

# Genome Biology

## Extensive genomic and transcriptomic variation defines the chromosome-scale assembly of *Haemonchus contortus*, a model gastrointestinal worm

--Manuscript Draft--

<b>Manuscript Number:</b>		
<b>Full Title:</b>	Extensive genomic and transcriptomic variation defines the chromosome-scale assembly of <i>Haemonchus contortus</i> , a model gastrointestinal worm	
<b>Article Type:</b>	Research	
<b>Funding Information:</b>	Wellcome Trust (067811)	Not applicable
	Wellcome Trust (WT206194)	Not applicable
	Biotechnology and Biological Sciences Research Council (BB/M003949/1)	Not applicable
	Biotechnology and Biological Sciences Research Council (BB/P024610/1)	Not applicable
	Biotechnology and Biological Sciences Research Council (BB/K020048/1)	Not applicable
<b>Abstract:</b>	<p><b>Background</b></p> <p><i>Haemonchus contortus</i> is a globally distributed and economically important gastrointestinal pathogen of small ruminants, and has become the key nematode model for studying anthelmintic resistance and other parasite-specific traits among a wider group of parasites including major human pathogens. Two draft genome assemblies for <i>H. contortus</i> were reported in 2013, however, both were highly fragmented, incomplete, and differed from one another in important respects. While the introduction of long-read sequencing has significantly increased the rate of production and contiguity of de novo genome assemblies broadly, achieving high quality genome assemblies for small, genetically diverse, outcrossing eukaryotic organisms such as <i>H. contortus</i> remains a significant challenge.</p> <p><b>Results</b></p> <p>Here, we report using PacBio long read and OpGen and 10X Genomics long-molecule methods to generate a highly contiguous 283.4 Mbp chromosome-scale genome assembly including a resolved sex chromosome. We show a remarkable pattern of almost complete conservation of chromosome content (synteny) with <i>Caenorhabditis elegans</i>, but almost no conservation of gene order. Long-read transcriptome sequence data has allowed us to define coordinated transcriptional regulation throughout the life cycle of the parasite, and refine our understanding of cis- and trans-splicing relative to that observed in <i>C. elegans</i>. Finally, we use this assembly to give a comprehensive picture of chromosome-wide genetic diversity both within a single isolate and globally.</p> <p><b>Conclusions</b></p> <p>The <i>H. contortus</i> MHco3(ISE).N1 genome assembly presented here represents the most contiguous and resolved nematode assembly outside of the <i>Caenorhabditis</i> genus to date, together with one of the highest-quality set of predicted gene features. These data provide a high-quality comparison for understanding the evolution and genomics of <i>Caenorhabditis</i> and other nematodes, and extends the experimental tractability of this model parasitic nematode in understanding pathogen biology, drug discovery and vaccine development, and important adaptive traits such as drug resistance.</p>	
<b>Corresponding Author:</b>	Stephen R Doyle, PhD Wellcome Trust Sanger Institute	

	Hinxton, Cambridgeshire UNITED KINGDOM
<b>Corresponding Author Secondary Information:</b>	
<b>Corresponding Author's Institution:</b>	Wellcome Trust Sanger Institute
<b>Corresponding Author's Secondary Institution:</b>	
<b>First Author:</b>	Stephen R Doyle, PhD
<b>First Author Secondary Information:</b>	
<b>Order of Authors:</b>	<p>Stephen R Doyle, PhD</p> <p>Alan Tracey</p> <p>Roz Laing</p> <p>Nancy Holroyd</p> <p>David Bartley</p> <p>Wojtek Bazant</p> <p>Helen Beasley</p> <p>Robin Beech</p> <p>Collette Britton</p> <p>Karen Brooks</p> <p>Kirsty Maitland</p> <p>Axel Martinelli</p> <p>Jennifer D Noonan</p> <p>Michael Paulini</p> <p>Michael A Quail</p> <p>Elizabeth Redman</p> <p>Faye H Rodgers</p> <p>Guillaume Sallé</p> <p>Geetha Sankaranarayanan</p> <p>Kevin L Howe</p> <p>Neil Sargison</p> <p>Eileen Devaney</p> <p>Matthew Berriman</p> <p>John S Gilliard</p> <p>James A Cotton</p>
<b>Order of Authors Secondary Information:</b>	
<b>Suggested Reviewers:</b>	<p>Neil Young, PhD Senior Research Fellow, The University of Melbourne, Australia <a href="mailto:nyoung@unimelb.edu.au">nyoung@unimelb.edu.au</a> Extensive helminth genomics experience/expertise</p> <p>Tim Geary, PhD Professor, Queen's University Belfast, UK <a href="mailto:T.Geary@qub.ac.uk">T.Geary@qub.ac.uk</a> Expert on parasitic nematode biology and anthelmintics</p> <p>Erik Andersen, PhD</p>

	<p>Associate Professor, Northwestern University, USA  Erik.Andersen@northwestern.edu  C. elegans population genomicist</p>
	<p>Christian Rödelsperger, PhD  Departmental Project Leader, Max Planck Institute for Developmental Biology,  Germany  christian.roedelsperger@tuebingen.mpg.de  Nematode evolutionary biologist focused on bioinformatics / genomics</p>
	<p>Erich Schwarz, PhD  Assistant Research Professor, Cornell University, USA  ems394@cornell.edu  Functional genomicist, with interests in C. elegans and other nematodes. Lead author on draft McMaster strain <i>Haemonchus contortus</i> genome assembly published in 2013.</p>
	<p>Kirsten Gunsalus, PhD  Professor of Biology; NYU-AD Faculty Director of Bioinformatics, New York University, USA  kcg1@nyu.edu  Uses functional genomics and computational approaches to study development in C. elegans and related nematodes</p>
	<p>Vicky Hunt, PhD  Group Leader, The University of Bath, UK  bs1vlh@bath.ac.uk  Comparative and functional genomics of Caenorhabditids and parasitic nematodes</p>
<b>Additional Information:</b>	
<b>Question</b>	<b>Response</b>
Has this manuscript been submitted to this journal before?	No
<b>Experimental design and statistics</b>	Yes
Full details of the experimental design and statistical methods used should be given in the Methods section, as detailed in our <a href="#">Minimum Standards Reporting Checklist</a> . Information essential to interpreting the data presented should be made available in the figure legends.	
Have you included all the information requested in your manuscript?	
<b>Resources</b>	Yes
A description of all resources used, including antibodies, cell lines, animals and software tools, with enough information to allow them to be uniquely identified, should be included in the Methods section. Authors are strongly encouraged to cite <a href="#">Research Resource</a>	

Identifiers (RRIDs) for antibodies, model organisms and tools, where possible.

Have you included the information requested as detailed in our [Minimum Standards Reporting Checklist](#) " target="\_blank">>Minimum Standards Reporting Checklist?

[Click here to view linked References](#)

1   **Extensive genomic and transcriptomic variation defines**  
2   **the chromosome-scale assembly of *Haemonchus contortus*,**  
3   **a model gastrointestinal worm**

4  
5   Stephen R. Doyle<sup>1\*</sup>, Alan Tracey<sup>1</sup>, Roz Laing<sup>2</sup>, Nancy Holroyd<sup>1</sup>, David Bartley<sup>3</sup>, Wojtek  
6   Bazant<sup>1</sup>, Helen Beasley<sup>1</sup>, Robin Beech<sup>4</sup>, Collette Britton<sup>2</sup>, Karen Brooks<sup>1</sup>, Kirsty Maitland<sup>2</sup>,  
7   Axel Martinelli<sup>1</sup>, Jennifer D. Noonan<sup>4</sup>, Michael Paulini<sup>5</sup>, Michael A. Quail<sup>1</sup>, Elizabeth  
8   Redman<sup>6</sup>, Faye H. Rodgers<sup>1</sup>, Guillaume Sallé<sup>7</sup>, Geetha Sankaranarayanan<sup>1</sup>, Kevin L. Howe<sup>5</sup>,  
9   Neil Sargison<sup>8</sup>, Eileen Devaney<sup>2</sup>, Matthew Berriman<sup>1</sup>, John S. Gilleard<sup>6</sup>, James A. Cotton<sup>1\*</sup>

10  
11   \*Corresponding authors

12  
13   1. **Wellcome Sanger Institute, Hinxton, Cambridgeshire, CB10 1SA, United Kingdom**

14   Stephen R. Doyle: [stephen.doyle@sanger.ac.uk](mailto:stephen.doyle@sanger.ac.uk); Alan Tracey: [alt@sanger.ac.uk](mailto:alt@sanger.ac.uk); Nancy Holroyd: [neh@sanger.ac.uk](mailto:neh@sanger.ac.uk);  
15   Wojtek Bazant: [wojtek.bazant@sanger.ac.uk](mailto:wojtek.bazant@sanger.ac.uk); Helen Beasley: [hb521@cam.ac.uk](mailto:hb521@cam.ac.uk); Karen Brooks: [kd1@sanger.ac.uk](mailto:kd1@sanger.ac.uk) ;  
16   Axel Martinelli: [axel.martinelli@gmail.com](mailto:axel.martinelli@gmail.com); Michael A. Quail: [mq1@sanger.ac.uk](mailto:mq1@sanger.ac.uk); Faye Rodgers: [fr7@sanger.ac.uk](mailto:fr7@sanger.ac.uk);  
17   Geetha Sankaranarayanan: [gs13@sanger.ac.uk](mailto:gs13@sanger.ac.uk); Matthew Berriman: [mb4@sanger.ac.uk](mailto:mb4@sanger.ac.uk); James A. Cotton:  
18   [jc17@sanger.ac.uk](mailto:jc17@sanger.ac.uk)

19  
20   2. **Institute of Biodiversity Animal Health and Comparative Medicine, College of Medical, Veterinary and Life  
21   Sciences, University of Glasgow, Garscube Campus, Glasgow, G61 1QH, United Kingdom**

22   Roz Laing: [Rosalind.Laing@glasgow.ac.uk](mailto:Rosalind.Laing@glasgow.ac.uk); Collette Britton: [Collette.Britton@glasgow.ac.uk](mailto:Collette.Britton@glasgow.ac.uk); Kirsty Maitland:  
23   [Kirsty.Maitland@glasgow.ac.uk](mailto:Kirsty.Maitland@glasgow.ac.uk); Eileen Devaney: [Eileen.Devaney@glasgow.ac.uk](mailto:Eileen.Devaney@glasgow.ac.uk)

24  
25   3. **Moredun Research Institute, Pentlands Science Park, Bush Loan, Penicuik EH26 0PZ, United Kingdom**

26   David Bartley: [Dave.Bartley@moredun.ac.uk](mailto:Dave.Bartley@moredun.ac.uk)

27  
28   4. **Institute of Parasitology, McGill University, 2111 Lakeshore Road, Sainte Anne-de-Bellevue, Québec, H9X3V9  
29   Canada**

30   Robin Beech: [robin.beech@mcgill.ca](mailto:robin.beech@mcgill.ca); Jennifer D. Noonan: [jennifer.noonan@mail.mcgill.ca](mailto:jennifer.noonan@mail.mcgill.ca)

31  
32   5. **European Molecular Biology Laboratory, European Bioinformatics Institute, Hinxton, Cambridgeshire, CB10 1SA,  
33   United Kingdom**

34   Michael Paulini: [mh6@ebi.ac.uk](mailto:mh6@ebi.ac.uk); Kevin Howe: [klh@ebi.ac.uk](mailto:klh@ebi.ac.uk)

35  
36   6. **Department of Comparative Biology and Experimental Medicine, Faculty of Veterinary Medicine, University of  
37   Calgary, Calgary, Alberta, Canada**

38   Elizabeth Redman: [elmredma@ucalgary.ca](mailto:elmredma@ucalgary.ca); John Gilleard: [jsgillea@ucalgary.ca](mailto:jsgillea@ucalgary.ca)

39  
40   7. **INRAE - U. Tours, UMR 1282 ISP Infectiologie et Santé Publique, Centre de recherche Val de Loire, Nouzilly, France**

41   Guillaume Salle: [guillaume.salle@inrae.fr](mailto:guillaume.salle@inrae.fr)

42  
43   8. **Royal (Dick) School of Veterinary Studies, University of Edinburgh, Edinburgh, EH25 9RG, United Kingdom**

44   Neil Sargison: [Neil.Sargison@ed.ac.uk](mailto:Neil.Sargison@ed.ac.uk)

## Abstract

**Background** *Haemonchus contortus* is a globally distributed and economically important gastrointestinal pathogen of small ruminants, and has become the key nematode model for studying anthelmintic resistance and other parasite-specific traits among a wider group of parasites including major human pathogens. Two draft genome assemblies for *H. contortus* were reported in 2013, however, both were highly fragmented, incomplete, and differed from one another in important respects. While the introduction of long-read sequencing has significantly increased the rate of production and contiguity of *de novo* genome assemblies broadly, achieving high quality genome assemblies for small, genetically diverse, outcrossing eukaryotic organisms such as *H. contortus* remains a significant challenge.

**Results** Here, we report using PacBio long read and OpGen and 10X Genomics long-molecule methods to generate a highly contiguous 283.4 Mbp chromosome-scale genome assembly including a resolved sex chromosome. We show a remarkable pattern of almost complete conservation of chromosome content (synteny) with *Caenorhabditis elegans*, but almost no conservation of gene order. Long-read transcriptome sequence data has allowed us to define coordinated transcriptional regulation throughout the life cycle of the parasite, and refine our understanding of *cis*- and *trans*-splicing relative to that observed in *C. elegans*. Finally, we use this assembly to give a comprehensive picture of chromosome-wide genetic diversity both within a single isolate and globally.

**Conclusions** The *H. contortus* MHco3(ISE).N1 genome assembly presented here represents the most contiguous and resolved nematode assembly outside of the *Caenorhabditis* genus to date, together with one of the highest-quality set of predicted gene features. These data provide a high-quality comparison for understanding the evolution and genomics of *Caenorhabditis* and other nematodes, and extends the experimental tractability of this model parasitic nematode in understanding pathogen biology, drug discovery and vaccine development, and important adaptive traits such as drug resistance.

## Keywords

*Haemonchus contortus*; chromosomal genome assembly; PacBio long-read sequencing; IsoSeq cDNA sequencing; genetic variation; transcriptomics; population genomics

1  
2   **Background**  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65

78 A complete sequence assembly that reflects the biologically defined chromosomal DNA  
79 sequence of an organism is the ultimate goal of a *de novo* genome assembly project. Until  
80 recently, this goal was unattainable for all but the smallest of genomes or for select  
81 organisms, where genome assembly projects have been supported by large consortia and  
82 extensive financial and logistical support, for example, the model nematode *Caenorhabditis*  
83 *elegans* [1]. However, high-throughput sequencing has significantly impacted the rate at  
84 which draft genomes are produced and, together with the recent introduction of long-read  
85 sequencing from Pacific Biosciences (PacBio) and Oxford Nanopore, the quality and  
86 contiguity of assemblies of both large and small genomes has completely changed.  
87 Now, the ability to generate highly contiguous and closed prokaryote genomes is becoming  
88 routine, even from metagenomic samples containing multiple strains or species [2,3], and  
89 high-quality chromosome-scale genomes for a number of eukaryotic species are rapidly  
90 becoming available [4–6]. These technologies generally require microgram quantities of  
91 high molecular weight DNA, and the success of a genome assembly is highly dependent on  
92 the degree of polymorphism present in the DNA that is sequenced. This can be minimised  
93 by sequencing DNA from a single individual – so that the only polymorphism is  
94 heterozygosity between homologous chromosomes in the case of a diploid or polyploid  
95 individual – or from a clonal population, where multiple but genetically identical individuals  
96 can be pooled. However, for many small organisms, this is not possible; the only option is to  
97 sequence DNA from a pool of genetically distinct organisms and derive a consensus  
98 assembly. As modern assembly algorithms are designed and trained on haploid or diploid  
99 genomes from single individuals [7–9], the excessive diversity presents a significant  
100 challenge to the assembly process, and consequently, typically results in fragmented

101 assemblies that contain misassembled and haplotypic sequences.  
1  
2  
3  
4  
5 103 Parasitic helminths are a diverse group of organisms for which significant efforts have been  
6  
7 104 made to develop genomic resources [10]. For some species, including *Onchocerca volvulus*  
8  
9 105 [11], *Strongyloides ratti* [12], *Taenia multiceps* [13], and *Echinococcus multilocularis* [14],  
10  
11 106 high-quality near-complete assemblies have been achieved using multiple complementary  
12  
13 107 sequencing and mapping technologies, and usually with extensive manual improvement  
14  
15 108 work. Similar quality assemblies for a few other species (*Schistosoma mansoni*, *Brugia*  
16  
17 109 *malayi*, *Hymenolepis microstoma*, and *Trichuris muris*) are available from WormBase  
18  
19 109 ParaSite [15], but await formal description in the peer-reviewed literature. However, for  
20  
21 110 many others, access to sufficient high-quality DNA, small organismal size, and high genetic  
22  
23 111 diversity are significant challenges that need to be overcome to achieve highly contiguous  
24  
25 112 assemblies. The nematode *Haemonchus contortus* is one such species. *H. contortus* is a  
26  
27 113 major pathogen of sheep and goats worldwide and is recognised as a model parasite due to  
28  
29 114 its experimental tractability for drug discovery [16,17], vaccine development [18,19] and  
30  
31 115 anthelmintic resistance research [20]. Draft genome assemblies were published in 2013 for  
32  
33 116 two anthelmintic-sensitive isolates, MHco3(ISE).N1 [21] and McMaster [22]; both  
34  
35 117 assemblies were produced with a combination of short-read, high-throughput sequencing  
36  
37 118 technologies (including both short- and long-insert libraries), resulting in estimated genome  
38  
39 119 sizes of 370 Mbp and 320 Mbp, respectively. However, direct comparison of these two  
40  
41 120 assemblies shows discordance, revealing clear differences in gene family composition [23]  
42  
43 121 and variation in assembly quality and gene content [24]. Although some of these differences  
44  
45 122 reflect true biological variation in this highly polymorphic species [25–27], many differences  
46  
47 123 are technical artefacts, as revealed by comparing fragmented genome assemblies from  
48  
49 124

125 sequencing DNA derived from pools of a genetically diverse organism. These fragmented  
126 genome assemblies are also less than ideal substrates for many downstream uses of  
127 reference genomes, for example, in interpreting the results of genetic experiments in which  
128 patterns of genome-wide variation are defined and interpreted [28].  
9

10 129  
11  
12  
13 130 Here we present the chromosome-scale assembly of the MHco3(ISE).N1 isolate of *H.*  
14  
15 131 *contortus*, representing the first high-quality assembly for any strongylid parasite, an  
16  
17 132 important clade of parasitic nematodes that includes parasitic species of major veterinary  
18  
19 133 and medical importance. Using a hybrid assembly approach incorporating short and long-  
20  
21 134 read sequencing followed by manual finishing, we demonstrate how a highly contiguous  
22  
23 135 assembly overcomes some of the critical limitations imposed by draft assemblies in  
24  
25  
26 136 interpreting large-scale genetic and functional genomics datasets. We provide insight into  
27  
28 137 patterns of transcriptional change and putative co-regulation using a significantly improved  
29  
30 138 genome annotation, derived *de novo* from both short- and long-read cDNA sequencing, and  
31  
32 139 describe within- and between-population genetic diversity that shapes the genome. This  
33  
34 140 manually refined chromosome-scale assembly now offers insight into genome evolution  
35  
36 141 among a broad group of important parasite species and offers a robust scaffold for genome-  
37  
38 142 wide analyses of important parasite traits such as anthelmintic resistance.  
39  
40  
41  
42  
43  
44  
45  
46 143  
47  
48  
49 144  
50  
51  
52 145  
53  
54  
55  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65

146 **Results**

1  
2  
3 **Chromosome structure of *Haemonchus contortus***

4  
5  
6 We have built upon our previous assembly (version 1 (V1); [21]) using a hybrid approach,  
7  
8 iteratively incorporating Illumina short-insert and 3 kbp libraries, PacBio long-read, OpGen  
9  
10 optical mapping, and 10X Genomics linked-read data (**Additional file 1: Figure S1**; see  
11  
12 **Additional file 2: Table S1** for new sequencing data generated) to generate a largely  
13  
14 complete, chromosomal-scale genome assembly. Consistent with the karyotype for *H.*  
15  
16 *contortus* [29,30], the core genome assembly consists of five autosomal scaffolds and one  
17  
18 sex-linked scaffold (**Additional file 1: Figure S2** shows completion of the X chromosome  
19  
20 from assembly V3 to V4, and the relative X-to-autosome genome coverage between a single  
21  
22 male and female parasite), each containing terminal telomeric sequences, and a single  
23  
24 mitochondrial contig. We assigned chromosome names based on synteny with *C. elegans*  
25  
26 chromosomes (**Figure 1 A**); over 80% of the 7,361 one-to-one orthologous genes are shared  
27  
28 between syntenic chromosomes between the two species demonstrating the high  
29  
30 conservation of genes per chromosome (**Figure 1 B**; top), however, vast rearrangements are  
31  
32 evident and very little conservation of gene order remains (**Figure 1 B**; bottom). We  
33  
34 determined the extent of microsynteny by comparing conserved, reorientated, and  
35  
36 recombined gene pairs between the two species (**Additional file 1: Figure S3 A**); the  
37  
38 distance between ortholog pairs in *H. contortus* and *C. elegans* is correlated up to ~ 100 kbp  
39  
40 (rho = 0.469, p-value = 2.2E-16) in which pair order is conserved (**Additional file 1: Figure S3**  
41  
42 **B**), likely representing selective constraint to maintain evolutionary conserved operons, but  
43  
44 is lost above 100 kbp (rho = 0.017, p-value = 0.6595) where a greater frequency of  
45  
46 recombination between pairs is evident (**Additional file 1: Figure S3 B,C**). Beyond pairs of  
47  
48 genes, synteny breaks down rapidly; although almost 50% of shared ortholog pairs are  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65

170 adjacent to each other, only a single group of 10 orthologs are colinear between the

171 genomes of the two species (**Additional file 1: Figure S3 D**).

