

RESEARCH ARTICLE

High-quality reference genome of *Fasciola gigantica*: Insights into the genomic signatures of transposon-mediated evolution and specific parasitic adaptation in tropical regions

Xier Luo^{1,2}, Kuiqing Cui², Zhiqiang Wang², Zhipeng Li², Zhengjiao Wu², Weiyi Huang², Xing-Quan Zhu³, Jue Ruan^{1,2*}, Weiyu Zhang^{2*}, Qingyou Liu^{1,2*}

1 Shenzhen Branch, Guangdong Laboratory of Lingnan Modern Agriculture, Genome Analysis Laboratory of the Ministry of Agriculture and Rural Affairs, Agricultural Genomics Institute at Shenzhen, Chinese Academy of Agricultural Sciences, Shenzhen, China, **2** State Key Laboratory for Conservation and Utilization of Subtropical Agro-bioresources, Guangxi University, Nanning, China, **3** College of Veterinary Medicine, Shanxi Agricultural University, Taigu, China

 These authors contributed equally to this work.

* ruanjue@caas.cn (JR); zweiyu@gxu.edu.cn (WZ); qyliu-gene@gxu.edu.cn (QL)



OPEN ACCESS

Citation: Luo X, Cui K, Wang Z, Li Z, Wu Z, Huang W, et al. (2021) High-quality reference genome of *Fasciola gigantica*: Insights into the genomic signatures of transposon-mediated evolution and specific parasitic adaptation in tropical regions. *PLoS Negl Trop Dis* 15(10): e0009750. <https://doi.org/10.1371/journal.pntd.0009750>

Editor: Neil David Young, The University of Melbourne, AUSTRALIA

Received: April 12, 2021

Accepted: August 23, 2021

Published: October 5, 2021

Peer Review History: PLOS recognizes the benefits of transparency in the peer review process; therefore, we enable the publication of all of the content of peer review and author responses alongside final, published articles. The editorial history of this article is available here: <https://doi.org/10.1371/journal.pntd.0009750>

Copyright: © 2021 Luo et al. This is an open access article distributed under the terms of the [Creative Commons Attribution License](https://creativecommons.org/licenses/by/4.0/), which permits unrestricted use, distribution, and reproduction in any medium, provided the original author and source are credited.

Data Availability Statement: The whole genome assembly (contig version) and gene annotation reported in this paper have been deposited in the

Abstract

Fasciola gigantica and *Fasciola hepatica* are causative pathogens of fascioliasis, with the widest latitudinal, longitudinal, and altitudinal distribution; however, among parasites, they have the largest sequenced genomes, hindering genomic research. In the present study, we used various sequencing and assembly technologies to generate a new high-quality *Fasciola gigantica* reference genome. We improved the integration of gene structure prediction, and identified two independent transposable element expansion events contributing to (1) the speciation between *Fasciola* and *Fasciolopsis* during the Cretaceous-Paleogene boundary mass extinction, and (2) the habitat switch to the liver during the Paleocene-Eocene Thermal Maximum, accompanied by gene length increment. Long interspersed element (LINE) duplication contributed to the second transposon-mediated alteration, showing an obvious trend of insertion into gene regions, regardless of strong purifying effect. Gene ontology analysis of genes with long LINE insertions identified membrane-associated and vesicle secretion process proteins, further implicating the functional alteration of the gene network. We identified 852 predicted excretory/secretory proteins and 3300 protein-protein interactions between *Fasciola gigantica* and its host. Among them, copper/zinc superoxide dismutase genes, with specific gene copy number variations, might play a central role in the phase I detoxification process. Analysis of 559 single-copy orthologs suggested that *Fasciola gigantica* and *Fasciola hepatica* diverged at 11.8 Ma near the Middle and Late Miocene Epoch boundary. We identified 98 rapidly evolving gene families, including actin and aquaporin, which might explain the large body size and the parasitic adaptive character resulting in these liver flukes becoming epidemic in tropical and subtropical regions.

Genome Warehouse in BIG Data Center, Beijing Institute of Genomics (China National Center for Bioinformatics), Chinese Academy of Sciences, under accession number GWHAZTT00000000 that is publicly accessible at <https://ngdc.cnbc.ac.cn/gwh/>. The AGP file for Hi-C was uploaded as supplement file. The Pacbio sequencing reads has been deposited into the genome sequence archive (GSA) in BIG under accession code CRA003783. The whole genome assembly also can be obtained in the National Center for Biotechnology Information (NCBI) under Bioproject PRJNA691688.

Funding: This work was supported by the National Natural Science Fund (grants No. U20A2051, and 31860638) to QL, National Natural Science Fund (grants No. 31760648) to KC, National Natural Science Fund (grants No. 31960706) to WZ, Guangxi Natural Science Foundation (grants No. AB18221120) to QL, Guangxi Distinguished Scholars Program (grants No. 201835) to QL, and Science and Technology Major Project of Guangxi (grants No. Guike AA17204057) to WZ. The funders had no role in study design, data collection and analysis, decision to publish, or preparation of the manuscript.

Competing interests: The authors have declared that no competing interests exist.

Author summary

Fascioliasis is a neglected zoonotic tropical disease of humans, which reduces the productivity of animal industries, and imposes an economic burden of at least 3.2 billion dollars annually. Although there are four assemblies for *F. hepatica* and two assemblies for *F. gigantica* at NCBI, the inherent limited ability of short reads based assemblies made the completeness of genome sequences and the quality of gene annotation challenging. Here, we report the Pacbio assembly of reference genome for *F. gigantica*, and the quality of assembly and gene annotation are significantly improved compared with previous assemblies. Besides, we found the evidence of transposon-mediated evolution, especially for LINE insertions into gene regions between 41 Ma and 62 Ma, contributing to the speciation and adaptation of the *Fasciola* ancestors. Furthermore, we identified *F. gigantica* specific gene duplication including 98 gene families, and 3300 protein-protein interactions between *F. gigantica* and the host in the small intestine and liver environment. These results illustrate the genomic and gene evolution of *F. gigantica* potentially shaping multiple parasitic characters.

Introduction

Fasciola gigantica and *Fasciola hepatica*, known as liver flukes, are two species in the genus *Fasciola*, which cause fascioliasis commonly in domestic and wild ruminants, but also are causal agents of fascioliasis in humans. Fascioliasis reduces the productivity of animal industries, imposes an economic burden of at least 3.2 billion dollars annually worldwide, and is a neglected zoonotic tropical disease of humans, according the World Health Organization's list [1]. *F. gigantica*, the major fluke infecting ruminants in Asia and Africa, has been a serious threat to the farming of domesticated animals, such as cows and buffaloes, and dramatically reduces their feed conversion efficiency and reproduction [2]. The prevalence of *F. gigantica* infection has greatly affected subsistence farmers, who have limited resources to treat their herds, and has hindered economic development and health levels, especially in developing countries.

The various omics technologies provide powerful tools to advance our understanding of the molecules that act at the host-parasite interface, and allow the identification of new therapeutic targets against fascioliasis [3]. To date, four assemblies for *F. hepatica* and two assemblies for *F. gigantica* have been deposited at the NCBI [4–7]. These assemblies reveal a large genome with a high percentage of repeat regions in *Fasciola* species, and provided valuable insights into features of adaptation and evolution. However, these assemblies are based on the short read Illumina sequencing or hybrid sequencing methods, with limited ability to span large families of repeats. Various limitations have led to the current assemblies in the genus *Fasciola* being fragmented (8 kb to 33 kb and 128 kb to 1.9 Mb for contig and scaffold N50s, respectively). Subsequent gene annotation analysis using current assemblies were also challenging, with abundant transposition events occurring over evolutionary history, which significantly increased the repeat components in intron regions, resulting in considerable fragmentation in gene annotation.

Infection by *Fasciola* causes extensive damage to the liver, and excretory/secretory (E/S) proteins play an important role in host-parasite interactions. Parasite-derived molecules interact with proteins from the host cell to generate a protein interaction network, and these proteins partly contribute to *Fasciola*'s striking ability to avoid and modulate the host's immune response [8]. Previous proteomics of E/S proteins have highlighted the importance of secreted

extracellular vesicles (EVs) and detoxification enzymes to modulate host immunity by internalizing with host immune cells [9,10]. The anthelmintic drug, triclabendazole (TCBZ), is currently the major drug available to treat fascioliasis at the early and adult stages, which acts by disrupting β -tubulin polymerization [11]; however, over-reliance on TCBZ to treat domesticated ruminants has resulted in selection for resistance to liver flukes [12]. Drug and vaccine targets for molecules associated with reactive oxygen species (ROS)-mediated apoptosis have recently been validated as an effective tools in multiple helminth parasites [13]. Increased understanding of host-parasite and drug-parasite interactions would facilitate the development of novel strategies to control fascioliasis.

In recent years, there have been increasing numbers of human cases of fascioliasis, becoming a major public health concern in many regions [14,15]. However, high quality genome assemblies for liver flukes are still insufficient. In the present study, we combined multiple sequencing technologies to assemble a chromosome-level genome for *F. gigantica* and provided integrated gene annotation. Protein-protein interactions were analyzed between the predicted *F. gigantica* secretome and host proteins expressed in the small intestine and liver. In addition, gene family analysis identified a series of genes expansions in *F. gigantica*. Interestingly, the distribution of repeat sequences in the genome exhibit an excess of long interspersed element (LINE) duplications inserted into intronic regions, potentially helping to explain the duplications of transposable element (TE) plasticizing gene structures and possibly acting as long-term agents in the speciation of *Fasciola*.

Results

Pacbio long reads-based *de novo* assembly and gene annotation

The *F. gigantica* genome contains abundant repeat sequences that are difficult to span using short read assembly methods, and the complex regions also hinder integrated gene annotation of the genome. Therefore, in the present study, multiple sequencing technologies, have been applied: (1) Single-molecule sequencing long reads ($\sim 91\times$ depth) using the Pacbio Sequel II platform; (2) paired-end reads ($\sim 66\times$ depth) using the Illumina platform; and (3) chromosome conformation capture sequencing (Hi-C) data ($\sim 100\times$ depth) (S1 Table). The initial assembly was performed using the Pacbio long reads, followed by mapping using single-molecule sequencing and Illumina sequencing reads to polish assembly errors and sequencing mistakes, resulting in a contig N50 size of 4.89 Mb (Fig 1A). The Hi-C data were used to build final super-scaffolds, resulting in a total length of 1.35 Gb with a scaffold N50 size of 133 Mb (Fig 1B and S1 Fig and Table 1 and S1–S3 Tables). The final assembly consists of 10 pseudo-chromosomes covering more than 99.9% of the *F. gigantica* genome, and the length distribution was approximate equal to the estimation by karyotype in previous research (S2 Fig and S4 Table) [16]. The assessment of nucleotide accuracy shows that the error rate was 5.7×10^{-6} in the genome. QUILT analysis [17] showed a high mapping and coverage rate using both Illumina short reads and Pacbio long reads, in which 99.73% of reads mapped to 99.85% of the genome with more than $10\times$ depth (S5 Table).

