

1 **Regulation of DNA methylation on key parasitism genes of** 2 **Cysticercus cellulosae revealed by integrative** 3 **epigenomic-transcriptomic analyses**

4
5 Shumin Sun^{*,1}, Xiaolei Liu^{*,1}, Guanyu Ji^{†,1}, Xuelin Wang^{*}, Junwen Wang[†], Xue Bai^{*},
6 Jing Xu^{*}, Jianda Pang^{*}, Yining Song^{*}, Xinrui Wang^{*}, Fei Gao^{*,§,2}, Mingyuan Liu^{*,‡,2}

7
8 ^{*}College of Animal Science and Technology, Inner Mongolia University for
9 Nationalities, Tongliao; Key Lab for Zoonoses Research, Ministry of Education,
10 Institute of Zoonoses, Jilin University, P. R. China

11 [†]E-GENE, Huahan Industrial Zone, Pingshan District, Shenzhen 518083, China

12 [§]Agricultural Genomics Institute at Shenzhen, Chinese Academy of Agricultural
13 Sciences, Pengfei Road 7, Dapeng District, Shenzhen 518112, China

14 [‡] Jiangsu Co-innovation Center for Prevention and Control of Important Animal
15 Infectious Diseases and Zoonoses, Yangzhou, Jiangsu, P. R. China

16
17 ¹ These authors contributed equally to this work.

18 ²Corresponding authors: Key Lab for Zoonoses Research, Ministry of Education,
19 Institute of Zoonoses, Jilin University, 5333 Xian Rd, Changchun. E-mail:
20 liumy@jlu.edu.cn; flys828@gmail.com

21 **Keywords:** Cysticercus cellulosae, DNA methylation, epigenetics, gene regulation

22

1 **Abstract**

2 **Background:** The life cycle of *Taenia solium* is characterized by different stages of
3 development, requiring various kinds of hosts that can appropriately harbor the eggs
4 (proglottids), the oncospheres, the larvae and the adults. Similar to other metazoan
5 pathogens, *T. solium* undergoes transcriptional and developmental regulation via
6 epigenetics during its complex lifecycle and host interactions.

7 **Result:** In the present study, we integrated whole-genome bisulfite sequencing and
8 RNA-seq technologies to characterize the genome-wide DNA methylation and its
9 effect on transcription of *Cysticercus cellulosae* of *T. solium*. We confirm that the *T.*
10 *solium* genome in the cysticercus stage is epigenetically modified by DNA
11 methylation in a pattern similar to that of other invertebrate genomes, i.e., sparsely or
12 moderately methylated. We also observed an enrichment of non-CpG methylation in
13 defined genetic elements of the *T. solium* genome. Furthermore, an integrative
14 analysis of both the transcriptome and the DNA methylome indicated a strong
15 correlation between these two datasets, suggesting that gene expression might be
16 tightly regulated by DNA methylation. Importantly, our data suggested that DNA
17 methylation might play an important role in repressing key parasitism-related genes,
18 including genes encoding excretion-secretion proteins, thereby raising the possibility
19 of targeting DNA methylation processes as a useful strategy in therapeutics of
20 cysticercosis.

21 **Conclusion:** Our study will provide a foundation for future studies to explore this key
22 epigenetic modification in development of *Cysticercus cellulosae* and in human

1 cysticercus disease.

2

3 **Introduction**

4 Cysticercus cellulosae, the larval stage of *T. solium*, resides in the central nervous
5 system, skeletal muscle, and other organs of both pigs and humans [1], resulting in the
6 high prevalence of cysticercosis worldwide. As a neglected tropical disease prioritized
7 by the World Health Organization, serious human disease burden [2] and annual
8 economic losses in livestock are caused by infection with this pork tapeworm. To
9 better control this disease, the mechanisms of transcriptional and developmental
10 regulation during its complex lifecycle and host interactions should be better
11 understood.

12 ADNA methylation, i.e., 5-methylcytosine (m5C) is an important epigenetic
13 mechanism that is present in the genomes of *Trichinella spiralis* [3] and
14 Platyhelminthes (*Schistosoma mansoni* [4]) parasitic nematodes. Via regulating gene
15 transcription, DNA methylation plays an important role in parasitism. Similar to *T.*
16 *solium*, *S. mansoni* belongs the phylum of Platyhelminthes. A previous study showed
17 that *S. mansoni* contains conserved DNA methyltransferase 2 (DNMT2) and
18 methyl-CpG binding proteins (MBD) [4]. Importantly, demethylation induced by
19 5-azacytidine can disrupt egg production and maturation, indicating an essential role
20 for DNA methylation in the normal development of this parasitic worms in this
21 phylum [4].

22 In the present study, we aimed to identify functional DNA methylation machinery and

1 detect cytosine methylation levels in the cysticercus cellulosae of *T. solium* based on a
2 draft genome that has been sequenced and annotated previously [5]. To achieve this
3 aim, we applied the whole-genome bisulfite sequencing (WGBS) method to
4 characterize the genome-wide DNA methylation pattern at single-base resolution [6].
5 Based on this unbiased characterization, our results confirm that in the cysticercus
6 stage, the *T. solium* genome [5] is epigenetically modified by DNA methylation in a
7 pattern similar to that of other invertebrate genomes, i.e., sparsely or moderately
8 methylated [7, 8]. We also observed an enrichment of non-CpG methylation in
9 defined genetic elements of *T. solium* genome, which is a pattern different from
10 mammalian methylomes [7, 8]. Furthermore, we applied RNA-seq technology to
11 profile gene expression. An integrative analysis on both the transcriptome and DNA
12 methylome indicated a strong correlation between these two datasets, suggesting that
13 gene expression might be tightly regulated by DNA methylation. Importantly, our
14 data suggested that DNA methylation might play an important role in repressing key
15 parasitism-related genes, including genes encoding excretion–secretion proteins. In
16 summary, for the first time, we provide data to characterize the DNA methylome and
17 the transcriptome of the *T. solium* cysticercus cellulosae. Our data will be valuable to
18 the community and will allow researchers to provide new insights into the mechanism
19 of methylation in cysticercosis in future studies.

20 **Materials and Methods**

21 **Sample collection and nuclei acid extraction**

22 Individual cysticerci were isolated from a single, naturally infected pig (Neimeigu

Province, China) and rinsed thoroughly several times with phosphate-buffered saline. The cysticerci were first frozen in liquid nitrogen and then finely ground to a powder-like texture. Genomic DNA was extracted using the phenol chloroform extraction method, and total RNA was purified using Trizol reagent (Invitrogen, CA, USA) according to the manufacturer's instructions. RNA was dissolved in diethylpyrocarbonate (DEPC)-treated water and treated with DNase I (Invitrogen, CA, USA). The quantity and quality of the DNA and RNA were tested by ultraviolet-Vis spectrophotometry with a NanoDrop 2000 (Thermo Scientific CA, USA).

BlastP searches and phylogenetic analysis of DNMTs

Reciprocal BlastP comparisons were first performed to identify DNMTs and MBD orthologs. Significant hits were defined as those satisfying the following criteria: E-value < 1e-5 and aligned segments covering at least 30% of the sequence length of the hit. For phylogenetic analysis, multiple sequence alignment was performed by Clustal W [9]. The MEGA7 with the neighbor-joining method [10, 11] based on the JTT+ G (Jones-Taylor-Thornton and Gamma Distribution) model was applied to reconstruct the phylogenetic tree.

MethylC-seq library construction and sequencing

Prior to library construction, 5 µg of genomic DNA extracted from a cysticercosis body was spiked with 25 ng unmethylated lambda DNA (Promega, Madison, WI, USA) and fragmented using a Covarias sonication system to a mean size of approximately 200 bp. After fragmentation, libraries were constructed according to

the Illumina Paired-End protocol with some modifications. Briefly, purified randomly fragmented DNA was treated with a mix of T4 DNA polymerase, Klenow fragment and T4 polynucleotide kinase to repair blunt ends and phosphorylate the ends. The blunt DNA fragments were subsequently 3' adenylated using Klenow fragment (3'-5' exo-), followed by ligation to adaptors synthesized with 5'-methylcytosine instead of cytosine using T4 DNA ligase. After each step, DNA was purified using a QIAquick PCR purification kit (Qiagen, Shanghai, China). Next, a ZYMO EZ DNA Methylation-Gold KitTM (ZYMO Research, Irvine, CA, USA) was employed to convert unmethylated cytosine to uracil, according to the manufacturer's instructions, and 220 to 250 bp converted products were size selected. Finally, PCR was carried out in a final reaction volume of 50 µl consisting of 20 µl of size selected fractions, 4 µl of 2.5 mM dNTPs, 5 µl of 10× buffer, 0.5 µl of JumpStartTM Taq DNA Polymerase, 2 µl of PCR primers and 18.5 µl water. The thermal cycling program was 94°C for 1 minute; 10 cycles of 94°C for 10 s, 62°C for 30 s, 72°C for 30 s; and then a 5-minute incubation at 72°C before holding the products at 12°C. The PCR products were purified using a QIAquick gel extraction kit (Qiagen). Before analysis with an Illumina HiSeq2500, the purified products were analyzed using a Bioanalyzer analysis system (Agilent, Santa Clara, CA, USA) and quantified by real-time PCR. Raw sequencing data were processed using the Illumina base-calling pipeline (Illumina Pipeline version 1.3.1). The sodium bisulfite non-conversion rate was calculated as the percentage of cytosines sequenced at cytosine reference positions in the lambda genome.