---

172

173

174 **Figure 1. Chromosomal synteny between *Haemonchus contortus* and *Caenorhabditis elegans***

175

176 **A.** Genome-wide comparison of chromosomal organisation between *H. contortus* V4 and *C. elegans*, analysed using *promer* and visualised using *Circos*. **B.** Comparison of shared orthologs, based on assignment to chromosomes (top plot) and their relative position on the chromosome (bottom plot). In the top plot, chromosome assignment of 7,361 one-to-one orthologs between *C. elegans* and *H. contortus* V4 annotations demonstrates that the majority of genes on any given *H. contortus* chromosome are found on the same chromosome in *C. elegans*. Values within each panel represent the percentage of shared orthologs found on the same chromosome of both *H. contortus* and *C. elegans*. In the bottom plot, the relative genomic position of the one-to-one orthologs between *C. elegans* (y-axis) and *H. contortus* V4 chromosomes (x-axis) is shown. Complete within-chromosome synteny would be expected to show a positive linear relationship between the genomic positions of orthologs between *C. elegans* and *H. contortus* chromosomes; this is not observed, with an almost complete reshuffling of orthologs between the two species. Box headers are coloured by *H. contortus* chromosomes, whereas data points within the boxes are coloured by *C. elegans* chromosomes.

---

191

192

193 The assembly is approximately 283.4 Mbp in length (scaffold N50 = 47.3 Mbp; contig N50 =

194 3.8 Mbp), representing only 75.6% of the 369.8 Mbp V1 draft assembly length (**Table 1**); the reduction in assembly size is largely due to the identification and separation of redundant haplotypic sequences present in the polymorphic draft V1 and preliminary PacBio assemblies generated. The V4 assembly is highly resolved with only 185 gaps, a significant

198 reduction from the 41,663 gaps present in the V1 assembly (**Table 1**). As a measure of  
1  
2 assembly completeness, we identified 242 of 248 universally conserved orthologs measured  
3  
4 by the Core Eukaryotic Genes Mapping Approach (CEGMA; **Additional file 2: Table S2**), and  
5  
6 859 of 982 metazoan Benchmarking Universal Single-Copy Orthologs (BUSCOs). The  
7  
8 remaining missing orthologs show phylogenetic structure among clade V nematodes  
9  
10 (Additional file 1: Figure S4 A), suggesting that many are truly missing from the genomes of  
11  
12 these species rather than due to assembly artefacts. These data contained 141 more full-  
13  
14 length single-copy orthologs and 159 fewer duplicated BUSCOs than identified in the V1  
15  
16 genome (**Additional file 2: Table S2**); considering that we identified a similar average  
17  
18 number of CEGMA orthologs per core gene in V4 relative to the *C. elegans* genome (1.1 in  
19  
20 V4, compared with 1.09 in *C. elegans*), these data represent a significant reduction in  
21  
22 erroneously-duplicated sequences in V4 compared with the draft *H. contortus* assemblies  
23  
24 (V1: 85.89% complete and 1.59 average orthologs per CEGMA gene; McMaster: 70.56%  
25  
26 complete and 1.43 average orthologs CEGMA gene; **Additional file 2: Table S2**).  
27  
28  
29 212  
30  
31 213 An assembly of a New Zealand (NZ) *H. contortus* strain was recently made available (ENA  
32  
33 accession: GCA\_007637855.2; [31]), which used our V4 chromosomal assembly to scaffold a  
34  
35 highly fragmented draft assembly (contigs = 172,899) into seven chromosome-scale  
36  
37 scaffolds. There are striking differences between the NZ and the V4 assemblies that are very  
38  
39 unlikely to represent genuine biological variation between strains. The significant expansion  
40  
41 of assembly size (465 Mbp vs 283 Mbp in the V4 assembly; **Table 1**) in the context of  
42  
43 elevated duplicated BUSCOs (7.3% vs 1.6% in the V4 assembly) and high average orthologs  
44  
45 per partial and complete CEGMA genes (1.42 and 1.71, respectively, vs 1.1 and 1.22 in the  
46  
47 V4 assembly; **Additional file 2: Table S2**) suggest that contaminating haplotypic sequences  
48  
49  
50 214  
51  
52 215  
53  
54 216  
55  
56 217  
57  
58 218  
59  
60 219  
61  
62 220  
63  
64 221  
65

1 have erroneously been incorporated into the New Zealand assembly. Similarly, although the  
2 coding sequences have not yet been made publically available, we suspect that the lower  
3 CEGMA statistics together with the reported low average exon number (4 exons per gene vs  
4 9 exons in the V4 assembly) and transcript lengths (666 bp vs 1,237 bp in the V4 assembly;  
5  
6 see below) suggest that frameshift errors due to INDELs from uncorrected PacBio reads  
7 and/or assembly gaps ( $n = 172,892$ ) are affecting coding sequences broadly throughout the  
8 NZ assembly. When compared against publically available genomes in WormBase ParaSite,  
9 the V4 genome assembly represents the most contiguous and resolved nematode assembly  
10 outside of the *Caenorhabditis* genus, and second-most complete Clade V parasitic nematode  
11 assembly based on the presence of universally conserved orthologs ( $n = 20$ ; *Ancylostoma*  
12  
13 *ceylanicum* is the most complete based on BUSCO analysis; **Additional file 1: Figure S4 A,B**).  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65

1  
2 242 **Table 1. Genome assembly summary statistics**  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65

	<i>H. contortus</i> V4 Chromosomes	<i>H. contortus</i> V4 Haplotypes	<i>H. contortus</i> V1 <sup>1</sup>	<i>H. contortus</i> McMaster <sup>2</sup>	<i>H. contortus</i> New Zealand <sup>3</sup>	<i>C. elegans</i> WB <sup>4</sup>
<b>Genome size (bp)</b>	283,439,308	248,771,548	369,846,877	319,640,208	465,720,674	100,286,401
<b>Scaffolds</b>	7	3,907	23,860	14,419	7	7
<b>Scaffold N50 (bp)</b>	47,382,676	167,270	83,287	56,328	83,970,805	17,493,829
<b>Scaffolds ≥ N50</b>	3	174	1,151	1,684	3	3
<b>Contigs</b>	192	4,226	65,523	55,322	172,899	7
<b>Contig N50 (bp)</b>	3,801,457	117,552	20,808	12,900	3,087	17,493,829
<b>Contigs ≥ N50</b>	20	439	4,943	6646	43,632	3
<b>Scaffold N content (bp)</b>	5,699,711	14,755,378	23,804,39	20,495,720	17,288,295	0
<b>Scaffold gaps</b>	185	323	41,663	40,903	172,892	0

243 1. Haem V1 [21]: haemonchus\_contortus.PRJEB506.WBPS10.genomic.fa

244 2. McMaster [22]: haemonchus\_contortus.PRJNA205202.WBPS11.genomic.fa

245 3. New Zealand [31]: ENA accession GCA\_007637855.2: GCA\_007637855.2\_ASM763785v2\_genomic.fna

246 4. *C. elegans* WBcel235: caenorhabditis\_elegans.PRJNA13758.WBPS9.genomic\_softmasked.fa

247

248

249 **Resolving haplotypic diversity and repeat distribution within the chromosomes**

250 Diverse haplotypic sequences were identified, particularly in our preliminary PacBio

251 assembly (**Additional file 1: Figure S1**; Initial assembly size: 516.6 Mbp), as a result of

252 sequencing DNA derived from a large pool of individuals necessary to obtain sufficient

material. We identified a dominant group of haplotypes in our PacBio assemblies which,  
1  
2 together with Illumina read coverage derived from a single worm, was used to partially  
3  
4 phase the chromosome assembly sequence by manually curating the haplotypes in the  
5 reference assembly. Mapping of the additional haplotypes to the reference-assigned  
6  
7 haplotype (**Figure 2 A**) highlighted both the diversity between the haplotype and  
8 chromosomal placed sequences (**Figure 2 B; Figure 2 C**, second circle) and coverage of the  
9  
10 haplotypes, spanning approximately  $66.8 \pm 24.1\%$  on average of the chromosome  
11 sequences (**Figure 2 C**, third circle; range  $85.1 \pm 12.3\%$  for chromosome I to  $51.8 \pm 16.6\%$  on  
12 the X chromosome, measured in 1 Mbp bins). We determined the relative usage of the  
13  
14 reference-assigned haplotype and those detected in the alternate haplotype assembly for  
15 individuals in the MHco3(ISE).N1 population from which the samples used in the assembly  
16  
17 were derived; whole-genome short-read sequencing of 11 single worms (10 male and 1  
18 female, identified by X chromosome coverage; female (XX) = full coverage, male (XO) = half  
19 coverage) followed by competitive read mapping to the reference genome and alternate  
20  
21 haplotypes revealed clear evidence of haplotype switching based on differences in the  
22 coverage profile throughout the genome and between samples (**Figure 2 D**). Haplotypic  
23 blocks are evident, including large regions on chromosome 2 and 5 that are devoid of  
24  
25 haplotype diversity (full coverage of the reference; blue), as well as regions in which  
26 individuals lack any read coverage due to the absence of the reference haplotype in their  
27 genome (e.g. red region on chromosome 1, samples 09 & 10). These data suggest that, in  
28  
29 spite of being semi-inbred through a limited number of single pair matings [32], complex  
30  
31 haplotypes are segregating among individuals in this population, and that the frequency and  
32  
33 diversity of these haplotypes are non-randomly distributed throughout the genome.  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65

1           Despite the reduction in genome size of the V4 assembly, the use of multiple long molecule  
2           technologies allowed a greater representation of repetitive sequences, accounting for  
3           approximately 36.43% (compared with 30.34% in V1) of the genome (**Additional file 2:**  
4           **Table S3**). Particular repeat classes, including LINEs, LTRs and DNA elements, were enriched  
5           towards the middle of the chromosomes and the chromosome ends (**Figure 2 C**, inner  
6           circle; **Additional file 1: Figure S5**). This distribution pattern is negatively correlated with  
7           observed recombination rate domains in both *H. contortus* [33] and *C. elegans* [34],  
8           suggesting that transposon sequences accumulate in regions of low recombination rate  
9           throughout the genome. On the X chromosome, this enrichment is striking in the sub-  
10          telomeric regions and, in particular, from approximately 8 to 35 Mbp along the  
11          chromosome (**Figure 2 C**, inner circle; **Additional file 1: Figure S5**). From a technical  
12          perspective, the high frequency of transposon sequences explains the increased frequency  
13          of shared sequences between X and the other autosomes, rather than between autosomes  
14          alone (**Figure 1 A**). Finally, the pattern of repetitive elements throughout the chromosomes  
15          and not the coding sequences (shown below) correlates with and perhaps is driving the  
16          nucleotide composition profile throughout the genome (**Figure 2 C**, outer circle).  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65

---

**Figure 2. Haplotype and repeat distribution within and between chromosomes**

**A.** Genome graph derived from the whole-genome alignment of chromosome and alternate haplotype assemblies visualised using *Bandage*. Colour density reflects alternative haplotype density frequency; dark red indicates longer reference only-haplotypes. **B.** Close view of genome graph structure (box), highlighting an example of alternate haplotypes from 7.36 to 7.42 Mbp on Chr X. Aligned regions are represented by thick red lines and paths between regions as thin black lines. DNA alignment of the example 56 kbp alternate

302 haplotype to the genome using *nucmer* and visualised using *Genome Ribbon* demonstrates  
1 the extent of the genetic diversity between these sequences. Two genes are present (red  
2 horizontal bars), of which their exons (red ticks) are predominantly found in regions of  
3 conservation (ribbons connecting the top and bottom sequences). **C.** Circos plot of genome-  
4 wide haplotype and repeat diversity. Outer: mean GC nucleotide frequency in 100 kbp  
5 windows. Dashed line represents genome-wide mean GC content (0.429); Second:  
6 Distribution (x-axis) and length (y-axis) of mapped haplotype sequences to the  
7 chromosome, with sequence similarity between the mapped haplotype and chromosomes  
8 indicated by the grey-scale gradient (high similarity = dark, low = light); third: proportion of  
9 chromosome covered by at least one haplotype, measured in 1 Mbp windows; inner: density  
10 of annotated repeats, including LINEs (green), SINEs (red), DNA (blue), and LTRs (orange).  
11 Data is visualised using the R package *circulize*. **D.** Haplotype switching between  
12 chromosome and additional haplotypes present in the genome assembly. Whole-genome  
13 sequencing was performed on individual worms ( $n = 11$ ; MHco3(ISE).N1 strain, used in the  
14 genome assembly), after which reads were mapped competitively against the chromosomes  
15 and haplotype sequences. Normalised read depth coverage is presented, which was  
16 calculated per 100 kbp relative to the “maximum genome-wide coverage” set at the 95%  
17 quantile to prevent high-coverage outlier regions from skewing coverage estimates. Blue =  
18 ~100% of maximum coverage; Yellow = ~50% of maximum coverage; Red = ~ 0% of  
19 maximum coverage.  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38

---

  
39 322  
40  
41 323  
42  
43  
44 324 **Generation of a high-quality transcriptome annotation incorporating short and long reads**  
45  
46 325 The genome improvements made from the draft V1 assembly to the chromosomal assembly  
47 resulted in a reduction of genome size and, at the same time, an increase in the number of  
48 core conserved genes. The parsimonious conclusion is that coding sequences were easier to  
49 detect in the new assembly and that the reduction of genome size removed previously  
50 annotated but redundant genes located on alternate haplotypes in the draft assembly.  
51  
52  
53  
54 328 Considering these broad-scale changes, we undertook a *de novo* gene finding strategy,  
55  
56  
57 329  
58  
59 330  
60  
61  
62  
63  
64  
65

1           331 rather than relying solely on the previous gene set, to annotate the V4 genome. To do so,  
2           332 we used our comprehensive Illumina short-read RNAseq collection sampled from seven  
3           333 stages of the parasite life cycle (previously described in the analysis of the V1 genome [21]),  
4           334 and supplemented this with PacBio Iso-Seq long-read sequencing of full-length cDNA from  
5           335 adult male, female, and juvenile L3 parasites. To experimentally enhance the recovery of  
6           336 low-abundant transcripts, a duplex-specific nuclease (DSN)-normalisation was performed on  
7           337 the full-length cDNA molecules prior to PacBio sequencing on both RSII and Sequel systems  
8           338 (**Additional file 2: Table S4**), resulting in 23,306 high-quality isoforms after processing.  
9  
10          339  
11  
12          340 The overall annotation strategy is presented in **Additional file 1: Figure S6** and consisted of  
13          341 four key stages: (i) RNAseq evidence was used to train *ab initio* gene finding using Braker  
14          342 [35]; (ii) full-length high-quality isoforms generated using Iso-Seq were used to identify  
15          343 candidate gene models using the annotation tool Program to Assemble Spliced Alignments  
16          344 (PASA; [36]; (iii) evidence from steps (i) and (ii), together with 110 manually curated genes  
17          345 and orthologs from the *C. elegans* and *H. contortus* V1 genomes, were incorporated using  
18          346 EvidenceModeller; and finally (iv) UTRs and multi-evidence gene models were updated  
19          347 using the full-length isoforms once again using PASA. We further curated the genome  
20          348 annotation using Apollo [37]; together, this strategy achieved a high level of sensitivity and  
21          349 precision based on comparison with a subset of manually curated genes (**Additional file 2:**  
22          350 **Table S5**).  
23  
24          351  
25  
26          352 Our V4 annotation comprises 19,489 nuclear genes encoding 20,987 transcripts with  
27  
28          353 approximately 56.4% (23,687 of 41,974) of UTRs annotated. The increase in genome  
29  
30          354 contiguity, together with the incorporation of full-length cDNA sequencing (**Figure 3 A**), has  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65

1 resulted in longer gene (40.5% increase), mRNA (41.9% increase), and exon (30.6% increase)  
2  
3 models on average than the previous V1 genome annotation (**Figure 3 B; Additional file 2:**  
4  
5 **Table S6**); our approach allowed better representation of the longest genes (447,146 bp  
6  
7 compared with 91,953 bp in V1), and intriguingly, identified 43 of the longest 50 (86%) and  
8  
9 73 of the longest 100 gene models on the X chromosome alone. The reduction in haplotypic  
10 sequences together with greatly improved gene models resolved 61% more one-to-one  
11  
12 orthologs with *C. elegans* (n = 7,361) than between *C. elegans* and V1 (n = 4,529)  
13  
14  
15 **(Additional file 2: Table S7)**. Surprisingly, a greater number of one-to-one orthologs are  
16  
17 identified between V4 and *H. placei* (n = 9,970;  
18  
19  
20 [https://parasite.wormbase.org/Haemonchus\\_placei\\_prjeb509/Info/Index/](https://parasite.wormbase.org/Haemonchus_placei_prjeb509/Info/Index/)) than V4 and V1  
21  
22  
23 (n = 9,595), likely reflecting the higher duplication rate and haplotypic nature of the V1  
24  
25  
26 relative to the *H. placei* assembly, although this also highlights the close relationship  
27  
28 between *H. contortus* and *H. placei*, a gastrointestinal pathogen commonly associated with  
29 cattle. The increased orthology between close (*H. placei*) and more distantly (*C. elegans*)  
30 related species provides support for a more complete and biologically representative gene  
31 set in the V4 annotation than the previously described V1 and McMaster annotations.  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47 **Figure 3. De novo transcriptome annotation incorporating full-length cDNA sequencing**  
48 **results in longer and more complete annotations**

49  
50 **A.** Example of an improved and now full-length gene model. HCON\_00001100, a nuclear  
51 hormone receptor family member orthologous to *C. elegans nhr-85*, was incomplete and  
52 fragmented in the V1 assembly (multiple colours representing CDSs from three distinct V1  
53 annotations, HCOI00573500, HCOI00573600, and HCOI00573700). The full-length gene  
54 model, comprised of six exons (black boxes) flanked by UTR sequences (white boxes) is  
55  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65

1           380 supported by full-length cDNA IsoSeq sequences (blue bars). **B.** Comparison of mean length  
2           381 and abundance of annotation features between *H. contortus* V1 (square), *H. contortus* V4  
3           382 (cross), *H. contortus* McMaster (triangle) and *C. elegans* (circle) transcriptome annotations.  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65

---

### 384 **Transcriptional dynamics throughout development and between sexes**

385 Both 2013 publications describing draft assemblies of *H. contortus* provide extensive  
386 descriptions of the key developmental transitions between sequential pairs of life stages  
387 throughout the life cycle [21,22]. We revisited this in the current annotation, and for visual  
388 comparison, we present the top 1,000 most variable transcripts across all life stage samples  
389 in the MHco3(ISE).N1 datasets (**Figure 4 A**; see **Additional file 2: Table S8** for all  
390 differentially expressed genes from pairwise comparisons of the life stages throughout the  
391 life cycle). Most striking was the transcriptional transition between the juvenile free-living  
392 life stages and the more mature parasitic stages, reflecting not only the development  
393 toward reproductive maturity but also significant changes in the parasite's environment. To  
394 explore these developmental transitions more explicitly, we determined the co-expression  
395 profiles for all genes throughout the lifecycle, revealing 19 distinct patterns of co-expression  
396 containing 8,412 genes (41% of all genes), with an average cluster size of 442.7 genes  
397 (**Additional file 1: Figure S7; Additional file 2: Table S9**). Consistent with the pattern of  
398 highly variable genes, two dominant clusters of genes were identified: (i) genes with higher  
399 expression in free-living stages that subsequently decreased in expression in parasitic stages  
400 (cluster 1; n = 1,550 genes; **Figure 4 B**) and (ii) genes that were lowly expressed in free-living  
401 stages that increased in expression in parasitic stages (cluster 8; n = 1,542 genes). Functional  
402 characterisation of the cluster 1 genes revealed 45 significantly enriched gene ontology (GO)  
403 terms (26 molecular function [MF], 13 biological process [BP], & 6 cellular compartment

1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65

405 (CC) terms; <https://biit.cs.ut.ee/gplink/l/8jECkGwQQS>) that predominantly described ion  
406 transport and channel activity, as well as transmembrane signalling receptor activity. For  
407 cluster 8, 22 GO terms were significantly enriched (11 MF, 8 BP, 3 CC;  
408 <https://biit.cs.ut.ee/gplink/l/oJoq6Zh0Re>) and included terms describing peptidase activity,  
409 as well as protein and organonitrogen metabolism. To provide further resolution, we tested  
410 *C. elegans* orthologs of this gene set, which revealed greater resolution on metabolic  
411 processes, including terms associated with glycolysis and gluconeogenesis, pyruvate and  
412 nitrogen metabolism (<https://biit.cs.ut.ee/gplink/l/b3Bx8rhJR5>). Finally, we explored  
413 whether co-expressed genes contained shared sequence motifs in their 5' UTR that may  
414 would be indicative of true transcriptional co-regulation; while broadly poorly represented,  
415 we did identify 70 enriched motif sequences of which 17 were associated with gene sets  
416 with significantly enriched GO terms (**Additional file 2: Table S10**), for example, cysteine-  
417 type peptidase activity in cluster 8 (GATAAGR; motif E-value = 2.90E-02; GO p-value =  
418 4.181E-13), DNA replication in cluster 15 (AAAAATVA; motif E-value = 3.50E-02; GO p-value  
419 = 6.269E-4), and molting cycle, collagen and cuticulin-based cuticle in cluster 16  
420 (GTATAAGM; motif E-value = 2.30E-02; GO p-value = 1.926E-6), suggestive of either direct or  
421 indirect co-regulation of expression.  
422  
423 The greatest differences in gene expression between life stages are between L4 (mixed male  
424 and females) and either adult males or adult females, or between adult males and adult  
425 females, with a strong sex-bias towards increased gene expression in adult males  
426 (**Additional file 2: Table S8**). Consistently with *C. elegans*, sex-biased expression is not  
427 randomly distributed throughout the genome [38,39]; only 5% (146/2,923) of *H. contortus*  
428 male-biased genes are found on the X chromosome (relative to 14.3% of all genes found on

