genome

102     sequence to confirm our observations.

103

## 104     Materials and methods

### 105     Source of parasite material

106     *F. buski* adult specimens were obtained from the intestine of freshly

107     slaughtered pig (*Sus scrofa* domestica) by screening for naturally infected pigs

108     among the animals routinely slaughtered for meat purpose at the local abattoirs.

109     The worms reported in this study represent geographical isolates from the

110     Shillong area (25.57°N, 91.88°E) in the state of Meghalaya, Northeast India.

111     Eggs were obtained by squeezing mature adult flukes between two glass slides.

112     Adult flukes collected from different pig hosts were processed singly for the

113    purpose of DNA extraction and eggs were recovered from each of these

114    specimens separately.

115    All sample were obtained from animals used for local meat consumption

116    and no live animals were handled or sacrificed for this study. Hence there was no

117    need to follow relevant guidelines as laid down for handing live vertebrates.

## 118    RNA and DNA extraction, Illumina sequencing

119    Briefly, the adult flukes were digested overnight in extraction buffer

120    [containing 1% sodium dodecyl sulfate (SDS), 25 mg Proteinase K] at 37°C, prior

121    to DNA recovery. Genomic DNA was then extracted from lysed individual worms

122    by ethanol precipitation technique or by FTA cards using Whatman's FTA

123    Purification Reagent [34]. Quality and quantity of the DNA was assessed by 0.8%

124    agarose gel electrophoresis (Sigma-Aldrich, USA), spectrophotometer (Nanodrop

125    1000, Thermofischer USA) and flurometer (Qubit, Thermofischer, USA). Whole

126    genome shot-gun libraries for Illumina sequencing were generated and the

127    genomic DNA was sheared (Covaris, USA), end-repaired, adenylated

128    phosphorylate ligated to Illumina adapters (TruSeq DNA, Illumina, USA), to

129    generate short (300-350 bp) fragment and long (500 – 550) bp insert libraries.

130    Both libraries were amplified using 10 cycles of PCR KapaHiFiHotstart PCR

131    ready mix (Kapa Biosystems Inc., USA) and adapter removal by Solid Phase

132    Reversible Immobilization (SPRI) beads (Ampure XP beads, Beckmann-coulter,

133    USA). Prior to sequencing, fragment analysis was performed (High Sensitivity

134    Bioanalyzer Chip, Agilent, USA), quantified by qPCR. 36cycle sequencing was

135    performed on Illumina Genome analyzer II (SBS kit v5, Illumina, USA) and

136  Illumina HiSeq1000 (Illumina, USA) to generate 16 million 100 paired-end reads

137  from the shot-gun library.

138  Total RNA was collected from a single frozen adult individual fluke (~100mg

139  tissue) using TRI reagent® (Sigma, USA). The integrity of total RNA was verified

140  using Bioanalyzer Agilent 2100 with RNA integrity number 9. Transcriptome

141  libraries for sequencing was constructed from poly A RNA isolated from 1ug total

142  RNA (OligoTex mRNA minikit, Qiagen Gmbh, Germany), that was fragmented,

143  reverse transcribed (Superscript II Reverse transcriptase, Invitrogen, USA) with

144  random hexamers as described in (TruSeq RNA Sample Preparation Kit,

145  Illumina, USA). The cDNA was end-repaired, adenlylated and cleaned up by

146  SPRI beads (Ampure XP, Beckman Coulter, USA) prior to ligation with Illumina

147  single index sequencing adapters (TruSeq RNA, Illumina, USA). The library was

148  amplified using 11 cycles of PCR and quantified using Nanodrop (Agilent). The

149  prepared library was validated for quality by running an aliquot on High

150  Sensitivity Bioanalyzer Chip (Agilent) that displayed a confident bioanalyzer

151  profile for transcriptome sequencing.

152  **De novo assembly, annotation and characterization of**

153  **transcriptome**

154  Illumina-generated paired end reads were filtered for low quality scores (PHRED

155  score < 30). From 32010511 (32.01 millions) high quality raw reads (>70% of

156  bases in a read with >20 phred score), 47957 contiguous sequences were

157  assembled. Reads containing ambiguous characters were also removed from the

158  dataset. The high-quality reads were then assembled using Trinity software [35]

159    designed for transcriptome assembly. Default parameters were used for this

160    purpose. While Trinity can identify different isoforms, our initial target was to

161    identify the unigene transcripts. Thus, transcripts were subjected to clustering

162    using CD-HIT-EST [36] at 90% similarity. To assess the quality of assembly, raw

163    reads were mapped back to unigenes using Bowtie [37]. Reads aligning in

164    multiple locations were randomly allocated to one of the unigenes. Different

165    properties, such as average read depth, coverage, and percent uniquely mapped

166    reads were calculated. The numbers of raw reads were normalized for length and

167    RPKM (reads per kilobase of exon per million of mapped reads) values were

168    calculated in order to quantify relative abundance of each unigene for estimating

169    the levels of expression.

170        All unigenes were annotated using both sequence and domain-based

171    comparisons that involved sequence similarity analysis with non-redundant (NR)

172    protein database of NCBI (http://www.ncbi.nlm.nih.gov/), Uniprot (Swissprot +

173    trEMBL) [38] and Nembase4 EST database [39] using BLASTx and tBLASTn

174    [40]. Significant threshold of E-value $\leq 10^{-5}$ was used as a cutoff. Since we were

175    not able to annotate large numbers of unigenes using sequence similarity,

176    conserved domains were further identified in InterPro database [41] (HMMpfam,

177    HMMsmart, HMMpanther, FPrintScan, BlastProDom) and Clusters of

178    Orthologous Groups (COG) database [42] using InterProScan and BLASTx

179    respectively. Functional categorization was assigned by finding Gene Ontology

180    terms of best BLASTx hits against NR database using Blast2GO software [43]

181    (http://www.blast2go.com/b2ghome). GO assignments were represented in a plot

182 generated by WEGO [44]. Comparative analysis of the transcriptome was

183 performed against other trematodes. Transcriptome datasets of trematode

184 parasites, viz. *Fasciola hepatica, Fasciola gigantica, Clonorchis sinensis,*

185 *Opisthorchis viverrini, Schistosoma mansoni* and *Schistosoma japonicum* were

186 downloaded from NCBI Short Read Archive (SRA) database

187 (http://www.ncbi.nlm.nih.gov/sra/). Unigenes were mapped to different KEGG

188 pathways using KAAS [45]. All unigenes were used against the KEGG database

189 by a Bi-directional best-hit method suitable for whole genome or transcriptome

190 data. KEGG orthologs were displayed in iPath2 [46]. ORFs were predicted using

191 ORFpredictor [47] (http://proteomics.ysu.edu/tools /OrfPredictor.html). Signal

192 peptides were identified using SignalP [48]

193 (http://www.cbs.dtu.dk/services/SignalP/) and transmembrane domains were

194 identified using TMHMM [49]. Kinases and peptidases were identified by

195 comparison against EMBL kinase database (http://www.sarfari.org/kinasesarfari)

196 and MEROPS database [50] consecutively. Orthologs of RNAi pathway proteins

197 of *Caenorhabditis elegans* in *F.buski* were identified by a reciprocal BLAST hit

198 strategy. Raw reads generated from *F. buski* transcriptome were submitted in

199 NCBI SRA (https://trace.ncbi.nlm.nih.gov/Traces/sra/?run=SRR941773) under

200 the accession number SRX326786.

201 ***De novo* assembly of draft genome and identification of**

202 **transposable elements:**

203 Short read sequences from the *F. buski* genome were assembled using a

204 number of assembly tools, such as SOAP denovo [51], Velvet [52] and Abyss

205 [53]. These algorithms use de Bruijn graph for genomic sequence assembly.

206 Default parameters were used to assemble 42425988 paired end reads of

207 different insert size libraries. Each assembly program was used with a range of

208 K-mers. Statistics of best assemblies in respect of N50 values by three programs

209 are given in S1 Table. These assembly results suggest Trinity assembler

210 generated higher N50 comparatively for this dataset. Further analysis of *F. buski*

211 genome was performed on this reference draft genome assembly. RepeatMasker

212 [54] was used to identify repetitive and transposable elements with default

213 parameters. It makes use of RepBase libraries [version 20120418]

214 (www.girinst.org) that are used as a reference for the identification of

215 transposable and repetitive elements in a query sequence. Transcripts obtained

216 from the transcriptome were matched to the draft assembly. The raw reads

217 generated from *F. buski* genome were submitted in NCBI SRA

218 (www.ncbi.nlm.nih.gov/sra/) under the accession numbers SRX3087506,

219 SRX327895 (genome data) and SRX326786 (transcriptome data).

220 # Results

221 ## *De novo* assembly and annotation of transcriptome

222 The whole genome and transcriptome data are sequenced from the adult

223 specimens of *F. buski* that were collected from pigs with a naturally acquired

224 infection at an abattoir in Shillong, Meghalaya. Detailed statistics of initial output

225 generated by Trinity is shown in S1 Table.  Initial output was then subjected to

226 CD-HIT-EST to identify clusters and the results showed 43111 clusters with 3298

227 singletons representing 46409 unigene transcripts.  After assessment of the

228  assembly (S2 Table), 30677 high-quality unigenes were used for subsequent

229  analysis. We were able to predict ORFs for 30628 unigenes with the mean length

230  of ORF being 185 amino acids. The methionine start codon (ATG) was present in

231  23759 (77.4%) ORFs.

232        *F. buski* unigenes matched with 12279 sequences of NCBI non-redundant

233  (NR) database with an E-value less than 1E-05 on BLAST analysis (S3 Table).

234  Best hit for each unigene was ranked according to E-values. We report many

235  significant BLAST hits to *C. sinensis* (55%)*, S. mansoni* (21%) and *S. japonicum*

236  (14%), with less similarity to *F. hepatica* or *F. gigantica*. This may be due to the

237  availability and skew of sequence databases for the two *Fasciola* species.