Combing *de novo*/homolog/RNA-seq prediction, a total of 12,503 protein coding genes were annotated in the *F. gigantica* genome. BUSCO assessment [18] indicated that the genome is 90.4% complete and 5.6% fragmented, underscoring the significant improvement of the genome continuity and gene-structure predictions compared with previous assemblies (S6 Table). Specifically, the average gene length in the annotated data is 28.8 kb, nearly twice the length of that in other digenean species, but contrasted with the similar average length of the coding sequences (CDSs). Through functional annotation, we found that 8569 of the genes could be characterized in the InterPro database [19,20], 7892 of them were mapped to the gene

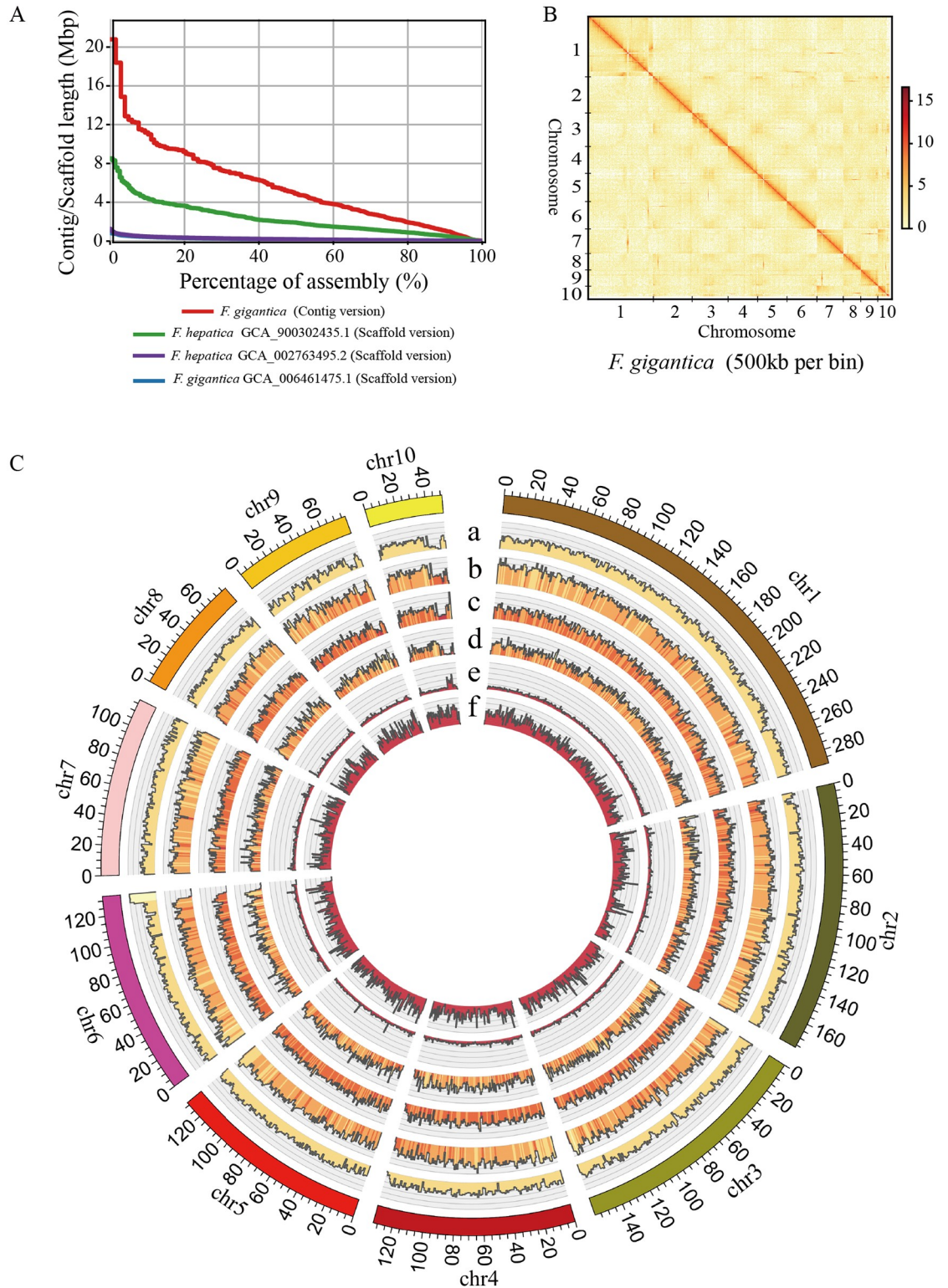


Fig 1. Landscape of the *Fasciola gigantica* genome. (A) Comparisons of the assembled contigs and scaffold lengths (y-axis) and tallies (x-axis) in *Fasciola* species. (B) Hi-C interactive heatmap of the genome-wide organization. The effective mapping read pairs between two bins were used as a signal of the strength of the interaction between the two bins. (C) Integration of genomic and annotation data using 1 Mb bins in 10 Hi-C assembled chromosomes. (a) Distribution of the GC content (GC content > 39% and < 52%); (b) distribution of the long interspersed element (LINE) percentage > 0% and < 50%; (c) distribution of the long

terminal repeat (LTR) percentage > 0% and < 50%; (d) distribution of the gene percentage > 0% and < 70%; (e) distribution of the heterozygosity density of our sample (percentage > 0% and < 1%); (f) distribution of the heterozygosity density of SAMN03459319 in the NCBI database. Hi-C, chromosome conformation capture sequencing;

<https://doi.org/10.1371/journal.pntd.0009750.g001>

ontology (GO) terms, and 5353 of them were identified by the Kyoto Encyclopedia of Genes and Genomes (KEGG) pathways database (S3–S4 Figs and S7 Table).

The unique repeat duplications in *Fasciola*

TEs are insertional mutagens and major drivers of genome evolution in eukaryotes, and replication of these sequences, resulting in variation of gene structure and expression, have been extensively documented [21,22]. Besides, TEs are molecular fossils, being remnants of past mobilization waves that occurred millions of years ago [23]. In the present study, we identified repeat sequences combined the analysis from RepeatModeler [24] and RepeatMasker [25], and detected a significant proportion of them neglected by previous studies. In the *F. gigantica* genome, we identified 945 Mb of repeat sequences, which was approximate 20% more than that identified in other assemblies in *Fasciola* species, while the lengths of non-repeat sequences were nearly identical. The most convincing explanation for the additional assembled repeat sequences was that the contigs constructed from Pacbio long reads spanned longer repeat regions, which were compressed in previous assemblies. Among these repeat sequences, there were 408 Mb of LINEs (corresponding to 30.3% of the assembled genome), 285 Mb of long terminal repeats (LTRs, corresponding to 21.2% of the assembled genome), and 162 Mb of unclassified interspersed repeats (corresponding to 12.0% of the assembled genome) (S5 Fig and S8 Table). According to the repeat landscapes, we found that there were two shared expansion events for LINEs and LTRs that occurred approximately 12 million years ago (Ma) and 65 Ma, and an additional expansion event at 33 Ma for LTRs (S6–S7 Figs). Our result confirmed previous study on family Fasciolidae [6], and the abundant repeat sequences in the *Fasciola* genomes aroused the interest concerning the role of repeats in evolution (Fig 2A), which implied a hypothesize that the expansion of TEs enlarged the genome size of an ancestor of

Table 1. Summary statistics for the genome sequences and annotation.

<i>F. gigantica</i>		
Genome	Total Genome Size (Mb)	1,348
	Chromosome Number	10
	Scaffold Number ^a	10+24
	Scaffold N50 (Mb)	133
	Scaffold L50	4
	Contig Number	1,022
	Contig N50 (Mb)	4.89
	Heterozygosity Rate (%)	1.9×10^{-3}
Annotation	Total Gene Number	12,503
	Average CDS Length (bp)	1552.7
	Average Gene Length (kb)	28.8
	Percentage of Genome Covered by CDSs (%)	1.5%
	BUSCO Assessment	90.4%
	Repeat Content	70.0%

^a number of chromosome level scaffolds and unplaced scaffolds. CDS, coding sequence.

<https://doi.org/10.1371/journal.pntd.0009750.t001>

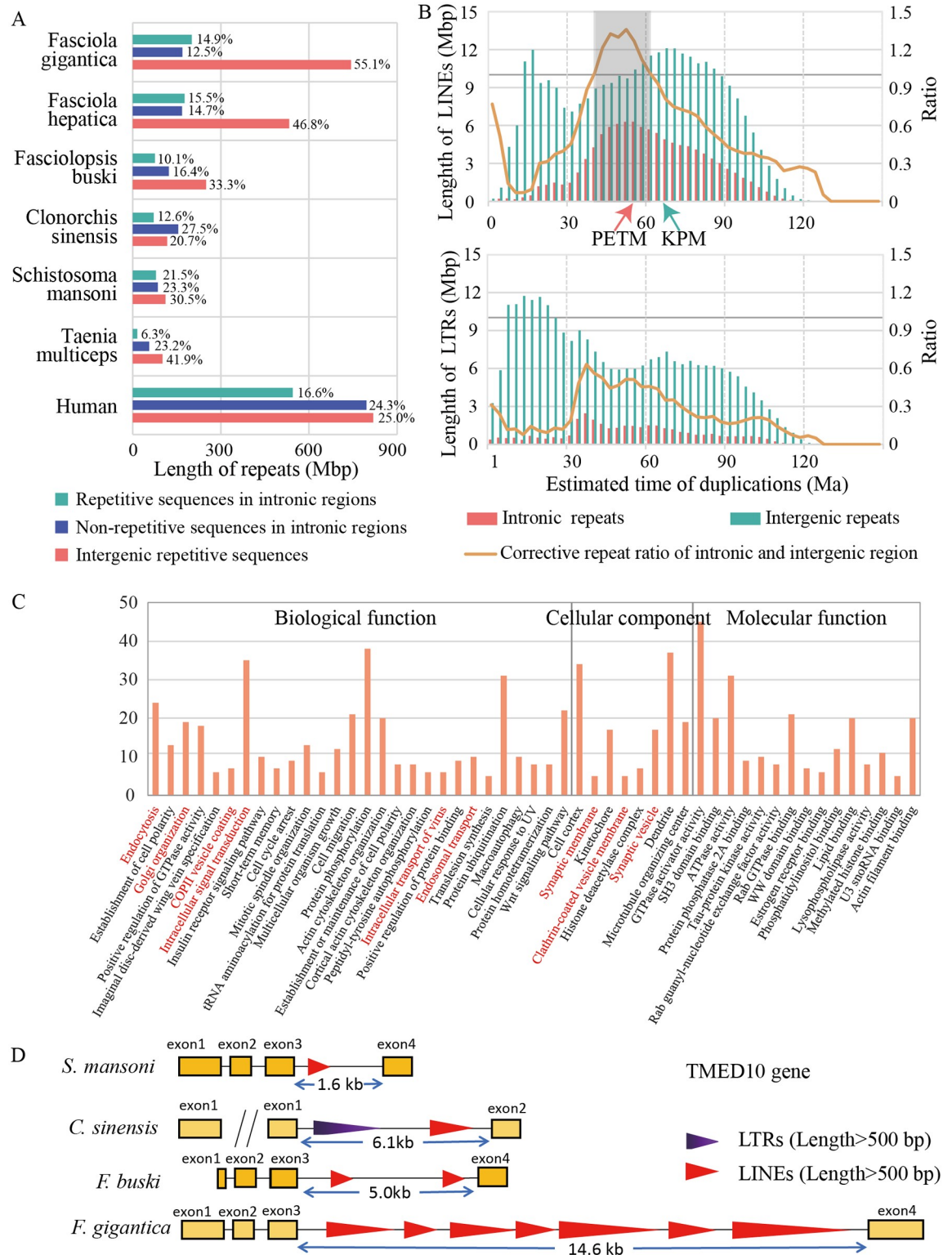


Fig 2. Identification of repeat expansion and alternative gene networks in the *Fasciola gigantica* genome. (A) The distribution of repetitive sequence length among the genomes of six flatworms and the human genome. (B) Landscape of LINES and LTRs distribution in the *Fasciola gigantica* genome. The x-axis shows the expansion time of TEs calculated by the divergence between repeat sequences. The mutation rate was set as 1.73×10^{-9} per year. The orange line represents the repeat length ratio, used to estimate the signatures of selection, which was corrected by the total length of intronic and intergenic regions in history. (C) The functional enrichment of genes

with more than 10 kb LINE insertions between 41 Ma and 62 Ma by Gene Ontology (GO) classification. The GO terms related to vesicle secretion are marked in red. (D) *TMED10* gene structure map. LINES original between 41 Ma and 62 Ma and longer than 500 bp identified by RepeatMasker were plotted. LTRs longer than 500 bp were plotted. Long interspersed element, LINE; long terminal repeat, LTR; TE, transposable element; TMED10, transmembrane P24 trafficking protein 10.

<https://doi.org/10.1371/journal.pntd.0009750.g002>

Fasciola to gain a new advantage by rewiring gene networks. To test this hypothesis, we focused on the genome-wide repeats distribution and test signatures of selection.