1 RNA-seq library construction and sequencing

2 Total RNA was extracted using the Invitrogen TRIzol Reagent and then treated with
3 RNase-free DNase I (Ambion, Guangzhou, China) for 30 minutes. The integrity of
4 total RNA was checked using an Agilent 2100 Bioanalyzer. cDNA libraries were
5 prepared according to the manufacturer's instructions (Illumina). The
6 poly(A)-containing mRNA molecules were purified using Oligo (dT) Beads
7 (Illumina) and 20 µg of total RNA from each sample. Tris-HCl (10 mM) was used to
8 elute the mRNA from the magnetic beads. To avoid priming bias when synthesizing
9 the cDNA, mRNA was fragmented before cDNA synthesis. Fragmentation was
10 performed using divalent cations at an elevated temperature. The cleaved mRNA
11 fragments were converted into double-stranded cDNA using SuperScript II, RNase H
12 and DNA Pol I, primed by random primers. The resulting cDNA was purified using a
13 QIAquick PCR Purification Kit (Qiagen). Then, the cDNA was subjected to end
14 repair and phosphorylation using T4 DNA polymerase, Klenow DNA polymerase and
15 T4 Polynucleotide Kinase (PNK). Subsequent purifications were performed using the
16 QIAquick PCR Purification Kit (Qiagen). These repaired cDNA fragments were
17 3'-adenylated using KlenowExo (Illumina) and purified using the MinElute PCR
18 Purification Kit (Qiagen), producing cDNA fragments with a single 'A' base
19 overhang at the 3' end for subsequent ligation to the adapters. Illumina PE adapters
20 were ligated to the ends of these 3'-adenylated cDNA fragments and then purified
21 using the MinElute PCR Purification Kit (Qiagen). To select a size range of templates
22 for downstream enrichment, the products of the ligation reaction were purified on 2%

1 TAE-Certified Low-Range Ultra Agarose (Bio-Rad, Hercules, CA, USA). cDNA
2 fragments (200 ± 20 bp) were excised from the gel and extracted using the QIAquick
3 Gel Extraction Kit (Qiagen). Fifteen rounds of PCR amplification were performed to
4 enrich the adapter-modified cDNA library using primers complementary to the ends
5 of the adapters (PCR Primer PE 1.0 and PCR Primer PE 2.0; Illumina).

6 **Transcriptome mapping**

7 RNA-seq reads were trimmed to a maximum length of 80 bp, and stretches of bases
8 having a quality score <30 at the ends of the reads were removed. Reads were mapped
9 using Tophat 2.0.11 [12]. As reference sequence for the transcriptome mapping we
10 used the current assembly of the *T. solium* database [13]. Expression was quantified
11 using cufflinks 2.1.1 [14]. RepeatMasker [15] were used to identify tandem repeats.

12 **Bisulfite mapping and methylation calling**

13 Reads were trimmed to a maximal length of 125 bp, and stretches of bases having a
14 quality score <30 at the ends of the reads were removed. Reads were mapped using
15 BSMAP 2.2.74 [16]. As a reference sequence for the bisulfite mapping we used the
16 current assembly of the *T. solium* genome [13]. Only reads mapping with both
17 partners of the read pairs at the correct distance were used. The CpG-specificity was
18 calculated by determining the number of cytosines called in all mapped reads at all
19 non-CpG positions and dividing by the number of all bases in all mapped reads at all
20 non-CpG positions. Methylation ratios were determined using a Python script
21 (methratio.py) distributed together with the BSMAP package for both the forward and
22 reverse strands.

1 **Protein network analyses**

2 The STRING online tool [17] was used with default parameters. Peptide sequences of
3 key genes were the input and were aligned to *Caenorhabditis elegans* protein
4 sequences.

5 **Data availability**

6 The *T. solium* methylome data have been deposited at NCBI/GEO/ under the
7 accession number GSE84086.

8 **Results**

9 **The presence of DNA methylation in the *T. solium* genome**

10 The methylation status of DNA is related to three types of enzymes, including DNA
11 methyltransferases, which affect maintenance methylation and de novo methylation.
12 To understand whether *T. solium* possesses the ability to methylate DNA, we first
13 conducted a reciprocal Blast alignment to identify genes that might be homologous to
14 known DNA (cytosine-5)-methyltransferases. As a result, two genes
15 (Scaffold00200.gene8095 and LongOrf.asmb1_16366) were identified that are
16 homologous to *DNMT3B* and *DNMT2*, respectively, with high sequence similarity
17 (e-value < 1e-10). Scaffold00067.gene4890 was aligned (e-value < 1e-5) with either
18 *DNMT3A* or *DNMT3B* from multiple species. In addition, more than one gene was
19 matched with *DNMT1*, among which Scaffold00068.gene4920 had the best hit
20 (e-value < 1e-10) (Table S1). Phylogenetic analyses by MEGA7 also supported these
21 results (Figure S1). Moreover, we searched for genes homologous to methyl-CpG
22 binding domain protein (MBD). Two candidate genes (LongOrf.asmb1_5021 and

LongOrf.assembl_14047) were homologous to *MBDs* in multiple species, including *Echinococcus granulosus*, which is closely related to *T. solium* evolutionarily (Table S1 and Figure S2). A full repertoire of functionally conserved amino acid residues was identified for both the potential DNMT2 and DNMT3 and the MBDs of *T. solium*, indicating that these proteins are functionally active (Table S2). However, a high level of divergence between *T. solium* and other species was observed for DNMT1 homologs (Table S2), which was in agreement with previous studies.

Given these results, we assessed the genome-wide DNA methylation profiles in *T. solium* using MethylC-Seq. There were 54.21 million raw reads generated (Table S3). BSMAP [16] was used to align the sequenced reads to the *T. solium* reference sequence, reaching an approximately 76.41% mapping rate. The average read depth was 11.32 per strand, while on average, over 50 Mb (90.52%) of each strand of the *T. solium* reference sequence was covered. Because of the potential occurrence of non-conversion and thymidine-cytosine sequencing errors, the false-positive rate was estimated by calculating the methylation level of lambda DNA, which is normally unmethylated (Materials and methods). We then applied the error rate (0.0041) to correct methylated cytosine sites (mC) identification according to the method described by Lister et al. [6], which is based on a binomial test and false discovery rate constraints. As a result, approximately 76.6 thousand mCs were estimated in the *T. solium* genome (accounting for 0.20% of the total cytosines sequenced with depth $\geq 5X$). Both symmetrical CpG methylation and asymmetrical non-CpG methylation were revealed.

1 **Characterization of overall methylation patterns**

2 We further characterized the global patterns of DNA methylation in the genomes of *T.*
3 *solium*. First, we showed the percentage of methylated cytosine of each sequence
4 context. Among the 76.6 thousand mCs across the entire genome, a majority (69.5%)
5 were in the context of CHH. In contrast, only 15.38% and 15.12% of the mCs were
6 located in the contexts of CHG and CpG, respectively (Figure S3 A). Furthermore,
7 most of the CpG and non-CpGs displayed a low methylation fraction (<30%) (Figure
8 S3B and C). These patterns are highly different from mammalian methylomes, in
9 which most 5mCs are located in CpG contexts and the majority of the CpGs are
10 highly methylated (>50%) [7]. Since most (69.5%) of the mCs in the *T. solium*
11 genome were in the CHH context, we further analyzed the sequence context of
12 mCHHs across the entire genome to further examine whether there is any sequence
13 bias in the enrichment of cytosine methylation in the CHH context. As a result, mCpA
14 was shown to be preferentially enriched within the methylated CHH dinucleotide
15 (Figure 1A). There were more than 21,000 methylated CpAs in each strand, meaning
16 that 55.73% of total CpAs were methylated in the entire genome (Figure 1B, D). This
17 result was consistent with reports that mCpA was predominantly found in another
18 tapeworm, *S. mansoni* [18, 19]. With regard to methylation levels, we did not observe
19 significant differences among different sequence contexts for mCs (Figure 1C).
20 We next examined whether there was any preference for the distance between
21 adjacent sites of DNA methylation in the *T. solium* genome. The relative distance
22 between mCs in each context within 50 nucleotides in introns was then analyzed

1 because of the steady methylation without any selective pressure by protein coding
2 genes in intron regions. Similar to the periodicity of 8–10 bases revealed in previous
3 studies on the Arabidopsis and human genomes [20], we also observed a strong
4 tendency of peaked enrichment of mCpA sites, which might be explained by a single
5 turn of the DNA helix (Figure 1E). Moreover, we found that mCpT revealed a similar
6 periodicity of 8-12 bases (Figure 1F), though the numbers of cytosines in the context
7 of CpG and CpC were too few to yield reliable results (Figure S3D and E). In
8 summary, our results indicated that the molecular mechanisms governing de novo
9 methylation at CpA sites may be similar among the *cysticercus* and the plant and
10 animal kingdoms.

