

The genomic basis of parasitism in the *Strongyloides* clade of nematodes

Vicky L Hunt^{1,15}, Isheng J Tsai^{2,3,15}, Avril Coghlan^{4,15}, Adam J Reid^{4,15}, Nancy Holroyd⁴, Bernardo J Foth⁴, Alan Tracey⁴, James A Cotton⁴, Eleanor J Stanley⁴, Helen Beasley⁴, Hayley M Bennett⁴, Karen Brooks⁴, Bhavana Harsha⁴, Rei Kajitani⁵, Arpita Kulkarni⁶, Dorothee Harbecke⁶, Eiji Nagayasu³, Sarah Nichol⁴, Yoshitoshi Ogura⁷, Michael A Quail⁴, Nadine Randle⁸, Dong Xia⁸, Norbert W Brattig⁹, Hanns Soblik⁹, Diogo M Ribeiro⁴, Alejandro Sanchez-Flores^{4,10}, Tetsuya Hayashi⁷, Takehiko Itoh⁵, Dee R Denver¹¹, Warwick Grant¹², Jonathan D Stoltzfus¹³, James B Lok¹³, Haruhiko Murayama³, Jonathan Wastling^{8,14}, Adrian Streit⁶, Taisei Kikuchi³, Mark Viney¹ & Matthew Berriman⁴

Soil-transmitted nematodes, including the *Strongyloides* genus, cause one of the most prevalent neglected tropical diseases. Here we compare the genomes of four *Strongyloides* species, including the human pathogen *Strongyloides stercoralis*, and their close relatives that are facultatively parasitic (*Parastrongyloides trichosuri*) and free-living (*Rhabditophanes* sp. KR3021). A significant paralogous expansion of key gene families—families encoding astacin-like and SCP/TAPS proteins—is associated with the evolution of parasitism in this clade. Exploiting the unique *Strongyloides* life cycle, we compare the transcriptomes of the parasitic and free-living stages and find that these same gene families are upregulated in the parasitic stages, underscoring their role in nematode parasitism.

More than 1 billion people are infected with intestinal nematodes^{1,2}. The World Health Organization has classified infections with soil-transmitted nematodes as one of the 17 most neglected tropical diseases and estimates that worldwide these infections cause an annual disease burden of 5 million years lost due to disability (YLD), greater than the annual disease burdens of malaria (4 million YLD) and HIV/AIDS (4.5 million YLD). Parasitic nematode infections can impair physical and educational development¹.

Strongyloides species are soil-transmitted gastrointestinal parasitic nematodes infecting a wide range of vertebrates³. Two species—*S. stercoralis* and *Strongyloides fuelleborni*—infect some 100–200 million people worldwide^{4,5}. Other *Strongyloides* species infect livestock, such as *Strongyloides papillosus* that infects sheep.

Strongyloides species are from a clade of nematodes^{6–8} that includes taxa with diverse lifestyles, including a free-living lifestyle (*Rhabditophanes*), parasitism of invertebrates, facultative parasitism of vertebrates (*Parastrongyloides*) and obligate parasitism of vertebrates (*Strongyloides*)^{6,7}. Nematodes have independently evolved parasitism of animals several times⁹, and thus understanding

the genomic adaptations to parasitism in one clade will help in understanding how parasitism has evolved across the phylum more widely.

The *Strongyloides* life cycle alternates between free-living and parasitic generations. The female-only, parthenogenetic¹⁰ parasitic stage lives in the small intestine of its host where it produces offspring that develop outside of the host, either directly into infective third-stage larvae (iL3s) or into a dioecious, sexually reproducing adult generation¹¹ whose progeny are iL3s. iL3s penetrate the skin of a host and migrate to its gut¹², where they develop into parasitic adults (Fig. 1). Therefore, this life cycle has two genetically identical adult female stages—one obligate and parasitic and one facultative and free-living; we have compared these stages at the transcriptome and proteome levels to identify the genes and gene products specifically present in the parasitic stage. The closely related genus *Parastrongyloides*^{3,13} is similar to *Strongyloides* species, except that its parasitic generation is dioecious and sexually reproducing and that it can have apparently unlimited cycles of its free-living adult generation^{3,14} (Fig. 1).

¹School of Biological Sciences, University of Bristol, Bristol, UK. ²Biodiversity Research Center, Academia Sinica, Taipei, Taiwan. ³Division of Parasitology, Faculty of Medicine, University of Miyazaki, Miyazaki, Japan. ⁴Wellcome Trust Sanger Institute, Wellcome Trust Genome Campus, Hinxton, UK. ⁵Graduate School of Bioscience and Biotechnology, Tokyo Institute of Technology, Meguro-ku, Tokyo, Japan. ⁶Max Planck Institute for Developmental Biology, Tübingen, Germany. ⁷Department of Bacteriology, Faculty of Medical Sciences, Kyushu University, Fukuoka, Japan. ⁸Department of Infection Biology, Institute of Infection and Global Health and School of Veterinary Science, University of Liverpool, Liverpool, UK. ⁹Department of Molecular Medicine, Bernhard Nocht Institute for Tropical Medicine, Hamburg, Germany. ¹⁰Unidad de Secuenciación Masiva y Bioinformática, Instituto de Biotecnología, Universidad Nacional Autónoma de México, Cuernavaca, Mexico. ¹¹Department of Integrative Biology, Oregon State University, Corvallis, Oregon, USA. ¹²Department of Animal, Plant and Soil Sciences, La Trobe University, Melbourne, Victoria, Australia. ¹³Department of Pathobiology, School of Veterinary Medicine, University of Pennsylvania, Philadelphia, Pennsylvania, USA. ¹⁴Faculty of Natural Sciences, University of Keele, Keele, UK. ¹⁵These authors contributed equally to this work. Correspondence should be addressed to T.K. (taisei_kikuchi@med.miyazaki-u.ac.jp), M.V. (mark.viney@bristol.ac.uk) or M.B. (mb4@sanger.ac.uk).

Received 17 June 2015; accepted 23 December 2015; published online 1 February 2016; doi:10.1038/ng.3495

Here we report the genome sequences for six nematodes from one clade: four species of *Strongyloides*, *S. stercoralis* (a parasite of humans and dogs), *Strongyloides ratti* and *Strongyloides venezuelensis* (both parasites of rats and important laboratory models of nematode infection), and *S. papillosus* (a parasite of sheep); *P. trichosuri* (which infects the brushtail possum *Trichosurus vulpecula*); and the free-living nematode *Rhabditophanes*⁶.

To investigate the genomic and molecular basis of parasitism in these nematodes, we compared (i) the genomes and gene families of the parasitic (*Strongyloides* and *Parastrongyloides*, the Strongyloididae) and free-living (*Rhabditophanes*) taxa (Fig. 1); (ii) the transcriptomes of parasitic adult females, free-living adult females and iL3s from *S. ratti* and *S. stercoralis*; and (iii) the proteomes of parasitic and free-living females of *S. ratti*. We have identified the genes present in the parasitic species and the genes and gene products uniquely upregulated in the parasitic stages of *S. stercoralis* and *S. ratti*; together, these are the major genomic and molecular adaptations to the parasitic lifestyle of these nematodes.

RESULTS

Chromosome biology

We have produced a high-quality 43-Mb reference genome assembly for *S. ratti* (Supplementary Note), with its two autosomes¹⁵ assembled into single scaffolds and the X chromosome¹⁵ assembled into ten scaffolds (Fig. 2 and Table 1). This assembly is the second most contiguous assembled nematode genome after the *Caenorhabditis elegans* reference genome¹⁶. We also produced high-quality draft assemblies of the 42- to 60-Mb genomes of *S. stercoralis*, *S. venezuelensis*, *S. papillosus*, *P. trichosuri* and *Rhabditophanes* sp. KR3021, which were 95.6–99.6% complete (Supplementary Table 1). With GC contents of 21% and 22%, respectively, the *S. ratti* and *S. stercoralis* genomes are the most AT rich reported thus far for nematodes (Supplementary Table 1). The ~43-Mb *S. ratti* and *S. stercoralis* genomes are small compared with the genomes of other nematodes. However, the total protein-coding content of each nematode genome is similar (18–22 Mb versus 14–30 Mb in eight outgroup species; Supplementary Table 1). Significant loss of introns, as well as shorter intergenic regions, accounts for the smaller genomes in the present study (Spearman's correlation between genome size and intron number $\rho = 0.91$, $P < 0.001$ and size of intergenic regions $\rho = 0.63$, $P = 0.02$; Supplementary Table 2). However, parsimony analysis of intronic positions conserved in two or more species showed that substantial intron losses occurred before the evolution of the *Rhabditophanes-Parastrongyloides-Strongyloides* clade (Supplementary Fig. 1) and are therefore not an adaptation associated with parasitism.

The canonical view of a nematode chromosome, defined nearly 20 years ago using *C. elegans* autosomes (and later confirmed in *Caenorhabditis briggsae*¹⁷) is of a gene-dense, repeat-poor 'center' of conserved genes (defined by homology with yeast genes¹⁶), flanked by two gene-poor, repeat-rich 'arms' in which most genes are less strongly conserved. *S. ratti* is the first non-*Caenorhabditis* nematode whose whole chromosomes have been assembled, and it presents a strikingly different organization, with relatively little variability in gene density, repeat density or gene conservation to yeast genes along its autosomes (Supplementary Figs. 2 and 3).

Synteny is highly conserved within the parasitic Strongyloididae but is much less conserved between this family and *Rhabditophanes* (Fig. 2). Scaffolds of the parasitic species largely correspond to blocks from a particular *S. ratti* chromosome but in a scrambled order. This suggests that intrachromosomal rearrangement is frequent but interchromosomal rearrangement is rare, a common phenomenon

in nematode chromosome evolution^{17–19}. The notable exceptions are the *S. papillosus* and *S. venezuelensis* scaffolds that have many blocks that are syntenic to both *S. ratti* chromosomes I and X (Supplementary Table 3). This pattern of synteny likely reflects the fusion event between chromosomes I and X in these species^{20–22}. Associated with this fusion is a change in the chromosome biology of sex determination in these species. *S. papillosus* undergoes chromatin diminution (where a chromosome fragments, after which part of the chromosome is eliminated during mitosis) to mimic the XX/XO sex-determining system of *S. ratti*²³ and *S. stercoralis*²⁰.

By analyzing the differential coverage of mapped sequence data from iL3s (which are all female) and adult males, we were able to identify regions of the *S. papillosus* X-I fusion chromosome that are eliminated from males during diminution (Supplementary Table 4). Six scaffolds were identified from the diminished region using existing genetic markers (Supplementary Table 5), but our read depth approach extended this map to 153 scaffolds (18% of the assembly; 10.9 Mb). Interestingly, some genes with orthologs on the X chromosome of *S. ratti* are not diminished in *S. papillosus*, so the dosage of these genes in males has changed since the species diverged, including three genes on *S. papillosus* chromosome II (confirming earlier work²⁰) and 33 genes that lie in non-diminished regions of the X-I fusion chromosome (Supplementary Table 6).