238  Although, an identity cutoff of E-value ≤ $1E^{-05}$ was applied for identification of

239  remote homologs, 68% of the unigenes were matched with similarity of E-values

240  less than $1E^{-30}$. Annotation of the unigenes depended on the length of the

241  unigenes, and shorter unigenes (13% of smaller than 500 bp) were not matched,

242  while the longer unigenes (91% of unigenes above 2000 bp) were annotated.

243  The high percentage of annotation suggested an accurate assembly, and

244  presence of functionally informative datasets for *F. buski*.

245        Unigenes were also annotated by sequence comparison against

246  Uniprot/Swissprot and Uniprot/trEMBL database. We found 8089 unigenes that

247  showed significant similarity with at least one Swissprot entry and 12231 with at

248  least one trEMBL entry, totaling 12242 unique unigenes. Analysis using

249  NEMBASE4 database helped us to find putative function of 6776 unigenes.

250 Overall 6752 unigenes displayed significant sequence matches against all three

251 databases used for sequence-based annotation.

252 Identification of conserved domains and motifs sometimes helps in

253 functional classification of genes. Interpro scan helped us to identify 3309 unique

254 domains/families in 6588 unigenes. Out of these, 1775 (53.6%) domains/families

255 were associated with only one unigene. Some of the best hits are given in S4

256 Table. Among these, protein kinases and ankyrin repeat domains were found to

257 be more prevalent. Protein kinases play a significant role in signaling; Ankyrin

258 repeat containing proteins are involved in a number of functions, such as integrity

259 of plasma membrane, and WD40 repeat is associated with signal transduction

260 and transcription regulation.

261 Unigenes were also classified using Clusters of Orthologous (COG)

262 database at http://www.ncbi.nlm.nih.gov/COG/. Each entry in COG is supported

263 by the presence of the sequence in at least three distinct lineages, thereby

264 representing an ancient conserved domain. A total of 4665 (15.2%) unigenes

265 were assigned with different functional terms (S5 Table). Though most of the

266 unigenes displayed more than one hit, a majority of them were from the same

267 COG functional class. Best hits were further processed for functional

268 classification. The results are shown in Fig 1. Large numbers of unigenes

269 mapped to "Post translational modification, protein turn over, chaperons" (354;

270 7.5%), "Translation, ribosomal structure and biogenesis" (289; 6.1%),

271 "Replication, recombination and repair" (216; 4.6%) and "Signal transduction

272 mechanisms" (195; 4.1%) in addition to "general function prediction" (938;

273    20.1%). In contrast, there were very few genes that mapped to functional classes

274    "cell motility" (7) and "nuclear structure" (2).

275    **Fig 1. COG classification of *F. buski* transcriptome.** The X-axis and Y-axis

276    represent different COG categories and number of unigenes associated with

277    each COG category respectively.

278

279           Since secretory proteins, such as those that are classified as excretory

280    and secretory proteins (ES), are involved in host-parasite interactions, we

281    identified proteins that belong to this category. Our analysis showed that 1501

282    unigenes contained signal peptides and there were 1406 ES proteins. Their

283    identification was done by the presence of a signal peptide domain and the

284    absence of any trans-membrane domain.

285           Sequences that matched significantly with NCBI nr database were further

286    processed using BLAST2GO for GO assignment. This analysis was done to

287    further classify expressed genes into different functional classes. A total 36433

288    GO terms were assigned to 6658 unigenes. The results of this analysis are

289    depicted in Fig 2. Majority of the GO terms were from "biological processes"

290    (2627; 56.2%), while others represented "molecular function" (1445; 30.9%) and

291    "cellular component" (602; 12.8%). The major sub-categories were found to be

292    mostly those associated with cell structure and function and metabolic

293    processes, such as "cell" (GO: 005623), "cell part" (GO: 0044464), "macro-

294    molecular complex" (GO: 0032991) and "organelle" (GO: 0043226), "catalytic

295    function" (GO: 0003824) and "metabolic process" (GO: 0008152). Considering all

296    levels of GO assignment, a significant number of sequences were found to be

297    associated with "Regulation of transcription, DNA-dependent" (GO: 0006355;

298    344), "integral to membrane" (GO: 0016021; 925) and "ATP binding" (GO:

299    0005524; 1086).

300    **Fig 2. Gene ontology classification of *F. buski* transcriptome.** GO terms

301    assigned to Unigenes were classified into three major functional classes; Cellular

302    component, Biological process and Molecular function.

303

304        A Venn diagram summarizing annotation analysis using five databases is

305    given in Fig 3. It was possible for us to annotate 12380 unigenes using both

306    sequence as well as domain-based comparisons. The mean length of annotated

307    unigenes was found to be 1727, which is higher than the mean length (1080) of

308    all the unigenes thus, showing the importance of assembly quality for annotation

309    (S1 Fig). It is expected that unannotated unigenes may participate in unique

310    biological functions specific to this organism.

311    **Fig 3. Venn diagram showing a number of annotated unigenes using**

312    **sequence based and domain based comparisons.** A) Unigenes were

313    compared with NR, Uniprot and Nembase4 for sequence-based annotation.

314    Circle intersections represent number of unigenes found in more than one

315    databases. B) Unigenes were compared against COG and INTERPRO

316    databases for domain-based annotation. C) Two circles A and B represent

317    number of sequences annotated using sequence based and domain based

318    comparisons respectively. The numbers outside circles represent number of

319    sequences remained unannotated after employing all these comparisons.

320

## Comparison with other trematodes

322    A comparative sequence analysis was performed with RNA-seq datasets of

323    different trematode families in order to understand functional divergence. The

324    results are summarized in Table 1 and in Venn diagram (Fig 4). The analysis

325    revealed that 54% of the unigenes displayed sequence similarity with one of the

326    trematodes and 19% with all trematodes analyzed at an E-value threshold of 1E$^{-05}$.

327    When sequence similarity values at different E-value cut-offs were compared

328    transcript sequences from *F. gigantica* appeared to be evolutionarily the closest

329    as it displayed highest number of matches at a given E-value among trematodes

330    studied. Further, to reconstruct evolutionary relationships of *F. buski* with other

331    trematodes, we computed a phylogenetic tree (Fig 5) concatenating all the

332    annotated 12 protein-coding genes (PCGs) from the *F. buski* mitochondrial DNA

333    (mtDNA). The mtDNA for the intestinal fluke was recovered from the *F. buski*

334    genome data and is available in NCBI SRA bearing accession no. SRX316736.

335    Previous studies confirm that alignments with > 10,000 nucleotides from

336    organelle genomes have ample phylogenetic signals for evolutionary

337    reconstruction of tapeworm phylogeny [55]. The phylogenetic tree is well

338    supported by very high posterior probabilities as the taxa grouped well into

339    distinct clades representing the different worm Families such as Opisthorchiidae,

340    Paragonimidae, Paramphistomidae and Fasciolidae (Trematoda); Ascarididae

341 (Nematoda) and Taeniidae (Cestoda).  *F. buski* claded well with *F. hepatica* and

342 *F. gigantica* with strong bootstrap support.

343 **Fig 4. Venn diagram showing number of homologs found in *F. buski***

344 **transcriptome against transcriptome/EST datasets of major trematode**

345 **families of Opisthorchiidae (*Opisthorchis viverrini* and *Clonorchis***

346 ***sinensis*), Fasciolidae (*Fasciola hepatica* and *Fasciola gigantica*) and**

347 **Schistosomatidae (*Schistosoma mansoni* and *Schistosoma japonicum*).**

348

349 **Fig 5. Bayesian phylogenetic relationship among the representative**

350 **helminth species based on 12 PCGs from their mitochondrial DNA.**

351 Phylogenetic analyses of concatenated nucleotide sequence datasets for all 12

352 PCGs were performed using four MCMC chains in bayesian analysis run for

353 1,000,000 generations, sampled every 1,000 generations. Bayesian posterior

354 probability (BPP) values were determined after discarding first 25% of trees as

355 burn-in**.** Posterior support values appear at nodes. Species representing

356 Nematoda (Ascaridida) were taken as outgroup.

357

358 **Table 1. Comparative sequence analysis of *F. buski* transcriptome with**

359 **different trematodes.**

| Organism | 1E-05 | 1E-15 | 1E-30 |
|---|---|---|---|
| *Fasciola gigantica* | 15538 (50.6%) | 12376 (40.3%) | 9880 (32.2%) |
| *Opisthorchis viverrini* | 11635 (37.9%) | 9250 (30.1%) | 7289 (23.7%) |
| *Clonorchis sinensis* | 11188 (36.4%) | 8858 (28.8%) | 6984 (22.7%) |
| *Fasciola hepatica* | 11004 (35.8%) | 7207 (23.4%) | 4695 (15.3%) |
| *Schistosoma mansoni* | 9742 (31.7%) | 7146 (23.2%) | 4984 (16.2%) |
| *Schistosoma* | 8634 (28.1%) | 6433 (20.9%) | 4685 (15.2%) |

| | | | |
|---|---|---|---|
| *japonicum* | | | |

360

361

362        Out of the homologs of *F. buski* unigenes found in trematodes but not in

363 model eukaryotes, 3552 could be annotated. These genes are likely to be

364 trematode specific and can be a useful resource in development of new

365 diagnostic and therapeutic tools. Our results suggest relatively high level of

366 sequence divergence among coding regions (ortholog divergence) of different

367 trematodes, probably a result of adaptation to different ecological conditions.

368

369

370

## 371 **Pathway analysis**

372        We used KEGG automatic annotation server (KAAS) for identification of

373 pathways associated with unigenes. About 3026 sequences were annotated with

374 2527 KO terms. In total, 287 different pathways, classified into six major groups,

375 were identified (Table 2). Some of the genes could not be tracked in the

376 transcriptome. This is presumably due to lack of expression in the adult flukes,

377 but could also represent transcripts that are not polyadenylated and thus were

378 not included in our transcript pools. In either case, a high quality genome

379 sequence data with better depth coverage would be required to confirm the

380 existence or absence of the genes. Major pathways deciphered are those

381 involved in metabolism (638 unigenes, 520 KO terms), such as carbohydrates,

382  amino acids, energy metabolism and genetic information processing (790

383  unigenes, 690 KO terms) transcription, translation, replication and repair, and

384  those involved in human diseases (530 unigenes, 440 KO terms) cancer and

385  neurodegenerative disorders. Overall, most highly represented KEGG terms by

386  virtue of number of unique KO identifiers are Spliceosome (33), RNA transport

387  (92) and protein processing in endoplasmic reticulum (83). Major components of

388  metabolic pathways found in *F. buski* transcriptome are shown in S2 Fig. We

389  could not detect most of the components of fatty acid biosynthesis pathway

390  except acetyl-CoA/propionyl-CoA carboxylase [EC: 6.4.1.2; EC: 6.4.1.3] and 3-

391  oxoacyl-[acyl-carrier-protein] synthase II [EC: 2.3.1.179]. We could not also

392  detect these genes in the current draft genome of *F. buski*. Probably this can be

393  better resolved with an improved genome assembly with higher depth coverage.