11 We then examined the distribution of methylation levels for the four categories of
12 methylated cytosines across the entire genome. In general, similar mosaic distribution
13 patterns were observed for methylation levels of all types of mCs, that is, relatively
14 highly methylated domains were interspersed within regions with low methylation
15 (Figure S4A). Furthermore, the distribution of mCs across the genome was also
16 uneven; dense mCs of specific categories were occasionally enriched in specific
17 scaffolds (Figure S4B). Such a pattern has been observed in previous studies on other
18 invertebrates. We also examined the patterns of methylation in annotated elements,
19 including genes, tandem repeats, and transposable elements. The methylation
20 percentage of each cytosine context in exons was higher than that in other annotated
21 elements, especially CpAs, which accounted for a more than 2-fold greater percentage
22 than the other contexts in exons (Figure 2A).

1 We then examined the average methylation level in each element, which showed that
 2 average CpG methylation levels were higher than those other types of methylated
 3 cytosines, similar to mammalian genomes. However, the genome-wide pattern was
 4 again divergent from mammalian genomes, as higher average methylation in exons
 5 and lower methylation in introns of CpG sites were observed (Figure 2B, C and
 6 Figure S5). The trend for the average methylation of CpC and CpT was similar to that
 7 of CpG. However, a uniform distribution of CpA methylation levels in each annotated
 8 element was displayed (Figure 2C and Figure S5). We also analyzed the methylation
 9 level of each cytosine context in repeat regions (Figure 2D, E). Previous studies have
 10 indicated that transposable elements are usually unmethylated in the honey bee *Apis*
 11 *mellifera* and silkworm *Bombyx mori* [21, 22]. In *cysticercus*, we observed a similar
 12 phenomenon as the above species except that relatively highly methylated rRNAs
 13 were observed in *T. solium*. Notably, CpAs were methylated at a higher level or
 14 frequency than other types in LINE/L1 (Figure 2D, E).

15 **The relationship between methylation and gene expression**

16 It was reported that DNA methylation plays an important role in regulating gene
 17 expression. We evaluated gene expression in *T. solium* using Illumina
 18 high-throughput RNA-seq technology. Most of the raw reads could be uniquely
 19 mapped to previously annotated genes (88.17%). A total of 9,718 annotated genes out
 20 of 11,903 could be aligned with at least one unique read. To characterize the
 21 relationship between DNA methylation and gene expression, we divided the
 22 expressed genes with at least one read into quartiles of expression levels. We then

1 examined the distribution of methylation levels for different quartiles of expressed
2 genes and genes exhibiting no expression. High CpG and CpC methylation levels
3 were observed in upstream and exon regions of genes with the lowest expression.
4 Moreover, a negative correlation could also be observed between CpA and CpT
5 methylation levels of upstream and exons and expression levels of these expressed
6 genes. However, for silent genes, mainly high CpG and CpC methylation levels of
7 downstream regions were observed (Figure 3). Taken together, methylation levels of
8 mCs from both CpG or non-CpG sequence contexts were correlated with gene
9 expression levels, though different regulation mechanisms might be involved.
10 Next, to infer whether methylated genes were enriched for specific molecular
11 functions, we filtered out a total of 1,647 of the genes with the lowest expression and
12 1,354 of the most highly expressed genes, based on the criteria that at least one mC
13 was present within their genic regions. Then, we applied the WEGO (Web Gene
14 Ontology Annotation Plotting) tool [23] to functionally categorize the gene ontology
15 (GO) terms of these genes. We found that these two sets of genes displayed similar
16 patterns of GO enrichment, specifically, “cell” and “cell part” in Cellular Component,
17 “binding” and “catalytic” Molecular Functions, and “cellular process” and “metabolic
18 process” in Biological Process were relatively enriched. This result suggested that the
19 genes heavily regulated by DNA methylation were more prone to signaling regulation
20 or interaction with environmental factors, e.g., diet or metabolism (Figure S6). In
21 summary, these results suggested the potential for the regulation of *T. solium* genes by
22 DNA methylation, especially those that function as regulators of cell-cell or

1 cell-environmental communication. Furthermore, different molecular mechanisms
2 might be involved depending on different mC contexts and genes.

3 **Regulation of DNA methylation on key parasitism genes of *T. solium***

4 To obtain further insight into the epigenetic regulation of parasite development,
5 survival and parasite-host interactions of *T. solium*, we next studied conserved genes
6 across tapeworm-species and genes encoding excretion–secretion proteins (ESPs) in *T.*
7 *solium*. For conserved genes, we applied a gene set that was reported previously in a
8 study by Bjorn Victor et al., in which 261 genes conserved between *Taenia* and
9 *Echinococcus* tapeworms were obtained by comparing the transcriptomes of five
10 important intestinal parasites, including *T. multiceps*, *T. solium*, *E. granulosus*, *E.*
11 *multilocularis* and *T. pisiformis* [24]. Based on their results, we further retrieved 216
12 genes with the best blastx hit for each contig ($e < 1e-10$) and studied their DNA
13 methylation status. A total of 190 of these genes contained at least one mC across
14 their genic regions. As indicated in Figure 4C, CpG and CpC methylation levels in
15 upstream and exon regions were higher than other types of methylation and in other
16 genic regions. A further examination of the 190 genes revealed that 71 genes
17 contained CpG or CpC methylation within their upstream or exon regions. Therefore,
18 we searched for extensively methylated genes based on the criterion that the CpG and
19 CpC methylation levels of the examined gene were significantly higher than the
20 average value of the 71 genes. As a result, we revealed 14 conserved genes that were
21 extensively methylated on CpG/CpC sites within their upstream regions and exons
22 ($p < 0.05$). Compared with those 26 genes without mC, we found these 14 genes were

expressed at a significantly lower level (Figure 4A), suggesting DNA methylation is a key mechanism for the transcriptional regulation of these conserved genes. For ESPs, we also applied a dataset containing 76 ESPs for *T. solium*, which was identified by Bjorn Victor et al. using a proteomics strategy [25]. We applied the BlastP algorithm to align these ESPs back to the genome and revealed 111 gene sequences that might encode these ESPs (Table S4). Using the same criterion for conserved genes, we found 13 extensively methylated genes. Similarly, gene expression comparisons again revealed that these 13 genes were expressed at significantly lower levels than the 26 non-methylated genes (Figure 4A). Using a similar strategy, we also looked into genes containing methylated CpAs and CpTs within their upstream or exon regions. However, no clear difference in gene expression levels was observed (Figure 4B). These results indicated that CpG/CpC methylation in upstream regions and exons played a major role in *T. solium* gene repression. Furthermore, we found a different distribution pattern between mCpG and mCpC for these repressed genes, in which mCpG were mostly distributed in upstream regions, while mCpC were more often in exons (Figure S5). Based on the above analyses, we revealed 27 key genes that might be repressed by DNA methylation mechanisms. Interestingly, a protein-protein interaction analysis using the STRING online tool [17] indicated strong mutual interactions among these conserved proteins and ESPs (Figure 5) based on annotation of the model organism *Caenorhabditis elegans*.

Discussion

The larval stage of the pork tapeworm *T. solium* is responsible for cysticercosis,

1 which represents an important public health problem that occurs mainly in developing
2 countries. *T. solium* cysticerci have developed diverse mechanisms to protect
3 themselves from host immune attack [26], among which epigenetics may play an
4 important role in gene regulation related to parasitism [24]. Recently, Geyer et al.
5 found that essential DNA methylation machinery components, such as DNMT2 and
6 MBD, are well conserved throughout the Platyhelminthes [4, 27]. Invertebrate
7 DNMT2s are believed to retain strong DNA methyltransferase activity [28], which is
8 different from vertebrate DNMT2s, which are considered tRNA methyltransferases
9 [29]. Our computational searches indicated that both DNMT2 and DNMT3 are found
10 in *T. solium*, which implies the potential existence of a more sophisticated DNA
11 methylation machinery. In addition, *MBD2/3* homologs were also identified in the *T.*
12 *solium* genome.