Extensive rearrangement of the mitochondrial gene order

The *S. stercoralis* mitochondrial genome is highly rearranged compared with the genomes of nematodes from clades I, III and V (ref. 24). Manual finishing of the mitochondrial genomes of the six species showed that the *Rhabditophanes* mitochondrial genome consists of two circular chromosomes, a feature of some other nematode species²⁵. Compared with eight outgroup species, *Rhabditophanes* sp. KR3021 has a conventional gene order but *Strongyloides* species and *P. trichosuri* have highly rearranged mitochondrial genomes (Fig. 2 and Supplementary Table 7). Similar observations have been reported in other clade IV parasitic nematodes^{25–28}, and there is evidence of mitochondrial recombination^{27,29}, which is rarely observed in animals³⁰. Consistent with published nematode mitochondrial genomes, the gene-based phylogeny of the mitochondrial genome (Fig. 2) conflicts with phylogenies based on nuclear genes^{27,31,32}, and the rearranged gene order of the mitochondrial genomes of *Strongyloides* species is accompanied by nucleotide divergence (Fig. 2).

Gene families associated with the evolution of parasitism

We predicted 12,451–18,457 genes across the six genomes, numbers comparable to those in other nematode species (Table 1 and Supplementary Fig. 4). We then used Ensembl Compara (Supplementary Note)³³ to identify orthologs and gene families (Supplementary Table 8) in these and eight outgroup species, encompassing four further nematode clades (Supplementary Fig. 4). By pinpointing when a new gene family arose and where a family has expanded or contracted, we could determine which gene families are associated with the evolution of parasitism. The largest acquisition of gene families (1,075 families) was found on the branch leading to the parasitic nematodes, including the *Strongyloides* species and *P. trichosuri* (Fig. 1 and Supplementary Fig. 4). Despite this highly dynamic pattern of gene gain and loss within each species' genome, the proportion of genes specific to *Strongyloides* (and Strongyloididae) is consistent across the phylogeny (Fig. 1). The branches leading to the five parasitic species also showed greater expansion of genes and families of genes as compared to that in the free-living *Rhabditophanes* sp.

KR3021. Gain and expansion of gene families in these parasitic species likely reflects the necessary adaptations required by these species to be able to parasitize vertebrate hosts while maintaining a free-living life cycle phase.

The two most expanded *Strongyloides* gene families encode astacin-like³⁴ and SCP/TAPS (SCP/Tpx-1/Ag5/PR-1/Sc7 (ref. 35),

also known as CAP-domain) proteins, present in multiple subfamilies (according to Ensembl Compara analysis (Supplementary Table 8) and protein domain combinations (Supplementary Table 9)). The astacin family of metallopeptidases was the most expanded, with 184–387 copies in the *Strongyloides*-*Parastrongyloides* species as compared with *Rhabditophanes*

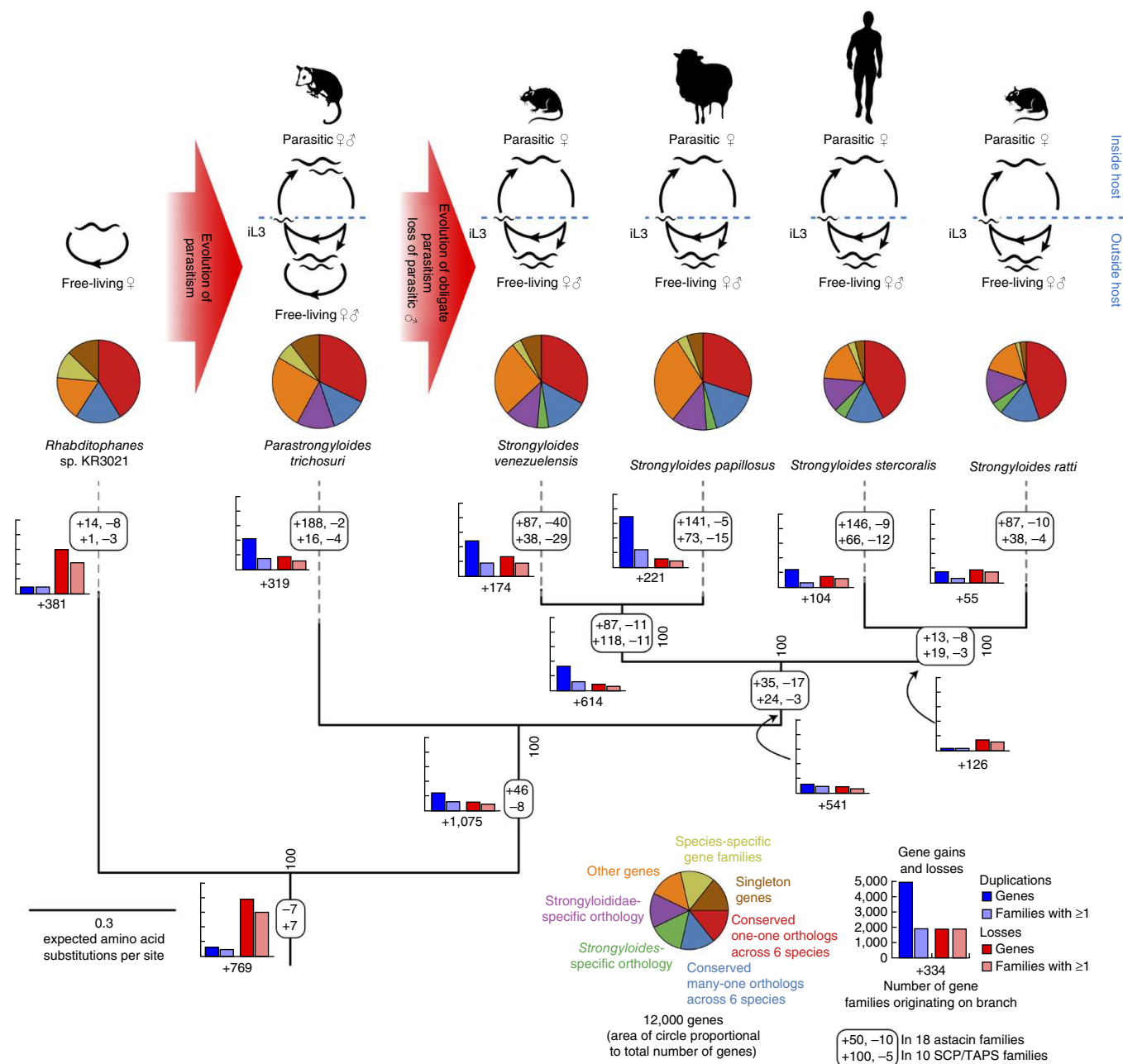


Figure 1 Evolution and comparative genomics of *Strongyloides* species and their relatives. The life cycles of six clade IV nematodes, showing the transition from a free-living lifestyle (in *Rhabditophanes* sp. KR3021) through facultative parasitism (*P. trichosuri*) to obligate parasitism (*Strongyloides* species) and the phylogeny of these species (maximum-likelihood phylogeny based on a concatenated alignment of 841,529 amino acid sites from 4,437 conserved single-copy orthologous genes). Values on nodes (all 100) are the number of bootstrap replicate trees, out of 100 bootstrap replicates, showing the split induced by the node. The phylogeny is annotated with the number of gene families appearing along each branch of the phylogeny ("+" values on each branch), and the histograms show the number of duplications (blue) and losses (red) for individual genes (dark color) and gene families (light color); the numbers of gene origins and gene losses in 18 astacin families (upper numbers in boxes) and ten SCP/TAPS families (lower numbers in boxes) as estimated by the Ensembl Compara pipeline are also shown. The pie charts summarize the evolutionary history of the genome of each species, defining genes shared by all six species, the five parasitic species (*Strongyloididae*, which includes all species except *Rhabditophanes*) and the four *Strongyloides* species or species-specific genes. The host species of the parasites are shown: *P. trichosuri*, brushtail possum; *S. ratti* and *S. venezuelensis*, rat; *S. stercoralis*, humans; *S. papillosus*, sheep.

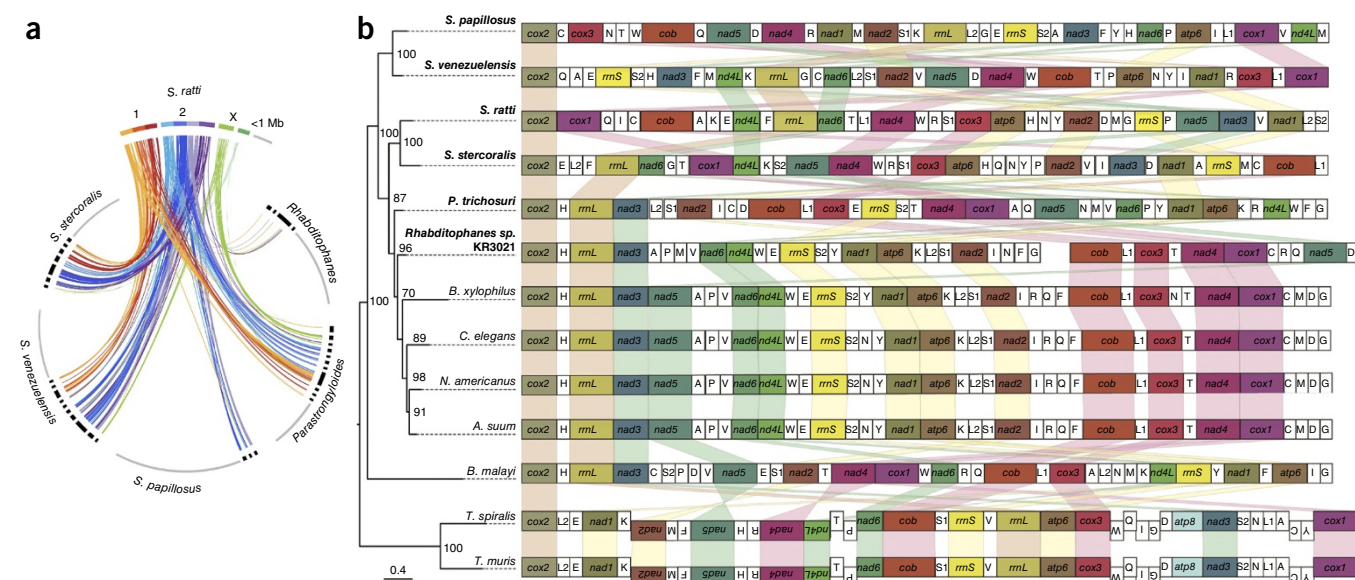


Figure 2 Nuclear genomic synteny and the mitochondrial genomes of four *Strongyloides* species, *P. trichosuri* and *Rhabditophanes* sp. KR3021. **(a)** The *S. ratti* genome, our best assembled genome, is used as the reference sequence; synteny is based on sequence matches. Gradation of color across the *S. ratti* chromosomes represents position along the chromosome for chromosome I (yellow-red), chromosome II (blue-purple) and chromosome X (green). Black boxes represent scaffolds >1 Mb in length; scaffolds <1 Mb in length are grouped together and shown in gray. **(b)** The mitochondrial gene order and phylogeny for our six species (highlighted in bold) and seven outgroup species that encompass four nematode clades. Our eighth outgroup species, *Meloidogyne hapla*, was excluded because of insufficient mitochondrial genome data. Inverted sequences are shown by gene boxes with inverted text. The maximum-likelihood tree (left) was constructed using 12 mitochondrial proteins. Amino acid sequences were aligned before concatenation, and then phylogenetic analysis was performed with RAXML v7.2.8 using the best-fitting empirical model of amino acid substitution with 1,000 bootstrap resampling replicates with the percentage support shown on the nodes (**Supplementary Note**). The scale bar shows the number of amino acid substitutions per site.

and with eight outgroup species, showing that this expansion accompanies the evolution of parasitism (Fig. 1 and **Supplementary Table 10**). Among the outgroup species, the hookworm *Necator americanus*³⁶ has 82 astacin-encoding genes and the free-living *C. elegans* has 40 (ref. 34).

