1    **Comprehensive analysis of human hookworm secreted proteins using a proteogenomic**

2    **approach**

3    Logan J[a], Manda SS[b,c], Choi YJ[d], Field M[a], Eichenberger RM[a], Mulvenna J[e], Nagaraj SH[f],

4    Fujiwara RT[g], Gazzinelli-Guimaraes P[g], Bueno L[g], Mati V[h], Mitreva M[d], Sotillo J[a#], Loukas A[a#]

5

6    [a]Centre for Biodiscovery and Molecular Development of Therapeutics, Australian Institute of Tropical

7    Health and Medicine, James Cook University, Cairns, QLD, Australia

8    [b]Cancer Data Science Group, ProCan, Children's Medical Research Institute, Faculty of Medicine and

9    Health, University of Sydney, Westmead, NSW, Australia

10   [c]LifeBytes India Pvt Ltd, Whitefield, Bangalore, India

11   [d]McDonnell Genome Institute, Washington University School of Medicine, St. Louis, Missouri, USA

12   [e]QIMR-Berghofer Medical Research Institute, Brisbane, QLD, Australia

13   [f]Institute of Health and Biomedical Innovation and Translational Research Institute, Queensland

14   University of Technology, Brisbane, QLD, Australia

15   [g]Department of Parasitology, Biological Sciences Institute, Federal University of Minas Gerais, Belo

16   Horizonte, Brazil

17   [h]Department of Health Sciences, Universidade Federal de Lavras, Lavras, Brazil.

18   [#] Both authors contributed equally to this work

19   Corresponding Authors:

20         Prof. Alex Loukas. Centre for Biodiscovery and Molecular Development of Therapeutics,

21   James Cook University, Cairns, 4878, Queensland, Australia. Email: alex.loukas@jcu.edu.au

22         Dr. Javier Sotillo. Centre for Biodiscovery and Molecular Development of Therapeutics,

23   James Cook University, Cairns, 4878, Queensland, Australia. Email: javier.sotillo@jcu.edu.au

24

25    **Running title:** Proteomic analysis of *Necator americanus* secreted proteins

26    **Abbreviations:**

27    ASP – Aspartyl

28    ASPR – Ancylostoma secreted protein related

29    CAP – Cysteine-rich secretory protein

30    CDS – Coding sequence

31    cRAP – Common repository of adventitious proteins

32    CRISP - Cysteine-rich secretory protein

33    CYS – Cysteine

34    DP – Declustering potential

35    DUF – Domain of unknown function

36    emPAI – Exponential modified protein abundance index

37    ES – Excretory/Secretory

38    FDR – False discovery rate

39    GO – Gene ontology

40    GSSPs – Genome search specific peptides

41    IAM – Iodoacetamide

42    iTOF – Interactive tree of life

43    MET – Metallo

44    ML – Maximum-likelihood

45    NK – Natural killer

46    NP – Number of peptides

47    RO – Reverse osmosis

48    SCP/TAPS – Sperm-coating protein

49    SER – Serine

50    SFG – Six-frame translated genome

51    SP – Signal peptide

52    SRA – Sequence read archive

53    TD – Transmembrane domain

54

Proteomic analysis of *Necator americanus* secreted proteins

55 **Summary**

56  The human hookworm *Necator americanus* infects more than 400 million people worldwide,

57  contributing substantially to the poverty in these regions. Adult stage *N. americanus* live in the

58  small intestine of the human host where they inject excretory/secretory (ES) products into the

59  mucosa. ES products have been characterized at the proteome level for a number of animal

60  hookworm species, but until now, the difficulty in obtaining sufficient live *N. americanus* has been

61  an obstacle in characterizing the secretome of this important human pathogen. Herein we describe

62  the ES proteome of *N. americanus* and utilize this information to conduct the first proteogenomic

63  analysis of a parasitic helminth, significantly improving the available genome and thereby