13 Based on these results, our present study focused on characterizing the DNA
14 methylome and transcriptome of *T. solium* cysticerci, aiming for providing
15 comprehensive omics profiles for this important parasitic stage of *T. solium*. We
16 revealed a mosaic methylation pattern in *T. solium* that is typical of other
17 invertebrates [3, 4, 21, 22, 30, 31]. Cytosine methylation was predominantly found in
18 the CpA dinucleotide context, similar to other invertebrate species, including
19 *Drosophila melanogaster* [32] and other *platyhelminths* such as *S. mansoni* [27],
20 which might be mediated by MBD2/3 proteins [33, 34]. These patterns in the DNA
21 methylome might be closely related to the activity of different DNMTs. As DNMT1
22 functions as a maintenance methylase by copying methylation after DNA replication

1 with the help of Uhrf1 [35], a lack of DNMT1 might help to explain why much
2 non-symmetrical methylation was observed in the *platyhelminth* genome. We also
3 found that a periodicity for two pairs of mCpA and mCpT sites spaced with 13 bases
4 between the pairs, corresponding to a single turn of the DNA helix, as previously
5 observed. A structural study of the mammalian de novo methyltransferase DNMT3A
6 and its partner protein DNMT3L found that two copies of each form a heterotetramer
7 that contains two active sites separated by a length of 8–10 nucleotides in a DNA
8 helix [36, 37]. Because we could not locate *DNMT3L* in the *T. solium* genome, the
9 consistent 8–10 nucleotide spacing we observed in the *T. solium* genome might be due
10 to *DNMT3A* alone or an unknown factor other than *DNMT3L*.

11 Gene methylation is believed to be an evolutionarily ancient means of transcriptional
12 control. Among plants, vertebrates and some invertebrates such as *T. spiralis*, the
13 notion that methylation in promoters primarily represses genes by impeding
14 transcriptional initiation has been widely accepted [8, 18, 19], whereas intermediate
15 levels of expression have been associated with genes experiencing the greatest extent
16 of methylation in the gene body, indicating a bell-shaped relationship [38–40].

17 However, in invertebrates, such as the fungus *Neurospora crassa* [41] and the
18 silkworm *Bombyx mori*, transcription initiation is unaffected. Thus, DNA methylation
19 shows remarkable diversity in its extent and function across eukaryotic evolution. In
20 our *T. solium* results, we also found that methylation levels of mCs were correlated
21 with gene expression levels. Depending on different sequence contexts, methylation
22 seemed to function differently in transcriptional regulation. Intriguingly, high CpG

1 and CpC methylation levels of downstream regions, but not of promoter regions, were
 2 observed for silent genes (Figure 4). In contrast, upstream methylation seemed to
 3 mostly affect genes with low expression. Currently, knowledge on methylation
 4 patterns and their effects on gene regulation in non-vertebrates are still limited, though
 5 species-specific diversity has been observed [42]. Therefore, more data should be
 6 collected for the DNA methylomes of each specific species to characterize their
 7 patterns and functions.

8 In addition to characterizing the general distribution pattern of genome-wide DNA
 9 methylation, we also focused on methylation status of important *T. solium* genes.

10 Based on previous studies, we looked into 27 extensively methylated genes that are
 11 important for *T. solium* development, survival and parasite-host interactions. We
 12 found 13 of these 27 genes mutually interacted based on annotations in the model
 13 organism *C. elegans*. Specifically, ESPs formed the core of the
 14 protein-protein-interaction network, while proteins encoded by conserved genes
 15 directly interacted with specific ESPs. Among the ESPs that were potentially
 16 regulated by DNA methylation, we found that two heat shock proteins (HSPs), hsp-90
 17 and hsp-1, were highlighted and mutually interacted. The heat shock response is a
 18 general homeostatic mechanism that protects cells and organisms from the deleterious
 19 effects of environmental stress [43]. Together with COX-2, these proteins were
 20 previously reported to be important parasitism-related proteins [44]. Furthermore, we
 21 also revealed two genes encoding diagnostic antigens that might be regulated by DNA
 22 methylation, including diagnostic antigen gp50 and an 8 kDa diagnostic protein. GP50

1 is a glycosylated and GPI-anchored membrane protein. In recent years, one
2 component of the lentil lectin purified glycoprotein (LLGP) antigens has been used
3 for antibody-based diagnosis of cysticercosis [45]. The 8 kDa family members are
4 metacestode excretory/secretory glycoproteins, which invoke strong antibody
5 reactions in infected individuals [46]. Importantly, our data suggest that DNA
6 methylation might play a key role in repressing their transcription, implying a
7 potential for drug development in the future that can target epigenetic modification
8 machinery to control this important neglected tropical disease.

9

10 **Acknowledgements**

11 This study was supported by the National Natural Science Foundation of China
12 (NSFC: 31460658, 31402185, 31440085, 31160504, 31030064, 31520103916) and
13 the China Postdoctoral Science Foundation (2012M520674).

14 Author contributions: M. Liu, F. Gao and S. Sun conceived and supervised the
15 project. X. Liu and J. Wang performed NGS library construction. G. Ji and H. Lu
16 conducted bioinformatic analysis. X. Bai, X. Wang, J. Pang, Y. Zhao, K. Yuan, X. Li
17 collected and prepared samples. S. Sun prepared the manuscript. All authors have read
18 and approved the manuscript for publication.

19

20

21

22 **References**

- 1 1. Schantz PM, Cruz M, Sarti E, Pawlowski Z: **Potential eradicability of taeniasis and**
2 **cysticercosis.** *Bull Pan Am Health Organ* 1993, **27**(4):397-403.
- 3 2. Sciutto E, Fragoso G, Fleury A, Laclette JP, Sotelo J, Aluja A, Vargas L, Larralde C: **Taenia**
4 **solium disease in humans and pigs: an ancient parasitosis disease rooted in developing**
5 **countries and emerging as a major health problem of global dimensions.** *Microbes Infect*
6 2000, **2**(15):1875-1890.
- 7 3. Gao F, Liu X, Wu XP, Wang XL, Gong D, Lu H, Xia Y, Song Y, Wang J, Du J *et al*: **Differential**
8 **DNA methylation in discrete developmental stages of the parasitic nematode *Trichinella***
9 **spiralis.** *Genome Biol* 2012, **13**(10):R100.
- 10 4. Geyer KK, Rodriguez Lopez CM, Chalmers IW, Munshi SE, Truscott M, Heald J, Wilkinson MJ,
11 Hoffmann KF: **Cytosine methylation regulates oviposition in the pathogenic blood fluke**
12 ***Schistosoma mansoni*.** *Nat Commun* 2011, **2**:424.
- 13 5. Aguilar-Diaz H, Bobes RJ, Carrero JC, Camacho-Carranza R, Cervantes C, Cevallos MA, Davila
14 G, Rodriguez-Dorantes M, Escobedo G, Fernandez JL *et al*: **The genome project of *Taenia***
15 ***solium*.** *Parasitol Int* 2006, **55 Suppl**:S127-130.
- 16 6. Lister R, Ecker JR: **Finding the fifth base: genome-wide sequencing of cytosine methylation.**
17 *Genome Res* 2009, **19**(6):959-966.
- 18 7. Feng S, Cokus SJ, Zhang X, Chen PY, Bostick M, Goll MG, Hetzel J, Jain J, Strauss SH, Halpern
19 ME *et al*: **Conservation and divergence of methylation patterning in plants and animals.**
20 *Proc Natl Acad Sci U S A* 2010, **107**(19):8689-8694.
- 21 8. Zemach A, McDaniel IE, Silva P, Zilberman D: **Genome-wide evolutionary analysis of**
22 **eukaryotic DNA methylation.** *Science* 2010, **328**(5980):916-919.
- 23 9. Saitou N, Nei M: **The neighbor-joining method: a new method for reconstructing**
24 **phylogenetic trees.** *Mol Biol Evol* 1987, **4**(4):406-425.
- 25 10. Jones DT, Taylor WR, Thornton JM: **The rapid generation of mutation data matrices from**
26 **protein sequences.** *Comput Appl Biosci* 1992, **8**(3):275-282.
- 27 11. Kumar S, Stecher G, Tamura K: **MEGA7: Molecular Evolutionary Genetics Analysis Version**
28 **7.0 for Bigger Datasets.** *Mol Biol Evol* 2016, **33**(7):1870-1874.
- 29 12. Trapnell C, Pachter L, Salzberg SL: **TopHat: discovering splice junctions with RNA-Seq.**
30 *Bioinformatics* 2009, **25**(9):1105-1111.
- 31 13. **The *Taenia solium* Genome Project** [<http://www.taeniasolium.unam.mx/taenia/>]. Accessed
32 03 Sept 2016.]
- 33 14. Trapnell C, Williams BA, Pertea G, Mortazavi A, Kwan G, van Baren MJ, Salzberg SL, Wold BJ,
34 Pachter L: **Transcript assembly and quantification by RNA-Seq reveals unannotated**
35 **transcripts and isoform switching during cell differentiation.** *Nat Biotechnol* 2010,
36 **28**(5):511-515.