1       429 the X chromosome;  $\chi^2 = 175.83$ , df = 1, p = 3.84E-40), whereas more female-biased genes  
2       430 are X linked than expected (24.6% [110/476];  $\chi^2 = 33.59$ , df = 1, p = 6.78E-9). Comparison of  
3       431 sex-biased genes on the X chromosome revealed that, while almost completely absent in  
4       432 males, X-linked female-biased expression identified *C. elegans* orthologs previously  
5       433 associated with organism development processes, including developmentally upregulated  
6       434 genes (e.g. *smad4*, transcription factor AP2, *vab-3*), genes involved in egg development (e.g.  
7       435 *egg-5*, *vit-1*, *cpg-2*, *egg-1*), as well as genes involved in developmental timing (e.g. *lin-14*,  
8       436 *kin-20* [40]) that are expressed downstream of the sex-determination regulator *tra-1*. Sex  
9       437 determination in *C. elegans* is mediated by the ratio of X chromosomes to autosomes,  
10      438 whereby organisms with a single X chromosome will develop as a male (X-to-autosome ratio  
11      439 = 0.5:1, i.e., XO) and those with two X chromosomes develop as a hermaphrodite (XX; X-to-  
12      440 autosome ratio = 1:1); similarly, *H. contortus* males are XO and females are XX [29,30]. The  
13      441 detection of the X to autosome copy number initiates a regulatory pathway that first  
14      442 activates a dosage compensation response that down-regulates gene expression in  
15      443 hermaphrodites by one-half, followed by sex determination mediated by the conserved  
16      444 transcription factor *tra-1*; in XX individuals *tra-1* is activated, whereas *tra-1* is repressed in  
17      445 XO individuals. Although the similarities in the XX/XO chromosome structures between *C.*  
18      446 *elegans* and *H. contortus* suggest that the mechanisms controlling sex determination may be  
19      447 the same, this has not been previously shown. Comparison of *H. contortus* RNAseq data  
20      448 from males and females showed no significant difference in the distribution of expression  
21      449 between autosomes and X-linked genes, and considering only 9.24% of genes (110 + 146 of  
22      450 2,770 X-linked genes) are differentially expressed on the X chromosome between sexes,  
23      451 these data suggest that X chromosome dosage compensation is active in XX females. To  
24      452 explore this further, we mapped *H. contortus* orthologs of *C. elegans* genes involved in the X  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65

453 chromosome sensing and dosage compensation response and subsequent sex  
1  
2 454 determination pathway (**Figure 4 C, Additional file 2: Table S11**; Adapted from WormBook  
3  
4 455 [41]). Five of the six *C. elegans* orthologs associated with the dosage compensation complex  
5  
6 456 (DCC) were found, however, almost all genes upstream that initiate the recruitment of the  
7  
8 457 DCC are absent, including *xol-1* and *sdc-1,-2, and -3*. In *C. elegans*, chromosome dosage is  
9  
10 458 determined by measuring the expression of key genes on the X and autosomes; in *H.*  
11  
12 459 *contortus*, only two of the *C. elegans* X-linked orthologs - *fox-1* and *sex-1* - are present,  
13  
14  
15 460 however, these are located on the autosomes. Thus, while the DCC may still play a role in  
16  
17 461 dosage compensation in *H. contortus*, the mechanism used to determine the X-to-autosome  
18  
19 462 ratio is likely different from *C. elegans*. Identification of most orthologs downstream of *her-*  
20  
21  
22 463 1, which receives X chromosome dosage signals to initiate sex determination, suggests  
23  
24  
25 464 greater conservation of this mechanism between species. Notable is the absence of *H.*  
26  
27 465 *contortus* orthologs of *tra-2* and *fem-3*, as they are key members of the pathway that  
28  
29 466 receive regulatory signals promoting male or female developmental pathways, respectively,  
30  
31  
32 467 in *C. elegans*.  
33  
34  
35  
36  
37  
38  
39 468  
40  
41 469  
42  
43  
44 470 **Figure 4. Transcriptional dynamics throughout development and between sexes**  
45  
46 471 **A.** Heatmap of the most variable transcripts across the life stages. Transcript abundance  
47 measured as transcripts per million (TPM) was determined using Kallisto, from which the  
48 TPM variance was calculated across all life stages. The top 1,000 most variable transcripts  
49  
50 473 with a minimum of 1 TPM are presented. Genes differentially expressed in each pairwise  
51  
52 474 comparison throughout the life cycle are presented in **Additional file 2: Table S8. B.**  
53  
54 475 Transcript co-expression profiles across life stages of *H. contortus*. Cluster 1, which contains  
55  
56 476 the most co-expressed genes, is presented and broadly reflects the transition between  
57  
58 477 juvenile free-living life stages and mature parasitic stages. The remaining profiles are  
59  
60 478  
61  
62  
63  
64  
65

479 presented in **Additional file 1: Figure S8**, with genes by cluster described in **Additional file 2:**  
1  
480 **Table S9. C.** Schematic of the chromosome dosage sensing and sexual development  
2  
481 pathway of *C. elegans*. Adapted from WormBook [41], we highlight the presence of *H.*  
3  
482 *contortus* orthologs (**bold**, underlined) of *C. elegans* genes previously demonstrated to  
4  
483 promote male (blue) or female (green) development. Missing orthologs from the *H.*  
5  
484 *contortus* genome are translucent. A full list of *C. elegans* genes and *H. contortus* orthologs  
6  
485 from the pathway are presented in **Additional file 2: Table S11.**  
7  
486

---

  
8  
487  
9  
488

### Transcriptome complexity is generated by extensive *cis*- and *trans*-splicing

20  
489 The significantly improved genome contiguity and transcriptome representation prompted  
21  
490 us to revisit the prevalence and characteristics of RNA splice variation throughout the  
22  
491 transcriptome. Spliced leader (SL) *trans*-splicing, which commonly functions to split  
23  
492 polycistronic mRNAs into monocistronic units by the addition of a common SL sequence to  
24  
493 unrelated pre-mRNAs to form mature mRNAs, has been extensively described in *C. elegans*  
25  
494 [42] and has been shown previously to operate in *H. contortus* [43–45] (**Figure 5 A**). To  
26  
495 explore this further in the V4 annotation (**Additional file 1: Figure S8 A,B**), we identified  
27  
496 7,293 (37.4%) and 1,139 (5.8%) genes with evidence of SL1 and SL2 *trans*-splicing,  
28  
497 respectively, within 200 bp upstream or 50 bp downstream of the first start codon (see  
29  
498 **Additional file 2: Table S12** for complete list of genes with SL1 and/or SL2). However,  
30  
499 additional evidence of *trans*-splicing was observed; 7,895 and 1,222 annotated transcripts  
31  
500 contained SL1/SL2 sequences, and a further 6,883 internal SL cut sites (overlapping internal  
32  
501 exons, excluding the 5' and 3' exons) indicated additional non-annotated isoforms were  
33  
502 present in the transcriptomes. While our results likely underestimate the frequency of  
34  
503 *trans*-splicing due to the modest RNA sequencing depth available, this approach reveals a  
35  
504 higher frequency of *trans*-splicing than previously described and presents additional  
36  
505

1 evidence for the prediction and annotation of novel gene isoforms. Further, we identify  
2 degeneracy in the nucleotide sequences of both SL1 and SL2 sequences (**Additional file 1:**  
3 **Figure S8 C**), and sequence conservation in the genome sequence immediately upstream of  
4 the splice site, particularly for SL2 (**Additional file 1: Figure S8 D**). In *C. elegans*, SL1  
5 sequences are typically associated with the non-operon genes, or the first coding sequence  
6 of an operon, whereas SL2 sequences are used for second and subsequent coding  
7 sequences in a polycistronic operon. The intercistronic regions between SL2-spliced genes  
8 and their upstream genes in *C. elegans* are typically short, usually approximately 50 to 200  
9 bp on average, with the majority less than 2 kbp [42]. In *H. contortus*, the intergenic  
10 distance upstream of genes with SL1 sequences is largely consistent with that of genes  
11 without SL1 sequences (**Figure 5 B**, left plot; median = 8,557 bp for SL1-spliced genes,  
12 8,603.5 bp for genes without a SL); however, a small proportion of SL1-spliced genes are  
13 found less than 100 bp downstream of the nearest gene. Conversely, while a larger  
14 proportion of genes with SL2 sequences are immediately downstream of the nearest gene,  
15 as expected (median = 566.5 bp), a substantial proportion (432 [35.4%]) are found with a  
16 large intergenic distance (> 2 kbp) to the nearest upstream gene (**Figure 5 B**, right plot).  
17 These observations suggest that while broadly consistent with the mechanism by which  
18 SL1/SL2 sequences are used in *C. elegans*, (i) *H. contortus* SL1-sequences may be found in  
19 downstream genes within a polycistron, (ii) SL2-spliced genes are not necessarily in a  
20 polycistron but may function independently, and/or (iii) that recognition of SL2 splicing  
21 within an operon remains effective over large physical distances, for example, the 2064 bp  
22 region between *deg-3* (*HCON\_00039370*) and *des-2* (*HCON\_00039380*) operon of *H.*  
23 *contortus* [46]. Although these variations have been described in *C. elegans* [47], they seem  
24 to be more of the norm rather than the exception in *H. contortus*.  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65

529  
1  
2 530 We extended the splice variation analysis to characterise *cis*-splicing of transcripts, with a  
3  
4 531 particular focus on differential isoform usage between life stages. Using a conservative  
5  
6 532 threshold to detect differential splicing, we identified 1,055 genes with differentially spliced  
7  
8 533 transcripts impacting at least one intron, representing on average 283 genes (FDR < 0.05;  
9  
10 534 range: 43-418) that switch between different isoforms in each pairwise comparison of life  
11  
12 535 stages throughout the lifecycle (**Figure 5 C**; see **Additional file 2: Table S13** for a summary  
13  
14 536 per pairwise comparison between life stages, and **Additional file 2: Table S14** for list of  
15  
16 537 genes that appear differentially spliced between pairwise life stages throughout the life  
17  
18 538 cycle). For example, 231 of 1,068 genes (21.63%; **Additional file 2: Table S13**) show  
19  
20 539 evidence of differential splicing between adult male and female samples. The most  
21  
22 540 differentially spliced intron is located in *HCON\_00073480*, a one-to-one ortholog of *C.*  
23  
24 541 *elegans daz-1* (**Figure 5 D**). *daz-1* is an RNA-binding protein involved in oogenesis but not  
25  
26 542 spermatogenesis in *C. elegans* [48]; in *H. contortus*, *HCON\_00073480* is predominantly  
27  
28 543 expressed in females (average TPM = 636.22; highest expression in hermaphrodites in *C.*  
29  
30 544 *elegans*), and is expressed in males, albeit approximately 23-fold less (average TPM =  
31  
32 545 26.71) (**Figure 5 E**). Only a single transcript is described in *C. elegans*; however, a total of five  
33  
34 546 isoforms are annotated in the *H. contortus* genome, with a clear difference in sex-specific  
35  
36 547 expression of exon 9 that is present in males and absent in females. Further, additional  
37  
38 548 differential splicing was found in the splice donor site of exon 8, and while both alternate  
39  
40 549 splice sites are present in both sexes, the sites are used equally in males but the longer  
41  
42 550 transcript is used in around 60% of transcripts in female worms (**Figure 5 F**).  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57 551  
58  
59  
60 552  
61  
62  
63  
64  
65

---

553

1

2

3 **Figure 5. Extensive *trans*- and *cis*-splicing of gene transcripts throughout the life cycle**

4

5 Example of a gene model, HCON\_00096900, visualised using *Apollo* that contains two SL1  
6 trans-spliced sites upstream of the transcriptional start sites of two distinct isoforms, each  
7 supported by full-length IsoSeq sequence reads. **B.** Distribution of intergenic distance  
8 (log<sub>10</sub>[bp]) between the start coordinates of genes containing SL1 (red), SL2 (blue)  
9 sequences, or with no SL sequence (grey), and the end coordinate of the gene immediately  
10 upstream. In the second plot, the relative frequency of SL1 or SL2 sequences is presented. **C.**  
11  
12 Heatmap of the top 100 differentially spliced introns identified in pairwise comparisons  
13 throughout the life cycle. Data are presented as the proportion of differentially excised  
14 intron, ranging from 0 (light) to 1 (dark), for each intron cluster for which differential splicing  
15 was identified. A complete list of genes with evidence of differential splicing is presented in  
16 **Additional file 2: Table S14.** **D.** Heatmap of the top 100 differentially spliced introns  
17 between adult male and adult female worms. The scale is the same as shown in (C). **E.**  
18 RNAseq expression data for adult female, adult male, all life stages combined, relative to the  
19 annotated isoforms for HCON\_00073480, a one-to-one ortholog of the *C. elegans* daz-1. **F.**  
20 Quantification of differential splicing of HCON\_00073480 between adult female (top) and  
21 adult male (bottom) samples, determined using leafcutter. Focused on exons 8, 9, & 10  
22 (black bars), the red lines depict the linking of splice junctions, with the thickness of the line  
23 relative to the proportion of reads that support that junction. In this example, two  
24 differential splicing events are observed; (i) alternate splice donors of exon 8, and (ii) exon  
25 skipping of exon 9.

26

27

28

29

30

31

32

33

34

35

36

37

38

39

40

41

42

43

44

45

46

47

48

49

50

51

52

53

54

55

56

57

58

59

60

61

62

63

64

65

**Distribution of global genetic diversity throughout the chromosomes**

H. contortus is a globally distributed parasite, characterised by large population sizes and in turn, high genetic diversity [25–27,49]. We have characterised the distribution of genetic diversity throughout the V4 chromosomes by exploiting a recent large scale analysis of global genetic diversity [25] that we have extended by sequencing an additional 74

1 individual worms from previously unsampled regions including Pakistan, Switzerland, USA,  
2  
3 and the United Kingdom (**Figure 6 A**;  
4  
5 [https://microreact.org/project/hcontortus\\_global\\_diversity](https://microreact.org/project/hcontortus_global_diversity)). Principal component analysis  
6  
7 (PCA) of mtDNA diversity (**Figure 6 B; Additional file 1: Figure S9**) among 338 individual  
8  
9 parasites showed that relationships between isolates are congruent with those described by  
10  
11 [25]; further, we show that Pakistani worms were genetically positioned between South  
12  
13 African and Indonesian individuals, whereas the US and the UK samples were found among  
14  
15 globally distributed samples that likely shared a common ancestral population that was  
16  
17 distributed globally via modern human movement [25]. Of note, UK parasites were largely  
18  
19 genetically distinct from the French worms despite being geographically close, and yet,  
20  
21 Swiss parasites were genetically closer to African worms than other European parasites. The  
22  
23 US and UK samples consisted of both laboratory-maintained strains and field isolates;  
24  
25 notably, apart from MHco3(ISE).N1, the laboratory strains still retained similar levels of  
26  
27 genetic diversity to the field samples. MHco3(ISE), from which the reference strain  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44 perhaps reflecting its African origins.  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65

600  
601 Genome-wide nucleotide diversity ( $\pi$ ), calculated in 100 kbp windows throughout the  
602 genome, is broadly consistent chromosome-wide between the autosomes (**Figure 6 C-G**).  
603 The high frequency of SNPs private to only a single population (relative to variants shared  
604 between two or more populations) represents a significant proportion of the overall  
605 variation (23.55% of total variation), consistent with previous reports describing high

1 regional genetic diversity [25–27]. On the X chromosome, nucleotide diversity is much lower  
2 than neutral expectation (i.e.,  $\pi_X/\pi_{\text{autosome}} \approx 0.75$  [50]) at a  $\pi_X/\pi_{\text{autosome}}$  median ratio of 0.363  
3 (Figure 6 H), which suggests a sex-ratio bias in favour of males. This observation, together  
4 with the distinct repeat patterns that were negatively correlated with recombination rate,  
5 suggests that recombination may be less frequent on the X chromosome than expected.  
6 Interestingly, and likely related to these observations, is that variation in X chromosome  
7 coverage (Additional file 1: Figure S11 A) increases in populations that are more genetically  
8 distinct from the reference strain (Additional file 1: Figure S11 B); these data suggest that  
9 while the sampled nucleotide diversity is observed to be low on the X chromosome, greater  
10 biological variation in chromosome content and structure is present among the global  
11 sample set that cannot be correctly accounted for by mapping to a linear reference alone.  
12  
13 Finally, we explored evidence of genetic differentiation between populations throughout  
14 the genome, revealing multiple discrete regions in each chromosome with high  $F_{ST}$  among  
15 populations (Figure 6 C-H; bottom plots). Consistent with previous reports, strong genetic  
16 selection for anthelmintic resistance has impacted genome-wide variation; the largest and  
17 most differentiated region is found on chromosome 1 surrounding the beta-tubulin isotype  
18 1 locus (HCON\_00005260; 7027492 to 7031447 bp), a target of benzimidazole-class  
19 anthelmintics for which known resistant mutations have been identified in these  
20 populations [25], as well as a broad region of elevated  $F_{ST}$  from approximately 35 to 45 Mbp  
21 on chromosome 5 previously shown to be associated with ivermectin resistance in  
22 controlled genetic crosses [51]. The top 1% of  $F_{ST}$  outlier regions per chromosome (280 10  
23 kbp regions in total) contains 268 genes enriched for 10 GO terms (7 MF & 3 BP;  
24 https://biit.cs.ut.ee/gplink/l/HKi5aHzqTa ), largely describing proteolytic and cysteine-type  
25

1       630 peptidase activity. Analysis of  $d_N/d_S$  from alignments of 9,970 one-to-one orthologs from *H.*  
2       631 *contortus* and the closely related *H. placei* identified 108 genes putatively under positive  
3       632 selection ( $d_N/d_S > 1$ ). Gene set analysis revealed no enrichment for specific GO terms among  
4       633 these genes; closer inspection revealed a large number of genes with no known function,  
5       634 despite most having orthologs in two or more clade V nematodes, highlighting both the  
6       635 challenge of interpreting genetic novelty in non-model organisms and the clear need for  
7       636 functional genomics to provide missing detail in the annotation of organisms like *H.*  
8       637 *contortus* [52,53]. We did, however, identify a number of genes likely associated with  
9       638 host/parasite interactions (e.g. OV39 ortholog [HCON\_00044830], multiple ShKT domain  
10      639 and signal peptide-containing genes [e.g. HCON\_00095770, HCON\_00124950,  
11      640 HCON\_00108210], BAT-2 domain [HCON\_00173800]), and broad development and  
12      641 neurological function (e.g. *nck-1*, *snt-2*, *flp-5*, *ned-8*, *eat-17*), which could reflect adaptation  
13      642 to different host species. We extended this analysis to *H. contortus* variation derived from  
14      643 the global diversity dataset; after filtering, 849 (4.8% of 17,607) genes were identified with a  
15      644  $d_N/d_S$  greater than 2. Similar to the between species analysis, no significant GO terms were  
16      645 enriched in this dataset. However, G-protein coupled receptors (e.g. *aex-2*, *npr-4*, *sdr-7*, *srb-*  
17      646 *17*, *npr-27*, *seb-4*, *seb-2*), and G-proteins (e.g. *goa-1*, *gpa-4*, *gpa-16*), ion and neuropeptide  
18      647 gated transmembrane channels (e.g. *acr-12*, *acr-19*, *acr-20*, *lgc-9*, *lgc-34*, *lgc-46*, *lgc-39*,  
19      648 *mod-1*, *ccb-2*, *kvs-4*) were among the most abundant *pfam* domains identified, highlighting  
20      649 the adaptive potential of sensory and response pathways to external stimuli.  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60  
61

654

1

2

3 **Figure 6. Genome- and population-wide genetic diversity**

---

4

5 **A.** Global distribution of populations from which genomic data were derived, including 264

6 samples from Sallé *et al.* (circles) and 74 new samples (triangles) presented. The insert

7 highlights the dense sampling from the United Kingdom (this study) and France [25]. **B.**

8

9 Principal component analysis of genetic diversity among the globally distributed samples,

10

11 using 2,401 mtDNA variant sites. The reference strain, MHco3(ISE).N1, is indicated. **C to H.**

12

13 Chromosome-wide comparison of genetic diversity. Each comparison contains three panels.

14

15 Top: Estimate of nucleotide diversity ( $\pi$ ) calculated in 100 kbp windows, accounting variants

16

17 found only in a single population (dark), or shared between two or more populations (light).

18

19 Middle: Estimate of genetic differentiation ( $F_{ST}$ ) between populations, measured in 100 kbp

20

21 windows. Bottom: Gene density per chromosome, where each gene is represented as a black

22

23 line on the sense (upper lines) or antisense (lower lines) DNA strand of the chromosome.

24

25 **667**

---

26

27

28

29 **668**

30

31 **669**

32

33

34 **670**

35

36

37

38

39

40

41

42

43

44

45

46

47

48

49

50

51

52

53

54

55

56

57

58

59

60

61

62

63

64

65

671    **Discussion**

1  
2  
3    672    The *H. contortus* MHco3(ISE).N1 V4 assembly presented here represents the largest *de novo*  
4  
5    673    chromosome-scale nematode genome to date and provides a critical genomic resource for a  
6  
7  
8    674    broad group of human and veterinary pathogens present among the Clade V nematodes.  
9

10  
11    675    This assembly is complemented with a high-quality *de novo* genome annotation that  
12  
13    676    incorporates full-length cDNA sequencing to provide a precise transcriptomic resource for  
14  
15  
16    677    within and between species comparison, improved UTR and isoform classification for  
17  
18  
19    678    functional genomics, and should help to improve the annotation of other nematode species.  
20

21    679    Finally, we provide a comprehensive description of genetic and transcriptomic diversity  
22  
23    680    within individuals and across the species, which informs not only our understanding of key  
24  
25  
26    681    developmental trajectories but also the parasite's capacity to adapt when faced with  
27  
28  
29    682    challenges such as drug exposure.