SCP/TAPS proteins are often immunomodulatory molecules in parasitic nematodes³⁵ and have been investigated as potential vaccine candidates against *N. americanus*^{37,38}. We found 89–205 SCP/TAPS-encoding genes in the *Strongyloides* genomes, including nine subfamilies not present in *P. trichosuri*, *Rhabditophanes* sp. KR3021 or the eight outgroup species (**Supplementary Tables 8 and 10**). In *N. americanus*, there are 137 SCP/TAPS-encoding genes³⁶, suggesting

that this gene family has independently expanded twice, in nematode clades IV and V.

Additional gene expansions included receptor-type protein tyrosine phosphatases, which have a putative role in signaling³⁹ and are expanded in *Strongyloides* and *Parastrongyloides* (52–75 genes) compared with *Rhabditophanes* (13 genes) and the eight outgroup species (up to 39 genes). Acetylcholinesterase-encoding genes were expanded in *Strongyloides* and *Parastrongyloides* (30–126 genes) compared with *Rhabditophanes* (1 gene) and our outgroup species (1–5 genes). Many parasitic nematodes secrete acetylcholinesterases, which are thought to facilitate their maintenance in hosts⁴⁰, and the expansion of this gene family in these parasitic species is consistent with this

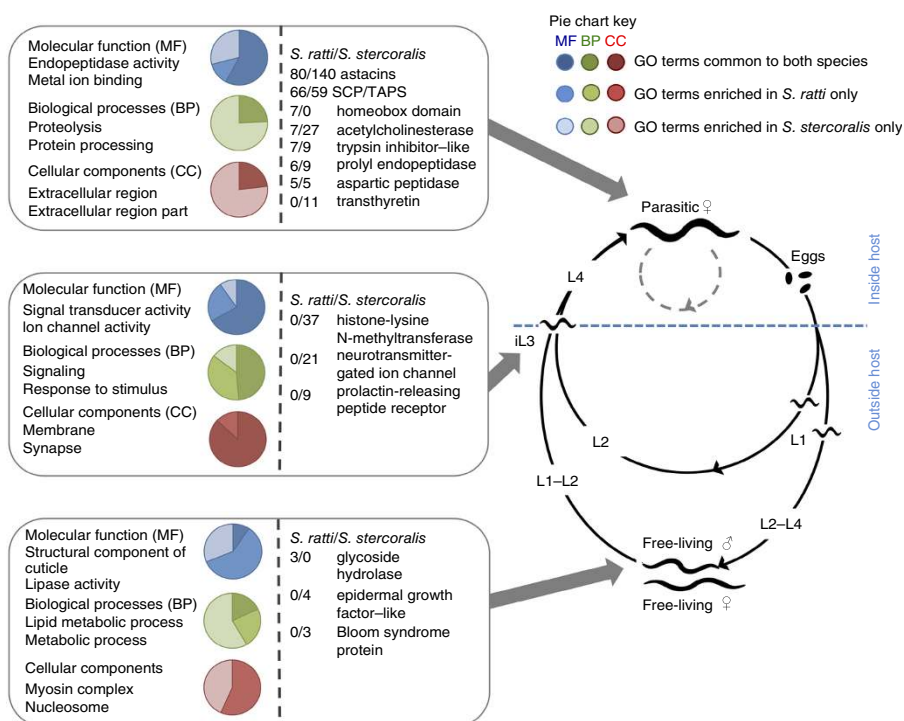
Table 1 Properties of the genome assemblies

	<i>S. ratti</i>	<i>S. stercoralis</i>	<i>S. papillosus</i>	<i>S. venezuelensis</i>	<i>P. trichosuri</i>	<i>Rhabditophanes</i> sp. KR3021	<i>C. elegans</i>
Clade	IV	IV	IV	IV	IV	IV	V
Number of chromosomes	3 (ref. 72)	3 (ref. 73)	2 (ref. 74)	2 (ref. 21)	3 (ref. 22)	5 ^a	6 (ref. 16)
Assembly version	V5.0.4	V2.0.4	V2.1.4	V2.0.4	V2.0.4	V2.0.4	WS244
Assembly size (Mb)	43.1	42.6	60.2	52.1	42.2	47.2	100.2
Number of scaffolds	115 ^b	675	4,353	520	1,391	380	6
N50 of scaffolds (kb)	11,700	431	86	715	837	537	17,500
N50 (number)	2	16	129	16	12	22	3
Maximum scaffold length (Mb)	16.8	5.0	1.7	5.9	6.2	7.3	20.9
GC content (%)	21	22	26	25	31	32	36
Number of genes	12,451	13,098	18,457	16,904	15,010	13,496	23,629
Number of exons	33,796	34,366	40,821	40,619	35,049	37,987	145,275
Exons, combined length (Mb)	17.5	17.9	22.4	20.3	20.8	17.8	30.1
Median exon length (bp)	263	265	304	261	348	276	146
Number of introns	21,345	21,268	22,364	23,715	20,039	24,491	169,506

Genome statistics are based on scaffolds, excluding scaffolds less than 1,000 bp in length. N50 is the size above which 50% of the assembled bases are distributed; N50 (number) is the number of scaffolds in which 50% of assembled bases exist.

^aSee **Supplementary Figure 7**. ^bTwelve scaffolds, covering 93% of the genome, are assigned to chromosomes; 103 scaffolds are not assigned to a chromosome.

Figure 3 The parasitic female, free-living female and iL3 transcriptomes of *Strongyloides* species. The progeny of the parasitic female pass out of the host (as larvae for *S. stercoralis* or as eggs and larvae for *S. ratti*), where iL3s can develop directly or free-living males and females develop whose progeny develop into iL3s; iL3s then infect hosts. The parasite of humans *S. stercoralis* can undergo internal autoinfection (gray dashed line) where iL3s develop and internally reinfect the same host. The transcriptome of the parasitic female, free-living female and iL3 were compared for *S. ratti* and *S. stercoralis*. Representative GO terms that were significantly enriched (left side of the box) and Ensembl Compara gene families significantly upregulated (right side of the box) for each of these three stages of the life cycle are summarized. The pie charts show the proportion of the GO terms common to *S. ratti* and *S. stercoralis* or unique to each species. Numbers on the right in the boxes represent the number of genes upregulated in each gene family for *S. ratti* and *S. stercoralis*. MF, molecular function; BP, biological process; CC, cellular component.



role. Some families show subclade-specific expansion; for instance, *S. papillosus* and *S. venezuelensis* have a paralogous expansion of genes encoding Speckle-type POZ domains⁴¹ (92–130 genes) compared with *S. ratti* and *S. stercoralis* (9 or 10 genes) (Fig. 1 and Supplementary Table 8).

No function or annotation could be assigned to approximately one-third (26–37%) of the genes present in the six species, but 50% of these genes could be assigned to novel gene families. The six largest of these families occurred only in *Strongyloides* and *Parastrongyloides*, comprising a total of 630 genes. We have named these *Strongyloides* genome project families *sgpf-1* to *sgpf-6*. Members of *sgpf-1* and *sgpf-5* are predicted to encode proteins with signal peptides that are highly glycosylated (Supplementary Table 11).

Expanded gene families are upregulated in parasitic stages

We identified genes and gene families that are likely to have a key role in the parasitic lifestyle of *S. ratti* and *S. stercoralis* by comparing the transcriptomes of parasitic and free-living female stages. We generated *S. ratti* transcriptome data and used previously published *S. stercoralis* data⁴². A total of 909 *S. ratti* and 1,188 *S. stercoralis* genes were upregulated in parasitic females as compared with free-living females (edgeR, fold change > 2, false discovery rate (FDR) < 0.01; Supplementary Tables 12 and 13), of which 423 *S. ratti* and 457 *S. stercoralis* orthologous genes were upregulated in the parasitic female stage of both species (Supplementary Table 14).

The two most expanded *Strongyloides* gene families—encoding SCP/TAPS³⁵ and astacin-domain^{43–46} proteins—dominated the list of genes differentially expressed by parasitic females. In *S. ratti* and *S. stercoralis*, respectively, 58 and 62% of putative astacin-like genes and 57 and 71% of SCP/TAPS genes were differentially expressed between parasitic and free-living females (Fig. 3 and Supplementary Tables 10 and 13). However, other paralogously expanded genes were not enriched among the upregulated genes, suggesting that they may not be important for parasitism. Both *Strongyloides* and *Parastrongyloides* infect their hosts by skin penetration; the larvae then migrate through the host, and

adult females in the host live in the mucosa of the small intestine^{47,48} where they feed on the host. Astacins are metallopeptidases that have previously been associated with a role in tissue migration by nematode infective larvae^{44,49}. Around half of the putative astacin-like proteins in *Strongyloides* species contain the canonical zinc-binding motif (HEXXHXXGXXH) of astacin active sites and likely have a role in penetrating the host mucosa in which parasitic females live. Teasing apart the role of different astacin gene family members in the migration and gut-dwelling phases of this life cycle could provide insights to allow new therapeutic interventions to be developed. For *S. ratti* and *S. stercoralis*, respectively, 63 and 53% of the SCP/TAPS genes upregulated in parasitic females encode a signal peptide, suggesting that the proteins may be secreted from the worm into the host. An immunomodulatory role for SCP/TAPS proteins has also been proposed on the basis of the inhibitory effect that these proteins have on neutrophil and platelet activity in hookworm infections^{35,50,51}.