394  While EC: 2.3.1.179 was identified in *F. hepatica*, *C. sinensis*, *O. viverrini*, *S.*

395  *mansoni* transcriptomes [27, 28]; it was reported to be missing in *F. gigantica*

396  [29]. It is not surprising as these parasites acquire fatty acids from the host [56].

397  Pathways encoding enzymes of different amino-acid-biosynthesis pathways were

398  represented only by one enzyme for valine, leucine and isoleucine biosynthesis

399  [map00290], and under-represented for pathways, such as lysine biosynthesis

400  [map00300] where no match was detected. In contrast, genes encoding enzymes

401  in fatty acid metabolism [map00071] and amino acid metabolism [map00280;

402  map00310; map00330] are well represented in the *F. buski* transcriptome. The

403  skewed representation of enzymes in the pathway analysis provided clues on the

404     gene regulation and parasite biology pointing towards catabolic process, with a

405     likely dependence in this stage of its life cycle to their host for nutrition.

406     **Table 2: KEGG pathways identified from *F. buski* transcriptome using**

407     **KAAS**.

| Metabolic Pathways | Unique KO terms | Top KEGG pathway terms |
|---|---|---|
| Carbohydrate metabolism | 196 | Glycolysis/Gluconeogenesis[ko00010] |
| Energy metabolism | 111 | Oxidative phosphorylation [ko00190] |
| Lipid metabolism | 105 | Glycerophospholipid metabolism [ko00564] |
| Nucleotide metabolism | 111 | Purine metabolism [ko00230] |
| Amino acid metabolism | 131 | Valine, leucine, isoleucine degradation [ko00280] |
| Metabolism of other amino acids | 37 | Glutathione metabolism [ko00480] |
| Glycan biosynthesis and metabolism | 91 | N-Glycan biosynthesis [ko00510] |
| Metabolism of cofactor and vitamins | 69 | Porphyrin and chlorophyll metabolism [ko00860] |
| Metabolism of terpenoids and polyketides | 18 | Terpenoid backbone biosynthesis [ko00900] |
| Biosynthesis of other secondary metabolites | 14 | Isoquinoline alkaloid biosynthesis [ko00950] |
| Xenobiotics biodegradation and metabolism | 35 | Drug metabolism – other enzymes [ko00983] |
| **Genetic information processing** | | |

| | | |
|---|---|---|
| Transcription | 142 | Spliceosome [ko03040] |
| Translation | 293 | RNA transport [ko3013] |
| Folding, sorting and degradation | 256 | Protein processing in endoplasmic reticulum [ko04141] |
| Replication and repair | 130 | Nucleotide excision repair [ko03420] |
| **Environmental information processing** | | |
| Membrane transport | 12 | ABC transporters [ko02010] |
| Signal transduction | 310 | PI3K-AKt signaling pathway |
| Signaling molecules and interaction | 21 | Neuroactive ligand-receptor interaction |
| **Cellular processes** | | |
| Transport and catabolism | 154 | Endocytosis [ko04144] |
| Cell motility | 39 | Regulation of actin cytoskeleton [ko04810] |
| Cell growth and death | 184 | Cell cycle [ko04681] |
| Cell communication | 103 | Focal adhesion [ko04510] |
| **Organismal systems** | | |
| Immune system | 167 | Chemokine signaling pathway [ko04062] |
| Endocrine system | 128 | Insulin signaling pathway [ko04910] |
| Circulatory system | 34 | Vascular smooth muscle contraction [ko04270] |
| Digestive system | 86 | Pancreatic secretion [ko04972] |
| Excretory system | 54 | Vasopressin-regulated water reabsorption [ko04962] |
| Nervous system | 209 | Neurotrophin signaling |

| | | pathway [ko04722] |
|---|---|---|
| Sensory system | 19 | Phototransduction - fly [ko04745] |
| Development | 46 | Axon guidance [ko04360] |
| Environmental adaptation | 37 | Circadian entrainment [ko04713] |
| **Human diseases** | | |
| Cancers | 334 | Pathways in cancer [ko05200] |
| Immune diseases | 25 | Rheumatoid arthritis [ko05323] |
| Neurodegenerative diseases | 226 | Huntington's disease [ko05016] |
| Substance dependence | 58 | Alcoholism [ko05034] |
| Cardiovascular diseases | 24 | Viral myocarditis [ko05416] |
| Endocrine and metabolic diseases | 10 | Type –II diabetes mellitus [ko04930] |
| Infectious diseases | 419 | Epstein –Barr virus infection [ko05169] |

408

409

# Highly expressed genes

411    Expression status of transcripts can provide clues on their likely biological

412    relevance [57]. To determine the relative expression levels RPKM values of each

413    unigene was evaluated (S6 Table). RNA-seq can provide sensitive estimates of

414    absolute gene expression variation [57]. We noted expression variation of

415    unigenes upto 4th order of magnitude. However, only a few unigenes (5)

416    displayed RPKM value >10,000 and 113 genes exceeding 1000. Among top 100

417    highly expressed genes, functions linked to translation (ribosomal genes-40S,

418  60S ribosomal proteins, Elongation Factor EF1α, Translation Initiation Factor

419  TIF1), protein stress and folding (heat shock proteins 20 and 90 and 10 kDa

420  chaperonins) constituted the majority. As expected, unigenes encoding

421  cytoskeletal proteins were expressed at relatively high levels. Since proteases

422  play an important role in parasite host-pathogen interaction, it was not surprising

423  to see highly expressed Cathepsin L and protease inhibitors as part of the

424  transcriptome. These observations are similar to that of *S. mansoni, F. gigantica,*

425  *C. sinensis and O. viverrini* [27, 29, 58]. In these trematodes stress response

426  genes, genes associated with ribosomes and translation, actin myosin complex

427  and proteolytic enzymes appear to be highly expressed. The results suggest that

428  the adult parasite is active in protein synthesis and is utilizing nutrients from the

429  host to function at a high metabolic load for reproduction and egg development.

## 430  Kinome

431  Eukaryotic or conventional protein kinases (ePKs) play important

432  regulatory roles in diverse cellular processes, such as metabolism, transcription,

433  cell cycle progression, apoptosis, and neuronal development [59]. These are

434  classified into eight groups, based on sequence similarity of their catalytic

435  domains, the presence of accessory domains, and their modes of regulation [60,

436  61, 62]. In addition to eight ePKs, a ninth group categorized as 'Other' and

437  consisting of a mixed collection of kinases that cannot be classified easily into

438  any of the groups is defined in the EMBL kinase database [41]. There are

439  extensive studies on kinome of the model organism *C. elegans*. Kinases from

440  this free-living nematode are deeply conserved in evolution, and the worm shares

441 family homologs as high as 80% with the human kinome. Out of a total of 438

442 worm kinases nearly half are members of worm-specific or worm-expanded

443 families. Studies point to recent evolution of such homologous genes in *C.*

444 *briggsae* involved in spermatogenesis, chemosensation, Wnt signaling and FGF

445 receptor-like kinases [63]. Our analysis of the *F. buski* transcriptome suggested

446 190 ePKs (250 unigenes) belonging to all the nine groups of protein kinases. The

447 results are shown in Fig 6 and S7 Table. The presence of multiple protein kinase

448 families, actively expressed in the adult stage suggests that *F. buski* encodes an

449 extensive signaling network and majority of these eukaryotic signal transduction

450 pathways are conserved. In *F. gigantica*, 308 sequences of protein kinases

451 belonging to eight ePK classes were found in the transcriptome dataset [29], in

452 contrast the free-living nematode *C. elegans* encoded nearly 400 pathways [64].

453 We ascribe this to the parasite biology or possible stage specific expression that

454 could not be determined from the transcriptome dataset of single stages.

455 **Fig 6. A pie chart displaying a number of significant matches found against**

456 **nine different protein kinase classes (according to EMBL kinase database)**

457 **in *F. buski* transcriptome.** The protein kinase groups are represented by: (i)

458 CMGC- cyclin-dependent, mitogen activated, glycogen synthase and CDK-like

459 serine/threonine kinases; (ii) CAMK - Calcium/Calmodulin-dependent

460 serine/threonine kinases; (iii) TK – Tyrosine kinases; (iv) TKL - Tyrosine kinase-

461 like; (v) AGC - cAMP-dependent, cGMP-dependent and protein kinase C

462 serine/threonine kinases; (vi) STE – serine/threonine protein kinases associated

463 with MAP kinase cascade; (vii) CK1 - Casein kinases and close relatives; (viii)

464 RGC - Receptor guanylate cyclase kinases: represented by a single protein,

465 receptor guanylate cyclase kinase and (ix) other unclassified kinases.