30  
31    683  
32  
33  
34    684    A key challenge to achieving this chromosomal assembly was the presence of extensive  
35  
36  
37    685    genetic diversity in the sequencing data. Although high-quality genome assemblies are  
38  
39  
40    686    becoming more prevalent through the use of long-read sequencing technologies [4,5,54], a  
41  
42    687    feature of these assemblies is that they are derived from a single, few, or clonal individual(s)  
43  
44  
45    688    in which the genetic diversity is low or absent; most of the genetic diversity present in a  
46  
47  
48    689    single diploid individual can be phased and separated into distinct haplotypes (for example,  
49  
50  
51    690    using software such as *Falcon-Unzip* [8] and *Supernova* [7] for PacBio and 10X Genomics  
52  
53  
54    691    data, respectively), and assembled separately. However, even under these scenarios,  
55  
56  
57    692    heterozygosity remains a key technical limitation towards achieving contiguous genome  
58  
59  
60    693    assemblies. Considering the small physical size of *H. contortus*, it was necessary to pool  
61  
62  
63    694    thousands of individuals to generate sufficient material for long-read sequencing, and while

1 we did so from a semi-inbred population, there was still significant genetic diversity present  
2  
3 in the population of worms sequenced. The generation of such diversity from a single  
4  
5 mating pair is likely influenced by reproductive traits including polygamous mating and high  
6  
7 fecundity of *H. contortus* female worms. How such haplotypic diversity is maintained in the  
8  
9 population or is tolerated within single individuals remains unclear, however, it is likely that  
10  
11 699 the holocentric chromosomes of *H. contortus* likely tolerate significant genetic diversity  
12  
13 700 during homologous pairing of chromosomes at meiosis [55]. In spite of this, we curated a  
14  
15 701 major haplotype in the ISE.N1 population to represent the chromosomes and demonstrated  
16  
17 702 how this haplotype and alternate haplotypes segregated in this semi-inbred population. The  
18  
19 703 absence in the literature to date of such chromosome-scale assemblies for physically small  
20  
21 704 and genetically diverse eukaryotic species likely reflects the significant investment required  
22  
23 705 to improve a draft assembly to chromosomes. However, it is likely that as DNA input  
24  
25 706 requirements for long-read sequencing decrease [56], the ability and ease with which highly  
26  
27 707 contiguous genomes can be generated for small and genetically diverse species will be  
28  
29 708 significantly enhanced.  
30  
31 709  
32  
33 710  
34  
35  
36 711 Completion of the X chromosome represented a key advancement in the current version of  
37  
38 712 the genome. In doing so, we revealed striking differences in: (i) the repeat distribution and  
39  
40 713 content, (ii) an enrichment for the longest genes of the genome, (iii) the lack of apparent  
41  
42 714 genetic diversity relative to the autosomes, and (iv) differences in the gene content that  
43  
44 715 regulates X chromosome dosage and sex determination relative to *C. elegans*. Although we  
45  
46 716 currently have little empirical data on recombination rates on the X chromosome, less  
47  
48 717 recombination is expected on the X chromosome relative to the autosomes due to  
49  
50 718 recombination on X only occurring during female gametogenesis. The domain structure  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65

1 revealed by the repeat distribution does suggest that recombination patterns may be  
2 distinct from the autosomes and warrants further investigation. The difference in nucleotide  
3 diversity between the X chromosome and autosomes further emphasises the impact of  
4 putative sex biases on the evolution of the X chromosome. *H. contortus* is polyandrous  
5 [29,33] and highly fecund, so that reproductive fitness is likely overdispersed between males  
6 and females, potentially resulting in different effective population sizes for male and female  
7 parasites. These differences may be exacerbated if males and females are differentially  
8 affected by anthelmintic exposure as has been demonstrated for some compounds [57].  
9  
10 Further, post-zygotic incompatibility between genetically divergent *H. contortus* strains has  
11 been previously demonstrated [58]; however, it is unclear if this occurs within naturally  
12 occurring populations and whether this would, in fact, impact effective sex ratios. The  
13 dichotomy between the lack of low nucleotide diversity and high variance between samples  
14 in the global cohort perhaps reflects a key limitation on the ability of a single linear  
15 reference genome to capture genome diversity in a species with high genetic diversity. From  
16 a technical perspective, if reads from regions of the genome that contain high diversity do  
17 not map to the reference, as demonstrated by differences in sequence read coverage  
18 between reference and haplotype sequences among ISE.N1 individuals, only a proportion of  
19 the total diversity will be captured. Long-read sequencing data (e.g. PacBio or Oxford  
20 Nanopore) from *Haemonchus* populations and bioinformatics approaches designed to  
21 capture more complex variation should help reveal the full extent of genetic variation in this  
22 species and shed light on the role of structural variation in shaping this genetic variation.  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65

1           743 the life cycle. These data identified gene sets that are putatively co-regulated, providing  
2           744 greater insight into coordinated developmental programs than has been described  
3           745 previously. The broad-scale coordinated gene expression associated with ion transport,  
4           746 channel activity, and transmembrane signalling receptor activity throughout the lifecycle  
5  
6           747 was of interest, particularly as ion channels and receptors are the primary targets of a  
7  
8           748 number of anthelmintic drugs used to control *H. contortus* and other helminth species of  
9  
10          749 human and veterinary importance [59]. A better understanding of these changes may  
11  
12          750 explain why some parasites at particular life stages are less sensitive to these drugs [60,61],  
13  
14          751 and may provide the rationale for selective treatment based on the known expression of  
15  
16          752 drug targets. We further explored transcriptome variation by defining and quantifying  
17  
18          753 differential isoform usage throughout the life cycle for the first time. The importance of  
19  
20          754 differential isoform usage is increasingly recognised in the sensitivity of drug targets in  
21  
22          755 nematodes such as *H. contortus*; for example, sex-dependent sensitivity of emodepside is  
23  
24          756 mediated by differential-splicing of the potassium channel *slo-1* [57], and similarly,  
25  
26          757 expression of truncated acetylcholine receptor subunits *acr-8* and *unc-63* isoforms are  
27  
28          758 associated with levamisole resistance [62,63]. Our results almost certainly under-estimate  
29  
30          759 the extent of *cis*-splice variation within and between life stages; stringent filtering resulted  
31  
32          760 in only a subset of genes tested (range: 759-1156 genes per pairwise comparison of life  
33  
34          761 stages), and yet, in each pairwise comparison, approximately 27.32% of genes on average  
35  
36          762 contained evidence of splice variation (**Additional file 2: Table S13**). Further, each pairwise  
37  
38          763 comparison identified cryptic splice sites, indicating that a proportion of splice junctions  
39  
40          764 from novel isoforms undergoing differential splicing were missing from the genome  
41  
42          765 annotation. Ongoing manual curation to the genome annotation to classify isoform variants,  
43  
44          766 together with more fine-grained quantification of transcript usage, will provide additional  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65

1       767 insight into the novelty that differential isoform usage provides, particularly for those genes  
2       768 for which isoform switching but not differential expression occurs between key life-stage  
3       769 transitions. Beyond the life cycle, the impact of high genetic diversity on transcriptomic  
4       770 variation within or between wild isolates is not yet clear, although both technical and  
5       771 biological variation between three divergent strains has recently been described [64]. A  
6       772 greater understanding of this genetic-to-transcriptomic variation is required, especially if  
7       773 phenotypic traits associated with transcriptomic differences, e.g. increased gene expression  
8       774 of a drug transporter correlated with resistance, are to be correctly attributed to genetic  
9       775 selection.

10      776

11      777 Finally, analyses of chromosome-scale assemblies allowed us to explore the distribution of  
12      778 genetic features throughout the genome and to speculate on the evolutionary events that  
13      779 have taken place to explain these distributions. For example, between species analyses have  
14      780 confirmed karyotype structures and/or revealed chromosomal rearrangements among  
15      781 rhabditine nematodes [65,66], filarial nematodes [11], and platyhelminths [67]. The V4  
16      782 assembly allowed the extent of genome rearrangements between *C. elegans* and *H.*  
17      783 *contortus* to be described for the first time, and similarly, we revealed numerous gene and  
18      784 repeat family expansions that were missed in the previous assembly and produced a more  
19      785 complete catalogue of coordinated transcription within and between life stages. As further  
20      786 genomes are completed from the Strongylid clade, determining the rates of structural  
21      787 evolution in this important group of parasitic organisms should become possible.

22      788 Investigating such evolutionary processes will be important for understanding genome-wide  
23      789 variation associated with trait variation, including important adaptive traits such as  
24      790 response to climate change and anthelmintic resistance. *H. contortus* is an established

1       791 model for understanding parasite response to anthelmintics and subsequent evolution of  
2       792 resistance to key drugs used to control parasitic nematodes [20]. Our chromosomal  
3       793 assembly has been instrumental in mapping genetic variation associated with ivermectin  
4       794 [51] and monepantel [68] resistance. In the case of ivermectin, the identification of the  
5       795 major genomic region on chromosome 5 associated with resistance would not have been  
6       796 possible without a chromosome-scale assembly, and highlights the potential limitations on  
7       797 the interpretation of such variation using less-contiguous genome assemblies as a reference  
8       798 [28].  
9  
10      799  
11  
12

13      800 **Conclusions**  
14  
15

16      801 Our chromosome-scale reference genome and comprehensive annotation of the  
17      802 MHco3(ISE).N1 isolate provides a significant step forward toward resolving the complete  
18      803 genome composition and transcriptional complexity of *H. contortus* and is an important tool  
19      804 for understanding how genetic variation is distributed throughout the genome in this  
20      805 species. The increasing availability of long-read sequencing together with technological  
21      806 advances in the sequencing of small diverse organisms holds great promise towards the  
22      807 rapid generation of high-quality contiguous genomes for organisms like *H. contortus*; future  
23      808 sequencing, including the generation of new *de novo* genomes and transcriptomes from  
24      809 individuals sampled throughout its geographic range, will provide great insight into how  
25      810 genetic variation is maintained and/or changes in response to persistent challenges,  
26      811 including anthelmintic exposure and climatic change, in this hyper-diverse globally  
27      812 important pathogen.  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65

813    **Materials and Methods**

814    **Parasite material**

815    All parasite material used in the genome assembly is derived from the semi-inbred strain of  
816    *H. contortus*, MHco3(ISE).N1. This strain was derived from a single round of single-worm  
817    mating using the MHco3(ISE) strain (protocol and procedure described in Sargison *et al.*  
818    [32]), and maintained as previously described [21].

820    **Iterative improvement of the genome assembly**

821    Our approach to improving the V1 genome to the final V4 chromosomal assembly is  
822    outlined in **Additional file 1: Figure S1**. The sequence data generated for this study are  
823    described in **Additional file 2: Table S1**.

825    Initial manual improvement on the V1 genome focused on iterative scaffolding with *SSPACE*  
826    [69] and gap-filling with *IMAGE* [70] using Illumina 500 bp and 3 kbp libraries, with  
827    additional low coverage data from 3, 8 and 20 kbp libraries generated using Roche 454  
828    sequencing. These improvements were performed alongside breaking of discordant joins  
829    using *Reapr* [71], and visual inspection using *Gap5* [72]. Substantial genetic variation was  
830    present in the sequence data due to the sequencing of DNA derived from a pool of  
831    individuals, resulting in a high frequency of haplotypes that assembled separately and  
832    therefore present as multiple copies of unique regions in the assembly. We surmised that  
833    much of the assembly fragmentation was due to the scaffolding tools not being able to deal  
834    with the changing rates of haplotypic variation so we attempted to solve this manually in  
835    *gap5*. We were aware that we did not have sufficient information to correctly phase these

1       836    haplotypes, so instead, we chose the longest scaffold paths available, accepting that our  
2       837    scaffolds would not represent single haplotypes but would rather be an amalgamation of  
3       838    haplotypes representing single chromosomal regions. This approach was initially difficult  
4  
5       839    and time-consuming, and was further confounded by a large number of repetitive  
6  
7       840    sequences present in each haplotype.  
8  
9  
10      841  
11  
12      842    Significant improvements in scaffold length were gained by the integration of OpGen  
13  
14      843    (<http://www.opgen.com/>) optical mapping data. Optical mapping was performed following  
15  
16      844    methods described previously [14] with the following exceptions: high molecular weight  
17  
18      845    DNA was prepared by proteinase K lysis of pools of ~500 *H. contortus* L3 embedded in  
19  
20      846    agarose plugs, after which one of three restriction enzymes (KpnI, AfII and KpnI) were used  
21  
22  
23      847    to generate three separate restriction map datasets. Initial attempts to generate a *de novo*  
24  
25  
26      848    assembly using optical molecules alone was unsuccessful, and therefore, optical contigs  
27  
28  
29      849    were generated using DNA sequence seeds from the genome assembly using  
30  
31  
32      850    *GenomeBuilder* (OpGen) and visualised and manually curated using *AssemblyViewer*  
33  
34  
35      851    (OpGen). Although this approach was successful, it was limited by the quality and integrity  
36  
37  
38      852    of the gap-dense scaffolds and arbitrary nature of the haplotype scaffolding.  
39  
40  
41  
42      853  
43  
44  
45  
46      854    Subsequent integration of PacBio long-read data alongside the optical mapping data  
47  
48  
49      855    resulted in major increases in contiguity. PacBio sequencing libraries were prepared and  
50  
51  
52      856    sequenced using the PacBio RSII system. A total of 32.3 Gbp raw subreads ( $n = 4,085,541$ ,  
53  
54      857    N50 = 10,299 bp) were generated from 33 flow cells, representing approximately 133.8×  
55  
56  
57      858    coverage (based on the estimated genome size of 283 Mbp). A *de novo* PacBio assembly was  
58  
59  
60      859    generated using *Sprai* (v0.9.9.18; <http://zombie.cb.k.u-tokyo.ac.jp/sprai/index.html>), which  
61  
62  
63  
64  
65

1 was mapped to the assembly and used to manually resolve many gaps and resolve some of  
2 the phasing between haplotypes. Using the *Sprai* PacBio *de novo* assembly, we were also  
3 able to incorporate contigs that were previously missing from the working assembly. The  
4 increase in contiguity of the PacBio assemblies, further improved using *canu* v1.3 [73],  
5 revealed two major but diverse haplotype groups segregating at approximately 65% and  
6 30% frequency in the pooled individuals sequenced; the presence of such diverse  
7 haplotypes resulted in a significantly expanded assembly over 500 Mbp. The major  
8 haplotype was more contiguous and therefore was chosen as the primary sequence to  
9 incorporate and improve the assembly. This approach was supported by competitive  
10 mapping of a single worm short read sequencing library (ENA: ERS086777), which was found  
11 to preferentially map to different scaffolds and/or different contigs within scaffolds that we  
12 inferred were different haplotypes in this single worm. Regions in the chromosome with no  
13 ERS086777 coverage, but for which an alternate haplotype was identified with coverage  
14 from the PacBio assembly, were manually removed from the main chromosomes and  
15 replaced. Further, this selective mapping correlated well with the optical contigs, and once  
16 these sequences were removed, much better optical alignments were produced further  
17 improving the assembly. Alternate haplotypes from the PacBio assembly, and those  
18 removed from the main assembly, were analysed as a separate haplotype assembly.  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65

The increase in contiguity and resolution of shorter repetitive regions using PacBio began to reveal chromosome-specific repetitive units. Although these repeats were highly collapsed in the assembly and were typically not spanned by optical molecules, we were able to iteratively identify and join contigs/scaffolds flanking large tandem repeats that had clear read-pair evidence that they only occurred in a single location in the genome (i.e. read pairs

1       884    were all mapped locally once the join was made). These were further supported by local  
2       885    assemblies of subsets of PacBio reads that contained copies of the repeat regions followed  
3       886    by *de novo* assembly using *canu* [73] to reconstruct the flanking unique regions surrounding  
4       887    a repeat. This iterative process resulted in the production of chromosome-scale scaffolds,  
5       888    each terminating with a 6 bp repeat consistent with being a telomeric sequence (sequence  
6       889    motif: TTAGGC).

7       890

8       891    The V3 assembly, used in the *H. contortus* genetic map [33], consisted of five autosomal  
9       892    chromosomes, with the X chromosome in two major scaffolds. The X chromosome-  
10      893    associated scaffolds were identified by synteny to the *C. elegans* X chromosome using  
11      894    *promer* [74], and by the expected ~0.75× coverage relative to the autosomes in pooled  
12      895    sequencing (due to a mixture of male and females present) and ~0.5× in single male  
13      896    sequencing libraries. We resolved the X chromosome using linked read long-range  
14      897    sequencing using the 10X Genomics Chromium platform ([www.10xgenomics.com](http://www.10xgenomics.com)). DNA  
15      898    prepared from pooled worms was used to generate 10X sequencing libraries, which were  
16      899    subsequently sequenced on a HiSeq 2500 using 250 bp PE reads. Reads were mapped using  
17      900    ***bwa mem*** (v0.7.17-r1188), followed by scaffolding using *ARCS* [75] and *LINKS* [76]. A single  
18      901    scaffold was generated using this approach, merging the two major X-linked scaffolds as  
19      902    well as two short sequences putatively identified as being X-linked but previously unplaced  
20      903    (**Additional file 1: Figure S2**). The completion of the X chromosome, together with the  
21      904    polishing of the genome first with *arrow*  
22      905    (<https://github.com/PacificBiosciences/GenomicConsensus>) followed by *Pilon* v1.22 [77],  
23      906    finalised the genome to produce the V4 version presented here.

907  
1  
2 908 **Genome Completeness**  
3  
4  
5 909 Assembly completeness was determined throughout the improvement process using Core  
6  
7  
8 910 Eukaryotic Genes Mapping Approach (CEGMA; v2.5) [78] and Benchmarking Universal  
9  
10 911 Single-Copy Orthologs (BUSCO; v3.0) [79] gene completeness metrics. For comparison,  
11  
12 912 these metrics were also determined for the complete set of nematode genomes from the  
13  
14  
15 913 WormBase ParaSite version 12 release.  
16  
17  
18 914  
19  
20 915 **Haplotype analyses**  
21  
22  
23 916 The identification and discrimination of the major and minor haplotypes in our assembly  
24  
25  
26 917 offered an opportunity to characterise the distribution of haplotypes throughout the  
27  
28 918 genome, and the relative diversity between them.  
29  
30  
31 919  
32  
33  
34 920 To visualise the broad-scale variation between the chromosome and haplotypes, we  
35  
36 921 generated an approximate genome graph using *minigraph*  
37  
38  
39 922 (<https://github.com/lh3/minigraph>; parameters: -xggs). The output GFA file was visualised  
40  
41 923 using *Bandage* [80]; a representative set of haplotypes were chosen for further comparison,  
42  
43  
44 924 which were compared using nucmer and visualised using *Genome Ribbon* [81].  
45  
46  
47 925  
48  
49 926 To quantify the distribution and diversity of haplotypes, we used *minimap* to align  
50  
51  
52 927 haplotypes per chromosome. Haplotype density was calculated in 1 Mbp windows  
53  
54  
55 928 throughout the genome using *bedtools coverage* [82] and visualised using the R package  
56  
57 929 *circlize* [83].  
58  
59  
60 930  
61  
62  
63  
64  
65

1       931 We determined the relative usage of detected haplotypes in the MHco3(ISE).N1 population  
2       932 by using Illumina sequencing of DNA from 11 individual worms and competitively mapping  
3       933 the sequencing data to both the reference and haplotype assemblies using *bwa mem* [84].  
4  
5       934 Relative sequence coverage was determined using *bedtools coverage* [82], which was  
6  
7       935 normalised per sample to enable equivalent comparison between samples.  
8  
9  
10      936  
11  
12      937  
13      938     **Transcriptome library preparation and sequencing**  
14  
15      939 We made use of the RNAseq data generated for each of the life stages as previously  
16      940 described [21], together with additional full-length cDNA (PacBio Iso-Seq) sequencing  
17  
18      941 performed with the specific aim of generating evidence for the annotation. Three life stages  
19  
20      942 were used in the Iso-Seq analyses: adult female, adult male, and L3. RNA was extracted by  
21  
22      943 homogenising the samples in Trizol in Fast prep (MP bio) followed by chloroform extraction  
23  
24      944 and ethanol precipitation of the RNA, after which 1.2 µg of total RNA was converted to  
25  
26      945 cDNA using SmartSeq2 protocol [85]. Size fractionation of equimolar pooled cDNA of the  
27  
28      946 different stages was performed using a SageELF electrophoresis system (Sage Sciences),  
29  
30      947 from which cDNA size fraction pools spanning 400 to 6,500 bp were collected. To improve  
31  
32      948 the coverage of transcripts that differed in abundance and also in length (and to help  
33  
34      949 minimise library preparation and sequencing biases towards shorter and highly abundant  
35  
36      950 transcripts), we used a cDNA normalisation method followed by size fractionation of the  
37  
38      951 pooled normalised cDNA to recover discrete size fractions of transcripts. The cDNA  
39  
40      952 normalisation, which exploits a duplex-specific nuclease (DSN), was performed using the  
41  
42      953 Trimmer-2 cDNA normalisation kit (Evrogen) following the manufacturer's instructions. Size  
43  
44      954 fractionated cDNA pools were sequenced on five individual flow cells on the RSII  
45  
46  
47      955 instrument. To increase the yield of the sequencing data, a subsequent sequencing run was  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65

1       955      performed after pooling the size fractions using a single flow cell on the Sequel instrument.  
2       956      A comparison of the sequences generated by the RSII and Sequel systems is presented in  
3  
4  
5       957      **Additional file 2: Table S4.**  
6  
7  
8       958  
9  
10      959      **Generation of a high-quality non-redundant Iso-Seq dataset**  
11  
12      960      Full-length Iso-Seq subreads were mapped to the V4 genome using *minimap2* (version 2.6-  
13  
14      961      r639-dirty; [86]) with splice-aware mapping enabled (parameter: *-ax splice*). The generation  
15  
16      962      of high-quality Iso-Seq reads was performed following the IsoSeq3 pipeline  
17  
18      963      (<https://github.com/PacificBiosciences/IsoSeq3>), which consisted of four steps: (i) the  
19  
20      964      generation of circular consensus sequencing (CCS) reads from raw subreads (*ccs*  
21  
22      965      *subreads.bam subreads.ccs.bam --noPolish --minPasses 1*), (ii) the removal of library  
23  
24      966      primers, and reorientation of reads (*lima subreads.ccs.bam primers.fa subreads.demux.bam*  
25  
26      967      *--isoseq --no-pbi*), (iii) clustering of transcripts, removal of concatemers and trimming of  
27  
28      968      polyA tails (*isoseq3 cluster subreads.demux.primer\_5p--primer\_3p.bam*  
29  
30      969      *subreads.isoseq3\_unpolished.bam*), and finally, (iv) polishing of the transcripts to generate a  
31  
32      970      high quality set of isoforms (*isoseq3 polish merged\_subreads.isoseq3\_unpolished.bam*  
33  
34      971      *merged\_subreads.bam merged\_subreads.isoseq3\_polished.bam* ).  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44      972  
45  
46  
47      973      **Gene model prediction and genome annotation**  
48  
49      974      A schematic of the strategy used for genome annotation is presented in **Additional file 1:**  
50  
51  
52      975      **Figure S6.** Initial attempts to transfer genome annotations from the V1 assembly to the  
53  
54      976      intermediate V3 assembly using tools such as *Ratt* [87] were unsuccessful, and with the  
55  
56  
57      977      generation of new Iso-Seq data, a new strategy was designed and implemented.  
58  
59  
60      978  
61  
62  
63  
64  
65

1 Existing RNAseq data were individually mapped to the V4 genome using *STAR* v2.5.2 [88],  
2 before combined into a single bam using *samtools merge*. The merged bam was used as  
3 input to *Braker* (version 2.0; [35,89], which can exploit mapped RNAseq information as hints  
4 to train the gene prediction tools *GeneMark-EX* and *Augustus*, substantially improving  
5 genome annotation. The use of exon, intron, and other hints derived from RNAseq or  
6 proteins mapped to the genome can be used to train *Augustus* directly; initial attempts to  
7 incorporate mapped Iso-Seq transcripts as an input to *Braker*, or by training *Augustus* using  
8 exon and intron hints, were successful and produced superior gene models compared with  
9 *Braker* using RNASeq alone. However, as Iso-Seq data encode full-length transcripts that  
10 include UTRs in a single sequence, and that in this data there was little discrimination  
11 between the depth of coverage between UTR and exon, many UTRs were incorrectly  
12 annotated as 5' or 3' exons. As such, only short-read RNAseq data were used as input to  
13 *Braker*, resulting in a primary annotation (GFF<sub>1</sub>; **Additional file 1: Figure S6**).  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65

To incorporate Iso-Seq data into the genome annotation, we used the *PASA* (v2.2.0) genome annotation tool [90] to map the high-quality full-length transcripts to the genome using both *blat* and *gmap*, akin to the use of a *de novo* transcriptome assembly produced from, for example, *trinity*, to generate a transcript database ([https://github.com/PASApipeline/PASApipeline/wiki/PASA\\_comprehensive\\_db](https://github.com/PASApipeline/PASApipeline/wiki/PASA_comprehensive_db)). The *Braker* GFF<sub>1</sub> was subsequently loaded into *PASA* and integrated with the mapped Iso-Seq transcripts to produce the first round merged annotation (GFF<sub>2</sub>; **Additional file 1: Figure S6**).  
  