Other gene families commonly upregulated in the parasitic females of both species, as compared with free-living females and iL3s, included ones encoding transthyretin-like proteins, prolyl endopeptidases, acetylcholinesterases, trypsin inhibitors and aspartic peptidases (Fig. 3 and Supplementary Table 15). The transthyretin-like genes had some of the highest fold changes in expression of genes upregulated in parasitic females (Supplementary Table 13). Transthyretin-like genes constitute a large, nematode-specific gene family⁵², are expressed in adult parasitic stages^{53–55} and are distant relatives of the vertebrate transthyretins that are involved in transporting thyroid hormones⁵⁶. While some aspartic peptidases are essential for the digestion of host hemoglobin in blood-borne parasites^{57,58}, it has been proposed that others are involved in digesting other host macromolecules⁵⁹.

Hypothetical protein-coding genes accounted for 20–37% of the differentially expressed genes from pairwise comparisons of parasitic females, free-living females and iL3s, and these included genes with the highest relative expression levels (Supplementary Table 13). These novel genes are likely to be important to these distinctive phases of

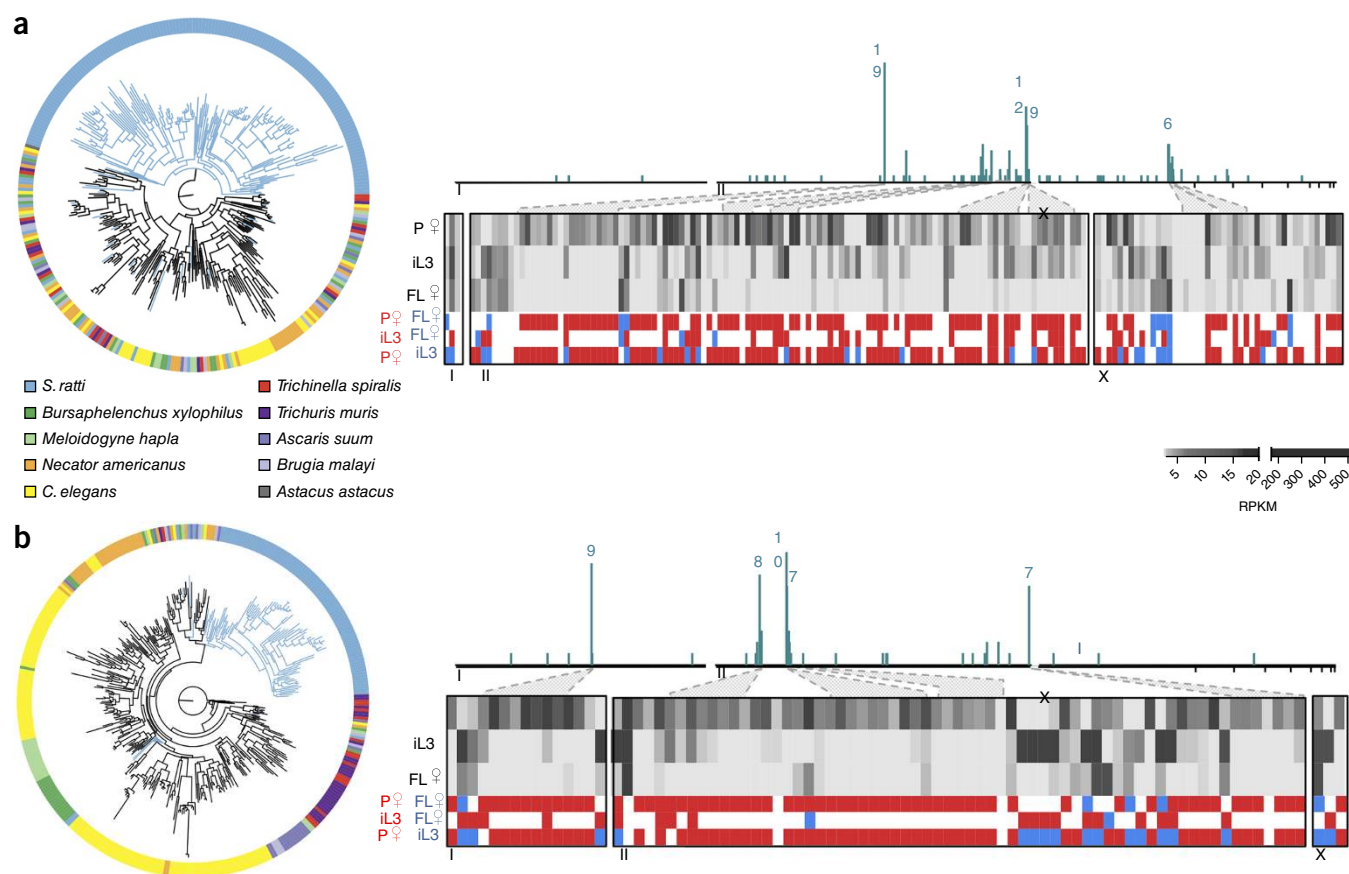


Figure 4 *Strongyloides*-specific expansion and chromosomal clustering of gene families. **(a,b)** Astacin-like **(a)** and SCP/TAPS **(b)** genes are the two major *S. ratti* gene families upregulated in the transcriptome of the parasitic female. Left, the phylogeny of each of these gene families for *S. ratti*, our eight outgroup species and the crayfish *Astacus astacus*. *S. ratti* genes are in light blue. Right, the distribution of these genes in the genome, plotted as clusters of physically adjacent genes in the genome. Numbers above the peaks are the number of genes in a cluster of physically neighboring genes; ticks below the axis denote scaffold boundaries for chromosome X. The transcriptomic expression of these genes (in RPKM, reads per kilobase per million mapped reads) for parasitic female (P), free-living female (FL) and iL3 is shown in grayscale, and the results of pairwise edgeR analyses of gene expression for these life cycle stages are shown in color where a gene is upregulated. The color representing upregulation (red or blue) relates to the color of the name of the life cycle stage for each pairwise comparison (fold change > 2, FDR < 0.01); no differential expression is shown as a white block.

the life cycle, including in parasitism. Three small novel gene families (*sgpf-7* to *sgpf-9*) were predominantly upregulated in *S. ratti* parasitic female, with two of the genes predicted to encode predominantly secretory or membrane-targeted proteins (**Supplementary Table 11**). In contrast, the largest hypothetical protein-coding gene families, *sgpf-1* to *sgpf-6*, accounted for only a small proportion (1% in both *S. ratti* and *S. stercoralis*) of all differentially expressed hypothetical protein-coding genes, suggesting that they do not have roles in parasitism.

Using gene ontology (GO) annotations to summarize the putative functions of the upregulated genes identified distinct differences between the life cycle stages of both species (**Fig. 3** and **Supplementary Table 16**). The genes upregulated in iL3s appear to be associated with sensing the environment and with signal transduction and were the most consistent between *S. ratti* and *S. stercoralis*. The products of genes expressed in free-living females have core metabolic and growth-related roles (such as in cytoskeleton and chromatin). In parasitic stages, the dominant functional categories were proteases, consistent with the abundant astacins (**Fig. 3** and **Supplementary Table 16**).

The products of putative parasitism genes are secreted

In parallel, we compared the somatic proteomes of parasitic and free-living females of *S. ratti*. Of 1,266 proteins detected overall, 569

were comparatively upregulated in parasitic females and 409 were comparatively upregulated in free-living females (**Supplementary Tables 12** and **17**). We found a modest overlap between the transcriptome and somatic proteome: 6% of genes upregulated in the parasitic female transcriptome were also upregulated in the proteome, and 10% of genes upregulated in the free-living female transcriptome were also upregulated in the proteome (**Supplementary Fig. 5** and **Supplementary Table 18**). A poor concordance between transcript and peptide abundance has been reported in many systems^{60–62} and likely reflects post-translational processes that decouple protein and mRNA abundance. In the present study, this may be compounded by the excretion and/or secretion of many gene products from parasitic stages to allow these proteins to interact with the host. Indeed, 43% of genes upregulated in the parasitic female transcriptome are predicted to encode signal peptides, compared with 26% of genes upregulated in the free-living female. Furthermore, while several of the putative parasitism gene families were highly upregulated in the somatic proteome (aspartic peptidases, prolyl endopeptidases and acetylcholinesterases; **Supplementary Table 17**), we found only five astacin-like and no SCP/TAPS proteins (**Supplementary Fig. 5**). To address this, we extended the analysis to the excretory/secretory (ES) proteome data of Soblik *et al.*⁶³.

In the ES proteome, we detected an additional 882 proteins and found greater consistency with the parasitic female transcriptome: 13% of the parasitic female ES proteins overlapped with the genes upregulated in the transcriptome (**Supplementary Table 18**). We also found 25 astacin and 14 SCP/TAPS gene products in the ES proteome. Other gene families highly upregulated in the parasitic female transcriptome were also dominant in the parasitic ES proteome, including prolyl endopeptidases, acetylcholinesterases and transthyretin-like proteins (**Supplementary Table 19**). Protein products of the novel gene families *sgpf-1* and *sgpf-5* were also identified in the ES products of both parasitic and free-living females (**Supplementary Table 11**). Other parasitic nematodes have been noted to have many protease-encoding genes, and different species appear to have expanded different protease families^{36,64–66}. Together, these and our findings suggest that expansion of protease-encoding genes and secretion of extensive quantities of proteases is likely to be an essential feature of nematode parasitism. These proteases are, presumably, used to penetrate host tissue, acquire resources from the host and protect the parasite from host-induced harm.

Parasitism-associated genes are in coexpressed clusters

We observed that genes upregulated in the parasitic females and iL3s were often physically clustered in the genome, more so than for genes upregulated in the free-living females (**Supplementary Table 20**). To test whether this clustering was significant, we asked whether clusters of three or more adjacent genes, upregulated in the same life cycle stage, occurred more often than would be expected by chance. We found that 31%, 4% and 26% of upregulated genes were in such clusters in *S. ratti* parasitic female, free-living female and iL3, respectively, whereas in *S. stercoralis* this was 34%, 2% and 34% (**Supplementary Table 20**). This clustering is more than would be expected by chance (**Supplementary Fig. 6** and **Supplementary Table 20**). The clusters in parasitic females were larger (19 and 16 genes in the largest *S. ratti* and *S. stercoralis* clusters, respectively) than those of the iL3 (9 and 14 genes) and free-living females (3 genes) (**Supplementary Table 20**). Although nematodes, including *S. ratti*⁶⁷, have operons, these clusters are unlikely to be operons because (i) the average intergenic distance among clustered genes does not differ from the genome-wide average (**Supplementary Fig. 6**) and (ii) cluster members include genes on both strands.