- 1 15. Tarailo-Graovac M, Chen N: **Using RepeatMasker to identify repetitive elements in genomic**
2 **sequences**. *Current protocols in bioinformatics* 2009, **25**:04.10.01-04.10.04.
- 3 16. Xi Y, Li W: **BMAP: whole genome bisulfite sequence MAPping program**. *BMC*
4 *Bioinformatics* 2009, **10**:232.
- 5 17. **STRING** [<http://string-db.org/>. Accessed 03 Sept 2016.]
- 6 18. Zhang X: **The epigenetic landscape of plants**. *Science* 2008, **320**(5875):489-492.
- 7 19. Weber M, Hellmann I, Stadler MB, Ramos L, Paabo S, Rebhan M, Schubeler D: **Distribution,**
8 **silencing potential and evolutionary impact of promoter DNA methylation in the human**
9 **genome**. *Nature genetics* 2007, **39**(4):457-466.
- 10 20. Cokus SJ, Feng S, Zhang X, Chen Z, Merriman B, Haudenschild CD, Pradhan S, Nelson SF,
11 Pellegrini M, Jacobsen SE: **Shotgun bisulphite sequencing of the Arabidopsis genome**
12 **reveals DNA methylation patterning**. *Nature* 2008, **452**(7184):215-219.
- 13 21. Xiang H, Zhu J, Chen Q, Dai F, Li X, Li M, Zhang H, Zhang G, Li D, Dong Y *et al*: **Single**
14 **base-resolution methylome of the silkworm reveals a sparse epigenomic map**. *Nat*
15 *Biotechnol* 2010, **28**(5):516-520.
- 16 22. Lyko F, Foret S, Kucharski R, Wolf S, Falckenhayn C, Maleszka R: **The honey bee epigenomes:**
17 **differential methylation of brain DNA in queens and workers**. *PLoS Biol* 2010,
18 **8**(11):e1000506.
- 19 23. Ye J, Fang L, Zheng H, Zhang Y, Chen J, Zhang Z, Wang J, Li S, Li R, Bolund L *et al*: **WEGO: a**
20 **web tool for plotting GO annotations**. *Nucleic Acids Res* 2006, **34**(Web Server
21 issue):W293-297.
- 22 24. Robert McMaster W, Morrison CJ, Kobor MS: **Epigenetics: A New Model for Intracellular**
23 **Parasite-Host Cell Regulation**. *Trends in parasitology* 2016, **32**(7):515-521.
- 24 25. Victor B, Kanobana K, Gabriel S, Polman K, Deckers N, Dorny P, Deelder AM, Palmblad M:
25 **Proteomic analysis of Taenia solium metacestode excretion-secretion proteins**. *Proteomics*
26 2012, **12**(11):1860-1869.
- 27 26. Hewitson JP, Grainger JR, Maizels RM: **Helminth immunoregulation: the role of parasite**
28 **secreted proteins in modulating host immunity**. *Molecular and biochemical parasitology*
29 2009, **167**(1):1-11.
- 30 27. Geyer KK, Chalmers IW, Mackintosh N, Hirst JE, Geoghegan R, Badets M, Brophy PM, Brehm
31 K, Hoffmann KF: **Cytosine methylation is a conserved epigenetic feature found throughout**
32 **the phylum Platyhelminthes**. *BMC Genomics* 2013, **14**:462.
- 33 28. Phalke S, Nickel O, Walluscheck D, Hortig F, Onorati MC, Reuter G: **Retrotransposon silencing**
34 **and telomere integrity in somatic cells of Drosophila depends on the cytosine-5**
35 **methyltransferase DNMT2**. *Nature genetics* 2009, **41**(6):696-702.

- 1 29. Goll MG, Kirpekar F, Maggert KA, Yoder JA, Hsieh CL, Zhang X, Golik KG, Jacobsen SE, Bestor
2 TH: **Methylation of tRNAAsp by the DNA methyltransferase homolog Dnmt2.** *Science* 2006,
3 **311**(5759):395-398.
- 4 30. Bird AP, Taggart MH, Smith BA: **Methylated and unmethylated DNA compartments in the**
5 **sea urchin genome.** *Cell* 1979, **17**(4):889-901.
- 6 31. del Gaudio R, Di Giaimo R, Geraci G: **Genome methylation of the marine annelid worm**
7 **Chaetopterus variopedatus: methylation of a CpG in an expressed H1 histone gene.** *FEBS*
8 *Lett* 1997, **417**(1):48-52.
- 9 32. Lyko F, Ramsahoye BH, Jaenisch R: **DNA methylation in Drosophila melanogaster.** *Nature*
10 2000, **408**(6812):538-540.
- 11 33. Raddatz G, Guzzardo PM, Olova N, Fantappie MR, Rampp M, Schaefer M, Reik W, Hannon GJ,
12 Lyko F: **Dnmt2-dependent methylomes lack defined DNA methylation patterns.** *Proc Natl*
13 *Acad Sci U S A* 2013, **110**(21):8627-8631.
- 14 34. Marhold J, Kramer K, Kremmer E, Lyko F: **The Drosophila MBD2/3 protein mediates**
15 **interactions between the MI-2 chromatin complex and CpT/A-methylated DNA.**
16 *Development* 2004, **131**(24):6033-6039.
- 17 35. Law JA, Jacobsen SE: **Establishing, maintaining and modifying DNA methylation patterns in**
18 **plants and animals.** *Nat Rev Genet* 2010, **11**(3):204-220.
- 19 36. Huff JT, Zilberman D: **Dnmt1-independent CG methylation contributes to nucleosome**
20 **positioning in diverse eukaryotes.** *Cell* 2014, **156**(6):1286-1297.
- 21 37. Jia D, Jurkowska RZ, Zhang X, Jeltsch A, Cheng X: **Structure of Dnmt3a bound to Dnmt3L**
22 **suggests a model for de novo DNA methylation.** *Nature* 2007, **449**(7159):248-251.
- 23 38. Zilberman D, Gehring M, Tran RK, Ballinger T, Henikoff S: **Genome-wide analysis of**
24 **Arabidopsis thaliana DNA methylation uncovers an interdependence between methylation**
25 **and transcription.** *Nature genetics* 2007, **39**(1):61-69.
- 26 39. Jjingo D, Conley AB, Yi SV, Lunyak VV, Jordan IK: **On the presence and role of human**
27 **gene-body DNA methylation.** *Oncotarget* 2012, **3**(4):462-474.
- 28 40. Nanty L, Carbajosa G, Heap GA, Ratnieks F, van Heel DA, Down TA, Rakyan VK: **Comparative**
29 **methylomics reveals gene-body H3K36me3 in Drosophila predicts DNA methylation and**
30 **CpG landscapes in other invertebrates.** *Genome Res* 2011, **21**(11):1841-1850.
- 31 41. Rountree MR, Selker EU: **DNA methylation inhibits elongation but not initiation of**
32 **transcription in Neurospora crassa.** *Genes Dev* 1997, **11**(18):2383-2395.
- 33 42. Schubeler D: **Function and information content of DNA methylation.** *Nature* 2015,
34 **517**(7534):321-326.
- 35 43. Ferrer E, Gonzalez LM, Foster-Cuevas M, Cortez MM, Davila I, Rodriguez M, Sciotto E,
36 Harrison LJ, Parkhouse RM, Garate T: **Taenia solium: characterization of a small heat shock**