466

## Protease and protease inhibitors

468 *F. buski* transcriptome suggests that a total of 478 proteases (MEROPS

469 terms), corresponding to 6 catalytic types, and 138 protease inhibitors (138

470 MEROPS terms) as defined in the MEROPS database are expressed in this

471 organism [50]. The results are shown in Figure 7. Some of the highly expressed

472 proteases are prolyl oligopeptidase family **of** serine proteases (206 MEROPS

473 terms; 16 families), lysostaphin subfamily **of** metalloproteases (126 MEROPS

474 terms; 21 families) and ubiquitin-specific peptidases **of the** cysteine proteases

475 family (81 MEROPS terms; 22 families) in addition to a significant number of

476 protease inhibitors (138 MEROPS terms; 14 families) (Fig 7 and S8 Table)

477 **Figure 7: Protease and protease inhibitors (MEROPS terms) found in *F.***

478 ***buski* transcriptome.**

479

## RNAi pathway genes

481 RNA interference (RNAi), a gene silencing process generally triggered by

482 double-stranded RNA (dsRNA) delivers gene-specific dsRNA to a competent

483 cell, thereby engaging the RNAi pathway leading to the suppression of target

484 gene expression. RNAi pathway has played a crucial role in our current

485 understanding of genotype to phenotype relations in many organisms, including

486 *C. elegans*. Presence of this pathway in *F. buski* may be useful in deciphering

487    functions of genes as these are not amenable to classical genetic approaches.

488    We tried to decipher from the transcriptome data if this pathway exists in *F. buski*

489    and the divergence from RNAi pathway of other related parasites.  The results

490    are summarized in Table 3. Orthologs were found for all the genes associated

491    with small RNA biosynthesis of *C. elegans* in *F. buski* except rde-4. This gene is

492    also absent in the genome of most of the parasitic nematodes except *Brugia*

493    *malayi* and *Ancylostoma caninum* [65]. The major differences from pathways

494    present in *Schistosoma mansoni* and *C. elegans*, are lack of genes related to

495    dsRNA uptake and spreading (sid-1,sid-2,rsd-6). Absence of these genes is

496    reported in most of the nematodes outside the genus *Caenorhabditis* [65]. In

497    contrast, phylogenetic neighbors of *Schistosoma* sp. possess ortholog of sid-1

498    which is responsible for dsRNA entry in soaked parasites, though this protein is

499    much larger than sid-1 in *C. elegans*. Another conserved protein rsd-3 is found in

500    almost all the nematodes studied, though its functional collaborators are missing

501    in most of the organisms. [65]. Other than small RNA biosynthesis, proteins

502    associated with nuclear effectors were mostly identified except the absence of

503    less conserved mes-3, rde-2, ekl-5 etc [65]. Overall, orthologs for most

504    conserved proteins from each functional category [65] were found, whereas, the

505    less conserved ones remained unknown.

506    **Table 3. Proteins of RNAi pathway identified from *F. buski transcriptome*.**

|  | *Caenorhabditis elegans* | *Fasciolopsis buski* |
|---|---|---|
| Small RNA biosynthesis | dcr-1, drh-1/3, drsh-1, pash-1, rde-4, xpo-1/2/3 | dcr-1, drh-1/3, drsh-1, pash-1, xpo-1/2/3 |
| dsRNA uptake and | sid-1/2, rsd-3/6 | rsd-3 |

| spreading | | |
|---|---|---|
| siRNA amplification | rsd-2, ego-1, rrf-1/3, smg-2/5/6 | smg-2/6 |
| Argonautes | alg-1/2/3/4, csr-1, ergo-1, nrde-3, ppw-1/2, prg-1/2, rde-1, sago-1/2 | alg-1/2 |
| RISC proteins | ain-1/2, tsn-1, vig-1 | tsn-1 |
| RNAi inhibitors | adr-1/2, eri-1/3/5/6/7, lin-15b, xrn-1/2 | eri-1/7, xrn-1/2 |
| Nuclear RNAi effectors | cid-1, ekl-1/4/5/6, gfl-1, mes-2/3/6, mut-2/7/16, rde-2, rha-1,zfp-1 | cid-1, ekl-4/6, gfl-1, mes-2/6, rha-1, zfp-1 |

507

508

## Transposable elements (TEs)

510    Repetitive elements are important structural features of a genome and

511 transcriptome as many of these are transcribed. Retrotransposable elements

512 (RTEs) constitute the main type of interspersed repeats particularly in higher

513 eukaryotic genomes [66]. Therefore, we analyzed the *F. buski* transcriptome for

514 identification of expressed repetitive elements, particularly interspersed

515 elements.  A total of 3720 retroelements that include 3477 LINES in the *F. buski*

516 transcriptome are outlined in Table 4.  All LINE elements are not expressed;

517 genomic LINE elements were also detected (S9 Table). In contrast to LINES, the

518 data was under-represented for SINE elements. T2#SINE/tRNA were identified

519 from the transcriptome sequences. These results are similar to reports from other

520 trematodes, such as *C. sinensis* [27, 29, 67] where SINE elements were also not

521 reported.  Additionally, several LTR retrotransposable elements and DNA

522 transposons (Table 4), whose roles in parasite biology and genome organization

523 are not yet well defined were identified.

524 **Table 4. Number of different repeat elements identified from genome and**

525 **transcriptome of *F. buski.***

526

| | Number of elements in Genome | Number of elements in Transcriptome |
|---|---|---|
| Retroelements | 3556 | 3720 |
| SINE | 0 | 2 |
| Penelope | 50 | 129 |
| LINEs: | 3155 | 3477 |
| L2/CR1/Rex | 94 | 96 |
| R2/R4/NeSL | 6 | 11 |
| RTE/Bov-B | 2808 | 2975 |
| LTR elements | 401 | 241 |
| BEL/Pao | 14 | 6 |
| Gypsy/DIRS1 | 387 | 235 |
| DNA transposons | 0 | 0 |
| Simple repeats | 51778 | 19541 |
| Low complexity | 425 | 1443 |

527

528

# Discussion

530     Human parasites are a diverse group of organisms, ranging from

531 unicellular protists to multicellular trematodes. Most well studied molecular data

532 on human host-parasite interactions are from parasites, such as *Plasmodium* and

533 *Leishmania*. Lack of information and complex ecological relationship of

534    multicellular flukes limits our understanding of trematode biology. Recently,

535    availability of new and cheap genome sequencing technologies (NGS) has

536    opened up research avenues for characterization of these organisms. This is

537    reflected in a number of studies that have been published recently on these

538    parasites based on genomic information.

539    North-east India is considered to be a hot bed of diversity and it is

540    believed that many organisms from this region have evolved interesting biology

541    due to their evolution in a unique ecological niche. In this report we present data

542    about characterization of a trematode parasite *F. buski* (the giant intestinal fluke

543    having a zoonotic potential) using genome sequencing and RNA-seq analysis.

544    Transcriptome sequencing helped us identify 30677 genes from the parasite, and

545    we annotated 12380 genes, driven by longer assembly, that enabled more

546    authentic functional annotation of the unigenes.

547    Majority of the highly expressed genes in *F. buski* have also been reported

548    in other organisms including trematodes. We note a high level expression of

549    transcripts encoding cytoskeletal elements, protein biosynthesis and folding in

550    the adult stage of the parasite. Parasites infect host and are reliant on the host

551    for their nutrient supply during their growth and development cycle [56]. The

552    parasite is studied at an adult infective stage, where development has switched

553    to the reproductive stage with a high metabolic load required for regular egg

554    production. Reproductive development and egg production require high levels of

555    transcript accumulation and associated protein synthesis, often in conjunction

556    with reorganization of cytoskeletal elements. High-energy costs of eggshell

557  production in *F. buski* require efficient energy generation and high expression of

558  these enzymes may explain this. Cytoskeleton proteins are involved in

559  intracellular transport and glucose uptake in larvae and adult parasites and serve

560  as outlets for their glycogen stores [68].  We noted the high expression of genes

561  encoding fatty acid binding proteins (FABP) that support the role of these

562  molecules in acquisition, storage, and transport of lipids. Since many of these

563  display unique properties, FABPs have been suggested as potential vaccine

564  candidates [69]. We also speculate the difference in the fatty acid metabolic

565  processes, with higher enrichment of gene targets for catabolism, compared to

566  biosynthesis, which can be critical to the parasite biology at the adult stage.

567  Empirical comparison with other data sets from trematodes suggest, absence of

568  biosynthesis in the adult stage, either by a gene-regulatory mechanism or lack of

569  components in the fatty acid biosynthesis. This would warrant further

570  investigation in the light that these parasites have been known to acquire fatty

571  acids from the host [56].

572      Different proteases are expressed in trematodes such as *F. hepatica, C.*

573  *sinensis* and *S. mansoni.* Proteases play a significant role in parasite physiology

574  as well as in host–parasite interactions [70-72]. In general genes encoding six

575  different 'catalytic type' of proteases, were identified in the *F. buski*

576  transcriptome, and were expressed at relatively high-levels. Amongst them, the

577  serine and cysteine proteases and serine protease inhibitors were highly

578  represented in unigenes expressed in the adult stage.  In the blood fluke, *S.*

579  *mansoni*, serine proteases were implicates in tissue invasion and host immune

580 evasion [73, 74]. Cysteine proteases have roles in nutrition, tissue/cell invasion,

581 excystment/encystment, exsheathment and hatching, protein processing and

582 immune-evasion [75]. Papain-like cysteine proteases, particularly the cathepsins,

583 facilitate skin and intestine infections, tissue migration, feeding and suppression

584 of host immune effector cell functions [76]. Serine protease inhibitors (serpins)

585 are a super family of structurally conserved proteins that inhibit serine proteases

586 and are involved in many important endogenous regulatory processes and

587 possible regulation of host immune modulation and/or evasion processes [77].

588 We report the expression of globin genes (Hemoglobin F2 and Myoglobin

589 I), which was also found to be abundant in *F. buski* transcriptome. The functional

590 role of trematode haemoglobin (Hbs) is still not clear, as adult parasitic

591 trematodes and also nematodes, such as *Ascaris suum*, are residents of semi-

592 anaerobic environment. It is believed that the function of hemoglobin is not only

593 restricted to $O_2$ transport but also as an oxygen scavenger, a heme reserve for

594 egg production, and as co-factor of NO deoxygenase [78-82].