For new TE insertions to persist through vertical inheritance, transposition events must be under strong purifying effect among gene loci to avoid disturbing their biological function. However, we observed many intronic repeat elements in *Fasciola*, resulting in a larger intron size per gene. If there are equal selection effects on newly inserted TEs in intronic and intergenic regions, there would be a high correlation between the distribution of insertion time and retained TE lengths between these two regions. By contrast, there would be fewer accumulated repeat sequences existing under purifying effect. In this study, we use the relative proportion of TEs between intronic and intergenic regions as a simple indicator, and use the inferred size of intronic and intergenic regions over evolutionary history as a control to estimate the signatures of selection. The results showed that TE insertions into intronic regions are under persistent intense purifying effect, except for LINES. There was an excess of persistent LINE insertions into intronic regions between 41 Ma and 62 Ma, indicating different modes of accumulating LINES into intronic regions compared with that in other periods (Fig 2B). Specifically, the time of the ancient intronic LINE expansion (~51.5 Ma) was different to the genome-wide LINE expansion time (~68.0 Ma), whereas the time was coincident with two important environmental change events, the Cretaceous-Paleogene boundary (KPB) mass extinction (~66.0 Ma) and the Paleocene-Eocene Thermal Maximum (PETM) (~55.8 Ma). Both the PETM and KPB events recorded extreme and rapid warming climate changes; however, rapid evolutionary diversification followed the PETM event, as opposed to near total mass extinction at the KPB [26]. Therefore, we selected genes with different LINE lengths, derived between 41 Ma and 62 Ma, and expected to identify a transposon-mediated alternative gene network contributing to the host switch and the shift from intestinal to hepatic habitats.

LINE-mediated alternative gene network

We identified a substantial proportion of genes with LINE insertions, derived between 41 Ma and 62 Ma, indicating a universal effect of the gene network. We selected 1288 genes with the LINE insertions of more than 10 kb, representing more than one third of the average gene length, and annotated the genes using Gene Ontology (GO) terms and processes and Kyoto encyclopedia of genes and genomes (KEGG) pathways (Fig 2C and S9–S11 Tables). These genes involve molecules internalizing substances from their external environment, including membrane-associated and vesicle secretion process proteins. Meanwhile, the gene network was likely adapted to the evolution of protein biosynthesis and modification of histones.

Enrichment analysis of GO terms showed that membrane and membrane-associated proteins are over-represented, involving “synaptic membrane” ($P = 3.52E-04$), “clathrin-coated vesicle membrane” ($P = 1.08E-03$), and “synaptic vesicle” ($P = 3.02E-03$), as well as vesicles secretion processes, such as “endocytosis” ($P = 7.06E-06$), “Golgi organization” ($P = 7.45E-05$), “COPII vesicle coating” ($P = 2.72E-04$), “intracellular signal transduction” ($P = 5.16E-04$), and “endosomal transport” ($P = 2.47E-03$). Besides, proteins relating to phosphorylation and GTPase activators were also enriched, such as “Protein phosphorylation” ($p = 1.73E-03$), “Regulation of small GTPase mediated signal transduction” ($P = 1.21E-03$), “GTPase activator activity” ($P = 1.13E-10$). The over-representation of genes involved in membrane transport and

signal transduction was particularly interesting because helminth parasites interfere with the host immune system by secreting molecules from surface tegument or gut. The *TMED10* gene in *F. gigantica* (encoding transmembrane P24 trafficking protein 10) was used as an example. *TMED10* is a cargo receptor involved in protein vesicular trafficking along the secretory pathway [27,28], and the genes have an 11.1 kb LINE insertion in the third intron, resulting in an over three-fold increment in the gene length (Fig 2D). The enrichment suggests that the gene network related to secretion could have experienced adaptive evolution during LINE transposition events. We further compared our dataset with the proteome result from *F. hepatica* extracellular vesicles (EVs) [9], and found 21 proteins that were also identified as surface molecules associated with EV biogenesis and vesicle trafficking (*IST1*, *VPS4B*, *TSG101*, *MYOF*, *ATG2B*, *STXBP5L*, and 15 Rho GTPase-activating related proteins). Specifically, *IST1*, *VPS4B*, and *TSG101* are members of the endosomal sorting complex required for transport (ESCRT) pathway, which promotes the budding and release of EVs. *TSG101*, a crucial member of the ESCRT-I complex, has an important role in mediating the biogenesis of multi-vesicular bodies, cargo degradation, and recycling of membrane receptors. Besides, the ESCRT pathway promotes the formation of both exosomal carriers for immune communication. During the formation of the immunological synapse between T-cells and antigen-presenting B cells, *TSG101* ensures the ubiquitin-dependent sorting of T-Cell Receptor (*TCR*) molecules to exosomes that undergo *VPS4*-dependent release into the synaptic cleft [29].

The most significant KEGG pathway was aminoacyl-tRNA biosynthesis ($P = 7.16E-04$), containing 15 out of 38 annotated aminoacyl tRNA synthetases (*AARSs*). *AARSs* are the enzymes that catalyze the aminoacylation reaction by covalently linking an amino acid to its cognate tRNA in the first step of protein translation. The large-scale insertion of LINEs reside in *AARS* genes suggested that the ancestor of *Fasciola* may have profited from the effect of transposition, with changes to protein biosynthesis and several metabolic pathways for cell viability. In addition, a significant number of genes are strongly associated with histone modulation, including “histone deacetylase complex” ($P = 1.89E-03$), “histone methyltransferase activity (H3-K36 specific)” ($P = 1.08E-03$), and “methylated histone binding” ($P = 2.37E-03$). Histone modifications play fundamental roles in the manipulation and expression of DNA. We found nine histone deacetylases and Histone methyltransferases in the gene set (*HDAC4*, *HDAC8*, *HDAC10*, *KMT2E*, *KMT2H*, *KMT3A*, *KDM8*, *NSD1*, and *NSD3*). Histone modifications can exert their effects by influencing the overall structure of chromatin and modifying and regulating the binding of effector molecules [30,31]; therefore, the variation of these genes might bring about evolution from a disturbed gene structure to a mechanism of genome stabilization to tackle a continuous genome amplification process in evolutionary history.

Genome-wide host-parasite interaction analysis

In the *Fasciola* genome, we predicted genes encoding 268 proteases, 36 protease inhibitors (PIs), and 852 predicted excretory/secretory (E/S) proteins that are commonly involved in interacting with hosts and modulating host immune responses (S8 Fig). The largest class of proteases was cysteine peptidases ($n = 113$), which was also identified in the *F. hepatica* genome (Fig 3A and S12 Table). The largest ($n = 19$, 52.8% of PIs) PI family was the I02 family of Kunitz-BPTI serine protease inhibitors, which bind to Cathepsin L with a possible immunoregulatory function [32] (S13 Table). GO enrichment analysis of E/S proteins showed that proteins related to “activation of cysteine-type endopeptidase activity” ($P = 6.14E-19$), “peroxidase activity” ($P = 3.79E-07$) and “protein disulfide isomerase activity” ($P = 3.75E-06$) are over-represented (Fig 3B, S14–S15 Tables). Indeed, there were 38 cysteine peptidases identified as E/S

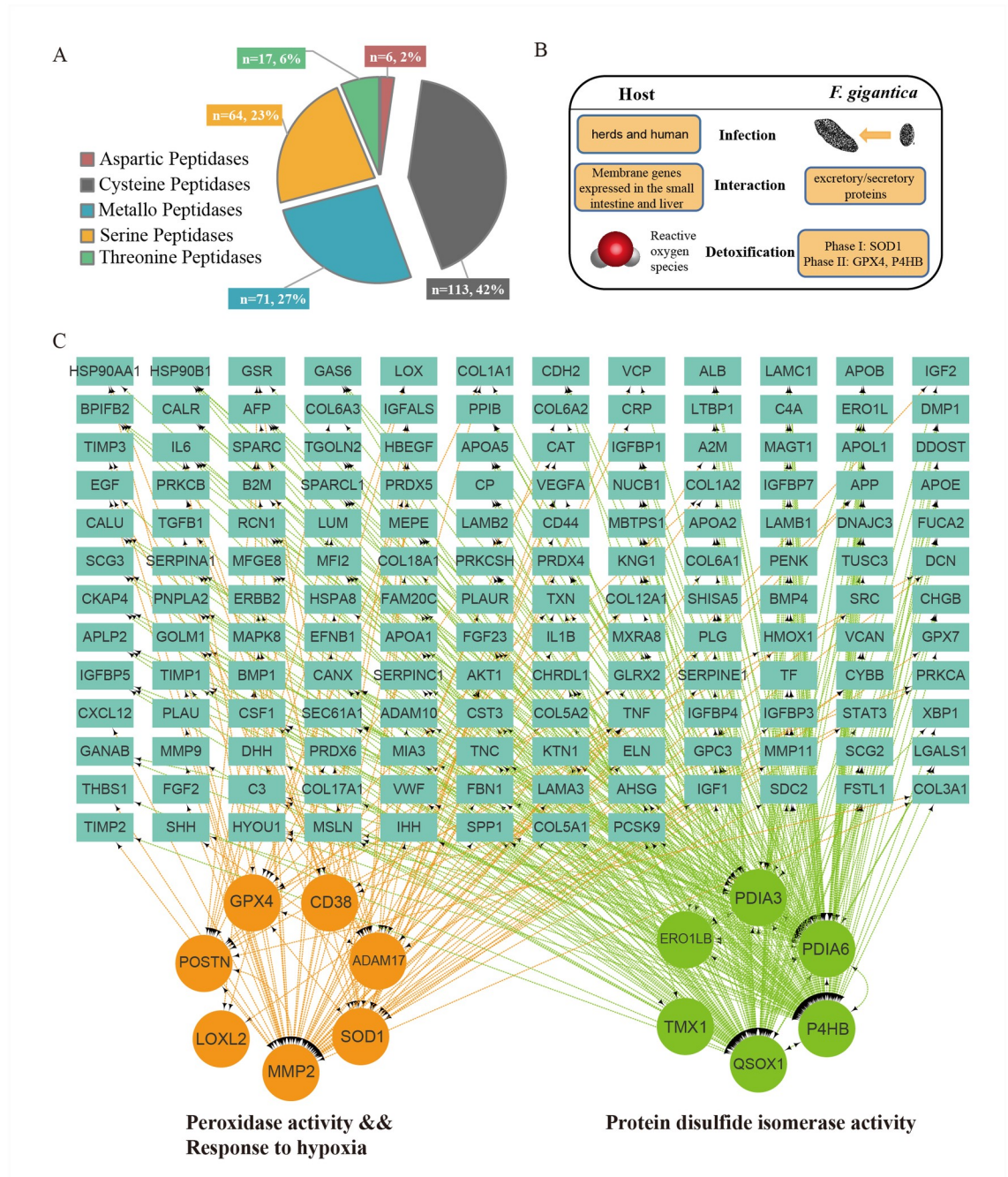


Fig 3. Genome-wide host-parasite interaction analysis. (A) Pie chart for proteases identified in *Fasciola gigantica*. (B) The interaction mode between the adult *Fasciola gigantica* and the host. (C) The protein-protein interaction (PPI) network of redox-related pathways in *Fasciola gigantica* with host proteins. The genes indicated in the three gene ontology (GO) terms were significantly enriched and have their encoded proteins have PPIs with excretory/secretory (E/S) proteins.

<https://doi.org/10.1371/journal.pntd.0009750.g003>

proteins, including cathepsin L-like, cathepsin B-like, and legumain proteins, which participate in excystment, migration through gut wall, and immune evasion [33].