- 1 **protein (Tsol-sHSP35.6) and its possible relevance to the diagnosis and pathogenesis of**
- 2 **neurocysticercosis. *Experimental parasitology* 2005, **110**(1):1-11.**
- 3 44. Choi W, Chu J: **The characteristics of the expression of heat shock proteins and COX-2 in the**
- 4 **liver of hamsters infected with *Clonorchis sinensis*, and the change of endocrine hormones**
- 5 **and cytokines. *Folia parasitologica* 2012, **59**(4):255-263.**
- 6 45. Hancock K, Patabhi S, Greene RM, Yushak ML, Williams F, Khan A, Priest JW, Levine MZ,
- 7 Tsang VC: **Characterization and cloning of GP50, a *Taenia solium* antigen diagnostic for**
- 8 **cysticercosis. *Molecular and biochemical parasitology* 2004, **133**(1):115-124.**
- 9 46. Ferrer E, Sanchez J, Milano A, Alvarez S, La Rosa R, Lares M, Gonzalez LM, Cortez MM, Davila
- 10 I, Harrison LJ *et al*: **Diagnostic epitope variability within *Taenia solium* 8 kDa antigen family:**
- 11 **implications for cysticercosis immunodetection. *Experimental parasitology* 2012,**
- 12 **130(1):78-85.**

13

14

15

16

17 **Figure legends**

18 **Figure 1. Cytosine DNA methylation in *T. solium*.** (A) Logo plots of the sequences
19 proximal to sites of cytosine DNA methylation in each sequence context in *T. solium*;
20 (B) Number of mCs for each type of dinucleotide; (C) Distribution of mCs; (D)
21 Percentage of each type of dinucleotide; (E, F) Prevalence of mCA/mCT sites (y-axis)
22 as a function of the number of bases between adjacent mCA/mCT sites (x-axis) based
23 on all non-redundant pair-wise distances up to 50 nt in all introns. The blue line
24 represents smoothing with cubic splines.

25

26 **Figure 2. Average methylation levels of different genomic regions.** (A, B) Average
27 density of methylation; (C) Average density of mC methylation distributed on the

1 genome. Two-kilobase regions upstream and downstream of each gene were divided
2 into 100-bp (bp) intervals. Each coding sequence or intron was divided into 20
3 intervals (5% per interval). (D) Average mC methylation level on repeat elements; (E)
4 Number of mCs on each repeat element.

5

6 **Figure 3. Relationship between mC DNA methylation and expression levels of**
7 **genes in *T. solium*.** Percentage of methylation within genes that were classified based
8 on expression levels. The first class includes silent genes with no sequencing reads
9 detected, and the second to fifth classes cover expressed genes from the lowest 25% to
10 the highest 25%. Regions of 2 kb upstream and downstream of each gene was divided
11 into 100-bp intervals, and each gene was divided into 20 intervals (5% per interval).

12

13 **Figure 4. Boxplots of gene expression levels of un-methylated genes, and genes**
14 **with conserved methylation and methylated ESP genes based on either (A)**
15 **CpG/CpC methylation or (B) CpA/CpT methylation.**

16

17 **Figure 5. Protein-protein interactions of conserved genes and ESP genes.**

18

19 **Supplemental Materials:**

20 **Figures S1: Phylogenetic tree of DNMT proteins.**

21 **Figures S2: Phylogenetic tree of mbd proteins.**

22 **Figures S3: Patterns and chromosomal distribution of DNA methylation in *T.***

1 *solium*.

2 **Figures S4: DNA methylation patterns and chromosomal distribution.**

3 **Figures S5: Average density of methylation levels of cytosine distributed on**
4 **genome.**

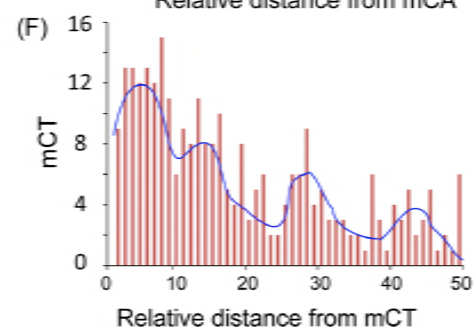
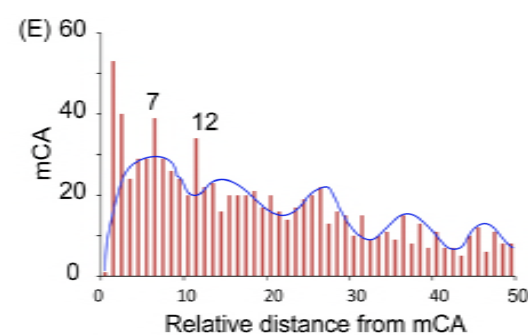
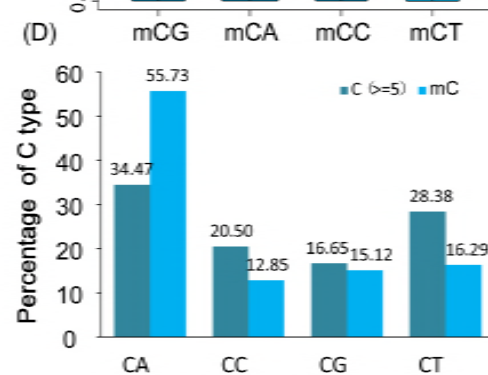
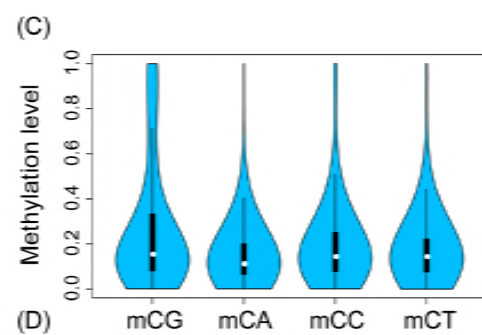
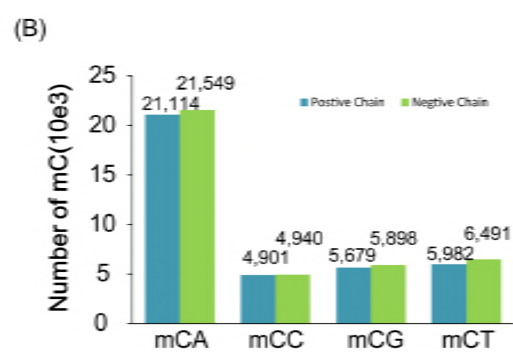
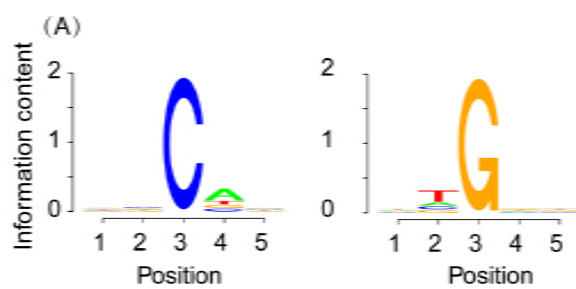
5 **Figures S6: Gene Ontology (GO) analysis for the genes with the lowest**
6 **expression (2nd) and the most highly expressed (5th) genes.**

7 **Table S1: Results of reciprocal BlastP searches of T.s DNMTs and MBD.**

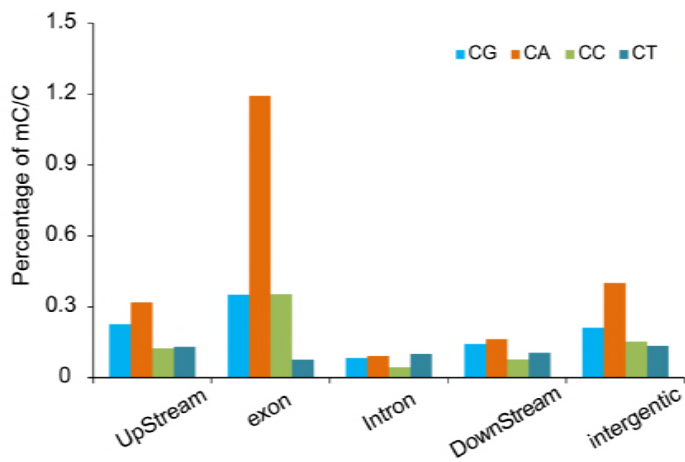
8 **Table S2: Results of ClustW on dnmt protein sequences.**

9 **Table S3: Data summary of MethylC-seq and RNA-seq.**

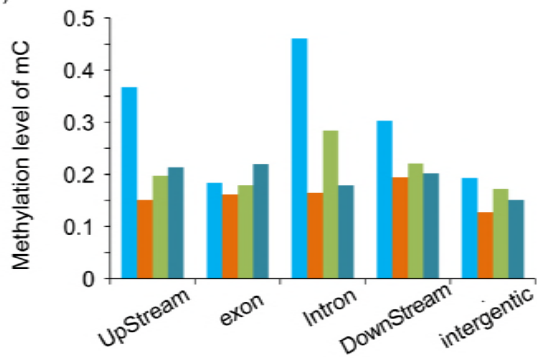
10 **Table S4: Summary of key parasitism genes that are methylated.**