Clusters of genes upregulated in the parasitic female were more likely to comprise genes from the same gene family. The majority (88 and 73% for *S. ratti* and *S. stercoralis*, respectively) of these parasitic female clusters were of genes belonging to the same *Compara* gene family; this is greater than that observed for iL3 (8–10%) (**Supplementary Tables 20–22**). Two gene families dominated parasitic female clusters: astacins (24 and 23% of parasitic female clusters for *S. ratti* and *S. stercoralis*, respectively) and SCP/TAPS (15 and 11%). Tandem expansions of astacin and SCP/TAPS genes could provide a plausible explanation for the preponderance of these gene families in the parasitic female expression clusters. However, even with exclusion of the astacin and SCP/TAPS families, most remaining parasitic female clusters still comprised genes from the same gene family (85 and 65% for *S. ratti* and *S. stercoralis*, respectively); fewer clusters from the same gene family occurred for iL3 (7 and 9%) compared to parasitic female (**Supplementary Table 21**).

Phylogenetic analysis of astacins, including those from the eight outgroup species, showed that 139 *S. ratti* genes form one distinct clade (**Fig. 4**), presumably derived from a single ancestral astacin gene. Similarly, the *S. ratti* SCP/TAPS gene family has almost exclusively expanded from one ancestral gene (**Fig. 4**). These gene clusters likely

arose by tandem duplication of genes, as has occurred for other large gene families, for example in *C. elegans*¹⁸. However, in contrast to *C. elegans*, physical adjacency of the duplicated genes has been maintained in *Strongyloides*, perhaps as a result of the expansions being recent and therefore not yet broken up by recombination. Alternatively, the adjacency may be functional, for example, if there is pressure to maintain a common regulatory environment. Clustering of gene families was relatively rare among *Rhabditophanes* sp. KR3102 and the eight outgroup species (**Supplementary Table 21**), meaning that this clustering is specific to the *Strongyloides-Parastrongyloides* lineage and thus to the parasitic lifestyle in this clade.

The clusters of genes upregulated in the parasitic females were themselves chromosomally clustered, forming ‘parasitism regions’ (**Fig. 4**). In *S. ratti*, one-third of genes upregulated in the parasitic female are concentrated in three regions of chromosome II, most notably in a 3.6-Mb region at one end of the chromosome, comprising 171 genes that were upregulated in the parasitic female transcriptome (**Supplementary Fig. 2**). A similar pattern is evident in *S. stercoralis*, where seven scaffolds and contigs with a high density of genes upregulated in the parasitic female also belong to chromosome II; 46% of the 171 genes upregulated in *S. ratti* belong to just eight different gene families, including those encoding aspartic peptidases, astacin-like proteins, SCP/TAPS proteins, transthyretin-like proteins and trypsin inhibitor-like proteins. This is the first report, to our knowledge, of chromosomal clustering of genes likely to be important in nematode parasitism, and this clustering hints at possible regulatory mechanisms for parasite development.

DISCUSSION

Understanding the molecular and genetic differences between parasitic and free-living organisms is of fundamental biological interest and is essential to identifying novel drug targets and other methods to control parasitic nematodes and the diseases that they cause. We have undertaken a comparative genomics study of six taxa from an evolutionary clade that transitions from a free-living to a parasitic lifestyle, which we combined with transcriptomic and proteomic analyses of parasitic and free-living female stages of *Strongyloides* species. Together, this is a powerful way to discover the molecular adaptations to parasitism among these nematodes. We find that a preponderance of the genes that are expanded in parasitic species are specifically used in the parasitic stages and are within genomic clusters, concentrated in regions of chromosome II. This is consistent with the idea that the within-host stages of parasitic nematodes deploy a specific biology that enables them to be successful parasites. The *Strongyloides* proteome and transcriptome have limited overlap, as has been observed in other systems. For the *Strongyloides* clade, we find that astacin- and SCP/TAPS-encoding genes are prominent among parasitism-associated genes. Other parasitic nematodes appear to have expanded the number of protease-encoding genes in their genome, which also appear to be used predominantly during the within-host stages. In *Strongyloides*, we have also found genomic clustering of these and other likely parasitism-associated genes, which is likely to have been initiated during the adaptation to parasitism, followed by subsequent repeated gene duplication, associated with adaptation to different hosts. This genomic arrangement may facilitate the expression of a parasitic transcriptional program by these parasites. Operons have been demonstrated in *Strongyloides*, and it will be important to determine whether these parasitism-associated genes are under operonic control.

Strongyloides is a particularly amenable laboratory system—both *S. ratti* and *S. venezuelensis* can be maintained in the laboratory in

their natural rat host, as well as in other rodents, and the parasite of humans *S. stercoralis* can also be maintained in the laboratory. In addition to providing a compelling model of the evolution of parasitism, transgenesis of *Strongyloides* and *Parastrongyloides* is possible^{68–71} uniquely among parasitic nematodes, which will allow functional genomic studies, directed by our findings, to further explore the genetic basis of nematode parasitism.

URLs. WormBase-ParaSite, <http://parasite.wormbase.org/>; World Health Organization soil-transmitted helminthiasis, http://www.who.int/gho/neglected_diseases/soil_transmitted_helminthiasis/en/; World Health Organization estimates of disease burden 2000–2012, http://www.who.int/healthinfo/global_burden_disease/estimates/en/index2.html; RepeatModeler, <http://www.repeatmasker.org/RepeatModeler.html>; TransposonPSI, <http://transposonpsi.sourceforge.net/>; SMALT, <http://www.sanger.ac.uk/resources/software/smalt/>; PHYLP package, <http://evolution.genetics.washington.edu/phylip.html>.

METHODS

Methods and any associated references are available in the [online version of the paper](#).

Accession codes. The *S. ratti*, *S. stercoralis*, *S. papillosus*, *S. venezuelensis*, *P. trichosuri* and *Rhabditophanes* genome assemblies, predicted transcripts, protein and annotation (.GFF) files are available from WormBase-ParaSite and are registered with the European Nucleotide Archive (ENA) under BioProject accessions [PRJEB125](#) (*S. ratti*_ED321_v5_0_4), [PRJEB528](#) (*S. stercoralis*_PV0001_v2_0_4), [PRJEB525](#) (*S. papillosus*_LIN_v2_1_4), [PRJEB530](#) (*S. venezuelensis*_HH1_v2_0_4), [PRJEB515](#) (*P. trichosuri*_KNP_v2_0_4) and [PRJEB1297](#) (*Rhabditophanes*_sp_KR3021_v2_0_4). The raw genomic data are available from the ENA via the accessions detailed in **Supplementary Table 23**. The transcriptomic data for *S. ratti* are available from ArrayExpress under accessions [E-ERAD-151](#) and [E-ERAD-92](#). For *S. venezuelensis*, transcriptomic data are available from the DNA Databank of Japan (DDBJ) under BioProject accession [PRJDB3457](#) (*S. venezuelensis*) (**Supplementary Table 24**).

Note: Any Supplementary Information and Source Data files are available in the online version of the paper.

ACKNOWLEDGMENTS

We thank, from the Wellcome Trust Sanger Institute, C. Griffiths, D. Willey, R. Rance and DNA Pipelines; J. Keane and D. Gordon for bioinformatics support; M. Dunn for the *S. venezuelensis* optical map; A. Babbage for laboratory support; and M. Zarowiecki for gene finding and functional annotation advice. We thank for technical help L. Hughes and L. Weldon (University of Bristol); H. Massey, Jr., X. Li and H. Shao (University of Pennsylvania); D.K. Howe and R.I. Wernick (Oregon State University); H. Denise (European Bioinformatics Institute); M. Yabana (Tokyo Institute of Technology); and A. Hino and R. Tanaka (University of Miyazaki) and A. Toyoda (National Institute of Genetics) for sequencing. The *S. ratti* transcriptome and proteome work was funded by Wellcome Trust grant 094462/Z/10/Z awarded to M.V., J.W. and M.B. The *S. ratti*, *S. stercoralis*, *S. papillosus*, *P. trichosuri* and *Rhabditophanes* sp. KR3021 genome sequencing and the *S. venezuelensis* optical mapping were funded by Wellcome Trust grant 098051. The *S. venezuelensis* work was supported by Japan Society for the Promotion of Science (JSPS) KAKENHI (24310142, 21590466 and 24780044), KAKENHI for Innovative Areas 'Genome Science' (221S0002) and the Integrated Research Project for Human and Veterinary Medicine of the University of Miyazaki. I.J.T. was supported by Academia Sinica. Work was funded by grants AI050668 and AI105856 from the US National Institutes of Health (NIH) to J.B.L. and by Resource-Related Research Grant RR02512 from the US NIH to M. Haskins, which provided research materials for the study. J.D.S. received support from US NIH training grant AI060516. A.K. was supported by a predoctoral stipend from the Max Planck Society. Work by A.K., D.H. and A.S. was funded by the Max Planck Society.

AUTHOR CONTRIBUTIONS

Cultivated and collected parasite material: V.L.H., D.R.D., W.G., J.B.L., E.N., H.M., A.S., A.K., J.D.S., D.H., T.K. and M.V. Prepared DNA, RNA and protein: V.L.H., N.R., H.M.B. and T.K. Prepared libraries: M.A.Q., H.M.B., E.N., Y.O. and J.D.S. Assembled the genomes: I.J.T., A.S.-F., R.K. and T.I. Quality checked the genomes: A.C. Provided genetic markers and mapping data: A.K., D.H. and A.S. Manually improved the genomes: A.T., H.B., K.B., S.N. and I.J.T. Predicted the genes: E.J.S., A.C., B.J.F. and I.J.T. Functionally annotated the genome: A.C., D.M.R. and B.H. Curated gene models: A.T., H.B., K.B., S.N., I.J.T. and V.L.H. Built a Compara database: B.H. Analyzed gene structure: I.J.T. Undertook proteomics and initial analysis: N.R., D.X., N.W.B., J.W. and V.L.H. Undertook ES work and analyzed the data: N.W.B., H.S., D.X. and V.L.H. Analyzed the transcriptome: V.L.H., I.J.T., B.J.F., A.J.R. and J.D.S. Analyzed the gene clusters: D.M.R. and V.L.H. Analyzed synteny and chromosome alignments: I.J.T. and A.J.R. Assembled and analyzed mitochondrial genomes: T.K. and I.J.T. Analyzed chromatin diminution: B.J.F., I.J.T. and A.C. Analyzed gene family clustering: A.J.R., J.A.C. and I.J.T. Coordinated the project, managed sequencing, assembly and finishing: N.H., T.H. and T.K. Wrote the manuscript: V.L.H., I.J.T., A.C., A.J.R., N.H., T.K., M.V. and M.B. Conceived the project: M.V., M.B., J.W., I.J.T. and T.K. Directed the project: M.V. and M.B.

COMPETING FINANCIAL INTERESTS

The authors declare no competing financial interests.

Reprints and permissions information is available online at <http://www.nature.com/reprints/index.html>.