In parasites, as in mammalian cells, ROS are produced as a by-product of cell metabolism and from the metabolism of certain pharmacological agents. The ability of a parasite to survive

in its host has been directly related to its antioxidant enzyme content [34]. To further analyze host-parasite interactions, we identified the protein-protein interactions (PPIs) between the *F. gigantica* secretome and human proteins expressed in the small intestine and liver [8]. In total, we identified 3300 PPIs, including rich interactions that directly or indirectly participated in the two phases of detoxification pathways (Fig 3C). Superoxide dismutase [Cu-Zn] (*SOD*, PPIs = 49) was first highlighted because of its important role on phase I detoxification against ROS, in which it catalyzes the dismutation of the superoxide radical to molecular oxygen and hydrogen peroxide (H_2O_2) [35]. Gene family analysis identified six *SOD* paralogs in *F. gigantica*, and two of them contained a signal peptide (Fig 4D). Previous enzyme activity assays also confirmed a significant difference between *SOD* activities and concentration in E/S proteins of two *Fasciola* species [36], suggesting an intense ability to resist superoxide radical toxicity. Meanwhile, the metabolite of phase I, H_2O_2 , can also damage parasites, which requires detoxification enzymes, including glutathione-dependent enzymes *GPx*, glutathione reductase, and other peroxidases. Protein disulfide-isomerase (*P4HB*, PPIs = 132) and phospholipid hydroperoxide glutathione peroxidase (*GPX4*, PPIs = 28) were as functioning in phase II detoxification. *GPx* catalyzes the reduction of hydroperoxides (ROOH) to water, using glutathione (*GSH*) as the reductant. *P4HB* also participates in the process by mediating homeostasis of the antioxidant glutathione [37]. However, we did not identify E/S proteins in the Cytochrome P450 (*CYP450*) family in phase III detoxification. Therefore, we speculated that successful parasite defense against *F. gigantica* mainly depends on the strong superoxide activity and efficient hydrogen peroxide detoxification.

Gene family analysis

Gene family analysis was performed using eight taxa (*F. gigantica*, *F. hepatica*, *Fasciolopsis buski* [38], *Clonorchis sinensis* [39], *Schistosoma mansoni* [40], *Taenia multiceps* [41], swamp buffalo [42], and human [43], which identified 17,992 gene families (Fig 4A). Phylogeny analysis of 559 single-copy orthologs showed that *F. gigantica* and *F. hepatica* shared a common ancestor approximately 11.8 million years ago (2.2–22.5 Ma, 95% highest posterior density [HPD]) near the Middle and Late Miocene Epoch boundary. The Miocene warming began 21 million years ago and continued until 14 million years ago, when global temperatures took a sharp drop at the Middle Miocene Climate Transition (MMCT). The divergence of the two *Fasciola* species may have resulted from the consequences of rapid climate changes, such as migration of the host causing geographic isolation. Our estimation is between the previously suggested date of 5.3 Ma based on 30 nuclear protein-coding genes [6], and 19 Ma based on cathepsin L-like cysteine proteases [44]. Although we used a more integrative gene dataset, the wide HPD interval could not be neglected, raising possible uncertainty from the complex process of speciation or inappropriate protein sequence alignment between members of the genus *Fasciola*.

The distribution of gene family size among different species is used to estimate which lineages underwent significant contractions or expansions. Compared with *F. hepatica*, *F. gigantica* shows more gene family expansion events (643 compared to 449) and a similar number of gene family contractions (713 compared to 672). The result emphasizes the general trend that, relative to the common ancestor of *Fasciola*, the ancestor of *F. gigantica* apparently underwent a higher extent of gene-expansion than did the ancestor of *F. hepatica*. Gene duplication is one of the primary contributors to the acquisition of new functions and physiology [45]. We identified 98 gene families, including 629 genes, as rapidly evolving families specific to *F. gigantica*. Family analysis showed a fascinating trend of gene duplication, with substantial enrichment for the “structural constituent of cytoskeleton” ($P = 3.52E-24$), “sarcomere organization”

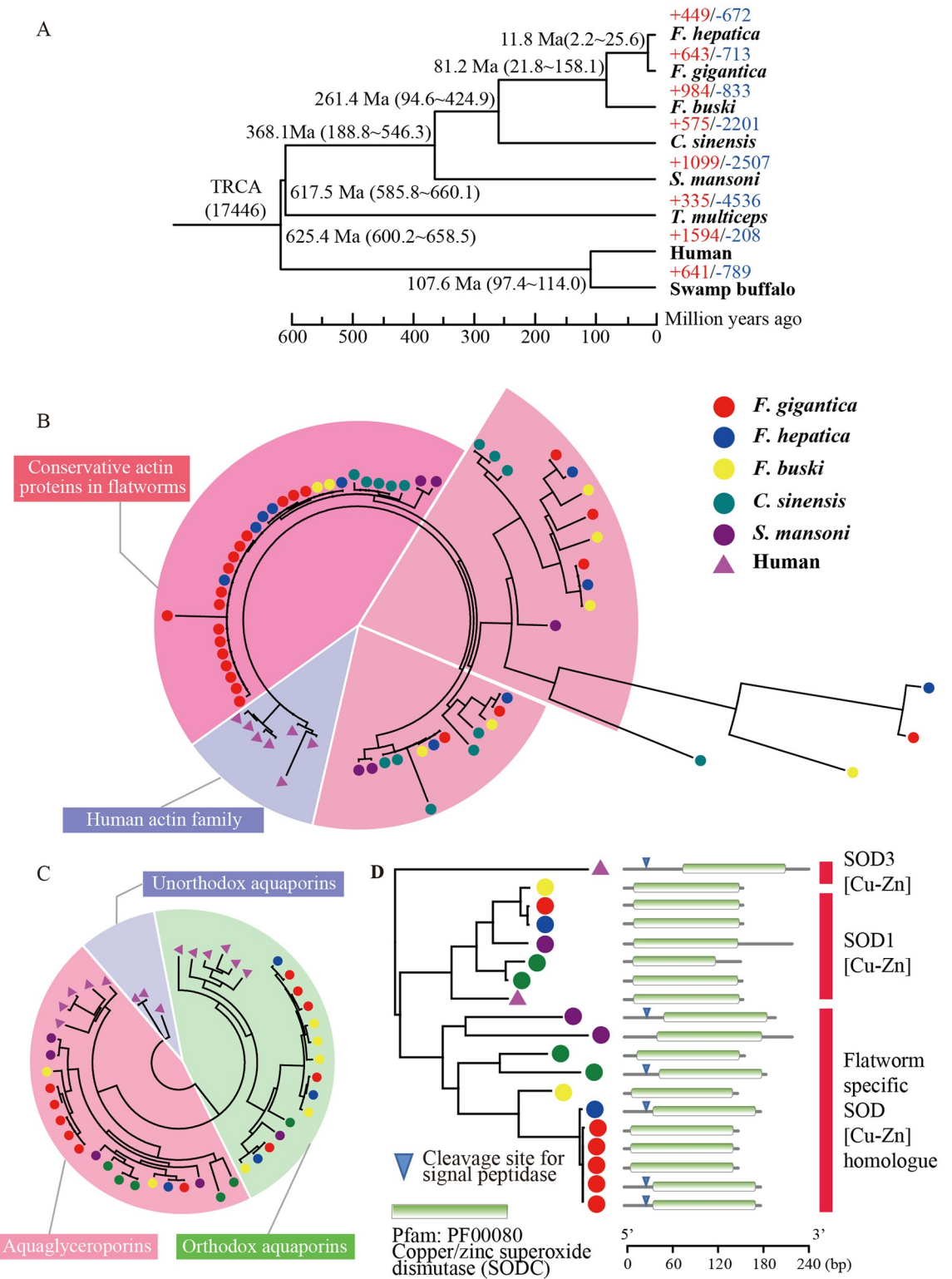


Fig 4. Phylogenetic tree and gene family analysis. (A) A phylogenetic tree generated using 559 single-copy orthologous genes. The numbers on the species names are the expanded (+) and contracted (-) gene families. The numbers on the nodes are the divergence time between species. (B) A phylogenetic tree of actin genes in flatworms and humans. All human homologue genes are selected as outgroup. (C) Phylogenetic tree of aquaglyceroporin (AQP) family genes in flatworms and humans. The human homologue genes (*AQP11*, *AQP12A*, and *AQP12B*) were selected as the outgroup. (D) A phylogenetic tree of copper/zinc superoxide dismutase (SOD) genes in flatworms and humans. The midpoint was selected as the root node.

<https://doi.org/10.1371/journal.pntd.0009750.g004>

($P = 2.29E-14$), “actin filament capping” ($P = 6.19E-13$), and “spectrin” ($P = 3.03E-11$) in *F. gigantica* (S16 Table). There were 24 actin paralogs in *F. gigantica*, in contrast to 8 actin paralogs in *F. hepatica*. Actin is one of the most abundant proteins in most cells, and actin filaments, one of the three major cytoskeletal polymers, provide structure and support internal movements of organisms [46]. They are also highly conserved, varying by only a few amino acids between algae, amoeba, fungi, and animals [47]. We observed three types of actin proteins in flukes, according to their identity from human actin family. Seventeen of the 24 actin proteins in *F. gigantica* are highly conserved (Identity > 95%) (Fig 4B). Consistent with the accepted role of the epidermal actin cytoskeleton in embryonic elongation [48,49], we speculated that the significant expansion of actin and spectrin genes increased the body size of *F. gigantica* via cell elongation or proliferation during morphogenesis. Another rapidly evolving family is the aquaglyceroporin subfamily in the membrane water channel family. We found six aquaglyceroporin paralogs in *F. gigantica*, which were over-represented in the GO term “water transport” ($P = 2.10E-06$) (Fig 4C). Aquaglyceroporins are highly permeated by glycerol and other solutes, and variably permeated by water, as functionally validated by several studies [50,51]. The mammalian aquaglyceroporins regulate glycerol content in epidermal, fat, and other tissues, and appear to be involved in skin hydration, cell proliferation, carcinogenesis, and fat metabolism. A previous study showed that *F. gigantica* could withstand a wider range of osmotic pressures compared with *F. hepatica* [52], and we speculated that a higher aquaglyceroporin gene copy number might help explain this observation.

It is worth mentioning that 57.6% of the rapidly evolving expansion genes specific to the *F. gigantica* genome were driven by tandem duplication, such that the newly formed duplicates preserved nearly identical sequences to the original genes. The newly formed genes would accumulate non-functionalizing mutations, or develop new functions over time. We found only few tandem duplicated genes that had non-functionalizing mutations, suggesting that adaptive evolution could have an important role in the consequences of these genes via a dosage effect or neo-functionalization.

Discussion

The genome of *Fasciola* species contains a large percentage of repeat sequences, making them the largest parasite genomes sequenced to date. Since the first assembly of *F. hepatica* was submitted in 2015 [5], several studies have aimed to improve the quality of assembly and gene annotation [4,6,7]. With advances in long read sequencing assembly and Hi-C scaffolding technologies, it is now viable to resolve the genomic “dark matter” of repetitive sequences, and other complex structural regions at relatively low cost [53]. Therefore, we present the highest quality genome and gene annotation for *F. gigantica* to date, and provide long-awaited integrated genome annotation for fascioliasis research.