595 The kinome of *F. buski* is similar to other eukaryotes, although the gene

596 family redundancy within protein kinases identified was fewer than reported from

597 free-living nematode. The functional repertoire suggest the presence of an

598 extensive signaling system, that may be required for sensing and modulating

599 signal response in the changing life-cycle and adult development/reproduction in

600 the intra-host environment. CMGC kinases were relatively the most abundant in

601 terms of number of sequences, followed by CAM kinases. This data correlates

602 well with the observation of liver parasite *F. gigantica* [29] and blood fluke,

603    *Schistosoma mansoni* [83]. CMGC kinases are known to regulate cell

604    proliferation and ensure correct replication and segregation of organelles in many

605    eukaryotic organisms including *Plasmodium falciparum* [84]. CAM kinases are

606    associated with calcium-mediated signaling. Among CAMKs in *F. buski*, the

607    majority of kinases belonged to the CAMK-like kinases (CAMK family 1 and 2),

608    death-associated protein kinases (DAPKs) and the myosin light chain kinases

609    (MLCKs). The primary function of MLCK is to stimulate muscle contraction

610    through phosphorylation of myosin II regulatory light chain (RLC), a eukaryotic

611    motor protein that interacts with filamentous actin [83]. The high abundance of

612    the above-mentioned kinase groups may reflect high mobility and muscle activity

613    in *F. buski* associated with feeding and egg sheading.

614      Expression of genes involved in RNAi pathway suggests partial evidence

615    of an active pathway. We noted absence of the protein required for RNA uptake,

616    and speculate that RNAi would be active if introduced into the cytoplasm

617    considering siRNA as a possible therapy.

618      *F. buski* genome also encodes different types of transposons as

619    interspersed repetitive sequences similar to other eukaryotes. Members of the

620    Phylum Platyhelminthes are also thought to contain diverse TEs, which comprise

621    up to 40% of their genomes. A total of 29 retrotransposons, belonging to one

622    non-long terminal-repeat (LTR) family (6 elements in CR1) and 3 LTR families (5

623    elements in Xena and Bel, and 13 elements in Gypsy) have been isolated from

624    the genomes of the digenean trematodes, *C. sinensis* and *Paragonimus*

625    *westermani*. CsRn1 of *C. sinensis* and PwRn1 of *P. westermani* are novel

626 retrotransposons, which are evenly distributed throughout their genomes and

627 expressed as full transcripts in high copy numbers. Phylogenetic studies have

628 revealed that the CsRn1 and PwRn1 elements formed a novel, tightly-conserved

629 clade, that has evolved uniquely in the metazoan genomes. Diverse

630 retrotransposon families are present in the lower animal taxa, and that some of

631 these elements comprise important intermediate forms marking the course of

632 evolution of the LTR retrotransposons. Retrotransposons in trematodes might

633 influence the remodeling of their host genomes [60].

634 Two novel families of tRNA-related SINEs have been described to be

635 widespread among all *Salmonoidei* genomes, with a role in human helminth

636 pathogen—*Schistosoma japonicum* (Trematoda: Strigeiformes) [85]. We could

637 identify only two SINE elements from the *F. buski* transcriptome data. Lack of

638 widespread presence of SINE elements in flukes raises interesting questions

639 about evolution of SINEs [85]. SINEs have been reported in early branching

640 protists, such as *Entamoeba histolytica* [86] The inability to detect SINE elements

641 in other trematodes and presence of two families in *F. buski* suggest that these

642 may have come by horizontal transfer possibly with some evolutionary

643 advantage. It is possible that SINEs may have been lost from the genome during

644 evolution in these organisms. We also found that not all LINE elements are

645 expressed in the adult stage. Overall the transcriptome of *F. buski* has shed new

646 light on the biology of this organism.

647 In conclusion, the rough draft genome and transcriptome characterized in

648 the present study will assist in future efforts to decode the entire genome of *F.*

649 *buski*. The transcriptome data of the adult stage of this giant intestinal fluke is

650 reported for the first time, and is archived in the NCBI SRA as well as in the

651 database (North-East India Helminth Parasite Information Database) developed

652 by our group (NEIHPID) [87]. We hope our study adds substantially to the public

653 information platform to achieve a fundamental understanding of the parasite

654 biology, which in turn would help in identification of potential drug targets and

655 host-pathogen interaction studies.

656

657 # Author Contributions

658 DKB and VT conceived and designed the experiments, analyzed the data,

659 contributed reagents/materials/analysis tools, wrote the paper. DKB, TR and PP

660 performed the experiments. DKB, TR and PP performed the bioinformatics

661 analysis. All authors reviewed the manuscript.

662

663 # Additional Information

664 The authors declare no competing financial interests.

665 # DNA Deposition

666 DNA sequences were deposited as follows:

667 National Centre for Biotechnology Information (NCBI) Bioproject Accession:

668 PRJNA212796 ID: 212796

669 NCBI Sequence Read Archive (SRA): **SRP028107,** SRX327895, SRX326786,

670 SRX316736

671

## **Acknowledgement**

678

# References

1. CDC (2009) "DPDx - Fasciolopsiasis". Laboratory Identification of Parasites of Public Health Concern. Available from: http://www.dpd.cdc.gov/dpdx/HTML/Fasciolopsiasis.htm.

2. Le T, Nguyen V, Phan B, Blair D, McManus D Case report: unusual presentation of *Fasciolopsis buski* in a Vietnamese child. Trans R Soc Trop Med Hyg. 2004; 98: 193-194.

3. Roy B, Tandon V. Morphological and microtopographical strain variations among *Fasciolopsis buski* originating from different geographical areas. Acta Parasitol. 1993; 38: 72-77.

4. Roy B, Tandon V *Fasciolopsis buski*. In: Miliotis MD, Bier JW, editors. International handbook of foodborne pathogens. New York, Basel: Marcel Decker, Inc; 2003. pp. 563-570.

5. Mas-Coma S, Bargues M, Valero M (2005) Fascioliasis and other plant-borne trematode zoonoses. Int J Parasitol 35: 1255-1278.

6. Chai JY, Shin EH, Lee SH, Rim HJ Foodborne Intestinal Flukes in Southeast Asia. Korean J Parasitol. 2009; 47(Suppl), S69.

7. Raymondo, D "Parasitology Training Manual - *Fasciolopsis buski*". Available from: http://www.practicalscience.com/fb.html.

8. Nakagawa K. The development of *Fasciolopsis buski* Lankester. J Parasitol. 1992; 8: 161-166.

9. Bunnag D, Radomyos P, Harinasuta T. Field trial on the treatment of fasciolopsiasis with praziquantel. Southeast Asian J Trop Med Public Health.

702  1983; 14: 216-219.

703 10. Harinasuta T, Bunnag D, Radomyos P. Efficacy of praziquantel on

704   fasciolopsiasis. Arzneimittelforschung 1984; 34: 1214-1215.

705 11. Handoyo I, Ismuljowono B, Darwis F, Rudiansyah. A survey of fasciolopsiasis

706   in Sei Papuyu village of Babirik subdistrict, Hulu Sungei Utara Regency,

707   South Kalimantan Province. Trop Biomed. 1986; 3: 113-118.

708 12. Lee HH. Studies on epidemiology and treatment of fasciolopsiasis in southern

709   Taiwan. Kaohsiung J Med Sci. 1986; 2: 21-27.

710 13. Taraschewski H, Mehlhorn H, Bunnag D, Andrews P, Thomas H. Effects of

711   praziquantel on human intestinal flukes (*Fasciolopsis buski* and *Heterophyes*

712   *heterophyes*). Zentralbl Bakteriol Mikrobiol Hyg. 1986; 262: 542-550.

713 14. Gupta, A. et al. *Fasciolopsis buski* (giant intestinal fluke) - a case report.

714   Indian J Pathol Microbiol. 1999; 42:59-60.

715 15. Datta S, Mukerjee GS, Ghosh JD. Comparative efficacy of triclabendazole,

716   oxyclozanide and rafoxanide against *Fasciolopsis buski* in naturally infected

717   pigs. Indian J Anim Health. 2004; 43: 53-56.

718 16. World Health Organization. Foodborne trematode infections. Bull WHO 1995;

719   73: 397-399.

720 17. World Health Organization Control of foodborne trematode infections. WHO

721   Tech Rep Ser. 1995; 849: 1-157.

722 18. Bhatti HS, Malla N, Mahajan RC, Sehgal R. Fasciolopsiasis - a re-emerging

723   infection in Azamgarh (Uttar Pradesh). Indian J Pathol Microbiol. 2000; 43:

724   73-76.

725   19. Muralidhar S, Srivastava L, Aggarwal P, Jain N, Sharma DK. Fasciolopsiasis -

726        a persisting problem in eastern U.P. - a case report. Indian J Pathol Microbiol.

727        2000; 43: 69-71.

728   20. Prasad PK, Tandon V, Chatterjee A, Bandyopadhyay S. PCR-based

729        determination of internal transcribed spacer (ITS) regions of ribosomal DNA

730        of giant intestinal fluke, *Fasciolopsis buski* (Lankester, 1857) Looss, 1899.

731        Parasitol Res. 2007; 101: 1581-1587.

732   21. Tandon V, Roy B, Prasad PK. Chapter 32. Fasciolopsis. In: Liu D, editor.

733        Molecular Detection of Human Parasitic Pathogens. Boca Raton, Florida:

734        CRC Press. 2012. pp. 353-364.

735   22. Bentley, D.R. et al. Accurate whole human genome sequencing using

736        reversible terminator chemistry. Nature. 2008; 456: 53-59.

737   23. Harris, T.D. et al. Single-molecule DNA sequencing of a viral genome. 2008;

738        Science 320: 106-109.

739   24. Margulies, M. et al. Genome sequencing in microfabricated high-density

740        picolitre reactors. Nature. 2005; 437: 376-380.

741   25. Pandey V, Nutter RC, Prediger E. Applied Biosystems SOLiD™ System:

742        Ligation-based sequencing. In: Michal Janitz, editor. Next-generation genome

743        sequencing: Towards personalized medicine. Wiley. 2008. pp. 29-41.

744   26. Wheat CW, Vogel H. Transcriptome sequencing goals, assembly, and

745        assessment. Methods Mol Biol. 2011; 772: 129-144.

746   27. Young, N.D. et al. Unlocking the transcriptomes of two carcinogenic

747        parasites, *Clonorchis sinensis* and *Opisthorchis viverrini*. PLoS Negl Trop Dis.

748    2010; 4:e719.