To incorporate further coding evidence into the annotation, we used *EvidenceModeller* v1.1.1 [36] to apply a differential weighting schema to multiple evidence sources, before

1003 merging them into a single genome annotation. Evidence used included the GFF<sub>1</sub> and GFF<sub>2</sub>  
1  
2  
3 annotations from *Braker* and *PASA* as *ab initio* and transcript evidence, respectively, and  
4  
5 protein sequences from *C. elegans* and *H. contortus* V1 genome assemblies (obtained from  
6  
7 Wormbase ParaSite release 9 caenorhabditis\_elegans.PRJNA13758.WBPS9.protein.fa.gz and  
8  
9 haemonchus\_contortus.PRJEB506.WBPS9.protein.fa.gz, respectively), as well as 110 curated  
10  
11 genes; each protein dataset was mapped to the genome using *exonerate* v2.2.0  
12  
13 (parameters: --model protein2genome --percent 50 --showtargetgff ; [91]) followed by GFF  
14  
15 extraction using *Exonerate\_to\_evm\_gff3.pl* from *EvidenceModeller*. The weighting was  
16  
17 applied as follows: PROTEIN (hc\_v1 \* curated models) = 2; PROTEIN (Ce) = 1;  
18  
19 ABINITIO\_PREDICTION (GFF<sub>1</sub> from Braker) = 2; TRANSCRIPT (GFF<sub>2</sub> from PASA) = 5.  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29 The merged output of *EvidenceModeller* (GFF<sub>3</sub>) was imported back into *PASA* to update  
30  
31 gene models with the full-length transcript information, including updates to UTRs and  
32  
33 evidence for alternative splicing. Consistent with the author's recommendations, two  
34  
35 rounds of iterative update and comparison were performed to maximise the use of the full-  
36  
37 length transcript, i.e., Iso-Seq, data (GFF<sub>4</sub>; **Additional file 1: Figure S6**). Gene IDs were  
38  
39 subsequently renamed with the prefix HCON, and numbered with an 8 digit identifier that  
40  
41 incremented by a value of 10 to allow identifiers to be subsequently added for new features  
42  
43  
44 if required (GFF<sub>5</sub>; **Additional file 1: Figure S6**).  
45  
46  
47  
48  
49  
50  
51  
52 The genome annotation, as well as available evidence used to generate the annotation  
53  
54 (including per life stage RNAseq, Iso-Seq [full-length high quality & CCS reads], coverage and  
55  
56 repeat tracks), was uploaded to the web portal *Apollo* [37] for further visualisation and  
57  
58 manual curation. This platform has been used to provide a means to continually improve  
59  
60  
61  
62  
63  
64  
65

1027 the *H. contortus* genome annotation via manual curation as new information becomes  
1  
2 available. These annotation improvements will be subsequently incorporated into  
3  
4  
5 1029 WormBase ParaSite.  
6  
7  
8 1030  
9  
10 1031 Annotation statistics were determined using *Genome Annotation Generator (GAG)* [92].  
11  
12  
13 1032 Annotation comparisons to a curated gene set to determine the sensitivity and specificity of  
14  
15 1033 the approaches used were performed using *gffcompare* v0.10.1  
16  
17  
18 1034 (<http://ccb.jhu.edu/software/stringtie/gffcompare.shtml>).  
19  
20  
21 1035  
22  
23 1036 The repetitive content of the genome was analysed using RepeatModeller (v1.0.11)  
24  
25  
26 1037 followed by RepeatMasker (v4.0.7). Further analysis of LTRs was performed using  
27  
28 1038 LTRharvest and LTRdigest [93] from the GenomeTools v1.5.9 [94] package.  
29  
30  
31 1039  
32  
33 1040 **Transcriptome analyses**  
34  
35  
36 1041 Quantitative transcriptome analyses were primarily performed using *kallisto* v0.43.1 [95],  
37  
38  
39 1042 which quantifies transcripts by pseudo-aligning raw sequencing reads to a transcriptome  
40  
41 1043 reference. The transcriptome reference was derived from the reference FASTA and  
42  
43  
44 1044 annotation GFF3 data using *gffread* (<https://github.com/gpertea/gffread>). Raw sequencing  
45  
46  
47 1045 reads from each of the life-stages were pseudo-aligned to the reference using *kallisto quant*  
48  
49 1046 (parameters: --bias --bootstrap-samples 100 --fusion), before the differential expression  
50  
51  
52 1047 between pairwise combinations of life-stages was determined using *sleuth* [96]. For each  
53  
54 1048 transcript, differential expression between life-stages was tested using a Wald test, for  
55  
56  
57 1049 which a q-value less than 0.05 was deemed significant.  
58  
59  
60 1050  
61  
62  
63  
64  
65

1051 To explore coordinated gene expression throughout the life-stages, we used *clust* v1.8.4  
1  
2  
3 1052 [97] to automatically define and group transcripts by expression level profiles. Transcripts  
4  
5 1053 per million (TPM) counts per life-stage were generated from the kallisto analysis described  
6  
7 1054 above and used as input to *clust*. To explore putative regulatory sequences that may be  
8  
9 1055 shared by co-expressed genes, we searched in each group of clustered transcripts for  
11  
12 1056 conserved motifs in the 100 bp region immediately upstream of the start codon for each  
13  
14 1057 transcript using *DREME* [98], after which motif-containing sequences were identified using  
15  
16 1058 *FIMO* [99] (MEME v5.0.4).  
17  
18 1059  
19  
20 21 1059  
22  
23 1060 To perform gene set analyses, subsets of *H. contortus* genes together with associated *C.*  
24  
25  
26 1061 *elegans* orthologs were obtained using WormBase ParaSite BioMart [15]. Functional  
27  
28 1062 enrichment analyses for both species were performed using g:Profiler (version  
29  
30  
31 1063 e97\_eg44\_p13\_d22abce), with a g:SCS multiple testing correction method applying a  
32  
33 1064 significance threshold of 0.05 [100].  
34  
35  
36 1065  
37  
38  
39 1066 *OrthoFinder* v2.2.7 [101] was used to determine orthologous relationships between protein  
40  
41 1067 sequences inferred from the *H. contortus* V1 and V4, the McMaster *H. contortus*, as well as  
42  
43  
44 1068 *H. placei* and *C. elegans* protein-coding gene annotations.  
45  
46  
47 1069  
48  
49 1070 **Trans-splicing: Characterisation of splice leader usage and distribution**  
50  
51  
52 1071 We characterised the number and distribution of genes for which mapped RNAseq  
53  
54 1072 transcripts contained the canonical *H. contortus* SL1 (GGTTTAATTACCCAAGTTGAG) and SL2  
55  
56  
57 1073 (GGTTTAACCCAGTATCTCAAG) sequences. For each splice leader sequence, *cutadapt* v1.13  
58  
59 1074 [102] was used to detect (allowing up to 10% divergence in sequence match [--error-rate =  
60  
61  
62  
63  
64  
65

1075 0.1] and minimum match length of 10 bp [--overlap=10]) and subsequently trim the leader  
1  
2 sequences from raw RNAseq reads from each of the life stages, after which all trimmed  
3  
4 reads were collated and mapped to the genome using *hisat2* v2.1.0 [103]. The relative  
5  
6 coverage of the sequence between the trimmed splice leader site and (i) the start codon  
7  
8 (using a window 200 bp upstream and 50 bp downstream of the start codon) and (ii)  
9  
10 (using a window 200 bp upstream and 50 bp downstream of the start codon) and (ii)  
11  
12 internal exons (excluding the 1st exon), was determined using *bedtools coverage*. *Weblogo*  
13  
14 [104] was used to visualise sequence conservation in the genomic sequence surrounding the  
15  
16 splice leader site, and in the splice leader sequences themselves. The spliced leader  
17  
18 sequence splice sites for all life stages were combined into a single dataset for plotting and  
19  
20 quantitative analyses.  
21  
22  
23  
24  
25  
26  
27  
28 **Cis-splicing: Differential isoform usage using Leafcutter**  
29  
30  
31 To analyse the extent of differential spliced isoform usage between life-stages, we used  
32  
33 *leafcutter* [105](see <http://davidaknowles.github.io/leafcutter/> for installation and scripts  
34  
35 described below) to quantify the proportion of reads that span introns and then  
36  
37 subsequently quantify differential intron usage across samples. First, raw RNAseq reads  
38  
39 were mapped to the reference genome using *STAR* [88](parameters: --*twopassMode basic*--  
40  
41 *outSAMstrandField intronMotif --outSAMtype BAM Unsorted*) to generate BAM files per life-  
42  
43 stage. Intron junctions (junc files) were generated from the BAM files, which were  
44  
45 subsequently clustered using *leafcutter\_cluster.py* (parameters: -m 30 -l 500000).  
46  
47  
48 Differential introns between pairwise life-stages using *leafcutter\_ds.R* (parameters: --  
49  
50  
51 min\_samples\_per\_intron=3 --min\_samples\_per\_group=3 --min\_coverage=10), which  
52  
53 required the clustered introns and an exon file (derived from the annotation GFF3) as input.  
54  
55  
56 Preparation of the data for visualisation of differentially spliced introns was performed using  
57  
58  
59  
60  
61  
62  
63  
64  
65

1099    *prepare\_results.R*; an annotation GTF was required as input, which was generated from the  
1  
2  
3 1100 annotation GFF3 using *gffread*. Visualisation was performed using *run\_leafviz.R*.  
4  
5 1101  
6  
7  
8 1102 **Population genetic analyses**  
9  
10 1103 To understand how genetic variation was distributed throughout the chromosomes, we  
11  
12 1104 collated sequencing data from our recent global analysis [25] together with in-house  
13  
14 1105 archival sequencing data to form the most complete sample cohort of *H. contortus* whole-  
15  
16 1106 genome data to date. Reanalysis of the global cohort data was necessary, as the previous  
17  
18 1107 analysis was performed using the V3 assembly. In total, 338 samples from 46 populations in  
19  
20 1108 17 countries were analysed; the samples together with geolocation, mitochondrial  
21  
22 1109 phylogeny describing genetic connectivity between samples, and ENA references to  
23  
24 1110 individual sequencing datasets are available to explore using *Microreact* [106] and can be  
25  
26 1111 visualised here: [https://microreact.org/project/hcontortus\\_global\\_diversity](https://microreact.org/project/hcontortus_global_diversity).  
27  
28  
29 1112  
30  
31 1113 Raw sequencing data were mapped to the V4 reference using *bwa mem* [84](parameters: -Y  
32  
33 1114 and -M options were used to use soft clipping for supplementary alignments, and to mark  
34  
35 1115 shorter hits as secondary, respectively), followed by identification of duplicate reads using  
36  
37 1116 *Picard* (v2.5.0; <https://github.com/broadinstitute/picard>). Samples for which multiple  
38  
39 1117 sequencing lanes were generated were mapped independently and subsequently combined  
40  
41 1118 into a single dataset using *samtools merge*. Mapping statistics were generated using  
42  
43 1119 *samtools* (v1.3) *flagstat* and *bamtools* (v 2.3.0) *stats*, and compared using *MultiQC* v1.3  
44  
45  
46 1120 [107].  
47  
48  
49  
50  
51  
52 1121  
53  
54 1122 Genetic variation within and between samples was determined using *GATK* (v 4.0.3.0).  
55  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65

1123 GVCFs were first created using *HaplotypeCaller*; to improve efficiency, this process was  
1  
2 performed in parallel for each chromosome for each sample separately, and subsequently  
3  
4 combined by chromosome using *CombineGVCFs*. Variants were determined per  
5  
6 chromosome using *GenotypeGVCFs*, after which the final raw variants were combined using  
7  
8  
9  
10 vcftools concat. A hard filter was applied using *GATK VariantFiltration* to mark variants using  
11  
12 the following criteria: QD < 2.0, FS > 60.0, MQ < 40.0, MQRankSum < -12.5, and  
13  
14 ReadPosRankSum < -8.0. Nucleotide diversity and F-stats within and between populations  
15  
16 was determined using *vcftools* v0.1.14[108] --window-pi and --weir-fst-pop, respectively,  
17  
18 using a window size of 100 kbp.  
19  
20  
21  
22  
23 1132  
24  
25  
26 1133 To analyse mitochondrial DNA diversity, bcftools view was used to select single nucleotide  
27  
28 1134 variants with a read depth > 10, homozygous variant, and with minor allele frequency >  
29  
30  
31 1135 0.01. To generate the phylogeny, variants were used to generate a consensus mitochondrial  
32  
33 genome per sample, which were aligned using *mafft* v7.205 [109], and from which  
34  
35 maximum likelihood phylogenies were generated using *iqtree* v1.6.5 [110](parameters: -alrt  
36  
37 1137 1000 -bb 1000 -nt 1) using automated ModelFinding [111]. Principal component analysis  
38  
39 1138 (PCA) of the mitochondrial DNA variation was also performed using the *poppr* [112] package  
40  
41 1139 in R.  
42  
43  
44 1140  
45  
46  
47 1141  
48  
49 1142 Sequencing coverage variation between the autosomes and X chromosome for each BAM  
50  
51  
52 1143 file was performed using *samtools* [113] bedcov, using a BED file of 100 kbp windows  
53  
54 1144 generated using *bedtools* [82] makewindows.  
55  
56  
57 1145  
58  
59 1146 Calculation of  $d_N/d_S$  between *H. contortus* and *H. placei* was performed using *codeml* in  
60  
61  
62  
63  
64  
65

1147 *PAML* v4.7 [114]. This ratio attempts to quantify selection by comparing the frequency of  
1  
2 substitutions at non-silent sites ( $d_N$ ) with the frequency of substitutions as silent sites ( $d_S$ );  
3  
4 ratios greater than one indicate positive selection for new protein-coding changes, ratios  
5 close to zero indicate neutral evolution, and ratios less than one indicate purifying selection.  
6  
7  
8 1150 Note that most genes will have ratios less than one due to the conservation of amino acid  
9  
10 sequences. One-to-one orthologs between *H. contortus* and *H. placei* were identified using  
11  
12 *Orthofinder* as described above, from which amino acid sequence alignments per  
13  
14 orthologous pair were obtained. Codon alignments were generated using *PAL2NAL* v14  
15  
16 1153 *codeml* was run using m1a and m2a models to examine nearly neutral and positive  
17  
18 selection scenarios, from which the log-likelihood (lnL) of each test were compared using a  
19  
20 likelihood ratio test. The D test statistic was calculated, where  $D = 2 * (\ln L_{m2a} - \ln L_{m1a})$ , and  
21  
22 compared against a chi-square distribution using  $df = 2$  and a significance threshold of 0.05.  
23  
24  
25 1156 The resulting gene set was filtered to only include genes with  $d_S < 2$ ,  $d_S > 0.02$ , and  $d_N/d_S > 1$ .  
26  
27  
28 1159  
29  
30  
31 1160  
32  
33  
34 To calculate  $d_N/d_S$  from the global *H. contortus* dataset, genome-wide variation was filtered  
35  
36 1161 using *vcftools* (parameter: --maf 0.01), from which missense and synonymous variants were  
37  
38 determined using *SnpEff* v4.3 (parameters: -no-downstream -no-intergenic -no-intron -no-  
39 1162 upstream -no-utr). For each gene, the frequency of missense and synonymous variants  
40  
41 1163 were each determined from counting variant alleles as a proportion of total genotypes  
42  
43  
44 1164 counted for each variant class. The resulting gene set was filtered to only include genes with  
45  
46  
47 1165  
48  
49 1166  
50  
51  
52 1167  $d_S < 2$ ,  $d_S > 0.02$ , and  $d_N/d_S > 2$ .  
53  
54  
55  
56  
57 1168  
58  
59 1169 **Statistical analysis and visualisation**  
60  
61  
62  
63  
64  
65 All statistical analyses, data exploration, and visualisation were performed using R (v3.5.0);

1171 https://www.R-project.org/ ) unless otherwise stated.