In previous study of Fasciolidae family, Choi et al. have discovered TE expansion in *Fasciola*, which also explained the large lineage-specific genome size and longer annotated gene [6]. We confirmed the result in *F. gigantica* genome and further identified signatures of selection based on unbalanced distribution of inserted TEs between intronic and intergenic regions in history. Especially, the strongest selection signal occurred in the speciation between the *Fasciola* and Fascioloides—a habitat switch from the small intestine to the liver in the host—during the PETM, which accompanied by LINE expansion biased toward intronic regions (Fig 5). This unexpected event provided a new evidence of adaptive evolution driven by transposition events and will prompt investigations of how such differences contribute mechanistically to the morphological phenotypes of liver flukes and related species. There are also many studies in other species supporting the hypothesis that TE invasions endured by organisms have

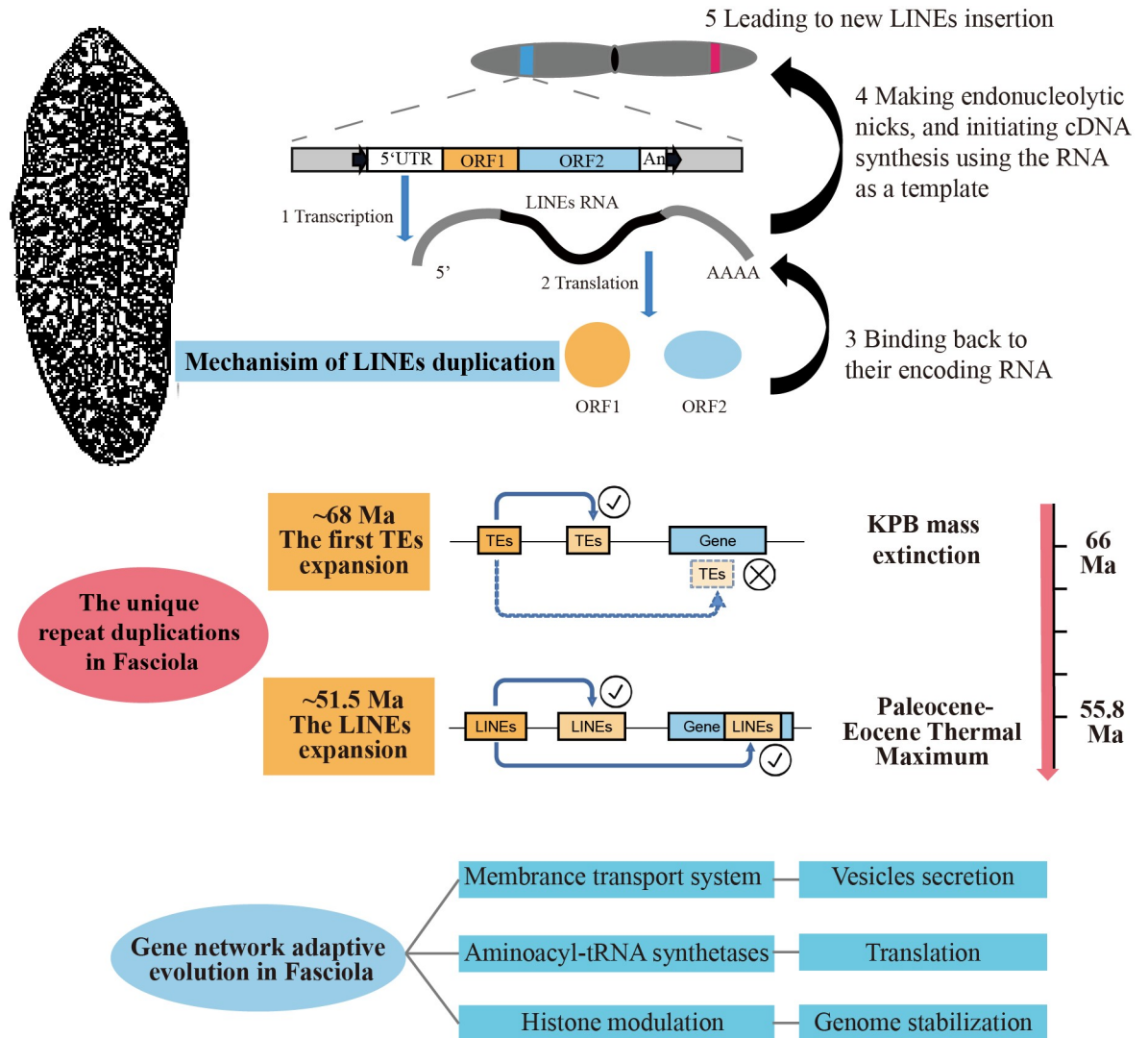


Fig 5. Schematic diagram of the process of *Fasciola*-specific repeat expansion during evolution.

<https://doi.org/10.1371/journal.pntd.0009750.g005>

catalyzed the evolution of gene-regulatory networks [54]. For example, Eutherian-specific TEs have the epigenetic signatures of enhancers, insulators, and repressors, and bind directly to transcription factors that are essential for pregnancy and coordinately regulate gene expression [55]. Similarly, genes with large-scale insertion of TEs in *Fasciola* species identified here, represent a signature of *Fasciola*-specific evolutionary gene network to distinguish other flukes of the family Fasciolidae. These genes overlap significantly with host-parasite interaction genes, including proteases and E/S proteins, and are enriched in the pathways of EV biogenesis and vesicle trafficking.

The data from genomic, transcriptomic, and proteomic studies can form a good complementary relationship to further our understanding of helminth parasites and their interaction with their hosts. Previous studies have identified a rich source of stage-specific molecules of interest using transcriptomic and proteomic analysis [56,57]. Here, we provided a comprehensive list of predicted E/S proteins in *F. gigantica* and predicted 3300 PPIs at the host-parasite interface, extending our understanding of how the phase I and phase II detoxification enzymes

counteract the effect of ROS. The ability of *Fasciola* species to infect and survive in different tissue environments is underpinned by several key E/S protein gene duplications. Both *Fasciola* species have a common expansion in the secretion of papain-like cysteine peptidase family (Clan A, family C1) [5]. Besides, *F. gigantica* has a specific variation in the *SOD* gene copy number, allowing it to regulate the catalytic activity of the superoxide radical released by the host. The effect of specific gene duplications can also be reflect in the increased body size of *F. gigantica*, which is an important morphometric character to distinguish *Fasciola* species and has a decisive influence on the final host species [58], although a gene level study of this phenotype is barely reported.

Overall, our study demonstrated that the combination of long-read sequencing with Hi-C scaffolding produced a very high-quality liver fluke genome assembly and gene annotation. Additionally, identification of the repeat distribution among the gene regions extended our understanding of the evolutionary process in *Fasciola* species. Further detailed functional studies of secretion might be of great scientific significance to explore their potential application in fascioliasis treatment.

Materials and methods

Ethics statement

This study was approved by the Research Ethics Committee of the Guangxi University (Permit code: GXU2019-029). In present study, experiment was performed by the Principle Guidance for the Use and Care of Laboratory Animals.

Sample collection and *de novo* sequencing

All animal work was approved by the Guangxi University Institutional Animal Care and Use Committee. For the reference genome sequencing, one *F. gigantica* at adult stage was derived from infected buffalo in the Guangxi Zhuang Autonomous Region. Nucleic acids were extracted using a QIAGEN DNeasy (DNA) kit (Qiagen Hilden, Germany). Three *de novo* genome sequencing methods were performed on the liver fluke: We generated (1) 122.4 Gb (~88× depth) PacBio Sequel II single-molecule long reads, with an average read length of 15.8 kb (PacBio, Menlo Park, CA, USA); (2) 89.5 Gb (~66× depth) Illumina HiSeq PE150 pair-end sequencing to correct errors (Illumina, San Diego, CA, USA); and (3) 134 Gb (~100× depth) chromosome conformation capture sequencing (Hi-C) data (sequenced by Illumina platform).

De novo assembly and assessment of the genome quality

A PacBio-only assembly was performed using Canu v2.0 [59,60] using new overlapping and assembly algorithms, including an adaptive overlapping strategy based on *tf-idf* weighted Min-Hash and a sparse assembly graph construction that avoids collapsing diverged repeats and haplotypes. To remove haplotigs and contig overlaps in the assembly, we used Purge_Dups based on the read depth [61]. Arrow (<https://github.com/PacificBiosciences/GenomicConsensus>) was initially used to reduce the assembly error in the draft assembly, with an improved consensus model based on a more straightforward hidden Markov model approach. Pilon [62] was used to improve the local base accuracy of the contigs via analysis of the read alignment information based on paired-end bam files (thrice). As a result, the initial assembly resulted had an N50 size of 4.89 Mb for the *F. gigantica* reference genome. ALLHiC was capable of building chromosomal-scale scaffolds for the initial genome using Hi-C paired-end reads containing putative restriction enzyme site information (S1 Text) [63]. The whole

genome assembly (contig version) have been deposited in the Genome Warehouse in BIG Data Center under accession number GWHAZTT00000000 and NCBI under Bioproject PRJNA691688.

Three methods were used to evaluate the quality of the genomes. First, we used QUality ASsessment Tool (QUAST) [64] to align the Illumina and PacBio raw reads to the *F. gigantica* reference genome to estimate the coverage and mapping rate. Second, all the Illumina paired-end reads were mapped to the final genome using BWA [65], and single nucleotide polymorphisms (SNPs) were called using Samtools and Bcftools. The predicted error rate was calculated by the homozygous substitutions divided by length of the whole genome, which included the discrepancy between assembly and sequencing data. Thirdly, we assessed the completeness of the genome assemblies and annotated the genes using BUSCO [18].

Genome annotation

Three gene prediction methods, based on *de novo* prediction, homologous genes, and transcriptomes, were integrated to annotate protein-coding genes. RNA-seq data of *F. gigantica* were obtained from the NCBI Sequence Read Archive, SRR4449208 [66]. RNA-seq reads were aligned to the genome assembly using HISAT2 (v2.2.0) [67] and subsequently assembled using StringTie (v2.1.3) [68]. PASA (v2.4) [69] was another tool used to assemble RNA-seq reads and further generated gene models to train *de novo* programs. Two *de novo* programs, including Augustus (v3.0.2) [70] and SNAP (v2006-07-28) [71], were used to predict genes in the repeat-masked genome sequences. For homology-based prediction, protein sequences from UniRef100 [72] (plagiorchiida-specific, $n = 75,612$) were aligned on the genome sequence using TBLASTn [73] (e-value $< 10^{-4}$), and GeneWise (version 2.4.1) [74] was used to identify accurate gene structures. All predicted genes from the three approaches were combined using MAKER (v3.1.2) [75] to generate high-confidence gene sets. To obtain gene function annotations, Interproscan (v5.45) [76] was used to identify annotated genes features, including protein families, domains, functional sites, and GO terms from the InterPro database. SwissProt and TrEMBL protein databases were also searched using BLASTp [77] (e-value $< 10^{-4}$). The best BLASTp hits were used to assign homology-based gene functions. BlastKOALA [78] was used to search the KEGG ORTHOLOGY (KO) database. The subsequent enrichment analysis was performed using clusterProfiler using total annotated genes as the background with the “enricher” function [79].

Repeat annotation and analysis

We combined *de novo* and homology approaches to identify repetitive sequences in our assembly and previous published assemblies, including *F. gigantica*, *F. hepatica*, and *Fasciolopsis buski*. RepeatModeler (v2.0.1) [24] was first used to construct the *de novo* identification and accurate compilation of sequence models representing all of the unique TE families dispersed in the genome. Then, RepeatMasker (v4.1.0) [25] was run on the genome using the combination of *de novo* libraries and a library of known repeats (Repbase-20181026). The relative position between a repeat and a gene was identified using bedtools [80], and the type of repeat was further divided to intronic and intergenic origin. The repeat landscape was constructed using sequence alignments and the complete annotations output from RepeatMasker, depicting the Kimura divergence (Kimura genetic distances between identified repeat sequences and their consensus) distribution of all repeats types. The most notable peak in the repeat landscapes was considered as the most convincing time of repeat duplication in that period. We inferred the time of LINEs insertion by transferring Kimura divergence in RepeatMasker to age ($t = d / 2\mu$). The distributions of TE elements were calculated with sliding windows ($n = 50$). In each

sliding window, we calculated the relative proportion of TE between intronic and intergenic regions, and further corrected them using the whole ratio between intronic and intergenic regions. To calculate mutation rate, we used 559 single-copy orthologs multiple sequence alignment among 8 species produced in the latter gene family analysis, and estimated the mutation rate using MCMCtree with global clock. A Markov chain Monte Carlo (MCMC) process was run for 2,000,000 iterations, with sample frequency of 100 after a burn-in of 1,000 iterations. The median of simulated data was selected as mutation rate ($\mu = 1.73 \times 10^{-9}$ per base per year).