749    28. Young ND, Hall RS, Jex AJ, Cantacessi C, Gasser RB. Elucidating the

750    transcriptome of *Fasciola hepatica* - a key to fundamental and

751    biotechnological discoveries for a neglected parasite. Biotechnol Adv. 2010;

752    28: 222-231.

753    29. Young, N.D. et al. A portrait of the transcriptome of the neglected trematode,

754    *Fasciola gigantica* - biological and biotechnological implications. PLoS Negl

755    Trop Dis. 2011; 5:e1004.

756    30. Almeida, G.T. et al. Exploring the *Schistosoma mansoni* adult male

757    transcriptome using RNA-seq. 2012; Exp Parasitol 132: 22-31.

758    31. Cantacessi, C. et al. A deep exploration of the transcriptome and

759    "excretory/secretory" proteome of adult *Fascioloides magna*. Mol Cell

760    Proteomics. 2012; 11: 1340-1353.

761    32. Garg, G. et al. The transcriptome of *Echinostoma caproni* adults: Further

762    characterization of the secretome and identification of new potential drug

763    targets.          J          Proteomics.          2013;          Available          from:

764    http://dx.doi.org/10.1016/j.jprot.2013.06.017

765    33. Biswal DK, Ghatani S, Shylla JA, Sahu R, Mullapudi N, Bhattacharya A,

766    Tandon V. An integrated pipeline for next generation sequencing and

767    annotation of the complete mitochondrial genome of the giant intestinal fluke,

768    Fasciolopsis buski (Lankester, 1857) Looss, 1899. PeerJ. 2013; 1:e207

769    34. Sambrook J, Fritsch EF, Maniatis T (1989) Molecular cloning: a laboratory

770    manual, 2nd edn. Cold Spring Harbor: Cold Spring Harbor Laboratory Press.

771    35. Grabherr MG, Haas BJ, Yassour M, Levin JZ, Thompson DA, et al. (2011)

772        Full-length transcriptome assembly from RNA-seq data without a reference

773        genome. Nat Biotechnol 29: 644-652.

774    36. Li W, Godzik A (2006) Cd-hit: a fast program for clustering and comparing

775        large sets of protein or nucleotide sequences. Bioinformatics 22: 1658-1659.

776    37. Langmead B, Trapnell C, Pop M, Salzberg SL (2009) Ultrafast and memory-

777        efficient alignment of short DNA sequences to the human genome. Genome

778        Biol 10: R25.

779    38. UniProt Consortium. UniProt: a hub for protein information. Nucleic Acids

780        Res. 2015; 43(Database issue): D204-12.

781    39. Elsworth B, Wasmuth J, Blaxter M NEMBASE4: the nematode transcriptome

782        resource. Int J Parasitol. 2011; 41: 881-894.

783    40. McGinnis S, Madden TL. BLAST: at the core of a powerful and diverse set of

784        sequence analysis tools. Nucleic Acids Res. 2004; 1:32 (Web Server issue):

785        W20-5.

786    41. Hunter, S. et al. InterPro: the integrative protein signature database. Nucleic

787        Acids Res. 2009; 37: D211-215.

788    42. Tatusov, R.L. et al. The COG database: an updated version includes

789        eukaryotes. BMC Bioinformatics. 2003; 4: 41.

790    43. Conesa, A. et al. Blast2GO: a universal tool for annotation, visualization and

791        analysis in functional genomics research. Bioinformatics. 2005; 21: 3674-

792        3676.

793    44. Ye, J. et al. WEGO: a web tool for plotting GO annotations. Nucleic Acids

794       Res. 2006; 34: W293-297.

795   45. Moriya Y, Itoh M, Okuda S, Yoshizawa A, Kanehisa M KAAS: an automatic

796       genome annotation and pathway reconstruction server. Nucleic Acids Res.

797       2007; 35: W182-185.

798   46. Letunic I, Yamada T, Kanehisa M, Bork P ipath: interactive exploration of

799       biochemical pathways and networks. Trends Biochem Sci. 2008; 33: 101-103.

800   47. Min XJ, Butler G, Storms R, Tsang A. OrfPredictor: predicting protein-coding

801       regions in EST-derived sequences. Nucleic Acids Res. 2005; 33. W677-680

802   48. Bendtsen JD, Nielsen H, von Heijne G, Brunak S. Improved prediction of

803       signal peptides: SignalP 3.0. J Mol Biol. 2004; 340: 783-795.

804   49. Krogh A, Larsson B, von Heijne G, Sonnhammer El. Predicting

805       transmembrane protein topology with a hidden Markov model: application to

806       complete genomes. J Mol Biol. 2001; 305: 567-580.

807   50. Rawlings ND, Barrett AJ, Bateman A. MEROPS: the database of proteolytic

808       enzymes, their substrates and inhibitors. Nucl Acids Res. 2012; 40: D343-

809       D350.

810   51. Zerbino DR, Birney E. Velvet: Algorithms for de novo short read assembly

811       using de Bruijn graphs. Genome Res. 2008; 18: 821–829.

812   52. Simpson, J.T. et al. ABySS: A parallel assembler for short read sequence

813       data. Genome Res. 2009; 19: 1117-1123.

814   53. Li, R. et al. De novo assembly of human genomes with massively parallel

815       short read sequencing. Genome Res. 2010; 20: 265-272.

816   54. Smit AFA, Hubley R & Green P. RepeatMasker Open-3.0. 1996-2010

817    <http://www.repeatmasker.org>

818    55. Waeschenbach A, Webster BL, Littlewood DT. Adding resolution to ordinal

819        level relationships of tapeworms (Platyhelminthes: Cestoda) with large

820        fragments of mtDNA. Mol Phylogenet Evol. 2012; 63:834-847

821    56. Barrett J. Biochemistry of Parasitic Helminths, University Park Press.

822        Baltimore. 1981; 308 p

823    57. Conesa A, Madrigal P, Tarazona S, Gomez-Cabrero D, Cervera A,

824        McPherson A, Szcześniak MW, Gaffney DJ, Elo LL, Zhang X, Mortazavi A. A

825        survey of best practices for RNA-seq data analysis. Genome Biol. 2016;

826        17:13. doi: 10.1186/s13059-016-0881-8. Review. Erratum in: Genome Biol.

827        2016; 17(1): 181.

828    58. Protasio, A.V. et al. A systematically improved high quality genome and

829        transcriptome of the human blood fluke Schistosoma mansoni. PLoS Negl

830        Trop Dis. 2012; 6: e1455.

831    59. Kannan N, Neuwald AF. Evolutionary constraints associated with functional

832        specificity of the CMGC protein kinases MAPK, CDK, GSK, SRPK, DYRK,

833        and CK2alpha. Protein Sci. 2004;13:2059-77

834    60. Hanks SK, Hunter T. Protein kinases 6. The eukaryotic protein kinase

835        superfamily: kinase (catalytic) domain structure and classification. FASEB J.

836        1995; 9: 576-596.

837    61. Manning G, Whyte DB, Martinez R, Hunter T, Sudarsanam S. The protein

838        kinase complement of the human genome. Science. 2002; 298: 1912-1934.

839    62. Miranda-Saavedra D, Barton GJ. Classification and functional annotation of

840        eukaryotic protein kinases. Proteins. 2007; 68: 893-914.

841    63. Eisenmann DM. Wnt signaling WormBook, ed. The C. elegans Research

842        Community, WormBook. 2005; doi/10.1895/wormbook.1.7.1. Available from:

843        http://www.wormbook.org.

844    64. Plowman G, Sudarsanam S, Bingham J, Whyte D, Hunter T. The Protein

845        kinases of *Caenorhabditis elegans:* A model for signal transduction in

846        multicellular organisms. Proc Natl Acad Sci USA. 1999; 96: 13603-13610

847    65. Dalzell, J.J. et al. RNAi effector diversity in nematodes. PLoS Negl Trop Dis.

848        2011; 5: e1176

849    66. Finnegan DJ. Retrotransposons. Curr Biol. 2012; 22(11):R432-7.

850    67. Brindley PJ, Laha T, McManus DP, Loukas A. Mobile genetic elements

851        colonizing the genomes of metazoan parasites. Trends Parasitol. 2003; 19:

852        79-87.

853    68. Robertson, A.P. et al. Antinematodal drugs - modes of action and resistance:

854        and worms will not come to Thee. In: Caffrey ER, editor. Parasitic Helminths,

855        Targets, Screens, Drugs and Vaccines. Weinheim, Germany: Wiley-

856        Blackwell. 2012; pp. 233-249.

857    69. Nie, H.M. et al. Cloning and characterization of the fatty acid-binding protein

858        gene from the protoscolex of *Taenia multiceps*. Parasitol Res. 2013; 112:

859        1833-1839.

860    70. Chen, M. et al. The anti-helminthic niclosamide inhibits Wnt/ Frizzled1

861        signaling. Biochemistry. 2009; 48: 10267-10274.

862    71. Horn, M. et al. Mapping the pro-peptide of the *Schistosoma mansoni*

863       cathepsin B1 drug target: modulation of inhibition by heparin and design of

864       mimetic inhibitors. ACS Chem Biol. 2011; 6: 609-617.

865   72. Jilkova, A. et al. Structural basis for inhibition of cathepsin B drug target from

866       the human blood fluke, *Schistosoma mansoni.* J Biol Chem. 2011; 286:

867       35770-35781.

868   73. Newport, G.R. et al. Cloning of the proteinase that facilitates infection by

869       schistosome parasites. J Biol Chem. 1998; 263: 13179-13184.

870   74. Fishelson Z. Novel mechanisms of immune evasion by *Schistosoma*

871       *mansoni*. Mem Inst Oswaldo Cruz. 1995; 90: 289-292.

872   75. Sajid M, McKerrow JH. Cysteine proteases of parasitic organisms. Mol

873       Biochem Parasitol. 2002; 120: 1-21.

874   76. Dalton, J.P. et al. Proteases in trematode biology. In: Maule AG, Marks NJ,

875       editors. Parasitic flatworms: Molecular biology, Biochemistry, Immunology

876       and Physiology. Oxford: CAB International. 2006. pp. 348-368.