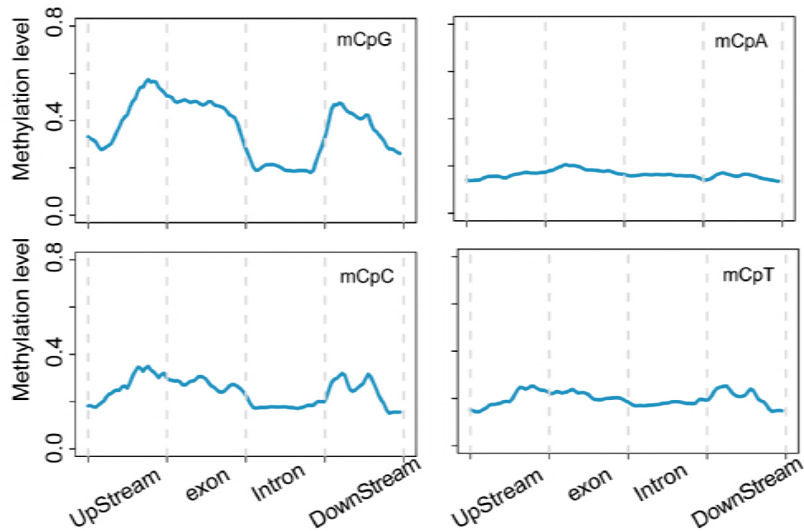
(A)



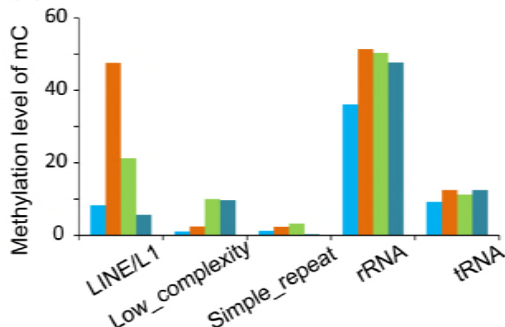
(B)



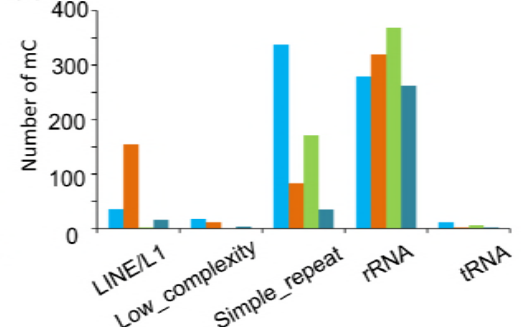
(C)

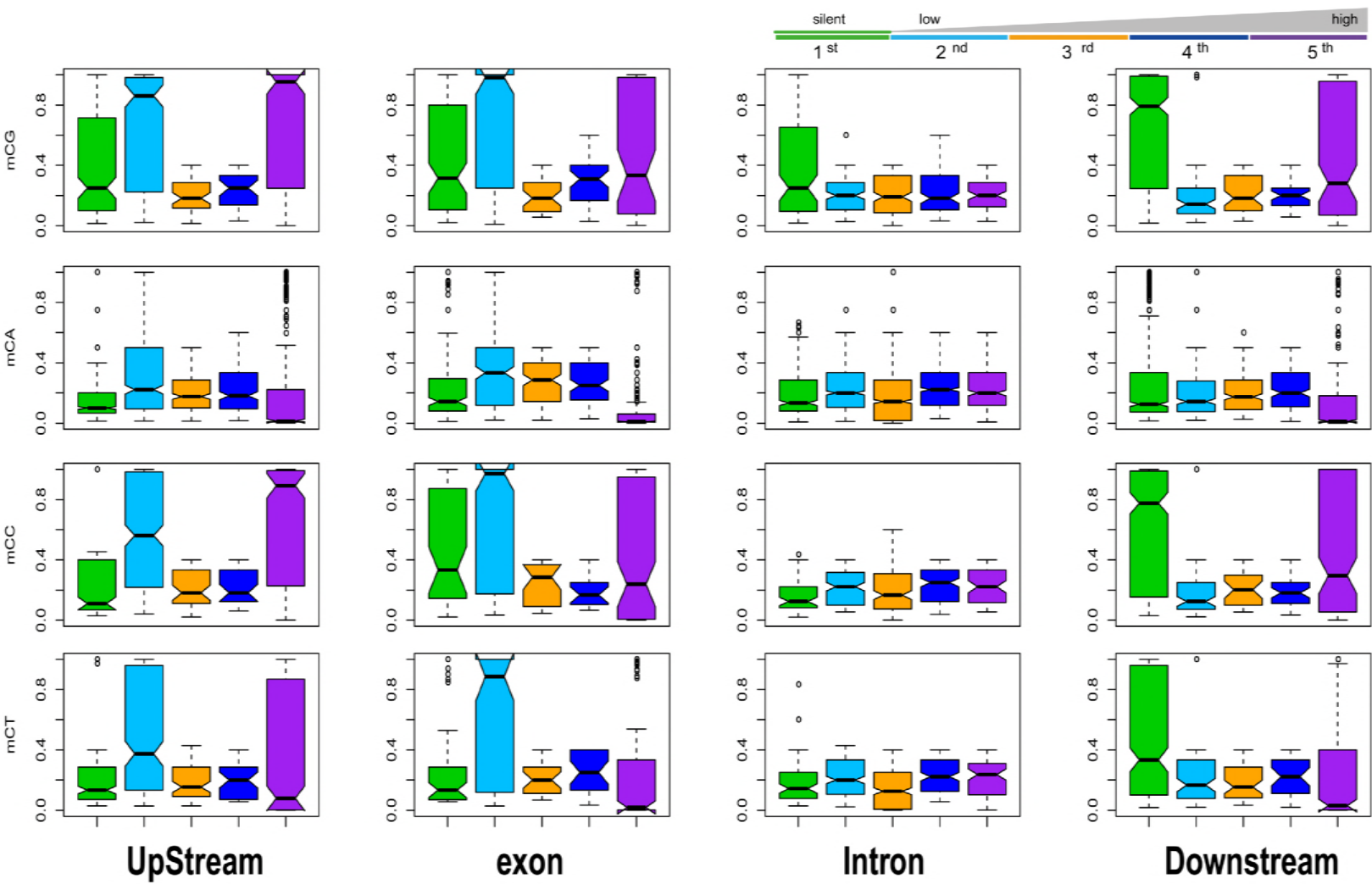


(D)