64  generating a robust description of the parasite secretome. The genome annotation resulted in a a

65  revised prediction of 3,425 fewer genes than initially reported, accompanied by a significant

66  increase in the number of exons and introns, total gene length and the percentage of the genome

67  covered by genes. Almost 200 ES proteins were identified by LC-MS/MS with SCP/TAPS

68  proteins, 'hypothetical' proteins and proteases among the most abundant families. These proteins

69  were compared to commonly used model species of human parasitic infections, including

70  *Ancylostoma caninum*, *Nippostrongylus brasiliensis* and *Heligmosomoides polygyrus*. Our

71  findings provide valuable information on important families of proteins with both known and

72  unknown functions that could be instrumental in host-parasite interactions, including protein

73  families that might be key for parasite survival in the onslaught of robust immune responses, as

74  well as vaccine and drug targets.

## Introduction

76    Hookworm infection is one of the most pertinent and life-limiting parasitic infections worldwide,

77    affecting more than 400 million people in tropical regions of Asia, Africa and South America (1,

78    2). Chronic infection with hookworms results in fatigue, abdominal pain, diarrhea, weight loss and

79    anemia (3). In children, it can cause growth retardation and impairments in cognitive development

80    (4), and in pregnant women these infections lead to poor birth outcomes including low birth

81    weight, increased perinatal morbidity and mortality (5, 6). Moreover, hookworm infections result

82    in 3.2 million disability-adjusted life years lost annually (7). Hookworm infection therefore

83    contributes significantly to widespread poverty in the majority of endemic regions (8).

84    The life cycle of the most widespread of the anthropophilic hookworms, *Necator americanus*, is

85    direct, with no intermediate hosts involved. Eggs are passed out of the body in the feces and, under

86    favorable conditions, hatch releasing first-stage (L1). Larvae undergo several molts to reach the

87    infective third (L3) stage which can penetrate human skin (9). Upon infection the parasite migrates

88    through the circulatory system to the lungs, where it moves up the trachea and is eventually

89    swallowed, thus commencing its final sojourn through the gastrointestinal tract to reside in the

90    small intestine where mature adult worms can live for up to 10 years (9). The adult stage of *N.*

91    *americanus* produces macromolecules known as excretory/secretory (ES) products, which consist

92    of a battery of proteins that have evolved to interact with human host tissues and facilitate

93    parasitism (10). These ES products have the potential to not only be targeted as potential vaccine

94    and diagnostic candidates, but also to shed light on how these parasites evade immune destruction

95    (11-14).

96    Despite their potential biotechnological utility, only a limited number of *N. americanus* ES

97    proteins have been described to date, and most of them have been identified as cDNAs based on

98      their homology to proteins from more readily accessible hookworm species from animals such as

99      *Ancylostoma caninum* (15). In terms of vaccine antigens, a handful of *N. americanus* ES products

100     including glutathione-S-transferases, aspartic proteases and sperm-coating proteins/Tpx-

101     1/Ag5/PR-1/Sc7 (SCP/TAPS), have been identified at the cDNA level, and vaccine efficacy of

102     recombinant proteins assessed in animal models and phase 1 clinical trials (13, 16, 17). SCP/TAPS,

103     also referred to as venom allergen-like (VAL) or Activation-associated Secreted Proteins (ASPs)

104     (Pfam accession number no. PF00188) have been reported from many helminths, but appear to be

105     significantly expanded in the genomes and secreted proteomes of gut-resident clade IV and V

106     parasitic nematodes, including the hookworms (11, 18-20).