Genome-wide host-parasite protein interaction analysis

In addition to the genome data that we generated for *F. gigantica*, we downloaded genome annotation information for human (GCA_000001405.28), swamp buffalo (GWHAAJZ000000000), *F. hepatica* (GCA_002763495.2), *Fasciolopsis buski* (GCA_008360955.1), *Clonorchis sinensis* (GCA_003604175.1), *Schistosoma mansoni* (GCA_000237925.2), and *Taenia multiceps* (GCA_001923025.3) from the NCBI database and BIG Sub (China National Center for Bioinformatics, Beijing, China). Proteases and protease inhibitors were identified and classified into families using BLASTp (e-value $< 10^{-4}$) against the MEROPS peptidase database (merops_scan.lib; (European Bioinformatics Institute (EMBL-EBI), Cambridge, UK)), with amino acids at least 80% coverage matched for database proteins. These proteases were divided into five major classes (aspartic, cysteine, metallo, serine, and threonine proteases). E/S proteins (i.e., the secretome) were predicted by the programs SignalP 5.0 [81], TargetP [82], and TMHMM [83]. Proteins with a signal peptide sequence but without a transmembrane region were identified as secretome proteins, excluding the mitochondrial sequences. Genome-wide host-parasite protein interaction analysis was performed by constructing the PPIs between the *F. gigantica* secretome and human proteins expressed in the tissues related to the liver fluke life cycle. For the hosts, we selected human proteins expressed in the small intestine and liver, and located in the plasma membrane and extracellular region. The gene expression and subcellular location information were obtained from the TISSUES [84] and Uniprot (EMBL-EBI) databases, respectively. For *F. gigantica*, secretome molecules were mapped to the human proteome as the reference, using the reciprocal best-hit BLAST method. These two gene datasets were used to construct host-parasite PPI networks. We downloaded the interaction files (protein.links.v11.0) in the STRING database [85], and only highly credible PPIs were retained by excluding PPIs with confidence scores below 0.7. The final STRING network was plotted using Cytoscape [86].

Gene family analysis

We chose the longest transcript in the downloaded annotation dataset to represent each gene, and removed genes with open reading frames shorter than 150 bp. Gene family clustering was then performed using OrthoFinder (v 2.3.12) [87], based on the predicted gene set for eight genomes, including *F. gigantica* (our assembly), *F. hepatica* (NCBI: GCA_002763495.2), *Fasciolopsis buski* (NCBI: GCA_008360955.1), *Clonorchis sinensis* (NCBI: GCA_003604175.1), *Schistosoma mansoni* (NCBI: GCF_000237925.1), *Taenia multiceps* (NCBI: GCA_001923025.3), swamp buffalo (BIG sub: GWHAAJZ000000000), and human (NCBI: GCF_000001405.39). This analysis yielded 17,992 gene families. To identify gene families that had undergone expansion or contraction, we applied the CAFE (v5.0.0) program [88], which inferred the rate and direction of changes in gene family size over a given phylogeny. Among the eight species, 559 single-copy orthologs were aligned using MUSCLE (v3.8.1551) [89], and we eliminated poorly aligned positions and divergent regions of the alignment using Gblock

0.91b [90]. RAxML (v 8.2.12) was then used with the PROTGAMMALGF model to estimate a maximum likelihood tree. Divergence times were estimated using PAML MCMCTREE [91]. A Markov chain Monte Carlo (MCMC) process was run for 2,000,000 iterations, with a sample frequency of 100 after a burn-in of 1,000 iterations under an independent rates model. Two independent runs were performed to check the convergence. The fossil-calibrated eukaryote phylogeny was used to set the root height for the species tree, taken from the age of Animals (602–661 Ma) estimated in a previous fossil-calibrated eukaryotic phylogeny [92] and the divergence time between the euarchontoglires and laurasiatheria: (95.3–113 Ma) [93].

To enhance the reproducibility of the results, we deposit the laboratory protocols in protocols.io (PROTOCOL DOI): <http://dx.doi.org/10.17504/protocols.io.bxatpien>.

Supporting information

S1 Fig. Genome-wide all-by-all chromosome conformation capture sequencing (Hi-C) interaction in *F. gigantica* (Bins = 500 K).

(TIF)

S2 Fig. Comparison of chromosome length between the chromosome conformation capture sequencing (Hi-C) assembly and estimates from published karyotype data by Jae Ku Rhee.

(TIF)

S3 Fig. Boxplot of average gene length.

(TIF)

S4 Fig. Boxplot of average coding sequence (CDS) length per gene.

(TIF)

S5 Fig. Divergence distribution of classified families of transposable elements. The classified transposon families in *F. gigantica*.

(TIF)

S6 Fig. Expansion time of long terminal repeats (LTRs) and long interspersed elements (LINEs). The mutation rate was 1.73×10^{-9} .

(TIF)

S7 Fig. Estimation of *F. gigantica* genome size based on the expansion time of repeat sequences during evolution. The mutation rate was 1.73×10^{-9} .

(TIF)

S8 Fig. Overlapping E/S proteins between this study and proteomic study by Di Maggio LS et al [94].

(TIF)

S1 Table. Genome sequencing strategy for buffaloes.

(XLSX)

S2 Table. Summary of the *Fasciola gigantica* genome assembly.

(XLSX)

S3 Table. Summary of different assemblies in *Fasciola* species.

(XLSX)

S4 Table. Summary of chromosome conformation capture sequencing (Hi-C) assembly of the chromosome length in *Fasciola gigantica*.

(XLSX)

S5 Table. Assessment of the completeness and accuracy of the genome.

(XLSX)

S6 Table. BUSCO assessment of the genome.

(XLSX)

S7 Table. Number of genes with functional classification gained using various methods.

(XLSX)

S8 Table. Transposable element content of *Fasciola gigantica* genome.

(XLSX)

S9 Table. The list of genes with more than 10 kb of long interspersed element (LINE) insertion between 41 Ma and 62 Ma.

(XLSX)

S10 Table. Gene ontology (GO) term category enrichment for genes with more than 10 kb of long interspersed element (LINE) insertion between 41 Ma and 62 Ma.

(XLSX)

S11 Table. Kyoto Encyclopedia of Genes and Genomes (KEGG) pathway enrichment for genes with more than 10 kb of long interspersed element (LINE) insertion between 41 Ma and 62 Ma.

(XLSX)

S12 Table. Kyoto Encyclopedia of Genes and Genomes (KEGG) pathway enrichment for genes with more than 10 kb of long interspersed element (LINE) insertion between 41 Ma and 62 Ma.

(XLSX)

S13 Table. Protein inhibitors in the *Fasciola gigantica* genome.

(XLSX)

S14 Table. Excretory/secretory (E/S) proteins in the *Fasciola gigantica* genome.

(XLSX)

S15 Table. Gene ontology (GO) term category enrichment for excretory/secretory (E/S) proteins.

(XLSX)

S16 Table. Gene ontology (GO) term category enrichment for rapidly evolving families specific to *F. gigantica*.

(XLSX)

S17 Table. Function annotation based on human uniprot gene using blastp with E-value < 10⁻⁴.

(XLSX)

S1 Text. AGP file for *Fasciola gigantica*.txt.

(DOC)

Author Contributions

Conceptualization: Kuiqing Cui, Zhiqiang Wang, Zhipeng Li.

Data curation: Kuiqing Cui, Zhiqiang Wang, Zhipeng Li.

Formal analysis: Kuiqing Cui, Zhiqiang Wang, Zhipeng Li.

Funding acquisition: Qingyou Liu.

Investigation: Xier Luo.

Methodology: Xier Luo, Jue Ruan.

Project administration: Xier Luo, Qingyou Liu.

Resources: Zhipeng Li, Zhengjiao Wu, Weiyi Huang.

Software: Xier Luo, Jue Ruan.

Supervision: Jue Ruan, Qingyou Liu.

Validation: Zhengjiao Wu, Weiyi Zhang.

Visualization: Weiyi Zhang, Qingyou Liu.

Writing – original draft: Xier Luo.

Writing – review & editing: Weiyi Huang, Xing-Quan Zhu, Jue Ruan, Weiyi Zhang, Qingyou Liu.

References

1. World Health O. Accelerating work to overcome the global impact of neglected tropical diseases—a roadmap for implementation. *Accelerating work to overcome the global impact of neglected tropical diseases—a roadmap for implementation*. 2012; 37.
2. Yadav SC, Sharma RL, Kalicharan A, Mehra UR, Dass RS, Verma AK. Primary experimental infection of riverine buffaloes with *Fasciola gigantica*. *Veterinary Parasitology*. 1999; 82(4):285–96. [https://doi.org/10.1016/s0304-4017\(99\)00005-9](https://doi.org/10.1016/s0304-4017(99)00005-9) PMID: 10384904.
3. Cwiklinski K, Dalton JP. Advances in *Fasciola hepatica* research using ‘omics’ technologies. *International Journal for Parasitology*. 2018; 48(5):321–31.
4. McNulty SN, Tort JF, Rinaldi G, Fischer K, Rosa BA, Smircich P, et al. Genomes of *Fasciola hepatica* from the Americas Reveal Colonization with *Neorickettsia* Endobacteria Related to the Agents of Potomac Horse and Human Sennetsu Fevers. *Plos Genetics*. 2017; 13(1). <https://doi.org/10.1371/journal.pgen.1006537> PMID: 28060841.
5. Cwiklinski K, Dalton JP, Dufresne PJ, La Course J, Williams DJL, Hodgkinson J, et al. The *Fasciola hepatica* genome: gene duplication and polymorphism reveals adaptation to the host environment and the capacity for rapid evolution. *Genome Biology*. 2015; 16. <https://doi.org/10.1186/s13059-015-0632-2> PMID: 25887684.
6. Choi Y-J, Fontenla S, Fischer PU, Thanh Hoa L, Costabile A, Blair D, et al. Adaptive Radiation of the Flukes of the Family Fasciolidae Inferred from Genome-Wide Comparisons of Key Species. *Mol Biol Evol*. 2020; 37(1):84–99. <https://doi.org/10.1093/molbev/msz204> PMID: 31501870.
7. Pandey T, Ghosh A, Todur VN, Rajendran V, Kalita P, Kalita J, et al. Draft Genome of the Liver Fluke *Fasciola gigantica*. *Acs Omega*. 2020; 5(19):11084–91. <https://doi.org/10.1021/acsomega.0c00980> PMID: 32455229.
8. Soyemi J, Isewon I, Oyelade J, Adebijoyi E. Inter-Species/Host-Parasite Protein Interaction Predictions Reviewed. *Current Bioinformatics*. 2018; 13(4):396–406. <https://doi.org/10.2174/1574893613666180108155851> PMID: 31496926.
9. de la Torre-Escudero E, Gerlach JQ, Bennett APS, Cwiklinski K, Jewhurst HL, Huson KM, et al. Surface molecules of extracellular vesicles secreted by the helminth pathogen *Fasciola hepatica* direct their internalisation by host cells. *PLoS Negl Trop Dis*. 2019; 13(1). <https://doi.org/10.1371/journal.pntd.0007087> PMID: 30657764.