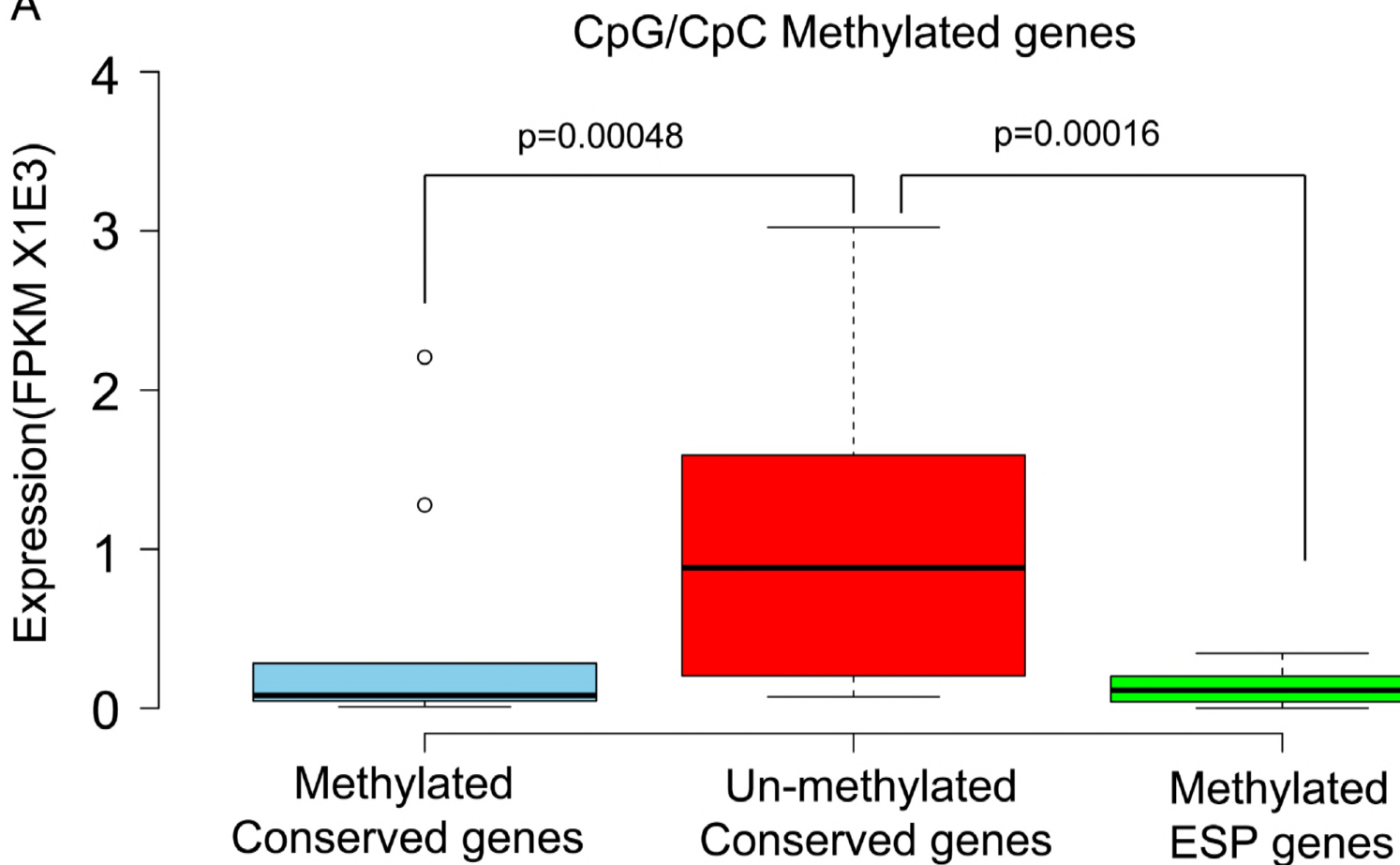


(E)

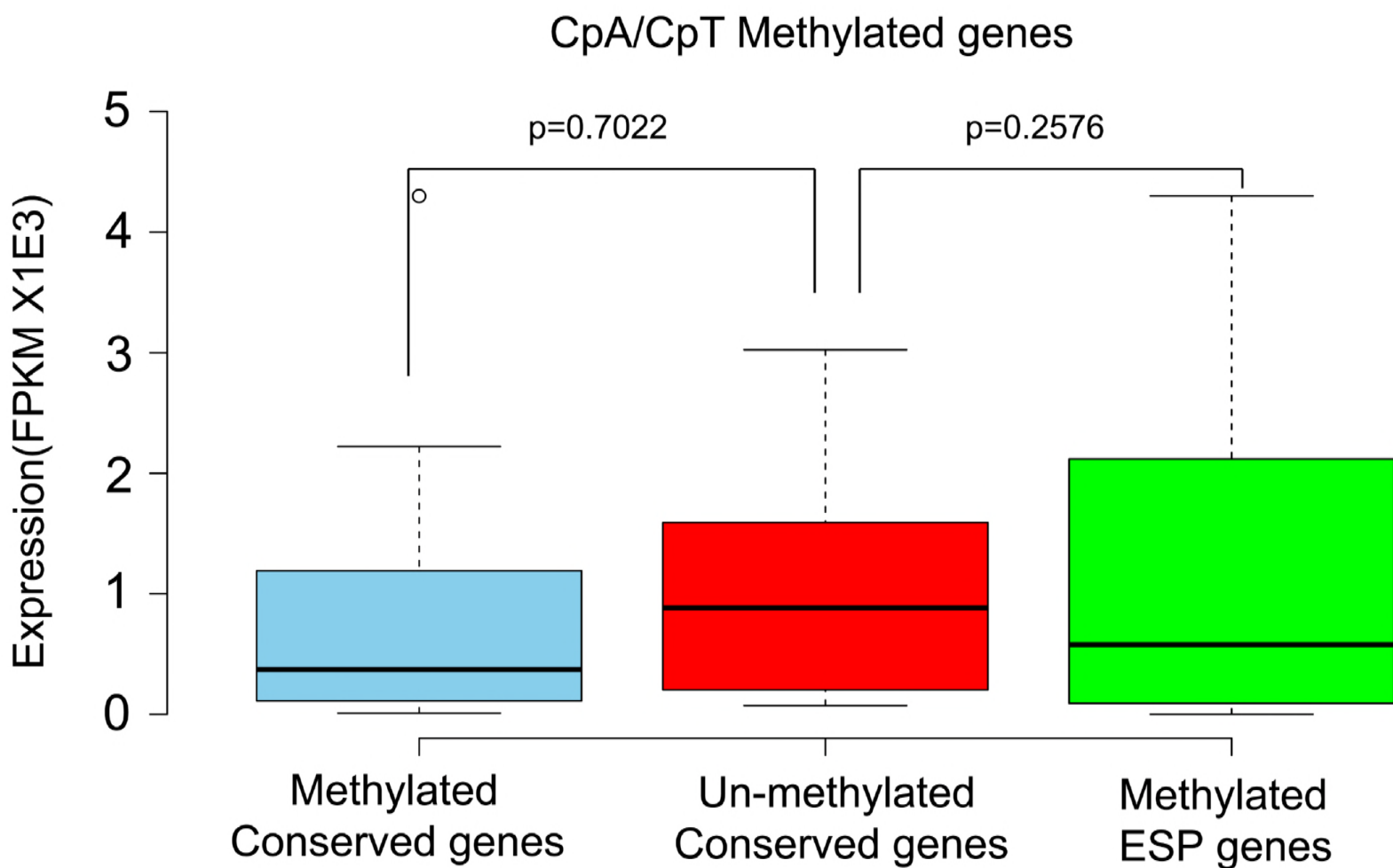


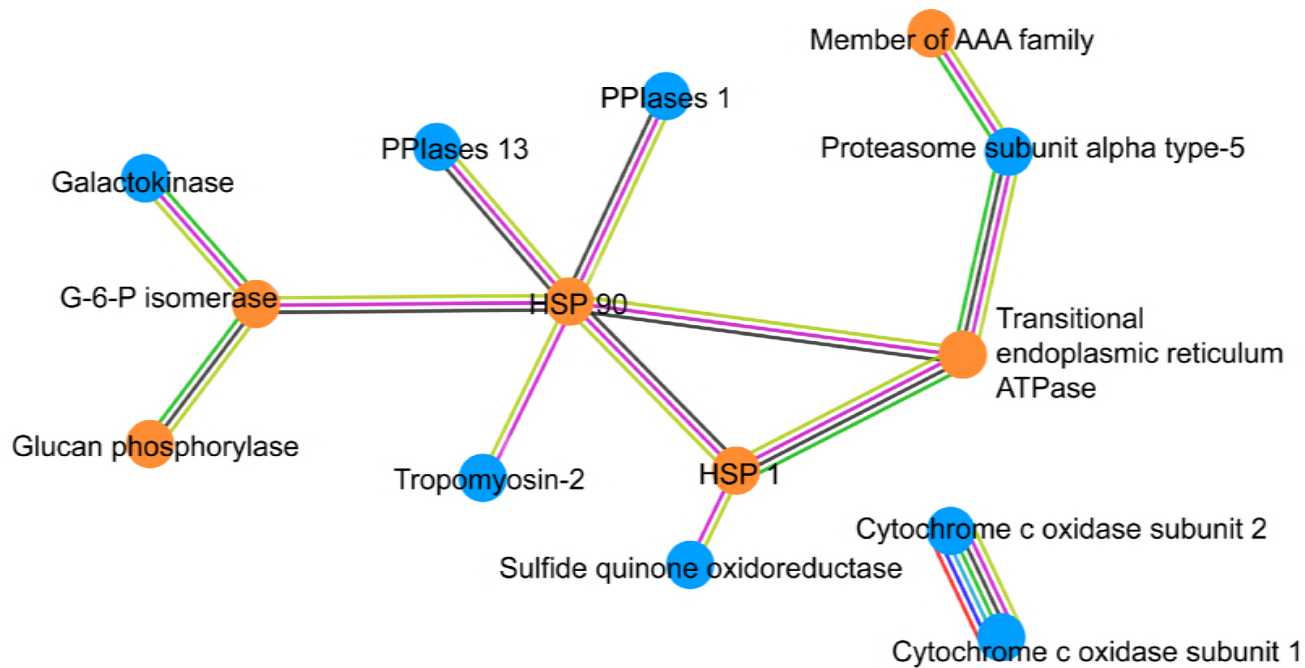


A



B





Node Color

- Conserved genes
- ESP

Known Interactions

- from curated databases
- experimentally determined

Predicted Interactions

- gene neighborhood
- gene fusions
- gene co-occurrence

Others

- textmining
- co-expression
- protein homology