107     Considerably little is known about the roles of *N. americanus* proteins in immunoregulation

108     compared with many of the parasites that serve as animal models for human infections. Hsieh and

109     colleagues reported an *N. americanus* protein(s) which selectively bound Natural Killer (NK) cells

110     resulting in IL-2- and IL-12-dependent IFN-γ production, although the identity of the protein was

111     not determined (21). Calreticulin, a protein from the ES products of *N. americanus*, was identified

112     as having an immunomodulatory role through inhibition of the hemolytic capacity of C1q, a human

113     complement protein (22). The relative paucity of functional information on *N. americanus* proteins

114     can be attributed, at least in part, to the difficulty in obtaining parasite material and the absence

115     until 2014 of the published *N. americanus* draft genome. Analysis of the 244 Mb draft genome

116     and the predicted 19,151 genes (11), provided important information about the molecular

117     mechanisms and pathways by which *N. americanus* interacts with its human host. In agreement

118     with published transcriptomes of *N. americanus* and genomes/transcriptomes of other hookworm

119     species, a select number of protein families were over-represented, including SCP/TAPS proteins

120     and different mechanistic classes of proteases with various functions including hemoglobinolysis,

121     and tissue penetration (23-28). Of the >19,000 predicted genes reported in the draft *N. americanus*

122     genome, 8,176 genes had no known InterPro domain. Additionally, more than half of the total

123     proteins had either no blast homology to any gene from the NCBI database (10,771 proteins) or

124     shared identity with a 'hypothetical' protein (3,043 proteins). These results highlight the

125     importance of further annotation and refinement of the *N. americanus* genome (29).

126     In general, the practicality of genomic sequence data is dependent on the accuracy of gene

127     annotation as well as the availability of functional, expression and localization information (30).

128     The high-throughput methods used when annotating a genome are prone to errors, therefore, to

129     validate the predicted protein-coding genes, an analysis of the proteome is essential. Mass

130     spectrometry provides useful data that can be used in a proteogenomic approach to improve

131     genome annotation and identify novel peptides containing predicted protein sequences (31).

132     Herein we perform the first proteogenomic analysis of a parasitic helminth, while also significantly

133     improving the genome annotation and comprehensively characterizing the ES proteome of adult

134     *N. americanus*. Our findings provide valuable information on important families of proteins with

135     both known and unknown functions that could be instrumental in host-parasite interactions,

136     including protein families that might be key for parasite survival and protection of the host against

137     excessive immunopathology. Characterization of these proteins will be useful for the identification

138     of vaccine and drug targets and diagnostics for hookworm infection.

139

140     **Experimental procedures**

141     *Ethics*

142     Ethical approval for hamster animal experimentation to obtain *N. americanus* adult worms was

143     obtained from the Federal University of Minas Gerais, Brazil (Protocol# 51/2013). Ethical

144 approval for human experimental infection with *N. americanus* and subsequent culturing of L3

145 was obtained from the James Cook University Human Research Ethics Committee (ID# H5936).

146

147 *Parasite material*

148 Adult *N. americanus* were manually isolated from the intestines of experimentally infected golden

149 hamsters (*Mesocricetus auratus*) upon euthanasia. For the isolation and purification of ES

150 products, the worms were washed 3 times in phosphate-buffered saline (PBS) before being

151 cultured overnight in a humidified incubator at 37°C, 5% $CO_2$, in RPMI 1640 (100U/ml penicillin,

152 100 μg/ml streptomycin sulphate, 0.25 μg/ml amphotericin B). The supernatant was collected the

153 following day and debris was removed by centrifuging the concentrated samples at 1,500 g for 3

154 minutes in a benchtop microfuge. ES products were concentrated with a 3 kDa cut-off Centricon

155 filter membrane (Merck Millipore) and samples were stored at -80°C until use. For somatic adult

156 extracts, data was used from Tang *et al* (11). In brief, whole worms were ground under liquid

157 nitrogen and solubilized using lysis buffer (1.0% (v/v) Triton X-100 in 40 mM Tris, 0.1% (w/v)

158 SDS, pH 7.4). The extract was filtered through 20μm filter before fractionation.

159 For isolation of *N. americanus* L3, stool samples were collected from infected human volunteers

160 and cultured as follows. Reverse osmosis (RO) water was added to the stool sample until a thick