This work is licensed under a Creative Commons Attribution-NonCommercial-ShareAlike 3.0 Unported License. The images or other third party material in this article are included in the article's Creative Commons license, unless indicated otherwise in the credit line; if the material is not included under the Creative Commons license, users will need to obtain permission from the license holder to reproduce the material. To view a copy of this license, visit <http://creativecommons.org/licenses/by-nc-sa/3.0/>.

- Savioli, L. & Albonico, M. Soil-transmitted helminthiasis. *Nat. Rev. Microbiol.* **2**, 618–619 (2004).
- Pullan, R.L. & Brooker, S.J. The global limits and population at risk of soil-transmitted helminth infections in 2010. *Parasit. Vectors* **5**, 81 (2012).
- Viney, M.E. & Lok, J.B. in *WormBook* (ed. The C. elegans Research Community) doi:10.1895/wormbook.1.141.2 (2015).
- Albonico, M., Crompton, D.W. & Savioli, L. Control strategies for human intestinal nematode infections. *Adv. Parasitol.* **42**, 277–341 (1999).
- Crompton, D.W.T. in *Bailliere's Clinical Tropical Medicine and Communicable Diseases* (ed. Pawlowski, Z.S.) 489–510 (Academic Press, 1987).
- Dorris, M., Viney, M.E. & Blaxter, M.L. Molecular phylogenetic analysis of the genus *Strongyloides* and related nematodes. *Int. J. Parasitol.* **32**, 1507–1517 (2002).
- Blaxter, M., Koutsovoulos, G., Jones, M., Kumar, S. & Elsworth, B. in *Next-Generation Systematics*. (ed. Cotton, J., Hughes, J. & Olson, P.) (in the press).
- Holovachov, O. et al. A phylogenetic tree of nematodes based on about 1200 full-length small subunit ribosomal DNA sequences. *Nematology* **11**, 927–950 (2009).
- Blaxter, M.L. et al. A molecular evolutionary framework for the phylum Nematoda. *Nature* **392**, 71–75 (1998).
- Viney, M.E. A genetic analysis of reproduction in *Strongyloides ratti*. *Parasitology* **109**, 511–515 (1994).
- Viney, M.E., Matthews, B.E. & Walliker, D. Mating in the nematode parasite *Strongyloides ratti*: proof of genetic exchange. *Proc. Biol. Soc.* **254**, 213–219 (1993).
- Tindall, N.R. & Wilson, P.A.G. An extended proof of migration routes of immature parasites inside hosts: pathways of *Nippostrongylus brasiliensis* and *Strongyloides ratti* in the rat are mutually exclusive. *Parasitology* **100**, 281–288 (1990).
- Mackerras, M. *Strongyloides* and *Parastrongyloides* (Nematoda: Rhabdiasoidea) in Australian marsupials. *Aust. J. Zool.* **7**, 87 (1959).
- Grant, W.N. et al. *Parastrongyloides trichosuri*, a nematode parasite of mammals that is uniquely suited to genetic analysis. *Int. J. Parasitol.* **36**, 453–466 (2006).
- Nemetschke, L., Eberhardt, A.G., Viney, M.E. & Streit, A. A genetic map of the animal-parasitic nematode *Strongyloides ratti*. *Mol. Biochem. Parasitol.* **169**, 124–127 (2010).
- C. elegans Sequencing Consortium. Genome sequence of the nematode *C. elegans*: a platform for investigating biology. *Science* **282**, 2012–2018 (1998).
- Hillier, L.W. et al. Comparison of *C. elegans* and *C. briggsae* genome sequences reveals extensive conservation of chromosome organization and synteny. *PLoS Biol.* **5**, e167 (2007).
- Foth, B.J. et al. Whipworm genome and dual-species transcriptome analyses provide molecular insights into an intimate host-parasite interaction. *Nat. Genet.* **46**, 693–700 (2014).
- Kikuchi, T. et al. Genomic insights into the origin of parasitism in the emerging plant pathogen *Bursaphelenchus xylophilus*. *PLoS Pathog.* **7**, e1002219 (2011).

20. Nemetschke, L., Eberhardt, A.G., Hertzberg, H. & Streit, A. Genetics, chromatin diminution, and sex chromosome evolution in the parasitic nematode genus *Strongyloides*. *Curr. Biol.* **20**, 1687–1696 (2010).
21. Hino, A. *et al.* Karyotype and reproduction mode of the rodent parasite *Strongyloides venezuelensis*. *Parasitology* **141**, 1736–1745 (2014).
22. Kulkarni, A., Dyka, A., Nemetschke, L., Grant, W.N. & Streit, A. *Parastrongyloides trichosuri* suggests that XX/XO sex determination is ancestral in Strongyloididae (Nematoda). *Parasitology* **140**, 1822–1830 (2013).
23. Harvey, S.C. & Viney, M.E. Sex determination in the parasitic nematode *Strongyloides ratti*. *Genetics* **158**, 1527–1533 (2001).
24. Hu, M., Chilton, N.B. & Gasser, R.B. The mitochondrial genome of *Strongyloides stercoralis* (Nematoda)—idiosyncratic gene order and evolutionary implications. *Int. J. Parasitol.* **33**, 1393–1408 (2003).
25. Armstrong, M.R., Blok, V.C. & Phillips, M.S. A multipartite mitochondrial genome in the potato cyst nematode *Globodera pallida*. *Genetics* **154**, 181–192 (2000).
26. Gibson, T. *et al.* The mitochondrial subgenomes of the nematode *Globodera pallida* are mosaics: evidence of recombination in an animal mitochondrial genome. *J. Mol. Evol.* **64**, 463–471 (2007).
27. Lunt, D.H. & Hyman, B.C. Animal mitochondrial DNA recombination. *Nature* **387**, 247 (1997).
28. Humphreys-Pereira, D.A. & Elling, A.A. Mitochondrial genome plasticity among species of the nematode genus *Meloidogyne* (Nematoda: Tylenchida). *Gene* **560**, 173–183 (2015).
29. Piganeau, G., Gardner, M. & Eyre-Walker, A. A broad survey of recombination in animal mitochondria. *Mol. Biol. Evol.* **21**, 2319–2325 (2004).
30. Ballard, J.W.O. & Whitlock, M.C. The incomplete natural history of mitochondria. *Mol. Ecol.* **13**, 729–744 (2004).
31. Sultana, T. *et al.* Comparative analysis of complete mitochondrial genome sequences confirms independent origins of plant-parasitic nematodes. *BMC Evol. Biol.* **13**, 12 (2013).
32. Sun, L., Zhuo, K., Lin, B., Wang, H. & Liao, J. The complete mitochondrial genome of *Meloidogyne graminicola* (Tylenchida): a unique gene arrangement and its phylogenetic implications. *PLoS One* **9**, e98558 (2014).
33. Vilella, A.J. *et al.* EnsemblCompara GeneTrees: complete, duplication-aware phylogenetic trees in vertebrates. *Genome Res.* **19**, 327–335 (2009).
34. Park, J.-O. *et al.* Characterization of the astacin family of metalloproteases in *C. elegans*. *BMC Dev. Biol.* **10**, 14 (2010).
35. Cantacessi, C. *et al.* A portrait of the “SCP/TAPS” proteins of eukaryotes—developing a framework for fundamental research and biotechnological outcomes. *Biotechnol. Adv.* **27**, 376–388 (2009).
36. Tang, Y.T. *et al.* Genome of the human hookworm *Necator americanus*. *Nat. Genet.* **46**, 261–269 (2014).
37. Bethony, J.M. *et al.* Randomized, placebo-controlled, double-blind trial of the Na-ASP-2 hookworm vaccine in unexposed adults. *Vaccine* **26**, 2408–2417 (2008).
38. Goud, G.N. *et al.* Expression of the *Necator americanus* hookworm larval antigen Na-ASP-2 in *Pichia pastoris* and purification of the recombinant protein for use in human clinical trials. *Vaccine* **23**, 4754–4764 (2005).
39. Lemmon, M.A. & Schlessinger, J. Cell signaling by receptor tyrosine kinases. *Cell* **141**, 1117–1134 (2010).
40. Selkirk, M.E., Lazari, O. & Matthews, J.B. Functional genomics of nematode acetylcholinesterases. *Parasitology* **131** (suppl.), S3–S18 (2005).
41. Kwon, J.E. *et al.* BTB domain-containing speckle-type POZ protein (SPOP) serves as an adaptor of Daxx for ubiquitination by Cul3-based ubiquitin ligase. *J. Biol. Chem.* **281**, 12664–12672 (2006).
42. Stoltzfus, J.D., Minot, S., Berriman, M., Nolan, T.J. & Lok, J.B. RNAseq analysis of the parasitic nematode *Strongyloides stercoralis* reveals divergent regulation of canonical dauer pathways. *PLoS Negl. Trop. Dis.* **6**, e1854 (2012).
43. Jing, Y., Toubarro, D., Hao, Y. & Simões, N. Cloning, characterisation and heterologous expression of an astacin metalloprotease, Sc-AST, from the entomoparasitic nematode *Steinernema carpocapsae*. *Mol. Biochem. Parasitol.* **174**, 101–108 (2010).
44. Williamson, A.L. *et al.* *Ancylostoma caninum* MTP-1, an astacin-like metalloprotease secreted by infective hookworm larvae, is involved in tissue migration. *Infect. Immun.* **74**, 961–967 (2006).
45. Lun, H.M., Mak, C.H. & Ko, R.C. Characterization and cloning of metalloproteinase in the excretory/secretory products of the infective-stage larva of *Trichinella spiralis*. *Parasitol. Res.* **90**, 27–37 (2003).
46. Semenova, S.A. & Rudenskaya, G.N. The astacin family of metalloproteinases. *Biochem. Suppl. Ser. B: Biomed. Chem.* **3**, 17–32 (2009).
47. Maruyama, H., El-Malky, M., Kumagai, T. & Ohta, N. Secreted adhesion molecules of *Strongyloides venezuelensis* are produced by oesophageal glands and are components of the wall of tunnels constructed by adult worms in the host intestinal mucosa. *Parasitology* **126**, 165–171 (2003).
48. Maruyama, H., Yabu, Y., Yoshida, A., Nawa, Y. & Ohta, N. A role of mast cell glycosaminoglycans for the immunological expulsion of intestinal nematode, *Strongyloides venezuelensis*. *J. Immunol.* **164**, 3749–3754 (2000).
49. Gomez Gallego, S. *et al.* Identification of an astacin-like metallo-proteinase transcript from the infective larvae of *Strongyloides stercoralis*. *Parasitol. Int.* **54**, 123–133 (2005).
50. Moyle, M. *et al.* A hookworm glycoprotein that inhibits neutrophil function is a ligand of the integrin CD11b/CD18. *J. Biol. Chem.* **269**, 10008–10015 (1994).
51. Del Valle, A., Jones, B.F., Harrison, L.M., Chadderdon, R.C. & Cappello, M. Isolation and molecular cloning of a secreted hookworm platelet inhibitor from adult *Ancylostoma caninum*. *Mol. Biochem. Parasitol.* **129**, 167–177 (2003).
52. Parkinson, J. *et al.* A transcriptomic analysis of the phylum Nematoda. *Nat. Genet.* **36**, 1259–1267 (2004).
53. Saverwyns, H. *et al.* Analysis of the transthyretin-like (TTL) gene family in *Ostertagia ostertagi*—comparison with other strongylid nematodes and *Caenorhabditis elegans*. *Int. J. Parasitol.* **38**, 1545–1556 (2008).
54. Jacob, J., Vanholme, B., Haegeman, A. & Gheysen, G. Four transthyretin-like genes of the migratory plant-parasitic nematode *Radopholus similis*: members of an extensive nematode-specific family. *Gene* **402**, 9–19 (2007).
55. Chehayeb, J.F., Robertson, A.P., Martin, R.J. & Geary, T.G. Proteomic analysis of adult *Ascaris suum* fluid compartments and secretory products. *PLoS Negl. Trop. Dis.* **8**, e2939 (2014).
56. Richardson, S.J., Henneby, S.C., Smith, B.J. & Wright, H.M. Evolution of the thyroid hormone distributor protein transthyretin in microbes, *C. elegans*, and vertebrates. *Ann. NY Acad. Sci.* **1040**, 448–451 (2005).
57. Williamson, A.L. *et al.* Cleavage of hemoglobin by hookworm cathepsin D aspartic proteases and its potential contribution to host specificity. *FASEB J.* **16**, 1458–1460 (2002).
58. Longbottom, D. *et al.* Molecular cloning and characterisation of a putative aspartate proteinase associated with a gut membrane protein complex from adult *Haemonchus contortus*. *Mol. Biochem. Parasitol.* **88**, 63–72 (1997).
59. Mello, L.V., O'Meara, H., Rigden, D.J. & Paterson, S. Identification of novel aspartic proteases from *Strongyloides ratti* and characterisation of their evolutionary relationships, stage-specific expression and molecular structure. *BMC Genomics* **10**, 611 (2009).
60. Foss, E.J. *et al.* Genetic basis of proteome variation in yeast. *Nat. Genet.* **39**, 1369–1375 (2007).
61. Haider, S. & Pal, R. Integrated analysis of transcriptomic and proteomic data. *Curr. Genomics* **14**, 91–110 (2013).
62. Ghazalpour, A. *et al.* Comparative analysis of proteome and transcriptome variation in mouse. *PLoS Genet.* **7**, e1001393 (2011).
63. Soblik, H. *et al.* Life cycle stage-resolved proteomic analysis of the excretome/secretome from *Strongyloides ratti*—identification of stage-specific proteases. *Mol. Cell. Proteomics* **10**, M111.010157 (2011).
64. Laing, R. *et al.* The genome and transcriptome of *Haemonchus contortus*, a key model parasite for drug and vaccine discovery. *Genome Biol.* **14**, R88 (2013).
65. Schwarz, E.M. *et al.* The genome and developmental transcriptome of the strongylid nematode *Haemonchus contortus*. *Genome Biol.* **14**, R89 (2013).
66. Schwarz, E.M. *et al.* The genome and transcriptome of the zoonotic hookworm *Ancylostoma ceylanicum* identify infection-specific gene families. *Nat. Genet.* **47**, 416–422 (2015).
67. Guiliano, D.B. & Blaxter, M.L. Operon conservation and the evolution of trans-splicing in the phylum Nematoda. *PLoS Genet.* **2**, e198 (2006).
68. Shao, H. *et al.* Transposon-mediated chromosomal integration of transgenes in the parasitic nematode *Strongyloides ratti* and establishment of stable transgenic lines. *PLoS Pathog.* **8**, e1002871 (2012).
69. Li, X. *et al.* Successful transgenesis of the parasitic nematode *Strongyloides stercoralis* requires endogenous non-coding control elements. *Int. J. Parasitol.* **36**, 671–679 (2006).
70. Li, X. *et al.* Transgenesis in the parasitic nematode *Strongyloides ratti*. *Mol. Biochem. Parasitol.* **179**, 114–119 (2011).
71. Grant, W.N. *et al.* Heritable transgenesis of *Parastrongyloides trichosuri*: a nematode parasite of mammals. *Int. J. Parasitol.* **36**, 475–483 (2006).
72. Bolla, R.I. & Roberts, L.S. Gametogenesis and chromosomal complement in *Strongyloides ratti* (Nematoda: Rhabdiosaidea). *J. Parasitol.* **54**, 849–855 (1968).
73. Hammond, M.P. & Robinson, R.D. Chromosome complement, gametogenesis, and development of *Strongyloides stercoralis*. *J. Parasitol.* **80**, 689–695 (1994).
74. Albertson, D.G., Nwaorgu, O.C. & Sulston, J.E. Chromatin diminution and a chromosomal mechanism of sexual differentiation in *Strongyloides papillosus*. *Chromosoma* **75**, 75–87 (1979).