877   77. Molehin AJ, Gobert GN, McManus DP. Serine protease inhibitors of parasitic

878       helminths. Parasitology. 2012; 139: 681-695.

879   78. Goldberg DE. The enigmatic oxygen-avid hemoglobin of *Ascaris.* Bioessays.

880       1995; 17: 177-182.

881   79. Rashid, A.K. et al. Trematode myoglobins, functional molecules with a distal

882       tyrosine. J Biol Chem 1997; 272: 2992-2999.

883   80. Minning, D.M. et al. *Ascaris* haemoglobin is a nitric oxide-activated

884       "deoxygenase." Nature; 1999. 401: 497-502.

885   81. Rashid AK, Weber RE Functional differentiation in trematode hemoglobin

886      isoforms. Eur J Biochem 1999. 260: 717-725.

887   82. de Guzman, J.V. et al. Molecular characterization of two myoglobins of

888      *Paragonimus westermani*. J Parasitol. 2007; 93: 97-103.

889   83. Andrade, L.F. et al. Eukaryotic protein kinases (ePKs) of the helminth parasite

890      *Schistosoma mansoni*. BMC Genomics. 2011; 12: 215.

891   84. Ward P, Equinet L, Packer J, Doerig C. Protein kinases of the human malaria

892      parasite *Plasmodium falciparum*: the kinome of a divergent eukaryote. BMC

893      Genomics. 2004; 5: 79.

894   85. Melamed P, Chong KL, Johansen MV. Evidence for lateral gene transfer from

895      salmonids to two schistosome species. Nat Genet. 2004; 36: 786-787.

896   86. Matveev V, Nishihara H, Okada N. Novel SINE families from salmons validate

897      *Parahucho* (Salmonidae) as a distinct genus and give evidence that SINEs

898      can incorporate LINE-related 3'-tails of other SINEs. Mol Biol Evol. 2007; 24:

899      1656-66.

900   87. Biswal DK, Debnath M, Kharumnuid G, Thongnibah W, Tandon V. Northeast

901      India Helminth Parasite Information Database (NEIHPID): Knowledge Base

902      for Helminth Parasites. PLoS One. 2016; 11(6):e0157459

903

904 # Supporting information

905 **S1 Fig. Histogram displaying length of the unigenes.** Black bars represent all

906 unigenes while red ones represent annotated unigenes only. Clearly, mean of

907 length of annotated unigenes is higher than overall mean length of all unigenes.

908 **S2 Fig. Major components of metabolic pathways present in F. buski**

909 **transcriptome.** Colored edges represent proteins homologous to any F.buski

910 unigene.

911 **S1 Table. Assembly statistics for *F. buski* genome using three denovo**

912 **assemblers and transcriptome using Trinity.**

913 **S2 Table. Assembly assessment of uigenes used for annotation**

914 **S3 Table.** BLASTx hits against NR database

915 **S4 Table. Top 30 conserved domain/families identified by InterProScan**

916 **from F. buski transcriptome**

917 **S5 Table. COG categories assigned to unigenes**

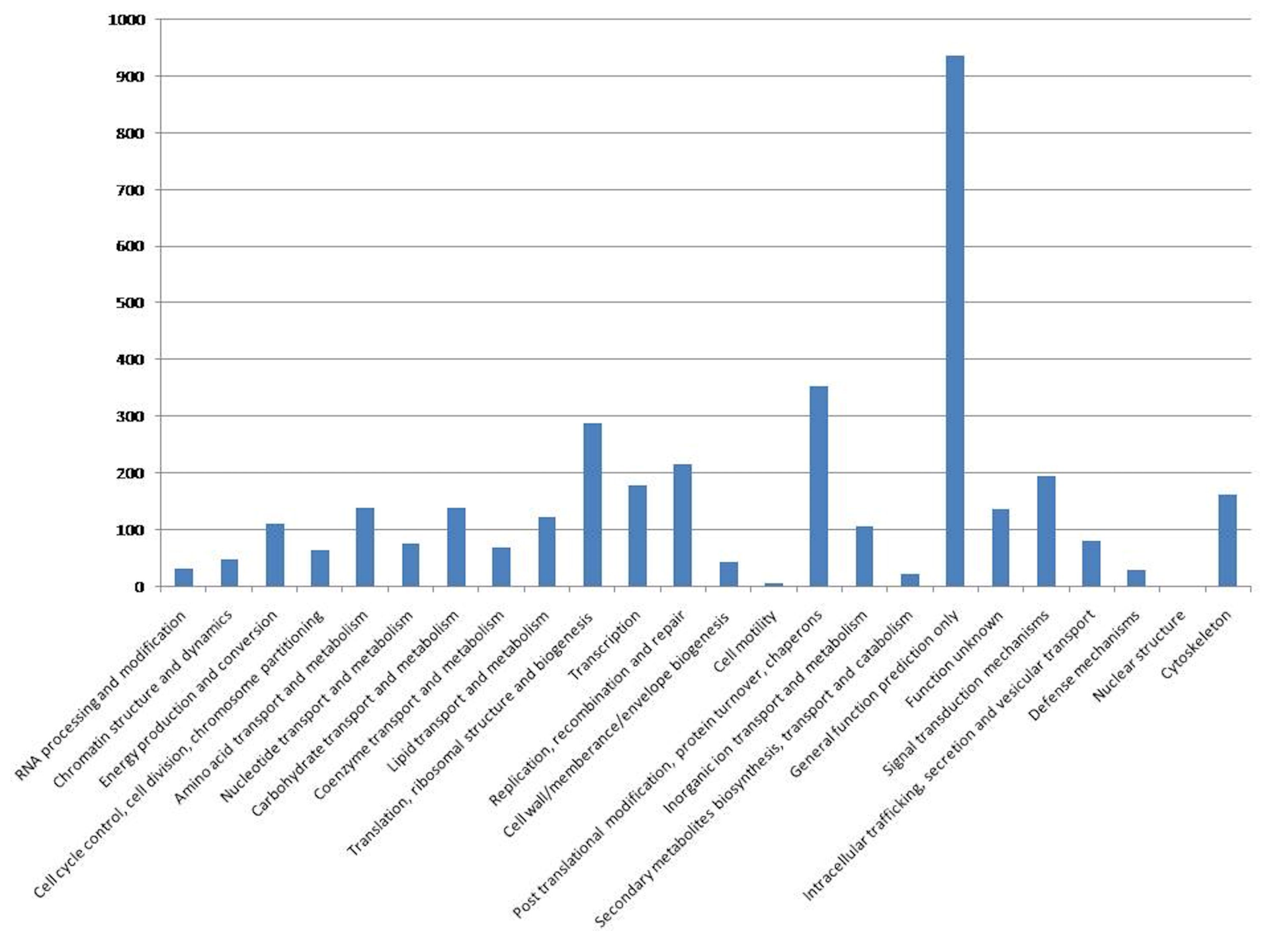918 **S6 Table. RPKM values of each of the unigenes**

919 **S7 Table. List of different kinases identified from *F. buski* transcriptome**

920 **S8 Table. List of different protease and protease inhibitors identified from**