1

2

3 **1172**

4

5

6

7

8

9

10

11

12

13

14

15

16

17

18

19

20

21

22

23

24

25

26

27

28

29

30

31

32

33

34

35

36

37

38

39

40

41

42

43

44

45

46

47

48

49

50

51

52

53

54

55

56

57

58

59

60

61

62

63

64

65

1173 **Abbreviations**

1  
2 1174 BUSCO: Benchmarking Universal Single-Copy Orthologs  
3  
4 1175 CCS: circular consensus sequencing  
5  
6 1176 CEGMA: Core Eukaryotic Genes Mapping Approach  
7  
8 1177 DCC: dosage compensation complex  
9  
10 1178 McMaster: draft genome assembly of the *H. contortus* Australian isolate McMaster,  
11  
12 1179 published in Schwarz *et al.* 2013  
13  
14 1180 DSN: duplex-specific nuclease  
15  
16 1181 PCA: Principal component analysis  
17  
18 1182 SL: spliced leader  
19  
20 1183 V1: version 1 of the *H. contortus* MHco3(ISE).N1 genome, published in Laing *et al.* 2013  
21  
22 1184 V3: version 3 of the *H. contortus* MHco3(ISE).N1 genome, unpublished  
23  
24 1185 V4: version 4 of the *H. contortus* MHco3(ISE).N1 genome, presented here for the first time  
25  
26 1186  
27  
28 1187  
29  
30 1188 **Declarations**

31  
32 1189 **Ethics approval and consent to participate**

33  
34 1190 The use of experimental animals to maintain parasite populations for the purposes  
35  
36 1191 described in this manuscript was approved by the Moredun Research Institute Experiments  
37  
38 1192 and Ethics Committee and were conducted under approved British Home Office licenses in  
39  
40 1193 accordance with the Animals (Scientific Procedures) Act of 1986. The Home Office licence  
41  
42 1194 numbers were PPL 60/03899 and 60/4421, and the experimental identifiers for these  
43  
44 1195 studies were E06/58, E06/75, E09/36 and E14/30.  
45  
46 1196  
47 1197 **Consent for publication**

48  
49 1198 Not applicable.  
50  
51 1199  
52  
53 1200 **Availability of data and materials**

54  
55 1201 The raw sequence data is available under the ENA accession PRJEB506, with reference to  
56  
57 1202 specific sequencing libraries described throughout the text, and/or in Table S3 of Laing et al.  
58  
59 1203 2013. RNAseq data is available from the ENA study ID PRJEB1360. The genome assembly has  
60  
61  
62  
63  
64  
65

1204 been made available at ENA (assembly accession: GCA\_000469685.2) and WormBase  
1  
1205 ParaSite ([https://parasite.wormbase.org/Haemonchus\\_contortus\\_prjeb506/Info/Index](https://parasite.wormbase.org/Haemonchus_contortus_prjeb506/Info/Index)). A  
2  
1206 static version of the genome annotation used in this paper is available at  
3  
1207 [ftp://ftp.sanger.ac.uk/pub/pathogens/sd21/HCON\\_V4\\_GENOME/](ftp://ftp.sanger.ac.uk/pub/pathogens/sd21/HCON_V4_GENOME/) (signoff date: 25th Jan  
4  
1208 2019), however, the most up-to-date version of the annotation can be accessed at  
5  
1209 WormBase ParaSite  
6  
1210 ([https://parasite.wormbase.org/Haemonchus\\_contortus\\_prjeb506/Info/Index/](https://parasite.wormbase.org/Haemonchus_contortus_prjeb506/Info/Index/); [15]).  
7  
1211 Genome variation data can be visualised using MicroReact  
8  
1212 ([https://microreact.org/project/hcontortus\\_global\\_diversity](https://microreact.org/project/hcontortus_global_diversity)), from which links to ENA  
9  
1213 accessions for individual sample sequencing data from the global collection can be obtained.  
10  
1214 Code and workflows used to analyse data are available at  
11  
1215 [https://github.com/stephenrdoyle/hcontortus\\_genome](https://github.com/stephenrdoyle/hcontortus_genome).  
12  
1216  
13  
1217  
14  
1218 **Competing interests**  
15  
1219  
16  
1220 The authors declare that they have no competing interests.  
17  
1221  
18  
1222  
19  
1223  
20  
1224  
21  
1225  
22  
1226  
23  
1227  
24  
1228  
25  
1229  
26  
1230  
27  
1231  
28  
1232  
29  
1233  
30  
1234  
31  
1235  
32  
1236  
33  
1237  
34  
1238  
35  
1239  
36  
1240  
37  
1241  
38  
1242  
39  
1243  
40  
1244  
41  
1245  
42  
1246  
43  
1247  
44  
1248  
45  
1249  
46  
1250  
47  
1251  
48  
1252  
49  
1253  
50  
1254  
51  
1255  
52  
1256  
53  
1257  
54  
1258  
55  
1259  
56  
1260  
57  
1261  
58  
1262  
59  
1263  
60  
1264  
61  
1265  
62  
1266  
63  
1267  
64  
1268  
65

1235 **References**

- 1  
2  
3 1236 1. C. elegans Sequencing Consortium. Genome sequence of the nematode *C. elegans*: a  
4 platform for investigating biology. *Science* [Internet]. 1998;282:2012–8. Available from:  
5  
6  
7 1238 <https://www.ncbi.nlm.nih.gov/pubmed/9851916>  
8  
9  
10 1239 2. Nicholls SM, Quick JC, Tang S, Loman NJ. Ultra-deep, long-read nanopore sequencing of  
11 mock microbial community standards. *Gigascience* [Internet]. 2019;8. Available from:  
12  
13  
14 1241 <http://dx.doi.org/10.1093/gigascience/giz043>  
15  
16  
17 1242 3. Stewart RD, Auffret MD, Warr A, Walker AW, Roehe R, Watson M. The genomic and  
18 proteomic landscape of the rumen microbiome revealed by comprehensive genome-  
19 resolved metagenomics [Internet]. *bioRxiv*. 2018. Available from:  
20  
21  
22 1245 <http://dx.doi.org/10.1101/489443>  
23  
24  
25 1246 4. Yoshimura J, Ichikawa K, Shoura MJ, Artiles KL, Gabdank I, Wahba L, et al. Recompleting  
26 the *Caenorhabditis elegans* genome [Internet]. *Genome Research*. 2019. p. 1009–22.  
27  
28  
29 1248 Available from: <http://dx.doi.org/10.1101/gr.244830.118>  
30  
31  
32 1249 5. Low WY, Tearle R, Bickhart DM, Rosen BD, Kingan SB, Swale T, et al. Chromosome-level  
33 assembly of the water buffalo genome surpasses human and goat genomes in sequence  
34 contiguity. *Nat Commun* [Internet]. 2019;10:260. Available from:  
35  
36 1251 <http://dx.doi.org/10.1038/s41467-018-08260-0>  
37  
38 1252  
39  
40  
41 1253 6. Belser C, Istace B, Denis E, Dubarry M, Baurens F-C, Falentin C, et al. Chromosome-scale  
42 assemblies of plant genomes using nanopore long reads and optical maps. *Nat Plants*  
43  
44  
45 1255 [Internet]. 2018;4:879–87. Available from: <http://dx.doi.org/10.1038/s41477-018-0289-4>  
46  
47  
48 1256 7. Weisenfeld NI, Kumar V, Shah P, Church DM, Jaffe DB. Direct determination of diploid  
49 genome sequences. *Genome Res* [Internet]. 2017;27:757–67. Available from:  
50  
51 1258 <http://dx.doi.org/10.1101/gr.214874.116>  
52  
53  
54  
55 1259 8. Chin C-S, Peluso P, Sedlazeck FJ, Nattestad M, Concepcion GT, Clum A, et al. Phased  
56 diploid genome assembly with single-molecule real-time sequencing. *Nat Methods*  
57  
58 1261 [Internet]. 2016;13:1050–4. Available from: <http://dx.doi.org/10.1038/nmeth.4035>

- 1262 9. Koren S, Rhie A, Walenz BP, Dilthey AT, Bickhart DM, Kingan SB, et al. De novo assembly  
1  
1263 of haplotype-resolved genomes with trio binning. *Nat Biotechnol* [Internet]. 2018;36:1174–  
2  
1264 82. Available from: <http://dx.doi.org/10.1038/nbt.4277>  
3  
1265 10. International Helminth Genomes Consortium. Comparative genomics of the major  
4  
1266 parasitic worms. *Nat Genet* [Internet]. 2018;51:163–74. Available from:  
5  
1267 <http://dx.doi.org/10.1038/s41588-018-0262-1>  
6  
1268 11. Cotton JA, Bennuru S, Grote A, Harsha B, Tracey A, Beech R, et al. The genome of  
7  
1269 *Onchocerca volvulus*, agent of river blindness. *Nat Microbiol* [Internet]. 2016;2:16216.  
8  
1270 Available from: <http://dx.doi.org/10.1038/nmicrobiol.2016.216>  
9  
1271 12. Hunt VL, Tsai IJ, Coghlan A, Reid AJ, Holroyd N, Foth BJ, et al. The genomic basis of  
10  
1272 parasitism in the *Strongyloides* clade of nematodes. *Nat Genet* [Internet]. 2016;48:299–307.  
11  
1273 Available from: <http://dx.doi.org/10.1038/ng.3495>  
12  
1274 13. Li W, Liu B, Yang Y, Ren Y, Wang S, Liu C, et al. The genome of tapeworm *Taenia*  
13  
1275 *multiceps* sheds light on understanding parasitic mechanism and control of coenurosis  
14  
1276 disease. *DNA Res* [Internet]. 2018;25:499–510. Available from:  
15  
1277 <http://dx.doi.org/10.1093/dnares/dsy020>  
16  
1278 14. Tsai IJ, Zarowiecki M, Holroyd N, Garciarrubio A, Sánchez-Flores A, Brooks KL, et al. The  
17  
1279 genomes of four tapeworm species reveal adaptations to parasitism. *Nature* [Internet].  
18  
1280 2013;496:57–63. Available from: <http://dx.doi.org/10.1038/nature12031>  
19  
1281 15. Howe KL, Bolt BJ, Shafie M, Kersey P, Berriman M. WormBase ParaSite – a  
20  
1282 comprehensive resource for helminth genomics. *Mol Biochem Parasitol* [Internet].  
21  
1283 2017;215:2–10. Available from: <http://dx.doi.org/10.1016/j.molbiopara.2016.11.005>  
22  
1284 16. Kaminsky R, Ducray P, Jung M, Clover R, Rufener L, Bouvier J, et al. A new class of  
23  
1285 anthelmintics effective against drug-resistant nematodes. *Nature* [Internet]. 2008;452:176–  
24  
1286 80. Available from: <http://dx.doi.org/10.1038/nature06722>  
25  
1287 17. Preston S, Jabbar A, Nowell C, Joachim A, Ruttkowski B, Baell J, et al. Low cost whole-  
26  
1288 organism screening of compounds for anthelmintic activity. *Int J Parasitol* [Internet].  
27  
1289 2015;45:333–43. Available from: <http://dx.doi.org/10.1016/j.ijpara.2015.01.007>

- 1290 18. Bassetto CC, Amarante AFT. Vaccination of sheep and cattle against haemonchosis. J  
1291 Helminthol [Internet]. 2015;89:517–25. Available from:  
1292 <http://dx.doi.org/10.1017/S0022149X15000279>
- 1293 19. Knox DP, Redmond DL, Newlands GF, Skuce PJ, Pettit D, Smith WD. The nature and  
1294 prospects for gut membrane proteins as vaccine candidates for *Haemonchus contortus* and  
1295 other ruminant trichostrongyloids. Int J Parasitol [Internet]. 2003;33:1129–37. Available  
1296 from: [http://dx.doi.org/10.1016/s0020-7519\(03\)00167-x](http://dx.doi.org/10.1016/s0020-7519(03)00167-x)
- 1297 20. Gilleard JS. *Haemonchus contortus* as a paradigm and model to study anthelmintic drug  
1298 resistance. Parasitology [Internet]. 2013;140:1506–22. Available from:  
1299 <http://dx.doi.org/10.1017/S0031182013001145>
- 1300 21. Laing R, Kikuchi T, Martinelli A, Tsai IJ, Beech RN, Redman E, et al. The genome and  
1301 transcriptome of *Haemonchus contortus*, a key model parasite for drug and vaccine  
1302 discovery. Genome Biol [Internet]. 2013;14:R88. Available from:  
1303 <http://dx.doi.org/10.1186/gb-2013-14-8-r88>
- 1304 22. Schwarz EM, Korhonen PK, Campbell BE, Young ND, Jex AR, Jabbar A, et al. The genome  
1305 and developmental transcriptome of the strongylid nematode *Haemonchus contortus*.  
1306 Genome Biol [Internet]. 2013;14:R89. Available from: <http://dx.doi.org/10.1186/gb-2013-14-8-r89>
- 1308 23. Laing R, Martinelli A, Tracey A, Holroyd N, Gilleard JS, Cotton JA. *Haemonchus contortus*:  
1309 Genome Structure, Organization and Comparative Genomics. Adv Parasitol [Internet].  
1310 2016;93:569–98. Available from: <http://dx.doi.org/10.1016/bs.apar.2016.02.016>
- 1311 24. Wintersinger J, Mariene GM, Wasmuth J. One species, two genomes: A critical  
1312 assessment of inter isolate variation and identification of assembly incongruence in  
1313 *Haemonchus contortus* [Internet]. bioRxiv. 2018. Available from:  
1314 <http://dx.doi.org/10.1101/384008>
- 1315 25. Sallé G, Doyle SR, Cortet J, Cabaret J, Berriman M, Holroyd N, et al. The global diversity  
1316 of *Haemonchus contortus* is shaped by human intervention and climate. Nature  
1317 Communications [Internet]. 2019;10:4811. Available from:

1318 https://www.nature.com/articles/s41467-019-12695-4  
1  
2  
3 1319 26. Troell K, Engström A, Morrison DA, Mattsson JG, Höglund J. Global patterns reveal  
4 strong population structure in *Haemonchus contortus*, a nematode parasite of  
5 domesticated ruminants. *Int J Parasitol* [Internet]. 2006;36:1305–16. Available from:  
6  
7 1321  
8 1322 http://dx.doi.org/10.1016/j.ijpara.2006.06.015  
9  
10  
11 1323 27. Blouin MS, Yowell CA, Courtney CH, Dame JB. Host movement and the genetic structure  
12 of populations of parasitic nematodes. *Genetics* [Internet]. 1995;141:1007–14. Available  
13 from: https://www.ncbi.nlm.nih.gov/pubmed/8582607  
14  
15 1325  
16  
17  
18 1326 28. Doyle SR, Cotton JA. Genome-wide Approaches to Investigate Anthelmintic Resistance.  
19  
20 1327 Trends Parasitol [Internet]. 2019;35:289–301. Available from:  
21  
22 1328 http://dx.doi.org/10.1016/j.pt.2019.01.004  
23  
24  
25 1329 29. Redman E, Grillo V, Saunders G, Packard E, Jackson F, Berriman M, et al. Genetics of  
26  
27 1330 mating and sex determination in the parasitic nematode *Haemonchus contortus*. *Genetics*  
28  
29 1331 [Internet]. 2008;180:1877–87. Available from:  
30  
31 1332 http://dx.doi.org/10.1534/genetics.108.094623  
32  
33  
34 1333 30. Bremner KC. Cytological polymorphism in the nematode *Haemonchus contortus*  
35  
36 1334 (Rudolphi 1803) Cobb 1898. *Nature* [Internet]. 1954;174:704–5. Available from:  
37  
38 1335 http://dx.doi.org/10.1038/174704b0  
39  
40  
41 1336 31. Palevich N, Maclean PH, Baten A, Scott RW, Leathwick DM. The Genome Sequence of  
42  
43 1337 the Anthelmintic-Susceptible New Zealand *Haemonchus contortus*. *Genome Biol Evol*  
44  
45 1338 [Internet]. 2019;11:1965–70. Available from: http://dx.doi.org/10.1093/gbe/evz141  
46  
47  
48 1339 32. Sargison ND, Redman E, Morrison AA, Bartley DJ, Jackson F, Naghra-van Gijzel H, et al. A  
49  
50 1340 method for single pair mating in an obligate parasitic nematode. *Int J Parasitol* [Internet].  
51  
52 1341 2018;48:159–65. Available from: http://dx.doi.org/10.1016/j.ijpara.2017.08.010  
53  
54  
55 1342 33. Doyle SR, Laing R, Bartley DJ, Britton C, Chaudhry U, Gillear JS, et al. A Genome  
56  
57 1343 Resequencing-Based Genetic Map Reveals the Recombination Landscape of an Outbred  
58  
59 1344 Parasitic Nematode in the Presence of Polyploidy and Polyandry. *Genome Biol Evol*  
60  
61 1345 [Internet]. 2018;10:396–409. Available from: http://dx.doi.org/10.1093/gbe/evx269  
62  
63  
64  
65

- 1346 34. Rockman MV, Kruglyak L. Recombinational landscape and population genomics of  
1  
1347 *Caenorhabditis elegans*. PLoS Genet [Internet]. 2009;5:e1000419. Available from:  
2  
1348 <http://dx.doi.org/10.1371/journal.pgen.1000419>  
3  
1349 35. Hoff KJ, Lange S, Lomsadze A, Borodovsky M, Stanke M. BRAKER1: Unsupervised RNA-  
4  
1350 Seq-Based Genome Annotation with GeneMark-ET and AUGUSTUS. Bioinformatics  
5  
1351 [Internet]. 2016;32:767–9. Available from: <http://dx.doi.org/10.1093/bioinformatics/btv661>  
6  
1352 36. Haas BJ, Salzberg SL, Zhu W, Pertea M, Allen JE, Orvis J, et al. Automated eukaryotic  
7  
1353 gene structure annotation using EVidenceModeler and the Program to Assemble Spliced  
8  
1354 Alignments. Genome Biol [Internet]. 2008;9:R7. Available from:  
9  
1355 <http://dx.doi.org/10.1186/gb-2008-9-1-r7>  
10  
1356 37. Dunn NA, Unni DR, Diesh C, Munoz-Torres M, Harris NL, Yao E, et al. Apollo:  
11  
1357 Democratizing genome annotation. PLoS Comput Biol [Internet]. 2019;15:e1006790.  
12  
1358 Available from: <http://dx.doi.org/10.1371/journal.pcbi.1006790>  
13  
1359 38. Albritton SE, Kranz A-L, Rao P, Kramer M, Dieterich C, Ercan S. Sex-biased gene  
14  
1360 expression and evolution of the x chromosome in nematodes. Genetics [Internet].  
15  
1361 2014;197:865–83. Available from: <http://dx.doi.org/10.1534/genetics.114.163311>  
16  
1362 39. Reinke V, Gil IS, Ward S, Kazmer K. Genome-wide germline-enriched and sex-biased  
17  
1363 expression profiles in *Caenorhabditis elegans*. Development [Internet]. 2004;131:311–23.  
18  
1364 Available from: <http://dx.doi.org/10.1242/dev.00914>  
19  
1365 40. Berkseth M, Ikegami K, Arur S, Lieb JD, Zarkower D. TRA-1 ChIP-seq reveals regulators of  
20  
1366 sexual differentiation and multilevel feedback in nematode sex determination. Proc Natl  
21  
1367 Acad Sci U S A [Internet]. 2013;110:16033–8. Available from:  
22  
1368 <http://dx.doi.org/10.1073/pnas.1312087110>  
23  
1369 41. Ellis R. Sex determination in the germ line [Internet]. WormBook. 2006. Available from:  
24  
1370 <http://dx.doi.org/10.1895/wormbook.1.82.1>  
25  
1371 42. Allen MA, Hillier LW, Waterston RH, Blumenthal T. A global analysis of *C. elegans* trans-  
26  
1372 splicing. Genome Res [Internet]. 2011;21:255–64. Available from:  
27  
1373 <http://dx.doi.org/10.1101/gr.113811.110>

- 1374 43. Bektesh S, Van Doren K, Hirsh D. Presence of the *Caenorhabditis elegans* spliced leader  
1  
2 1375 on different mRNAs and in different genera of nematodes [Internet]. *Genes & Development*.  
3  
4 1376 1988. p. 1277–83. Available from: <http://dx.doi.org/10.1101/gad.2.10.1277>  
5  
6  
7 1377 44. Redmond DL, Knox DP. *Haemonchus contortus* SL2 trans-spliced RNA leader sequence.  
8  
9 1378 *Mol Biochem Parasitol* [Internet]. 2001;117:107–10. Available from:  
10  
11 1379 <https://www.ncbi.nlm.nih.gov/pubmed/11551637>  
12  
13  
14 1380 45. Laing R, Hunt M, Protasio AV, Saunders G, Mungall K, Laing S, et al. Annotation of Two  
15 1381 Large Contiguous Regions from the *Haemonchus contortus* Genome Using RNA-seq and  
16  
17 1382 Comparative Analysis with *Caenorhabditis elegans* [Internet]. *PLoS ONE*. 2011. p. e23216.  
18  
19 1383 Available from: <http://dx.doi.org/10.1371/journal.pone.0023216>  
20  
21  
22 1384 46. Rufener L, Mäser P, Roditi I, Kaminsky R. *Haemonchus contortus* acetylcholine receptors  
23  
24 1385 of the DEG-3 subfamily and their role in sensitivity to monepantel. *PLoS Pathog* [Internet].  
25  
26 1386 2009;5:e1000380. Available from: <http://dx.doi.org/10.1371/journal.ppat.1000380>  
27  
28  
29 1387 47. Blumenthal T. Operon and non-operon gene clusters in the *C. elegans* genome  
30  
31 1388 [Internet]. *WormBook*. 2014. p. 1–20. Available from:  
32  
33 1389 <http://dx.doi.org/10.1895/wormbook.1.175.1>  
34  
35  
36 1390 48. Maruyama R, Endo S, Sugimoto A, Yamamoto M. *Caenorhabditis elegans* DAZ-1 is  
37  
38 1391 expressed in proliferating germ cells and directs proper nuclear organization and  
39  
40 1392 cytoplasmic core formation during oogenesis. *Dev Biol* [Internet]. 2005;277:142–54.  
41  
42 1393 Available from: <http://dx.doi.org/10.1016/j.ydbio.2004.08.053>  
43  
44  
45 1394 49. Gilleard JS, Redman E. Genetic Diversity and Population Structure of *Haemonchus*  
46  
47 1395 *contortus* [Internet]. *Haemonchus contortus and Haemonchosis – Past, Present and Future*  
48  
49 1396 Trends. 2016. p. 31–68. Available from: <http://dx.doi.org/10.1016/bs.apar.2016.02.009>  
50  
51 1397 50. Wilson Sayres MA. Genetic Diversity on the Sex Chromosomes. *Genome Biology and*  
52  
53 1398 *Evolution* [Internet]. 2018;10:1064–78. Available from:  
54  
55 1399 <http://dx.doi.org/10.1093/gbe/evy039>  
56  
57  
58 1400 51. Doyle SR, Illingworth CJR, Laing R, Bartley DJ, Redman E, Martinelli A, et al. Population  
59  
60 1401 genomic and evolutionary modelling analyses reveal a single major QTL for ivermectin drug

- 1402 resistance in the pathogenic nematode, *Haemonchus contortus*. BMC Genomics [Internet].  
1  
1403 2019;20:218. Available from: <http://dx.doi.org/10.1186/s12864-019-5592-6>
- 1404 52. Viney M. The failure of genomics in biology [Internet]. Trends in Parasitology. 2014. p.  
1405 319–21. Available from: <http://dx.doi.org/10.1016/j.pt.2014.04.010>
- 1406 53. Wasmuth JD. Realizing the promise of parasite genomics. Trends Parasitol [Internet].  
1407 2014;30:321–3. Available from: <http://dx.doi.org/10.1016/j.pt.2014.04.008>
- 1408 54. Jain M, Koren S, Miga KH, Quick J, Rand AC, Sasani TA, et al. Nanopore sequencing and  
1409 assembly of a human genome with ultra-long reads. Nat Biotechnol [Internet].  
1410 2018;36:338–45. Available from: <http://dx.doi.org/10.1038/nbt.4060>
- 1411 55. Mandrioli M, Manicardi GC. Unlocking holocentric chromosomes: new perspectives from  
1412 comparative and functional genomics? Curr Genomics [Internet]. 2012;13:343–9. Available  
1413 from: <http://dx.doi.org/10.2174/138920212801619250>
- 1414 56. Kingan SB, Heaton H, Cudini J, Lambert CC, Baybayan P, Galvin BD, et al. A High-Quality  
1415 De novo Genome Assembly from a Single Mosquito Using PacBio Sequencing. Genes  
1416 [Internet]. Multidisciplinary Digital Publishing Institute; 2019 [cited 2019 Jan 22];10:62.  
1417 Available from: <https://www.mdpi.com/2073-4425/10/1/62/htm>
- 1418 57. Kashyap SS, Verma S, Voronin D, Lustigman S, Kulke D, Robertson AP, et al. Emodepside  
1419 has sex-dependent immobilizing effects on adult *Brugia malayi* due to a differentially spliced  
1420 binding pocket in the RCK1 region of the SLO-1 K channel. PLoS Pathog [Internet].  
1421 2019;15:e1008041. Available from: <http://dx.doi.org/10.1371/journal.ppat.1008041>
- 1422 58. Sargison ND, Redman E, Morrison AA, Bartley DJ, Jackson F, Hoberg E, et al. Mating  
1423 barriers between genetically divergent strains of the parasitic nematode *Haemonchus*  
1424 *contortus* suggest incipient speciation. Int J Parasitol [Internet]. 2019;49:531–40. Available  
1425 from: <http://dx.doi.org/10.1016/j.ijpara.2019.02.008>
- 1426 59. Wolstenholme AJ. Ion channels and receptor as targets for the control of parasitic  
1427 nematodes [Internet]. International Journal for Parasitology: Drugs and Drug Resistance.  
1428 2011. p. 2–13. Available from: <http://dx.doi.org/10.1016/j.ijpddr.2011.09.003>
- 58