ONLINE METHODS

Parasite material, sequencing and assembly. *S. ratti*, *S. stercoralis*, *S. venezuelensis* and *S. papillosus* larvae were obtained from fecal cultures of infected laboratory animals; for *P. trichosuri* and *Rhabditophanes* sp. KR3021, material was obtained from stages grown on agar plates (**Supplementary Fig. 7**) (full details on ethical approval are in the **Supplementary Note**). To produce the *S. ratti* reference genome, a combination of Sanger capillary, 454 and Illumina-derived sequence data was used, whereas data for the other species were generated using Illumina technology. The *S. ratti* genome was initially assembled using Newbler v.2.3 (ref. 75) (for the capillary and 454 sequence data) and AbySS v.1.3.1 (ref. 76) (for the Illumina-derived data); Illumina paired-end reads were mapped to this assembly with SMALT (H. Ponstingl, personal communication). The genomes of the other species, except *S. venezuelensis*, were assembled using a combination of the SGA assembler⁷⁷ and Velvet⁷⁸ from 100-bp paired-end Illumina reads, produced from short-fragment (~500-bp)⁷⁹ and 3-kb mate-pair libraries⁸⁰. Illumina reads were used in IMAGE⁸¹ and Gapfiller⁸² software to fill gaps and in iCORN⁸³ to correct base errors. Gap5 (ref. 84) was used to manually extend and link scaffolds using Illumina read pairs. Genetic markers²⁰ were mapped to the *S. ratti* assembly to order and orient scaffolds and to the *S. papillosus* assembly to assign scaffolds to chromosomes and regions of putative chromosomal diminution. The *S. venezuelensis* genome was assembled using the Platanus assembler⁸⁵ and improved as described above for the other species. The resulting v2 *S. venezuelensis* assembly was further scaffolded using an optical map produced by an Argus optical mapping platform (Opgen). CEGMA v2 (ref. 86) was used to assess the completeness of each assembly.

Assembled sequences were scanned for contamination from other species using a series of BLASTX and BLASTP⁸⁷ searches against vertebrate and invertebrate sequence databases. Repeat sequences in the assemblies were characterized using RepeatModeler and TransposonPSI.

Mitochondrial genomes were assembled using the MITObim assembler⁸⁸ with the *C. elegans* mitochondrial genes as seeds. The gene order of each assembly was confirmed by PCR. A mitochondrial protein-coding gene sequence phylogeny was constructed using RAXML v7.2.8 (ref. 89).

Identifying regions that undergo chromatin diminution or belong to the X chromosome. To identify chromosomal regions that undergo chromatin diminution in *S. papillosus* and scaffolds that belong to the X chromosome in *S. ratti*, *S. stercoralis* and *P. trichosuri*, DNA of males and females from each species was sequenced and mapped to the appropriate reference genome using SMALT v0.7.4 (H. Ponstingl, personal communication). The read depth was calculated for each scaffold using the BedTools function `genomecov`⁹⁰, and all scaffolds were classified as diminished/X-chromosome or non-diminished/autosomal on the basis of differences in read coverage. Because males are hemizygous for the diminished region in *S. papillosus*²⁰ and for the X chromosome in the other species, a male:female read depth ratio of 0.5:1 was expected in diminished or X-chromosome scaffolds relative to autosomes, whereas in non-diminished/autosomal regions the ratio would be expected to be close to 1:1.

Gene prediction and functional annotation. Genes were predicted using Augustus⁹¹—with a training set of approximately 200–400 manually curated genes per species, aligned transcript data and *S. ratti* protein sequences as hints—supplemented with non-overlapping predictions from MAKER⁹². If there was more than one alternative splice pattern for a gene prediction in the combined Augustus and MAKER gene set, we only kept the transcript corresponding to the longest predicted protein. Astacin gene models and a subset of SCP/TAPS gene models from *S. ratti*, *S. venezuelensis* and *S. stercoralis* were manually curated before phylogenetic analyses.

A protein name was assigned to each predicted protein on the basis of manually curated orthologs in UniProt⁹³ from selected species (human, zebrafish, *Drosophila melanogaster*, *C. elegans* and *Schistosoma mansoni* orthologs), where possible. If a predicted protein was not assigned a protein name on the basis of its orthologs, then a protein name was assigned using InterPro⁹⁴ domains in the protein.

GO terms were assigned by transferring GO terms from human, zebrafish, *C. elegans* and *D. melanogaster* orthologs using an approach based on the Ensembl Compara approach for transferring GO terms to orthologs in vertebrate

species³³ but modified for improved accuracy in transferring GO terms across phyla. Manually curated GO annotations were downloaded from the GO Consortium website⁹⁵, and, for a particular predicted protein in the present study, the manually curated GO terms were obtained for all its human, zebrafish, *C. elegans* and *D. melanogaster* orthologs. From this set, the last common ancestor term (in the GO hierarchy) was found for each pair of GO terms from orthologs of two different species (for example, a *C. elegans* ortholog and a zebrafish ortholog) and then transferred to our predicted protein. GO terms of the three possible types (molecular function, cellular component and biological process) were assigned to predicted proteins in this way. Additional GO terms were identified using InterProScan⁹⁶.

Gene orthology and species tree reconstruction. Eight outgroup species were used, encompassing four previously defined nematode clades⁹ (clade I, *Trichinella spiralis* and *Trichuris muris*; clade III, *Ascaris suum* and *Brugia malayi*; clade IV, *Bursaphelenchus xylophilus* and *Meloidogyne hapla*; clade V, *Necator americanus* and *C. elegans*), together with the six species from the present study to construct a Compara database using the Ensembl Compara pipeline³³. The database was used to identify orthologs and paralogs, gene duplications and gene losses, as well as gene families shared among the species or subsets of the species or specific to one species.