10. Jaikua W, Kueakhai P, Chaithirayanon K, Tanomrat R, Wongwairo S, Riengrojpitak S, et al. Cytosolic superoxide dismutase can provide protection against *Fasciola gigantica*. *Acta Tropica*. 2016; 162:75–82. <https://doi.org/10.1016/j.actatropica.2016.06.020> PMID: 27338185.
11. Brennan GP, Fairweather I, Trudgett A, Hoey E, McCoy, McConville M, et al. Understanding triclabendazole resistance. *Experimental and Molecular Pathology*. 2007; 82(2):104–9. <https://doi.org/10.1016/j.yexmp.2007.01.009> PMID: 17398281.
12. Kelley JM, Elliott TP, Beddoe T, Anderson G, Skuce P, Spithill TW. Current Threat of Triclaenzazole Resistance in *Fasciola hepatica*. *Trends in Parasitology*. 2016; 32(6):458–69. <https://doi.org/10.1016/j.pt.2016.03.002> PMID: 27049013.
13. Rehman A, Ullah R, Gupta D, Khan MAH, Rehman L, Beg MA, et al. Generation of oxidative stress and induction of apoptotic like events in curcumin and thymoquinone treated adult *Fasciola gigantica* worms. *Experimental Parasitology*. 2020; 209:107810. <https://doi.org/10.1016/j.exppara.2019.107810> PMID: 31801691.
14. Le TH, De NV, Agatsuma T, Thi Nguyen TG, Nguyen QD, McManus DP, et al. Human fascioliasis and the presence of hybrid/introgressed forms of *Fasciola hepatica* and *Fasciola gigantica* in Vietnam. *International Journal for Parasitology*. 2008; 38(6):725–30. <https://doi.org/10.1016/j.ijpara.2007.10.003> PMID: 18031748.
15. Ashrafi K, Valero MA, Panova M, Periago MV, Massoud J, Mas-Coma S. Phenotypic analysis of adults of *Fasciola hepatica*, *Fasciola gigantica* and intermediate forms from the endemic region of Gilan, Iran. *Parasitology International*. 2006; 55(4):249–60. <https://doi.org/10.1016/j.parint.2006.06.003> PMID: 16901748.
16. Rhee JK, Eun GS, Lee SB. Karyotype of *Fasciola* sp. obtained from Korean cattle. *Kisaengch'unghak chapchi The Korean journal of parasitology*. 1987; 25(1):37–44.
17. Gurevich A, Saveliev V, Vyahhi N, Tesler G. QUAST: quality assessment tool for genome assemblies. *Bioinformatics*. 2013; 29(8):1072–5. <https://doi.org/10.1093/bioinformatics/btt086> PMID: 23422339.
18. Waterhouse RM, Seppey M, Simao FA, Manni M, Ioannidis P, Klioutchnikov G, et al. BUSCO Applications from Quality Assessments to Gene Prediction and Phylogenomics. *Mol Biol Evol*. 2018; 35(3):543–8. <https://doi.org/10.1093/molbev/msx319> PMID: 29220515.
19. Apweiler R, Attwood TK, Bairoch A, Bateman A, Birney E, Biswas M, et al. The InterPro database, an integrated documentation resource for protein families, domains and functional sites. *Nucleic Acids Research*. 2001; 29(1):37–40. <https://doi.org/10.1093/nar/29.1.37> PMID: 11125043.
20. Jones P, Binns D, Chang H-Y, Fraser M, Li W, McAnulla C, et al. InterProScan 5: genome-scale protein function classification. *Bioinformatics*. 2014; 30(9):1236–40. <https://doi.org/10.1093/bioinformatics/btu031> PMID: 24451626.
21. Lanciano S, Cristofari G. Measuring and interpreting transposable element expression. *Nature Reviews Genetics*. 2020; 21(12):721–36. <https://doi.org/10.1038/s41576-020-0251-y> PMID: 32576954.
22. Richardson SR, Doucet AJ, Kopera HC, Moldovan JB, Garcia-Perez JL, Moran JV. The Influence of LINE-1 and SINE Retrotransposons on Mammalian Genomes. *Microbiol Spectr*. 2015; 3(2):MDNA3-2014. <https://doi.org/10.1128/microbiolspec.MDNA3-0061-2014> PMID: 26104698.
23. Bejerano G, Lowe C, Ahituv N, King B, Siepel A, Salama S, et al. A distal enhancer and an ultraconserved exon are derived from a novel retroposon. *Nature*. 2006; 441:87–90. <https://doi.org/10.1038/nature04696> PMID: 16625209.
24. Flynn JM, Hubley R, Goubert C, Rosen J, Clark AG, Feschotte C, et al. RepeatModeler2 for automated genomic discovery of transposable element families. *Proceedings of the National Academy of Sciences of the United States of America*. 2020; 117(17):9451–7. Epub 2020/04/16. <https://doi.org/10.1073/pnas.1921046117> PMID: 32300014.
25. Smit AFA. Interspersed repeats and other mementos of transposable elements in mammalian genomes. *Current Opinion in Genetics & Development*. 1999; 9(6):657–63. [https://doi.org/10.1016/S0959-437X\(99\)00031-3](https://doi.org/10.1016/S0959-437X(99)00031-3) PMID: 10607616.
26. Keller G, Mateo P, Punekar J, Khozyem H, Gertsch B, Spangenberg J, et al. Environmental changes during the Cretaceous-Paleogene mass extinction and Paleocene-Eocene Thermal Maximum: Implications for the Anthropocene. *Gondwana Research*. 2018; 56:69–89.
27. Pastor-Cantizano N, Montesinos JC, Bernat-Silvestre C, Marcote MJ, Aniento F. p24 family proteins: key players in the regulation of trafficking along the secretory pathway. *Protoplasma*. 2016; 253(4):967–85. <https://doi.org/10.1007/s00709-015-0858-6> PMID: 26224213.
28. Montesinos JC, Sturm S, Langhans M, Hillmer S, Marcote MJ, Robinson DG, et al. Coupled transport of Arabidopsis p24 proteins at the ER-Golgi interface. *J Exp Bot*. 2012; 63(11):4243–61. Epub 2012/05/10. <https://doi.org/10.1093/jxb/ers112> PMID: 22577184.

29. Juan T, Fürthauer M. Biogenesis and function of ESCRT-dependent extracellular vesicles. *Seminars in Cell & Developmental Biology*. 2018; 74:66–77. <https://doi.org/10.1016/j.semcdb.2017.08.022> PMID: 28807885
30. Bannister AJ, Kouzarides T. Regulation of chromatin by histone modifications. *Cell Research*. 2011; 21(3):381–95. <https://doi.org/10.1038/cr.2011.22> PMID: 21321607.
31. Oda H, Okamoto I, Murphy N, Chu J, Price SM, Shen MM, et al. Monomethylation of histone H4-lysine 20 is involved in chromosome structure and stability and is essential for mouse development. *Mol Cell Biol*. 2009; 29(8):2278–95. Epub 2009/02/17. <https://doi.org/10.1128/MCB.01768-08> PMID: 19223465.
32. Muiño L, Perteguer MJ, Gárate T, Martínez-Sernández V, Beltrán A, Romarís F, et al. Molecular and immunological characterization of *Fasciola* antigens recognized by the MM3 monoclonal antibody. *Molecular and Biochemical Parasitology*. 2011; 179(2):80–90. <https://doi.org/10.1016/j.molbiopara.2011.06.003> PMID: 21723885.
33. Dalton JP, Neill SO, Stack C, Collins P, Walshe A, Sekiya M, et al. *Fasciola hepatica* cathepsin L-like proteases: biology, function, and potential in the development of first generation liver fluke vaccines. *International Journal for Parasitology*. 2003; 33(11):1173–81. [https://doi.org/10.1016/s0020-7519\(03\)00171-1](https://doi.org/10.1016/s0020-7519(03)00171-1) PMID: 13678633.
34. Batra S, Chatterjee RK, Srivastava VML. Antioxidant system of *Litomosoides carinii* and *Setaria cervi*: effect of a macrofilaricidal agent. *Veterinary Parasitology*. 1992; 43(1):93–103. [https://doi.org/10.1016/0304-4017\(92\)90052-b](https://doi.org/10.1016/0304-4017(92)90052-b) PMID: 1496806.
35. McGonigle S, Dalton JP. Isolation of *Fasciola hepatica* haemoglobin. *Parasitology*. 1995; 111:209–15. <https://doi.org/10.1017/s003118200064969> PMID: 7675536.
36. Farahnak A, Golestani A, Eshraghian M. Activity of Superoxide Dismutase (SOD) Enzyme in the Excretory-Secretory Products of *Fasciola hepatica* and *F. gigantica* Parasites. *Iran J Parasitol*. 2013; 8(1):167–70. PMID: 23682275.
37. Okada K, Fukui M, Zhu B-T. Protein disulfide isomerase mediates glutathione depletion-induced cytotoxicity. *Biochemical and Biophysical Research Communications*. 2016; 477(3):495–502. <https://doi.org/10.1016/j.bbrc.2016.06.066> PMID: 27317486.
38. Biswal DK, Roychowdhury T, Pandey P, Tandon V. De novo genome and transcriptome analyses provide insights into the biology of the trematode human parasite *Fasciolopsis buski*. *Plos One*. 2018; 13(10). <https://doi.org/10.1371/journal.pone.0205570> PMID: 30325945.
39. Wang X, Chen W, Huang Y, Sun J, Men J, Liu H, et al. The draft genome of the carcinogenic human liver fluke *Clonorchis sinensis*. *Genome Biology*. 2011; 12(10). <https://doi.org/10.1186/gb-2011-12-10-r107> PMID: 22023798.
40. Berriman M, Haas BJ, LoVerde PT, Wilson RA, Dillon GP, Cerqueira GC, et al. The genome of the blood fluke *Schistosoma mansoni*. *Nature*. 2009; 460(7253):352–U65. <https://doi.org/10.1038/nature08160> PMID: 19606141.
41. Li W, Liu B, Yang Y, Ren Y, Wang S, Liu C, et al. The genome of tapeworm *Taenia multiceps* sheds light on understanding parasitic mechanism and control of coenurosis disease. *DNA Research*. 2018; 25(5):499–510. <https://doi.org/10.1093/dnares/dsy020> PMID: 29947776.
42. Luo X, Zhou Y, Zhang B, Zhang Y, Wang X, Feng T, et al. Understanding divergent domestication traits from the whole-genome sequencing of swamp- and river-buffalo populations. *National Science Review*. 2020; 7(3):686–701.
43. de Jong P, Catanese JJ, Osoegawa K, Shizuya H, Choi S, Chen YJ, et al. Initial sequencing and analysis of the human genome (vol 409, pg 860, 2001). *Nature*. 2001; 412(6846):565–6. <https://doi.org/10.1038/35057062> PMID: 11237011.
44. Irving JA, Spithill TW, Pike RN, Whisstock JC, Smooker PM. The Evolution of Enzyme Specificity in *Fasciola* spp. *Journal of Molecular Evolution*. 2003; 57(1):1–15. <https://doi.org/10.1007/s00239-002-2434-x> PMID: 12962301.
45. Näsvall J, Sun L, Roth JR, Andersson DI. Real-time evolution of new genes by innovation, amplification, and divergence. *Science*. 2012; 338(6105):384–7. <https://doi.org/10.1126/science.1226521> PMID: 23087246.
46. Pollard TD. Actin and Actin-Binding Proteins. *Cold Spring Harb Perspect Biol*. 2016; 8(8):a018226. <https://doi.org/10.1101/cshperspect.a018226> PMID: 26988969.
47. Dominguez R, Holmes KC. Actin Structure and Function. *Annual Review of Biophysics*. 2011; 40(1):169. <https://doi.org/10.1146/annurev-biophys-042910-155359> PMID: 21314430.
48. Priess JR, Hirsh DI. *Caenorhabditis elegans* morphogenesis: the role of the cytoskeleton in elongation of the embryo. *Developmental biology*. 1986; 117(1):156–73. [https://doi.org/10.1016/0012-1606\(86\)90358-1](https://doi.org/10.1016/0012-1606(86)90358-1) PMID: 3743895.