921 ***F. buski* transcriptome**

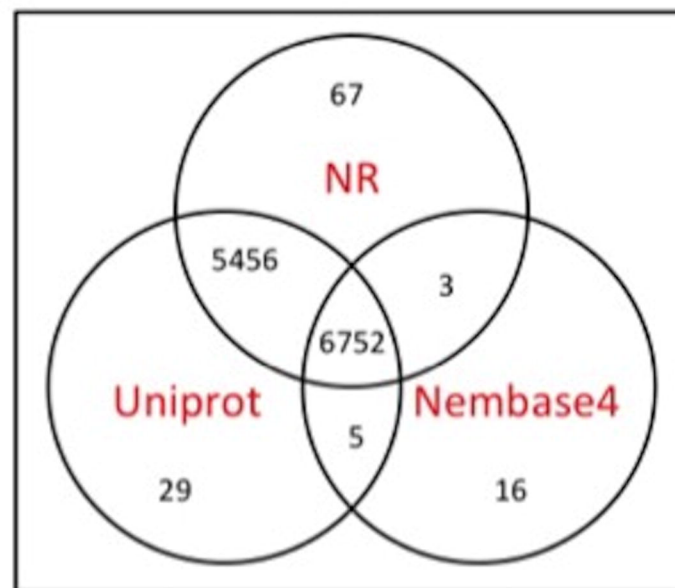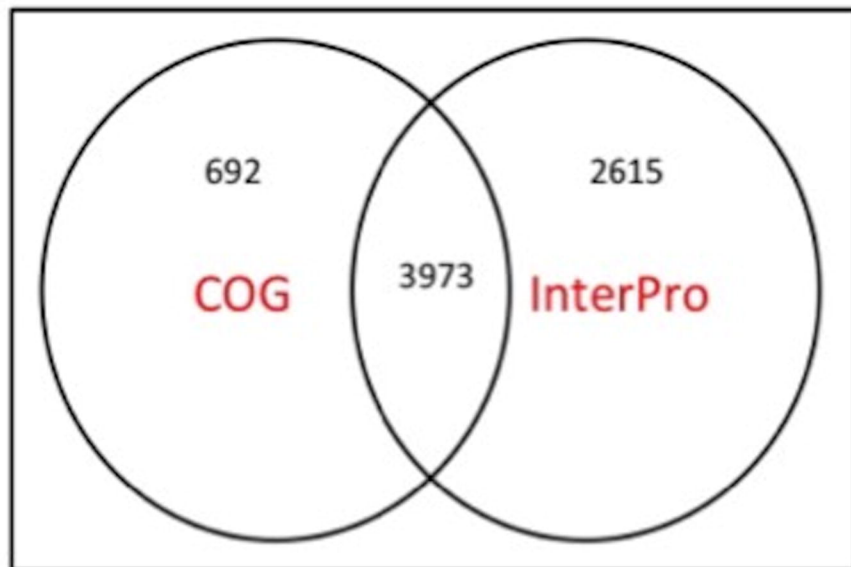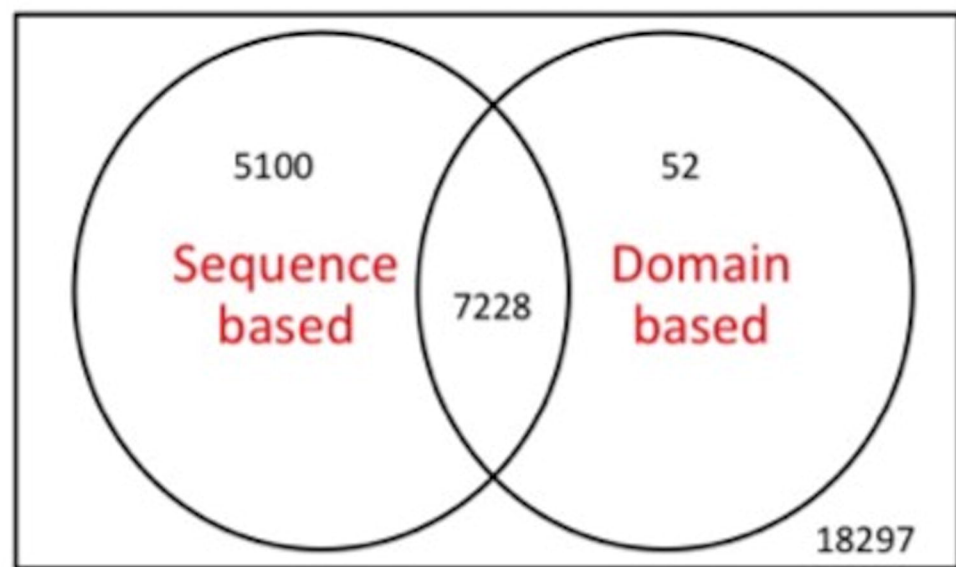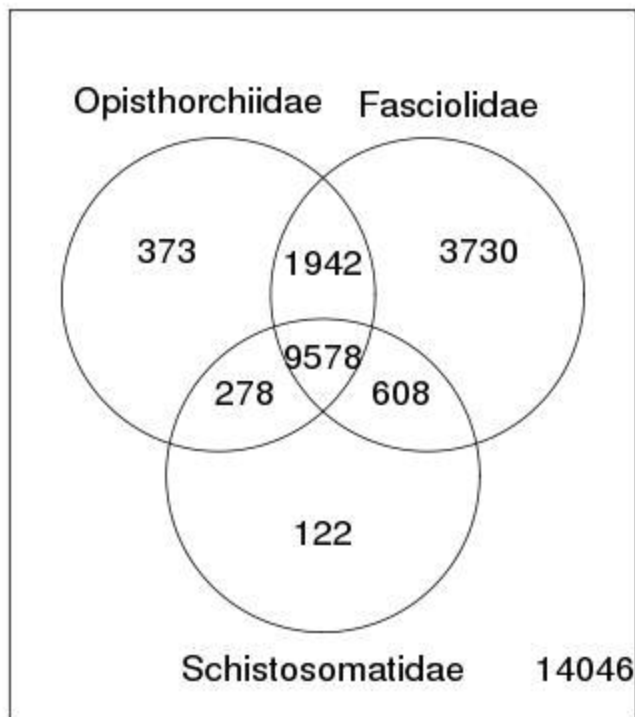922 **S9 Table. List of three LINE elements**

923

924

GO classification

A

| | |
|---|---|
| 67 | NR |
| 5456 | 3 |
| 6752 | |
| Uniprot | Nembase4 |
| 5 | |
| 29 | 16 |

B

| COG | InterPro |
|---|---|
| 692 | 2615 |
| 3973 | |

C

| Sequence based | Domain based |
|---|---|
| 5100 | 52 |
| 7228 | |

18297

```
            ┌─ 100 ─┌── Opisthorchis viverrini
            │       └── Clonorchis sinensis
         54 │
      100 ──┤──── Paragonimus westermani
            │──── Paramphistomum cervi
            │     ┌── Fasciolopsis buski
            │ 100 │── Fasciola hepatica
            │     └── Fasciola gigantica
  100 ──────┤          100
            │     ┌── Taenia solium
            │ 100 └── Taenia saginata
            │     ┌── Schistosoma japonicum
            │  98 └── Schistosoma mansoni
            └──── Ascaris lumbricoides
            └──── Ascaris suum
```
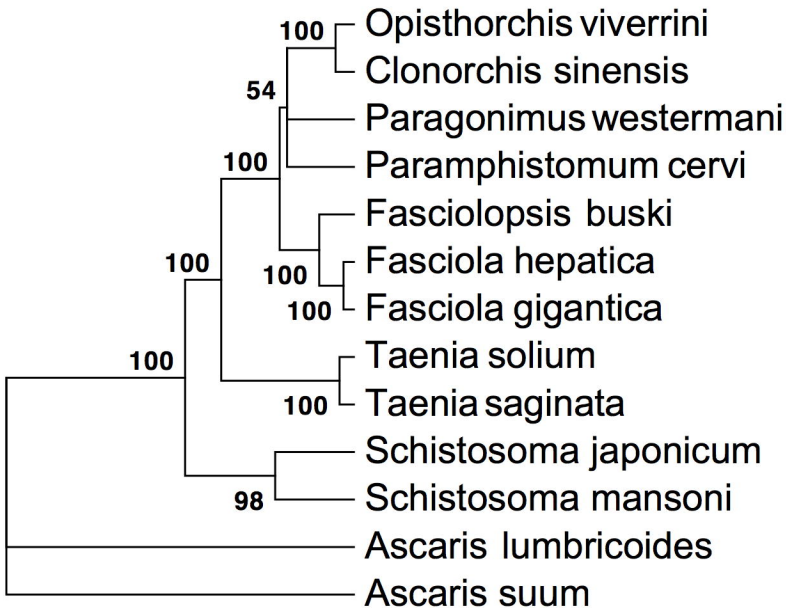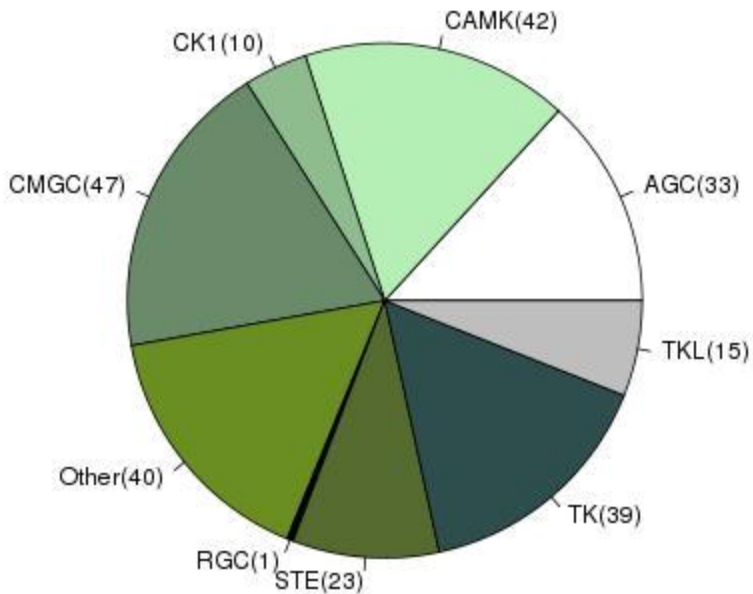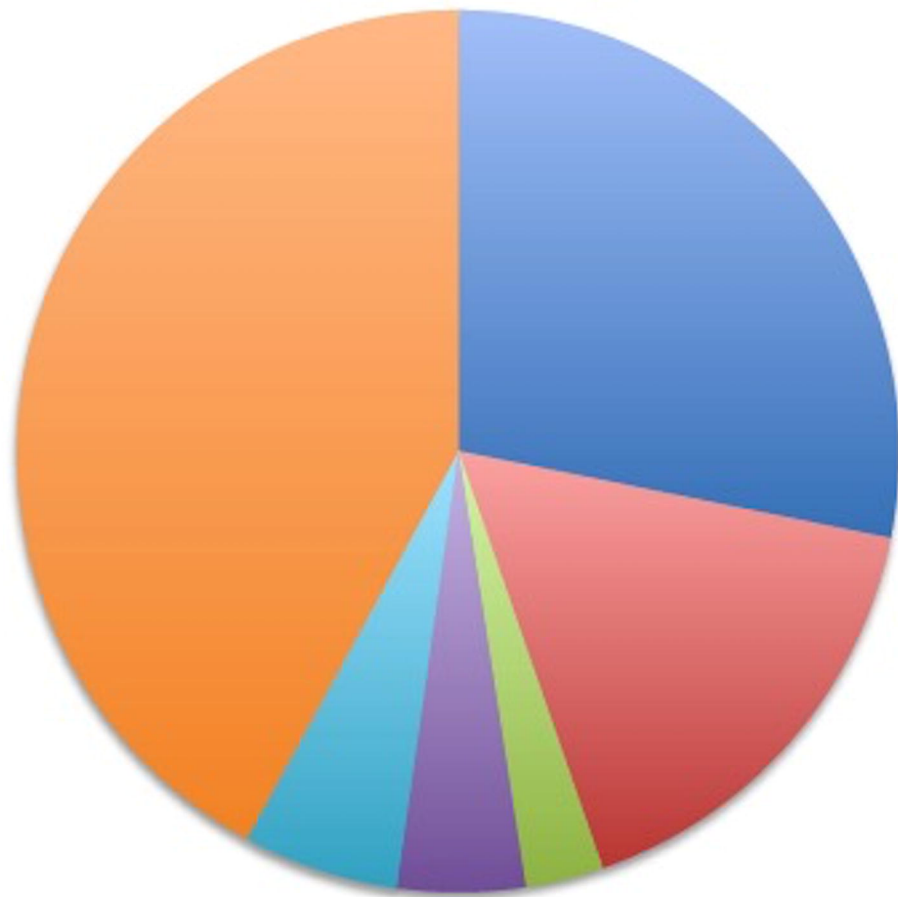
- Protease inhibitor (138)
- Cysteine (81)
- Aspartic (14)
- Unknown (23)
- Threonine (28)
- Serine (206)