In total, 4,437 gene families were identified that contained just one gene from each species and that were present in at least ten species out of the six species and the eight outgroups. An alignment for the proteins in each family was built using MAFFT version v6.857 (ref. 97), poorly aligning regions were trimmed using GBLOCKS v0.91b and the remaining columns were concatenated. For each alignment, the best-fitting amino acid substitution model was identified as that minimizing the Akaike Information Criterion from the set of models available in RAXML v8.0.24 (ref. 89), testing models with both predefined amino acid frequencies and observed frequencies in the data, and all with the CAT model of rate variation across sites. A maximum-likelihood phylogenetic tree was constructed on the basis of the concatenated alignment, with each protein alignment an independent partition of these data, applying the best-fitting substitution model identified above to each partition. This inference used RAXML v8.0.24 with ten random addition-sequence replicates and 100 bootstrap replicates and otherwise used default heuristic search settings.

Analysis of intron-exon structure and syntenic analysis. Introns that were present in two or more species were identified from gene structures and full-gene nucleotide alignments of 208 single-copy orthologs using Scipio⁹⁸ and GenePainter⁹⁹. The output from GenePainter was parsed into DOLLOP (PHYLIP package; see URLs) to infer intron gain and loss on every node of the species tree using maximum parsimony.

Whole-assembly nucleotide alignments were produced between *S. ratti* and the other five species using nucmer¹⁰⁰. Each scaffold from the other species was assigned a chromosome on the basis of its nucmer alignment to an *S. ratti* chromosome. To identify syntenic regions, conserved blocks of three consecutive orthologous genes or more in the same order and orientation were defined by DAGchainer¹⁰¹, between the *S. ratti* reference and each of the other five species. To gain a high-level view of syntenic, PROmer¹⁰² was used to identify very highly conserved sequence matches, on the basis of translated sequence, after which scaffolds from a particular species were ordered by matching to *S. ratti* chromosome and position in that chromosome and the matches were plotted using Circos¹⁰³.

Transcriptome and proteome analyses. For *S. ratti* and *S. stercoralis*, the transcriptomes were compared from the parasitic female, free-living female and iL3s; we note that parasitic and free-living adult females will have eggs *in utero*. For *S. ratti*, free-living females were picked individually from cultures of *S. ratti*-infected rat feces, from which iL3s were also collected; parasitic females were collected by dissection of *S. ratti*-infected rats¹⁰⁴. Two biological replicates were collected for parasitic and free-living females. These samples were divided approximately equally and used for both transcriptomic and proteomic analysis. A single biological sample was used for iL3 transcriptomic analysis. RNA was prepared from TRIzol and selected for poly(A) RNA with Dynabeads, acoustically sheared and reverse transcribed to construct Illumina libraries that were sequenced. For *S. stercoralis*, we used previously published

data⁴². RNA-seq data were analyzed using R v.3.0.2 and the Bioconductor package edgeR¹⁰⁵ to identify genes differentially expressed in all pairwise combinations of the three life cycle stages.

For *S. ratti*, the proteome was also compared between the parasitic and free-living females. Equivalent samples of the material collected for the transcriptome analyses were used. Protein was extracted by freeze-thaw cycles, mechanical grinding and chemical extraction and digested with trypsin. The resulting peptide mixture was analyzed by liquid chromatography–mass spectrometry. Proteins were identified and quantified using Progenesis. For downstream analyses, at least two unique peptides were required to identify proteins. Protein abundance (iBAQ) was calculated from Progenesis.

For both the transcriptome and proteome data, GO analysis was performed in R using TopGo v.2.16.0 and Fisher's exact test.

For the analysis of the ES proteome⁶³, converted raw spectral files were analyzed by the Mascot search engine, where an FDR <1% and a minimum of two significant peptides were required to identify proteins. Protein abundance was calculated from the Mascot algorithm empAI.

Astacins and SCP/TAPS. Genes encoding astacins and SCP/TAPS were identified using InterProScan. For these gene families, we aligned amino acid sequences of the members from all *S. ratti* and eight outgroup species using MAFFT⁹⁷. The alignments were edited with TCS¹⁰⁶ using the weighted option, and the distance matrix of the new alignment was calculated using ProtTest¹⁰⁷. The phylogenetic tree was constructed by maximum likelihood using RAxML⁸⁹ with 100 bootstrap replicates.

Gene clusters. Clusters of genes were identified as three or more adjacent genes upregulated in the same stage of the life cycle. The members of a cluster were considered to share a common gene family where ≥50% of the genes belonged to the same Compara gene family. To investigate the number of clusters expected by chance for a particular life cycle stage, for *n* genes upregulated in a particular stage, we randomly selected *n* genes from the genome and calculated the number of clusters seen for the *n* random genes; this was repeated 1,000 times and the mean value was calculated.

75. Margulies, M. *et al.* Genome sequencing in microfabricated high-density picolitre reactors. *Nature* **437**, 376–380 (2005).
76. Simpson, J.T. *et al.* ABySS: a parallel assembler for short read sequence data. *Genome Res.* **19**, 1117–1123 (2009).
77. Simpson, J.T. & Durbin, R. Efficient *de novo* assembly of large genomes using compressed data structures. *Genome Res.* **22**, 549–556 (2012).
78. Zerbino, D.R. Using the Velvet *de novo* assembler for short-read sequencing technologies. *Curr. Protoc. Bioinformatics* Chapter 11, 11–15 (2010).
79. Kozarewa, I. *et al.* Amplification-free Illumina sequencing-library preparation facilitates improved mapping and assembly of (G+C)-biased genomes. *Nat. Methods* **6**, 291–295 (2009).
80. Park, N. *et al.* An improved approach to mate-paired library preparation for Illumina sequencing. *Methods Next Gener. Seq.* **1**, 2084–7173 (2013).
81. Tsai, I.J., Otto, T.D. & Berriman, M. Improving draft assemblies by iterative mapping and assembly of short reads to eliminate gaps. *Genome Biol.* **11**, R41 (2010).
82. Nadalin, F., Vezzi, F. & Policriti, A. GapFiller: a *de novo* assembly approach to fill the gap within paired reads. *BMC Bioinformatics* **13** (suppl. 14), S8 (2012).
83. Otto, T.D., Sanders, M., Berriman, M. & Newbold, C. Iterative Correction of Reference Nucleotides (iCORN) using second generation sequencing technology. *Bioinformatics* **26**, 1704–1707 (2010).
84. Bonfield, J.K. & Whitwham, A. Gap5—editing the billion fragment sequence assembly. *Bioinformatics* **26**, 1699–1703 (2010).
85. Kajitani, R. *et al.* Efficient *de novo* assembly of highly heterozygous genomes from whole-genome shotgun short reads. *Genome Res.* **24**, 1384–1395 (2014).
86. Parra, G., Bradnam, K., Ning, Z., Keane, T. & Korf, I. Assessing the gene space in draft genomes. *Nucleic Acids Res.* **37**, 289–297 (2009).
87. Altschul, S.F. *et al.* Gapped BLAST and PSI-BLAST: a new generation of protein database search programs. *Nucleic Acids Res.* **25**, 3389–3402 (1997).
88. Hahn, C., Bachmann, L. & Chevreux, B. Reconstructing mitochondrial genomes directly from genomic next-generation sequencing reads—a baiting and iterative mapping approach. *Nucleic Acids Res.* **41**, e129 (2013).
89. Stamatakis, A. RAxML version 8: a tool for phylogenetic analysis and post-analysis of large phylogenies. *Bioinformatics* **30**, 1312–1313 (2014).
90. Quinlan, A.R. BEDTools: the Swiss-Army tool for genome feature analysis. *Curr. Protoc. Bioinformatics* **47**, 11.12.1–11.12.34 (2014).
91. Stanke, M. *et al.* AUGUSTUS: *ab initio* prediction of alternative transcripts. *Nucleic Acids Res.* **34**, W435–W439 (2006).
92. Holt, C. & Yandell, M. MAKER2: an annotation pipeline and genome-database management tool for second-generation genome projects. *BMC Bioinformatics* **12**, 491 (2011).
93. UniProt Consortium. UniProt: a hub for protein information. *Nucleic Acids Res.* **43**, D204–D212 (2015).
94. Mitchell, A. *et al.* The InterPro protein families database: the classification resource after 15 years. *Nucleic Acids Res.* **43**, D213–D221 (2015).
95. Gene Ontology Consortium. Gene Ontology Consortium: going forward. *Nucleic Acids Res.* **43**, D1049–D1056 (2015).
96. Jones, P. *et al.* InterProScan 5: genome-scale protein function classification. *Bioinformatics* **30**, 1236–1240 (2014).
97. Katoh, K. & Standley, D.M. MAFFT: iterative refinement and additional methods. *Methods Mol. Biol.* **1079**, 131–146 (2014).
98. Keller, O., Odronitz, F., Stanke, M., Kollmar, M. & Waack, S. Scipio: using protein sequences to determine the precise exon/intron structures of genes and their orthologs in closely related species. *BMC Bioinformatics* **9**, 278 (2008).
99. Hammesfahr, B., Odronitz, F., Mühlhausen, S., Waack, S. & Kollmar, M. GenePainter: a fast tool for aligning gene structures of eukaryotic protein families, visualizing the alignments and mapping gene structures onto protein structures. *BMC Bioinformatics* **14**, 77 (2013).
100. Delcher, A.L., Phillippy, A., Carlton, J. & Salzberg, S.L. Fast algorithms for large-scale genome alignment and comparison. *Nucleic Acids Res.* **30**, 2478–2483 (2002).
101. Haas, B.J., Delcher, A.L., Wortman, J.R. & Salzberg, S.L. DAGchainer: a tool for mining segmental genome duplications and synteny. *Bioinformatics* **20**, 3643–3646 (2004).
102. Kurtz, S. *et al.* Versatile and open software for comparing large genomes. *Genome Biol.* **5**, R12 (2004).
103. Krzywinski, M. *et al.* Circos: an information aesthetic for comparative genomics. *Genome Res.* **19**, 1639–1645 (2009).
104. Thompson, F.J., Barker, G.L.A., Hughes, L. & Viney, M.E. Genes important in the parasitic life of the nematode *Strongyloides ratti*. *Mol. Biochem. Parasitol.* **158**, 112–119 (2008).
105. Robinson, M.D., McCarthy, D.J. & Smyth, G.K. edgeR: a Bioconductor package for differential expression analysis of digital gene expression data. *Bioinformatics* **26**, 139–140 (2010).
106. Chang, J.-M., Di Tommaso, P. & Notredame, C. TCS: a new multiple sequence alignment reliability measure to estimate alignment accuracy and improve phylogenetic tree reconstruction. *Mol. Biol. Evol.* **31**, 1625–1637 (2014).
107. Abascal, F., Zardoya, R. & Posada, D. ProtTest: selection of best-fit models of protein evolution. *Bioinformatics* **21**, 2104–2105 (2005).