- 1429 60. Basáñez M-G, Pion SDS, Boakes E, Filipe JAN, Churcher TS, Boussinesq M. Effect of  
1 single-dose ivermectin on *Onchocerca volvulus*: a systematic review and meta-analysis  
2 [Internet]. *The Lancet Infectious Diseases*. 2008. p. 310–22. Available from:  
3  
4 1431 [http://dx.doi.org/10.1016/s1473-3099\(08\)70099-9](http://dx.doi.org/10.1016/s1473-3099(08)70099-9)  
5  
6 1432  
7  
8 1433 61. Sabah AA, Fletcher C, Webbe G, Doenhoff MJ. *Schistosoma mansoni*: Chemotherapy of  
9 infections of different ages [Internet]. *Experimental Parasitology*. 1986. p. 294–303.  
10  
11 1434 Available from: [http://dx.doi.org/10.1016/0014-4894\(86\)90184-0](http://dx.doi.org/10.1016/0014-4894(86)90184-0)  
12  
13 1435  
14  
15 1436 62. Neveu C, Charvet CL, Fauvin A, Cortet J, Beech RN, Cabaret J. Genetic diversity of  
16 levamisole receptor subunits in parasitic nematode species and abbreviated transcripts  
17 1437 associated with resistance. *Pharmacogenet Genomics* [Internet]. 2010;20:414–25. Available  
18  
19 20 1438 from: <http://dx.doi.org/10.1097/FPC.0b013e328338ac8c>  
21  
22  
23  
24 1440 63. Fauvin A, Charvet C, Issouf M, Cortet J, Cabaret J, Neveu C. cDNA-AFLP analysis in  
25 levamisole-resistant *Haemonchus contortus* reveals alternative splicing in a nicotinic  
26 1441 acetylcholine receptor subunit [Internet]. *Molecular and Biochemical Parasitology*. 2010. p.  
27  
28 1442 30 1443 105–7. Available from: <http://dx.doi.org/10.1016/j.molbiopara.2009.11.007>  
29  
30  
31  
32  
33 1444 64. Rezansoff AM, Laing R, Martinelli A, Stasiuk S, Redman E, Bartley D, et al. The  
34 confounding effects of high genetic diversity on the determination and interpretation of  
35 1445 differential gene expression analysis in the parasitic nematode *Haemonchus contortus*. *Int J  
36  
37 1446 Parasitol* [Internet]. 2019;49:847–58. Available from:  
38  
39 1447  
40  
41 1448 <http://dx.doi.org/10.1016/j.ijpara.2019.05.012>  
42  
43  
44 1449 65. Rödelsperger C, Meyer JM, Prabh N, Lanz C, Bemm F, Sommer RJ. Single-Molecule  
45 1450 Sequencing Reveals the Chromosome-Scale Genomic Architecture of the Nematode Model  
46  
47 1451 Organism *Pristionchus pacificus*. *Cell Rep* [Internet]. 2017;21:834–44. Available from:  
48  
49 1452 <http://dx.doi.org/10.1016/j.celrep.2017.09.077>  
50  
51  
52 1453 66. Tandonnet S, Koutsovoulos GD, Adams S, Cloarec D, Parihar M, Blaxter ML, et al.  
53  
54 1454 Chromosome-Wide Evolution and Sex Determination in the Three-Sexed Nematode  
55  
56 1455 *Auanema rhodensis*. *G3* [Internet]. 2019;9:1211–30. Available from:  
57  
58 1456 <http://dx.doi.org/10.1534/g3.119.0011>

- 1457 67. Picard MAL, Cosseau C, Ferré S, Quack T, Grevelding CG, Couté Y, et al. Evolution of gene  
1  
1458 dosage on the Z-chromosome of schistosome parasites. *eLife* [Internet]. 2018;7. Available  
2  
1459 from: <http://dx.doi.org/10.7554/eLife.35684>  
3  
1460 68. Niciura SCM, Tizioto PC, Moraes CV, Cruvinel GG, de Albuquerque ACA, Santana RCM, et  
4  
1461 al. Extreme-QTL mapping of monepantel resistance in *Haemonchus contortus*. *Parasit  
5  
1462 Vectors* [Internet]. 2019;12:403. Available from: [http://dx.doi.org/10.1093/bioinformatics/btq683](http://dx.doi.org/10.1186/s13071-019-<br/>6<br/>1463 3663-9</a><br/>7<br/>1464 69. Boetzer M, Henkel CV, Jansen HJ, Butler D, Pirovano W. Scaffolding pre-assembled<br/>8<br/>1465 contigs using SSPACE. <i>Bioinformatics</i> [Internet]. 2011;27:578–9. Available from:<br/>9<br/>1466 <a href=)  
10  
1467 70. Tsai IJ, Otto TD, Berriman M. Improving draft assemblies by iterative mapping and  
11  
1468 assembly of short reads to eliminate gaps. *Genome Biol* [Internet]. 2010;11:R41. Available  
12  
1469 from: <http://dx.doi.org/10.1186/gb-2010-11-4-r41>  
13  
1470 71. Hunt M, Kikuchi T, Sanders M, Newbold C, Berriman M, Otto TD. REAPR: a universal tool  
14  
1471 for genome assembly evaluation. *Genome Biol* [Internet]. 2013;14:R47. Available from:  
15  
1472 <http://dx.doi.org/10.1186/gb-2013-14-5-r47>  
16  
1473 72. Bonfield JK, Whitwham A. Gap5—editing the billion fragment sequence assembly.  
17  
1474 *Bioinformatics* [Internet]. 2010;26:1699–703. Available from:  
18  
1475 <http://dx.doi.org/10.1093/bioinformatics/btq268>  
19  
1476 73. Koren S, Walenz BP, Berlin K, Miller JR, Bergman NH, Phillippy AM. Canu: scalable and  
20  
1477 accurate long-read assembly via adaptive -mer weighting and repeat separation. *Genome  
21  
1478 Res* [Internet]. 2017;27:722–36. Available from: <http://dx.doi.org/10.1101/gr.215087.116>  
22  
1479 74. Kurtz S, Phillippy A, Delcher AL, Smoot M, Shumway M, Antonescu C, et al. Versatile and  
23  
1480 open software for comparing large genomes. *Genome Biol* [Internet]. 2004;5:R12. Available  
24  
1481 from: <http://dx.doi.org/10.1186/gb-2004-5-2-r12>  
25  
1482 75. Yeo S, Coombe L, Warren RL, Chu J, Birol I. ARCS: scaffolding genome drafts with linked  
26  
1483 reads. *Bioinformatics* [Internet]. 2018;34:725–31. Available from:  
27  
1484 <http://dx.doi.org/10.1093/bioinformatics/btx675>

- 1485 76. Warren RL, Yang C, Vandervalk BP, Behsaz B, Lagman A, Jones SJM, et al. LINKS: Scalable,  
1  
1486 alignment-free scaffolding of draft genomes with long reads. *Gigascience* [Internet].  
2  
1487 2015;4:35. Available from: <http://dx.doi.org/10.1186/s13742-015-0076-3>  
3  
1488 77. Walker BJ, Abeel T, Shea T, Priest M, Abouelliel A, Sakthikumar S, et al. Pilon: an  
4  
1489 integrated tool for comprehensive microbial variant detection and genome assembly  
5  
1490 improvement. *PLoS One* [Internet]. 2014;9:e112963. Available from:  
6  
1491 <http://dx.doi.org/10.1371/journal.pone.0112963>  
7  
1492 78. Parra G, Bradnam K, Korf I. CEGMA: a pipeline to accurately annotate core genes in  
8  
1493 eukaryotic genomes. *Bioinformatics* [Internet]. 2007;23:1061–7. Available from:  
9  
1494 <http://dx.doi.org/10.1093/bioinformatics/btm071>  
10  
1495 79. Simão FA, Waterhouse RM, Ioannidis P, Kriventseva EV, Zdobnov EM. BUSCO: assessing  
11  
1496 genome assembly and annotation completeness with single-copy orthologs. *Bioinformatics*  
12  
1497 [Internet]. 2015;31:3210–2. Available from:  
13  
1498 <http://dx.doi.org/10.1093/bioinformatics/btv351>  
14  
1499 80. Wick RR, Schultz MB, Zobel J, Holt KE. Bandage: interactive visualization of de novo  
15  
1500 genome assemblies. *Bioinformatics* [Internet]. 2015;31:3350–2. Available from:  
16  
1501 <http://dx.doi.org/10.1093/bioinformatics/btv383>  
17  
1502 81. Nattestad M, Chin C-S, Schatz MC. Ribbon: Visualizing complex genome alignments and  
18  
1503 structural variation [Internet]. bioRxiv. 2016. Available from:  
19  
1504 <http://dx.doi.org/10.1101/082123>  
20  
1505 82. Quinlan AR, Hall IM. BEDTools: a flexible suite of utilities for comparing genomic  
21  
1506 features. *Bioinformatics* [Internet]. 2010;26:841–2. Available from:  
22  
1507 <http://dx.doi.org/10.1093/bioinformatics/btq033>  
23  
1508 83. Gu Z, Gu L, Eils R, Schlesner M, Brors B. circlize implements and enhances circular  
24  
1509 visualization in R. *Bioinformatics* [Internet]. 2014;30:2811–2. Available from:  
25  
1510 <http://dx.doi.org/10.1093/bioinformatics/btu393>  
26  
1511 84. Li H. Aligning sequence reads, clone sequences and assembly contigs with BWA-MEM  
27  
1512 [Internet]. arXiv. 2013. Available from: <http://arxiv.org/abs/1303.3997>

- 1513 85. Picelli S, Björklund ÅK, Faridani OR, Sagasser S, Winberg G, Sandberg R. Smart-seq2 for  
1  
2 sensitive full-length transcriptome profiling in single cells. *Nat Methods* [Internet].  
3  
4 2013;10:1096–8. Available from: <http://dx.doi.org/10.1038/nmeth.2639>  
5  
6  
7 1516 86. Li H. Minimap2: pairwise alignment for nucleotide sequences. *Bioinformatics* [Internet].  
8  
9 2018;34:3094–100. Available from: <http://dx.doi.org/10.1093/bioinformatics/bty191>  
10  
11  
12 1518 87. Otto TD, Dillon GP, Degrave WS, Berriman M. RATT: Rapid Annotation Transfer Tool.  
13  
14 Nucleic Acids Res [Internet]. 2011;39:e57. Available from:  
15  
16 <http://dx.doi.org/10.1093/nar/gkq1268>  
17  
18  
19 1521 88. Dobin A, Davis CA, Schlesinger F, Drenkow J, Zaleski C, Jha S, et al. STAR: ultrafast  
20 universal RNA-seq aligner. *Bioinformatics* [Internet]. 2013;29:15–21. Available from:  
21  
22 <http://dx.doi.org/10.1093/bioinformatics/bts635>  
23  
24  
25 1524 89. Lomsadze A, Burns PD, Borodovsky M. Integration of mapped RNA-Seq reads into  
26 automatic training of eukaryotic gene finding algorithm. *Nucleic Acids Res* [Internet].  
27 2014;42:e119. Available from: <http://dx.doi.org/10.1093/nar/gku557>  
28  
29  
30  
31  
32 1527 90. Haas BJ, Delcher AL, Mount SM, Wortman JR, Smith RK Jr, Hannick LI, et al. Improving  
33 the *Arabidopsis* genome annotation using maximal transcript alignment assemblies. *Nucleic  
34 Acids Res* [Internet]. 2003;31:5654–66. Available from:  
35  
36  
37 <https://www.ncbi.nlm.nih.gov/pubmed/14500829>  
38  
39  
40  
41 1531 91. Slater GSC, Birney E. Automated generation of heuristics for biological sequence  
42 comparison. *BMC Bioinformatics* [Internet]. 2005;6:31. Available from:  
43  
44  
45 <http://dx.doi.org/10.1186/1471-2105-6-31>  
46  
47  
48 1534 92. Geib SM, Hall B, Derego T, Bremer FT, Cannole K, Sim SB. Genome Annotation  
49 Generator: a simple tool for generating and correcting WGS annotation tables for NCBI  
50 submission. *Gigascience* [Internet]. 2018;7:1–5. Available from:  
51  
52  
53 <http://dx.doi.org/10.1093/gigascience/giy018>  
54  
55  
56 1538 93. Ellinghaus D, Kurtz S, Willhoeft U. LTRharvest, an efficient and flexible software for de  
57 novo detection of LTR retrotransposons. *BMC Bioinformatics* [Internet]. 2008;9:18. Available  
58  
59 from: <http://dx.doi.org/10.1186/1471-2105-9-18>  
60  
61  
62  
63  
64  
65

- 1541 94. Gremme G, Steinbiss S, Kurtz S. GenomeTools: a comprehensive software library for  
1  
1542 efficient processing of structured genome annotations. IEEE/ACM Trans Comput Biol  
2  
1543 Bioinform [Internet]. 2013;10:645–56. Available from:  
3  
1544 <http://dx.doi.org/10.1109/TCBB.2013.68>  
4  
1545 95. Bray NL, Pimentel H, Melsted P, Pachter L. Near-optimal probabilistic RNA-seq  
5  
1546 quantification. Nat Biotechnol [Internet]. 2016;34:525–7. Available from:  
6  
1547 <http://dx.doi.org/10.1038/nbt.3519>  
7  
1548 96. Pimentel H, Bray NL, Puente S, Melsted P, Pachter L. Differential analysis of RNA-seq  
8  
1549 incorporating quantification uncertainty. Nat Methods [Internet]. 2017;14:687–90. Available  
9  
1550 from: <http://dx.doi.org/10.1038/nmeth.4324>  
10  
1551 97. Abu-Jamous B, Kelly S. Clust: automatic extraction of optimal co-expressed gene clusters  
11  
1552 from gene expression data. Genome Biol [Internet]. 2018;19:172. Available from:  
12  
1553 <http://dx.doi.org/10.1186/s13059-018-1536-8>  
13  
1554 98. Bailey TL. DREME: motif discovery in transcription factor ChIP-seq data. Bioinformatics  
14  
1555 [Internet]. 2011;27:1653–9. Available from:  
15  
1556 <http://dx.doi.org/10.1093/bioinformatics/btr261>  
16  
1557 99. Grant CE, Bailey TL, Noble WS. FIMO: scanning for occurrences of a given motif.  
17  
1558 Bioinformatics [Internet]. 2011;27:1017–8. Available from:  
18  
1559 <http://dx.doi.org/10.1093/bioinformatics/btr064>  
19  
1560 100. Raudvere U, Kolberg L, Kuzmin I, Arak T, Adler P, Peterson H, et al. g:Profiler: a web  
20  
1561 server for functional enrichment analysis and conversions of gene lists (2019 update)  
21  
1562 [Internet]. Nucleic Acids Research. 2019. p. W191–8. Available from:  
22  
1563 <http://dx.doi.org/10.1093/nar/gkz369>  
23  
1564 101. Emms DM, Kelly S. OrthoFinder: phylogenetic orthology inference for comparative  
24  
1565 genomics. Genome Biol [Internet]. 2019;20:238. Available from:  
25  
1566 <http://dx.doi.org/10.1186/s13059-019-1832-y>  
26  
1567 102. Martin M. Cutadapt removes adapter sequences from high-throughput sequencing  
27  
1568 reads. EMBnet.journal [Internet]. 2011;17:10. Available from:  
28  
1569 <http://dx.doi.org/10.1484/j.ej.11.10000>  
29  
1570 103. Zerbino DR, Birney E. FASTQC: a quality control tool for high throughput sequence data.  
30  
1571 Bioinformatics [Internet]. 2010;26:1292–3. Available from:  
31  
1572 <http://dx.doi.org/10.1093/bioinformatics/btp373>  
32  
1573 104. Jansson P, Almås I, Lundeberg J, Nyrén P. Quality assessment of DNA samples by  
33  
1574 capillary electrophoresis. Anal Chem [Internet]. 1999;71:1633–7. Available from:  
34  
1575 [http://dx.doi.org/10.1002/\(SICI\)1097-0231\(19990515\)71:11<1633::AID-ANAL1633>3.0.CO;2-1](http://dx.doi.org/10.1002/(SICI)1097-0231(19990515)71:11<1633::AID-ANAL1633>3.0.CO;2-1)

1569 http://dx.doi.org/10.14806/ej.17.1.200  
1  
2  
3 1570 103. Kim D, Langmead B, Salzberg SL. HISAT: a fast spliced aligner with low memory  
4 requirements. *Nat Methods* [Internet]. 2015;12:357–60. Available from:  
5  
6 1571 http://dx.doi.org/10.1038/nmeth.3317  
7  
8  
9 1572 104. Crooks GE, Hon G, Chandonia J-M, Brenner SE. WebLogo: a sequence logo generator.  
10  
11 1573 Genome Res [Internet]. 2004;14:1188–90. Available from:  
12  
13 1574 http://dx.doi.org/10.1101/gr.849004  
14  
15  
16 1575 105. Li YI, Knowles DA, Humphrey J, Barbeira AN, Dickinson SP, Im HK, et al. Annotation-free  
17 quantification of RNA splicing using LeafCutter. *Nat Genet* [Internet]. 2017;50:151–8.  
18  
19 1576 Available from: http://dx.doi.org/10.1038/s41588-017-0004-9  
20  
21  
22  
23 1577 106. Argimón S, Abudahab K, Goater RJE, Fedosejev A, Bhai J, Glasner C, et al. Microreact:  
24 visualizing and sharing data for genomic epidemiology and phylogeography. *Microb Genom*  
25  
26 1578 [Internet]. 2016;2:e000093. Available from: http://dx.doi.org/10.1099/mgen.0.000093  
27  
28  
29  
30 1579 107. Ewels P, Magnusson M, Lundin S, Käller M. MultiQC: summarize analysis results for  
31 multiple tools and samples in a single report. *Bioinformatics* [Internet]. 2016;32:3047–8.  
32  
33 1580 Available from: http://dx.doi.org/10.1093/bioinformatics/btw354  
34  
35  
36  
37 1581 108. Danecek P, Auton A, Abecasis G, Albers CA, Banks E, DePristo MA, et al. The variant call  
38 format and VCFtools. *Bioinformatics* [Internet]. 2011;27:2156–8. Available from:  
39  
40 1582 http://dx.doi.org/10.1093/bioinformatics/btr330  
41  
42  
43  
44 1583 109. Nakamura T, Yamada KD, Tomii K, Katoh K. Parallelization of MAFFT for large-scale  
45 multiple sequence alignments. *Bioinformatics* [Internet]. 2018;34:2490–2. Available from:  
46  
47 1584 http://dx.doi.org/10.1093/bioinformatics/bty121  
48  
49  
50  
51 1585 110. Nguyen L-T, Schmidt HA, von Haeseler A, Minh BQ. IQ-TREE: a fast and effective  
52 stochastic algorithm for estimating maximum-likelihood phylogenies. *Mol Biol Evol*  
53  
54 1586 [Internet]. 2015;32:268–74. Available from: http://dx.doi.org/10.1093/molbev/msu300  
55  
56  
57 1587 111. Kalyaanamoorthy S, Minh BQ, Wong TKF, von Haeseler A, Jermiin LS. ModelFinder: fast  
58 model selection for accurate phylogenetic estimates. *Nat Methods* [Internet]. 2017;14:587–  
59  
60 1588  
61  
62  
63  
64  
65

- 1596 9. Available from: <http://dx.doi.org/10.1038/nmeth.4285>
- 1  
2  
3 1597 112. Kamvar ZN, Tabima JF, Grünwald NJ. Poppr: an R package for genetic analysis of  
4 populations with clonal, partially clonal, and/or sexual reproduction. PeerJ [Internet].  
5 1598 2014;2:e281. Available from: <http://dx.doi.org/10.7717/peerj.281>  
6  
7  
8  
9 1600 113. Li H, Handsaker B, Wysoker A, Fennell T, Ruan J, Homer N, et al. The Sequence  
10 Alignment/Map format and SAMtools. Bioinformatics [Internet]. 2009;25:2078–9. Available  
11 1601 from: <http://dx.doi.org/10.1093/bioinformatics/btp352>  
12  
13 1602 114. Yang Z. PAML 4: phylogenetic analysis by maximum likelihood. Mol Biol Evol [Internet].  
14 1603 2007;24:1586–91. Available from: <http://dx.doi.org/10.1093/molbev/msm088>  
15  
16 1604 115. Suyama M, Torrents D, Bork P. PAL2NAL: robust conversion of protein sequence  
17 alignments into the corresponding codon alignments. Nucleic Acids Res [Internet].  
18 1605 2006;34:W609–12. Available from: <http://dx.doi.org/10.1093/nar/gkl315>  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65