49. McKeown C, Praitis V, Austin J. sma-1 encodes a beta(H)-spectrin homolog required for *Caenorhabditis elegans* morphogenesis. *Development*. 1998; 125(11):2087–98. PMID: [9570773](#).
50. de Almeida A, Martins AP, Mosca AF, Wijma HJ, Prista C, Soveral G, et al. Exploring the gating mechanisms of aquaporin-3: new clues for the design of inhibitors? *Molecular Biosystems*. 2016; 12(5):1564–73. <https://doi.org/10.1039/c6mb00013d> PMID: [26960817](#).
51. Soveral G, Casini A. Aquaporin modulators: a patent review (2010–2015). *Expert Opinion on Therapeutic Patents*. 2016; 27(1):49. <https://doi.org/10.1080/13543776.2017.1236085> PMID: [27622909](#).
52. Geadkaew A, von Bülow J, Beitz E, Grams SV, Viyanant V, Grams R. Functional analysis of novel aquaporins from *Fasciola gigantica*. *Molecular and Biochemical Parasitology*. 2011; 175(2):144–53. <https://doi.org/10.1016/j.molbiopara.2010.10.010> PMID: [21073907](#).
53. Sedlazeck FJ, Lee H, Darby CA, Schatz MC. Piercing the dark matter: bioinformatics of long-range sequencing and mapping. *Nature Reviews Genetics*. 2018; 19(6):329–46. <https://doi.org/10.1038/s41576-018-0003-4> PMID: [29599501](#).
54. Chuong EB, Elde NC, Feschotte C. Regulatory activities of transposable elements: from conflicts to benefits. *Nature Reviews Genetics*. 2017; 18(2):71–86. <https://doi.org/10.1038/nrg.2016.139> PMID: [27867194](#).
55. Lynch VJ, Leclerc RD, May G, Wagner GP. Transposon-mediated rewiring of gene regulatory networks contributed to the evolution of pregnancy in mammals. *Nature Genetics*. 2011; 43(11):1154–9. <https://doi.org/10.1038/ng.917> PMID: [21946353](#).
56. Zhang F-K, Zhang X-X, Elsheikha HM, He J-J, Sheng Z-A, Zheng W-B, et al. Transcriptomic responses of water buffalo liver to infection with the digenetic fluke *Fasciola gigantica*. *Parasites & Vectors*. 2017; 10(1):56. <https://doi.org/10.1186/s13071-017-1990-2> PMID: [28143561](#).
57. Zhang F-K, Hu R-S, Elsheikha HM, Sheng Z-A, Zhang W-Y, Zheng W-B, et al. Global serum proteomic changes in water buffaloes infected with *Fasciola gigantica*. *Parasites & Vectors*. 2019; 12(1):281. <https://doi.org/10.1186/s13071-019-3533-5> PMID: [31159882](#).
58. Valero MA, Darce NAn, Panova M, Mas-Coma S. Relationships between host species and morphometric patterns in *Fasciola hepatica* adults and eggs from the northern Bolivian Altiplano hyperendemic region. *Veterinary Parasitology*. 2001; 102(1):85–100. [https://doi.org/10.1016/s0304-4017\(01\)00499-x](https://doi.org/10.1016/s0304-4017(01)00499-x) PMID: [11705655](#).
59. Koren S, Walenz BP, Berlin K, Miller JR, Bergman NH, Phillippy AM. Canu: scalable and accurate long-read assembly via adaptive k-mer weighting and repeat separation. *Genome Research*. 2017; 27(5):722–36. <https://doi.org/10.1101/gr.215087.116> PMID: [28298431](#).
60. Berlin K, Koren S, Chin C-S, Drake JP, Landolin JM, Phillippy AM. Assembling large genomes with single-molecule sequencing and locality-sensitive hashing. *Nature Biotechnology*. 2015; 33(6):623–30. <https://doi.org/10.1038/nbt.3238> PMID: [26006009](#).
61. Guan D, McCarthy SA, Wood J, Howe K, Wang Y, Durbin R. Identifying and removing haplotypic duplication in primary genome assemblies. *Bioinformatics* (Oxford, England). 2020; 36(9):2896–8. <https://doi.org/10.1093/bioinformatics/btaa025> PMID: [31971576](#).
62. Walker BJ, Abeel T, Shea T, Priest M, Abouelliel A, Sakthikumar S, et al. Pilon: An Integrated Tool for Comprehensive Microbial Variant Detection and Genome Assembly Improvement. *Plos One*. 2014; 9(11). <https://doi.org/10.1371/journal.pone.0112963> PMID: [25409509](#).
63. Zhang X, Zhang S, Zhao Q, Ming R, Tang H. Assembly of allele-aware, chromosomal-scale autopolyploid genomes based on Hi-C data. *Nature Plants*. 2019; 5(8):833–45. <https://doi.org/10.1038/s41477-019-0487-8> PMID: [31383970](#).
64. Mikheenko A, Pribelski A, Saveliev V, Antipov D, Gurevich A. Versatile genome assembly evaluation with QUASt-LG. *Bioinformatics*. 2018; 34(13):i142–i50. Epub 2018/06/29. <https://doi.org/10.1093/bioinformatics/bty266> PMID: [29949969](#).
65. Li H, Durbin R. Fast and accurate long-read alignment with Burrows-Wheeler transform. *Bioinformatics*. 2010; 26(5):589–95. <https://doi.org/10.1093/bioinformatics/btp698> PMID: [20080505](#).
66. Zhang X-X, Cwiklinski K, Hu R-S, Zheng W-B, Sheng Z-A, Zhang F-K, et al. Complex and dynamic transcriptional changes allow the helminth *Fasciola gigantica* to adjust to its intermediate snail and definitive mammalian hosts. *BMC Genomics*. 2019; 20(1):729-. <https://doi.org/10.1186/s12864-019-6103-5> PMID: [31606027](#).
67. Kim D, Paggi JM, Park C, Bennett C, Salzberg SL. Graph-based genome alignment and genotyping with HISAT2 and HISAT-genotype. *Nature biotechnology*. 2019; 37(8):907–15. Epub 2019/08/02. <https://doi.org/10.1038/s41587-019-0201-4> PMID: [31375807](#).
68. Pertea M, Pertea GM, Antonescu CM, Chang T-C, Mendell JT, Salzberg SL. StringTie enables improved reconstruction of a transcriptome from RNA-seq reads. *Nature biotechnology*. 2015; 33(3):290–5. Epub 2015/02/18. <https://doi.org/10.1038/nbt.3122> PMID: [25690850](#).

69. Haas BJ, Delcher AL, Mount SM, Wortman JR, Smith RK Jr., Hannick LI, et al. Improving the Arabidopsis genome annotation using maximal transcript alignment assemblies. *Nucleic acids research*. 2003; 31(19):5654–66. <https://doi.org/10.1093/nar/gkg770> PMID: 14500829.
70. Stanke M, Steinkamp R, Waack S, Morgenstern B. AUGUSTUS: a web server for gene finding in eukaryotes. *Nucleic Acids Research*. 2004; 32(suppl_2):W309–W12. <https://doi.org/10.1093/nar/gkh379> PMID: 15215400.
71. Korf I. Gene finding in novel genomes. *BMC Bioinformatics*. 2004; 5(1):59. <https://doi.org/10.1186/1471-2105-5-59> PMID: 15144565.
72. The UniProt Consortium. UniProt: the universal protein knowledgebase. *Nucleic acids research*. 2017; 45(D1):D158–D69. Epub 2016/11/29. <https://doi.org/10.1093/nar/gkw1099> PMID: 27899622.
73. Gertz EM, Yu Y-K, Agarwala R, Schäffer AA, Altschul SF. Composition-based statistics and translated nucleotide searches: improving the TBLASTN module of BLAST. *BMC Biol*. 2006; 4:41-. <https://doi.org/10.1186/1741-7007-4-41> PMID: 17156431.
74. Birney E, Clamp M, Durbin R. GeneWise and Genomewise. *Genome research*. 2004; 14(5):988–95. <https://doi.org/10.1101/gr.1865504> PMID: 15123596.
75. Holt C, Yandell M. MAKER2: an annotation pipeline and genome-database management tool for second-generation genome projects. *BMC bioinformatics*. 2011; 12:491-. <https://doi.org/10.1186/1471-2105-12-491> PMID: 22192575.
76. Quevillon E, Silventoinen V, Pillai S, Harte N, Mulder N, Apweiler R, et al. InterProScan: protein domains identifier. *Nucleic acids research*. 2005; 33(Web Server issue):W116–20. <https://doi.org/10.1093/nar/gki442> PMID: 15980438.
77. Altschul SF, Madden TL, Schäffer AA, Zhang J, Zhang Z, Miller W, et al. Gapped BLAST and PSI-BLAST: a new generation of protein database search programs. *Nucleic acids research*. 1997; 25(17):3389–402. <https://doi.org/10.1093/nar/25.17.3389> PMID: 9254694.
78. Kanehisa M, Sato Y, Morishima K. BlastKOALA and GhostKOALA: KEGG Tools for Functional Characterization of Genome and Metagenome Sequences. *Journal of Molecular Biology*. 2016; 428(4):726–31. <https://doi.org/10.1016/j.jmb.2015.11.006> PMID: 26585406.
79. Yu G, Wang L-G, Han Y, He Q-Y. clusterProfiler: an R Package for Comparing Biological Themes Among Gene Clusters. *Omics-a Journal of Integrative Biology*. 2012; 16(5):284–7. <https://doi.org/10.1089/omi.2011.0118> PMID: 22455463.
80. Quinlan AR. BEDTools: The Swiss-Army Tool for Genome Feature Analysis. *Curr Protoc Bioinformatics*. 2014; 47:11.2.1–.2.34. <https://doi.org/10.1002/0471250953.bi1112s47> PMID: 25199790.
81. Petersen TN, Brunak S, von Heijne G, Nielsen H. SignalP 4.0: discriminating signal peptides from transmembrane regions. *Nature Methods*. 2011; 8(10):785–6. <https://doi.org/10.1038/nmeth.1701> PMID: 21959131.
82. Emanuelsson O, Nielsen H, Brunak S, von Heijne G. Predicting subcellular localization of proteins based on their N-terminal amino acid sequence. *Journal of Molecular Biology*. 2000; 300(4):1005–16. <https://doi.org/10.1006/jmbi.2000.3903> PMID: 10891285.
83. Krogh A, Larsson B, von Heijne G, Sonnhammer ELL. Predicting transmembrane protein topology with a hidden Markov model: Application to complete genomes. *Journal of Molecular Biology*. 2001; 305(3):567–80. <https://doi.org/10.1006/jmbi.2000.4315> PMID: 11152613.
84. Santos A, Tsafou K, Stolte C, Pletscher-Frankild S, O'Donoghue SI, Jensen LJ. Comprehensive comparison of large-scale tissue expression datasets. *Peerj*. 2015; 3. <https://doi.org/10.7717/peerj.1054> PMID: 26157623.
85. Jensen LJ, Kuhn M, Stark M, Chaffron S, Creevey C, Muller J, et al. STRING 8—a global view on proteins and their functional interactions in 630 organisms. *Nucleic Acids Research*. 2009; 37:D412–D6. <https://doi.org/10.1093/nar/gkn760> PMID: 18940858.
86. Doncheva NT, Morris JH, Gorodkin J, Jensen LJ. Cytoscape StringApp: Network Analysis and Visualization of Proteomics Data. *Journal of Proteome Research*. 2019; 18(2):623–32. <https://doi.org/10.1021/acs.jproteome.8b00702> PMID: 30450911.
87. Emms DM, Kelly S. OrthoFinder: phylogenetic orthology inference for comparative genomics. *Genome Biology*. 2019; 20(1):238. <https://doi.org/10.1186/s13059-019-1832-y> PMID: 31727128.
88. Han MV, Thomas GWC, Lugo-Martinez J, Hahn MW. Estimating Gene Gain and Loss Rates in the Presence of Error in Genome Assembly and Annotation Using CAFE 3. *Mol Biol Evol*. 2013; 30(8):1987–97. <https://doi.org/10.1093/molbev/mst100> PMID: 23709260.
89. Edgar RC. MUSCLE: a multiple sequence alignment method with reduced time and space complexity. *Bmc Bioinformatics*. 2004; 5:1–19. <https://doi.org/10.1186/1471-2105-5-1> PMID: 14706121.

90. Talavera G, Castresana J. Improvement of phylogenies after removing divergent and ambiguously aligned blocks from protein sequence alignments. *Systematic Biology*. 2007; 56(4):564–77. <https://doi.org/10.1080/10635150701472164> PMID: 17654362.
91. Yang Z, Rannala B. Bayesian Estimation of Species Divergence Times Under a Molecular Clock Using Multiple Fossil Calibrations with Soft Bounds. *Mol Biol Evol*. 2006; 23(1):212–26. <https://doi.org/10.1093/molbev/msj024> PMID: 16177230.
92. Parfrey LW, Lahr DJG, Knoll AH, Katz LA. Estimating the timing of early eukaryotic diversification with multigene molecular clocks. *Proceedings of the National Academy of Sciences*. 2011; 108(33):13624. <https://doi.org/10.1073/pnas.1110633108> PMID: 21810989.
93. Benton MJ, Donoghue PCJ. Paleontological evidence to date the tree of life (vol 24, pg 26, 2007). *Mol Biol Evol*. 2007; 24(3):889–91. <https://doi.org/10.1093/molbev/msl150> PMID: 17047029.
94. Di Maggio LS, Tirloni L, Pinto AFM, Diedrich JK, Yates JR Iii, Benavides U, et al. Across intra-mammalian stages of the liver fluke *Fasciola hepatica*: a proteomic study. *Scientific Reports*. 2016; 6(1):32796. <https://doi.org/10.1038/srep32796> PMID: 27600774.