Figure 1

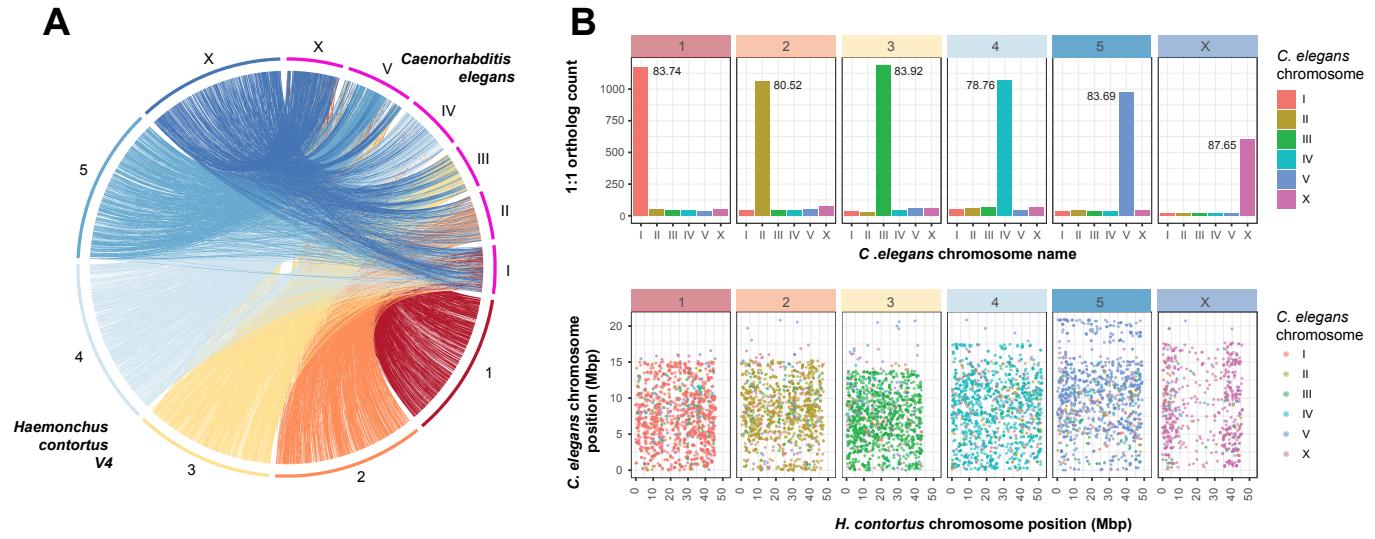
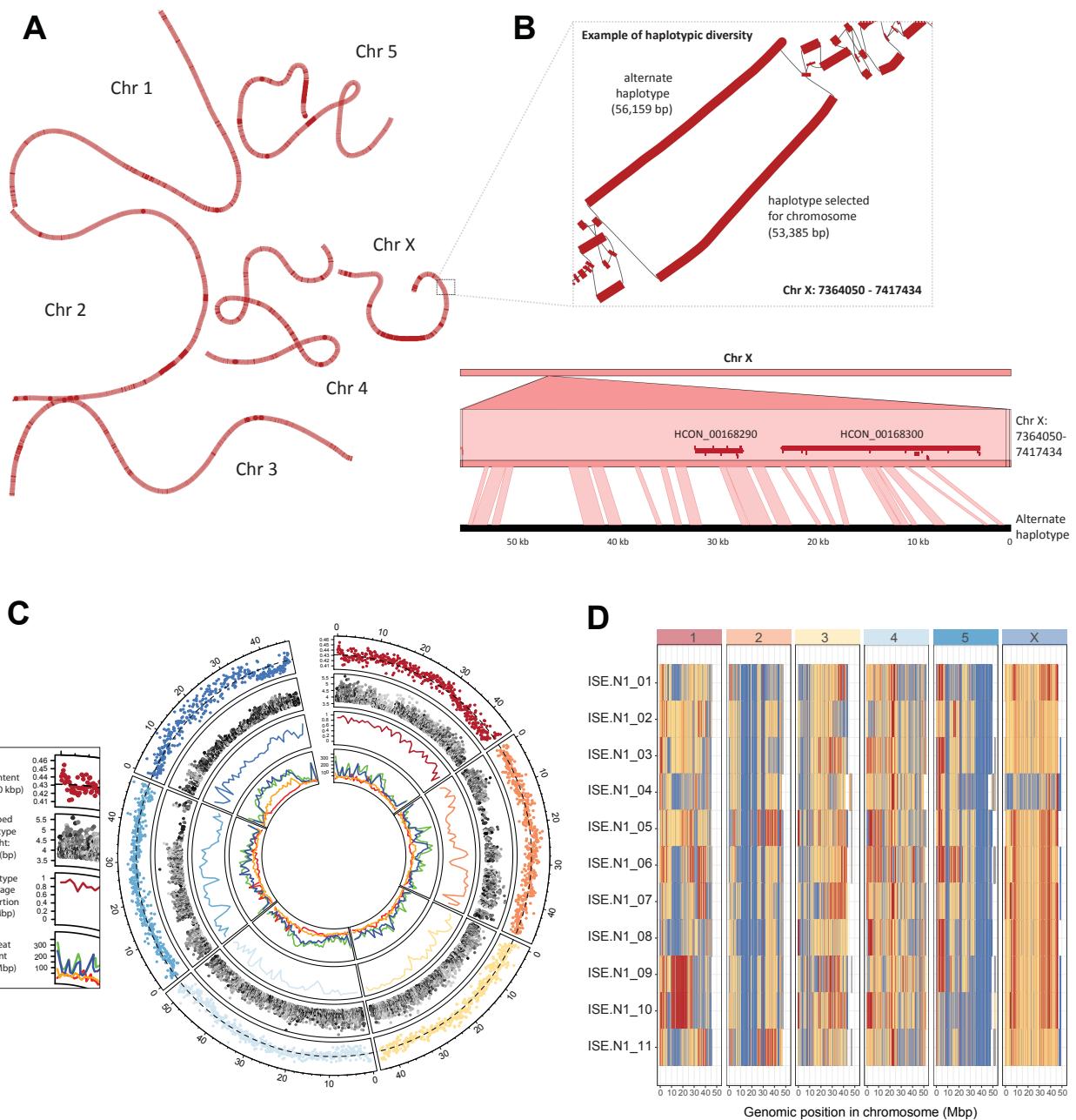
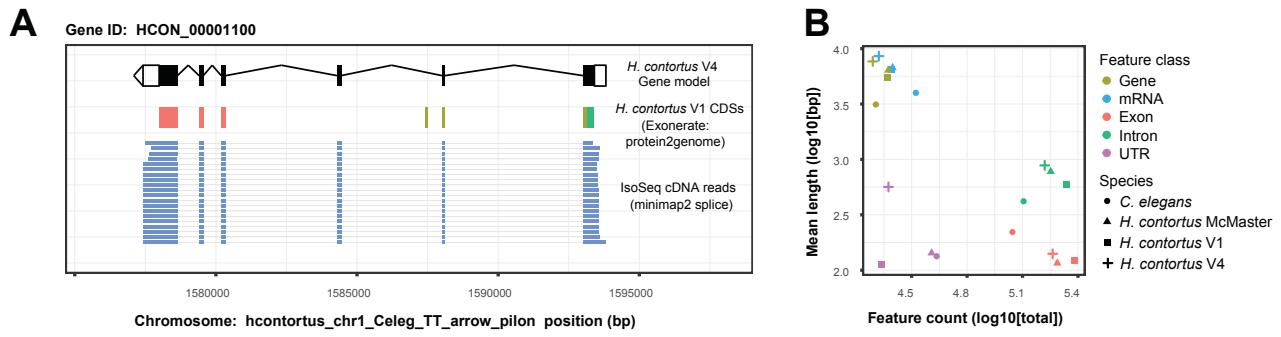
Click here to access/download;Figure;Figure\_1.ai 

Figure 2

Click here to access/download;Figure;Figure\_2.ai





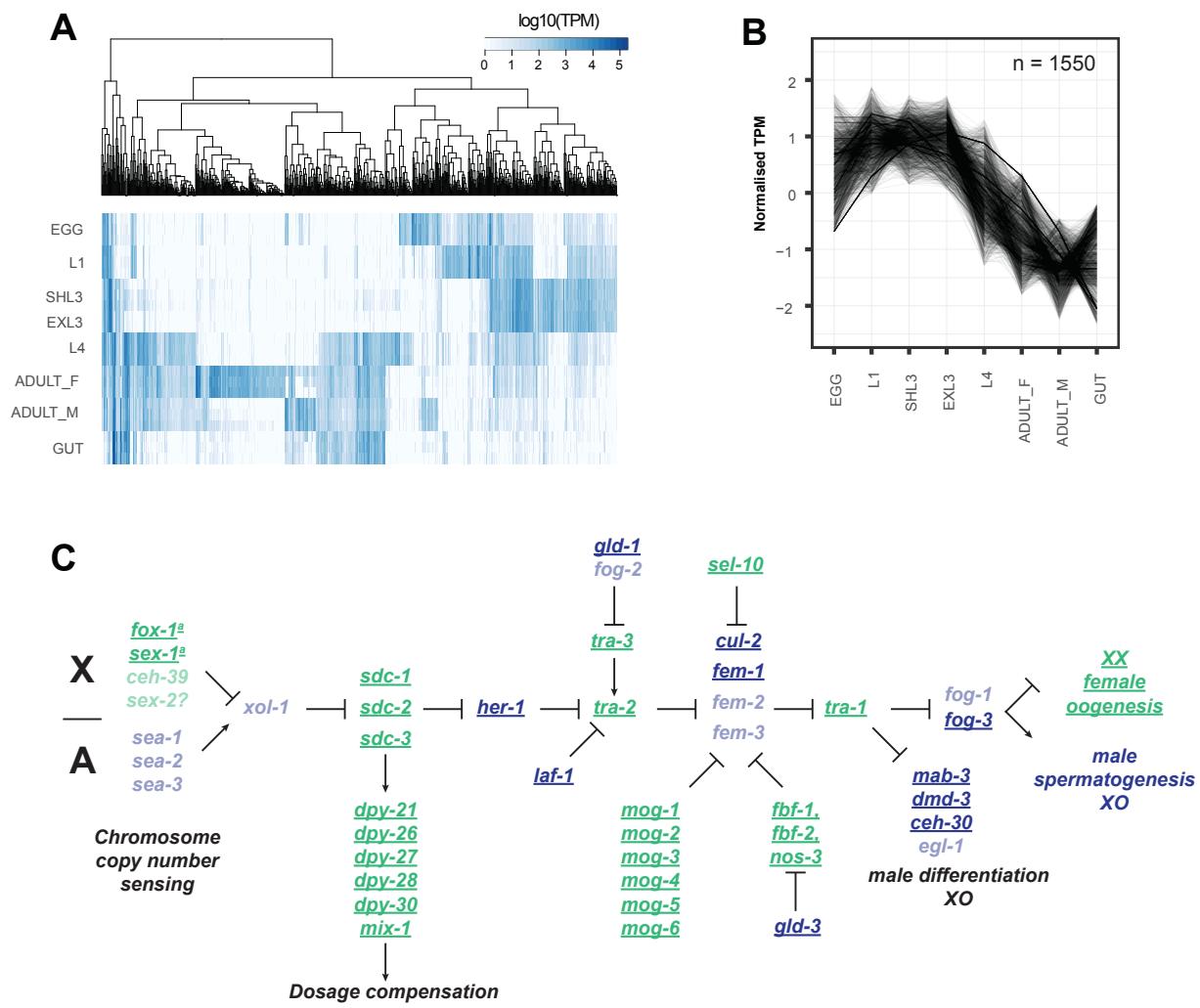


Figure 5

Click here to access/download;Figure;Figure\_5.ai

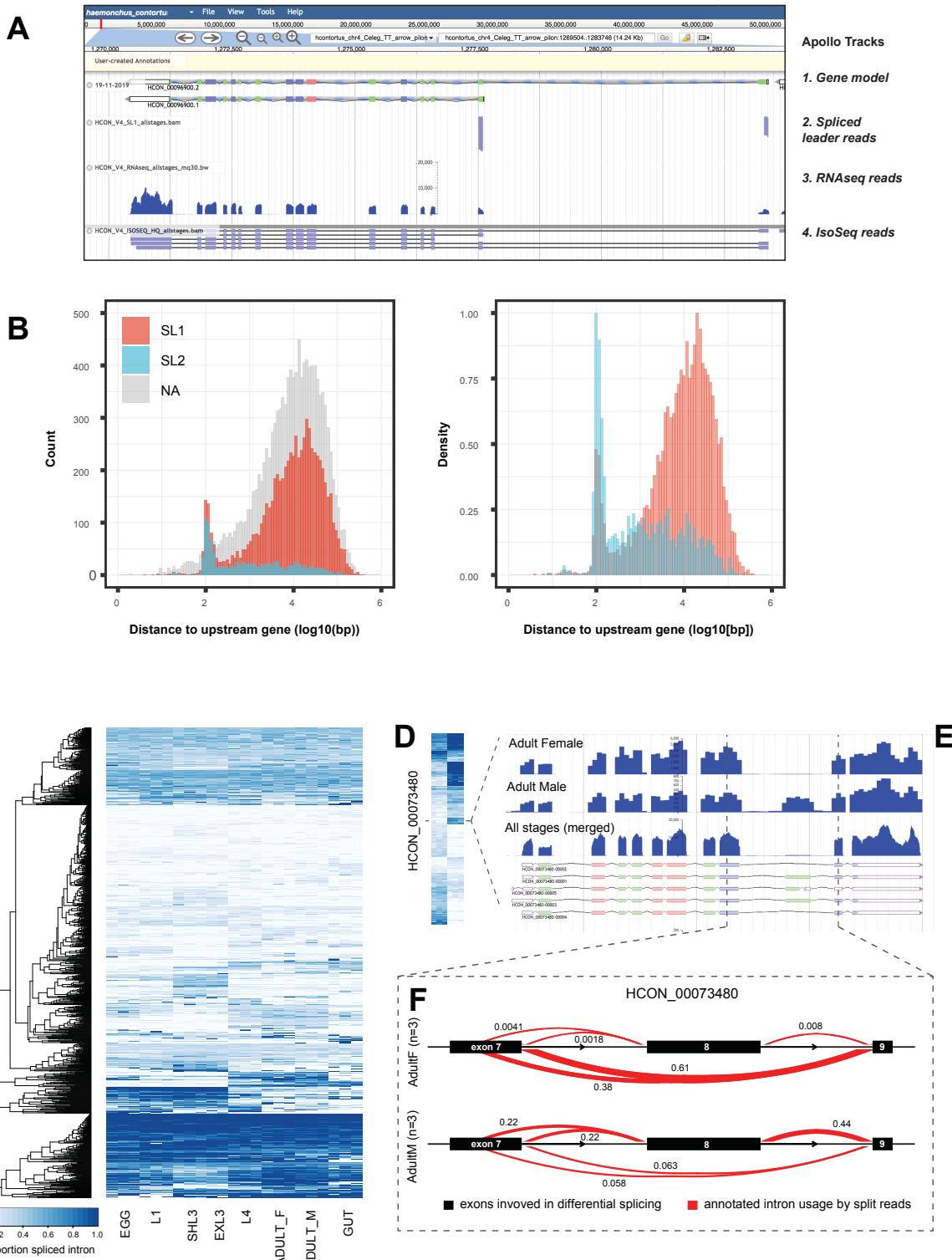
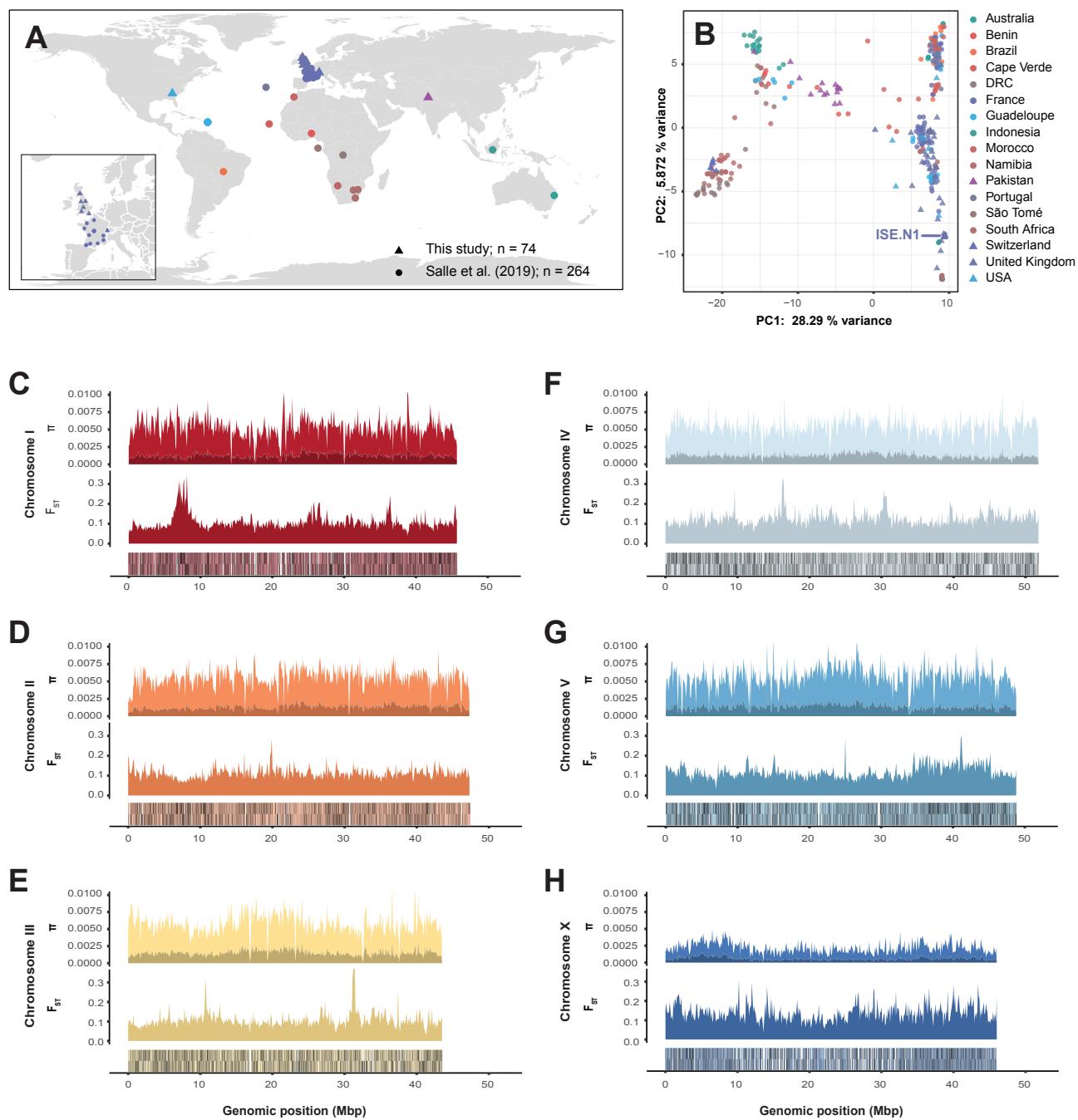


Figure 6

Click here to access/download;Figure;Figure\_6.ai





Click here to access/download  
**Supplementary Material**

Additional File 1 - Supplementary Figures.pdf



Click here to access/download  
**Supplementary Material**

Additional File 2 - Supplementary Tables.xlsx



18th February 2020

Dear Editor

We would like to submit the manuscript "**Extensive genomic and transcriptomic variation defines the chromosome-scale assembly of *Haemonchus contortus*, a model gastrointestinal worm**" for publication as a research article in *Genome Biology*. In our study, we report a highly resolved chromosomal genome assembly for *Haemonchus contortus*, a major gastrointestinal pathogen responsible for significant animal health and economic burdens worldwide. We describe chromosome-scale structural variation, patterns of transcriptional change and putative co-regulation throughout its life cycle using a comprehensive genome annotation derived from short- and long-read cDNA sequencing, and within- and between-population genetic diversity from a globally sampled cohort.

Our work represents a step-change in the genetic resources for *H. contortus*. A 370 Mbp highly fragmented and haplotypic draft assembly [ $n = 23,860$ ;  $N50 = 0.083$  Mbp] for this parasite strain was published in *Genome Biology* in 2013 (Laing et al.; gb-2013-14-8-r88). The utility and impact of this resource is highlighted by the fact that the publication describing the draft genome has been accessed over 20,000 times and cited at least 160 times in the last 7 years. Since completion in early 2019, we have made our chromosomal assembly [283 Mbp;  $n = 7$ ;  $N50 = 47.4$  Mbp] and all associated resources publicly available before publication via the main community portal for nematode research, WormBase Parasite, and since that time, has been used and referred to in at least 14 publications.

We believe that this work will be of great interest to communities that work on parasitic nematodes. *H. contortus* is the key model representing a group of parasites that includes most of the nematode species of relevance to veterinary health and the hookworms, which are important and neglected human pathogens that infect more than a billion people worldwide. However, we feel our work is particularly suited to *Genome Biology*, rather than a more specialised parasitology journal, as this assembly represents the most complete and accurate reference genome for any nematode species outside *Caenorhabditis*. Our analysis of resolved chromosomes reveals a remarkable pattern of almost complete conservation of chromosome content between the evolutionarily related *H. contortus* and *C. elegans*, but almost no conservation of gene order, and yet, important contrasts with this model, for example, in the patterns of *cis*- and *trans*-splicing, and chromosome dosage compensation components. Further, this work will be broadly applicable to readers of *Genome Biology* with an interest in chromosome-scale biology and potential insights gained from genetic and transcriptomic analyses derived from long molecule sequencing, especially for species with high genetic diversity.

We look forward to hearing from you.

Yours faithfully,

Two handwritten signatures are shown side-by-side. The signature on the left is "Stephen Doyle" and the signature on the right is "James Cotton". They are written in black ink on a light-colored background.

Stephen Doyle & James Cotton  
Corresponding authors, on behalf of all authors

**Wellcome Sanger Institute**  
Wellcome Genome Campus  
Hinxton, Cambridge  
CB10 1SA

**Genome Research Limited**  
Registered Office:  
215 Euston Road  
London NW1 2BE

A company registered  
in England (No. 2742969)  
and a charity registered  
in England (No. 1021457)

Wellcome Sanger Institute