161 paste was formed. This paste was then distributed onto moistened Filter Paper (VWR, Standard

162 Grade, 110 mm) in Petri dishes (Sarstedt, 150 mm) and placed in a 25°C incubator for 8 days.

163 Following incubation, the edges of each plate were gently rinsed with RO water to obtain clean L3

164 preparations. Somatic extracts were prepared by adding 100 μl lysis buffer (3 M Urea, 0.2% SDS,

165 1% Triton X, 50 mM Tris-HCl) to approximately 6,000 larvae before repeated vortexing and

166 sonication (4°C, probe sonicator, pulse setting) to digest the larvae. The extract was passed through

167 a 0.45 µm filter (Millipore) before in-gel fractionation.

168

169 *SDS-PAGE fractionation and in-gel trypsin digestion*

170 A total of 30 µg of *N. americanus* adult ES products was buffer exchanged into 50 mM $NH_4HCO_3$,

171 freeze-dried and resuspended in Laemmli buffer. The sample was boiled at 95°C for 5 minutes and

172 electrophoresed on a 12% SDS-PAGE gel for 40 minutes at 150V. The gel was stained with

173 Coomassie Blue and 20 pieces (approximately 1 mm thick) were cut and placed into Eppendorf

174 tubes. For the in-gel digestion, slices were de-stained and freeze-dried before incubating them at

175 65°C for 1 h in reduction buffer (20 mM dithiothreitol (DTT), Sigma, 50 mM $NH_4HCO_3$). Samples

176 were then alkylated in 50 mM iodoacetamide (IAM, Sigma), 50 mM $NH_4HCO_3$ for 40 minutes at

177 37°C, washed three times with 25 mM $NH_4HCO_3$ and digested with 20 µg/ml of trypsin (Sigma)

178 by incubating them for 16 h at 37°C with gentle agitation. Digestion was stopped and peptides

179 were released from the slices by adding 0.1% TFA, 70% acetonitrile. This step was repeated 3

180 times with pooling of the corresponding supernatants to maximise peptide recovery for each

181 sample. Finally, each sample was desalted with a ZipTip (Merck Millipore) and stored at -80°C

182 until use.

183

184 *In-solution trypsin digestion and off-gel fractionation*

185 A total of 70 µg of each *N. americanus* extract (L3 somatic and adult ES products) was buffer

186 exchanged into 50 mM $NH_4HCO_3$ before adding DTT to 20 mM and incubating for 10 minutes at

187 65°C. Alkylation was carried out by adding IAM to 55 mM and incubating for 45 minutes at room

188 temperature in the dark. Samples were digested with 2 µg of trypsin by incubating for 16 h at 37°C

8

189    with gentle agitation. Following trypsin digestion, peptides were fractionated using a 3100

190    OFFGEL Fractionator (Agilent Technologies) according to the manufacturer's protocol with a 24-

191    well format as described previously (19). In brief, Immobiline DryStrip pH 3-10 (24 cm) gel strips

192    were rehydrated in rehydration buffer in the assembled loading frame. Digested peptides were

193    diluted in dilution buffer to a volume of 3.6 ml and loaded equally across the 24 well cassette. The

194    sample was run at a current of 50 μA until 50 kilovolt hours (kVh) had elapsed. Upon completion

195    of the fractionation, samples were collected, desalted with ZipTip and stored at -80°C until use.

196

197    *Mass Spectrometry*

198    The extracts were analyzed by LC-MS/MS on a Shimadzu Prominence Nano HPLC (Japan)

199    coupled to a Triple Tof 5600+ mass spectrometer (SCIEX, Canada) equipped with a nano

200    electrospray ion source. Fifteen (15) μl of each extract was injected onto a 50 mm x 300 μm C18

201    trap column (Agilent Technologies, Australia) at 60 μl/min. The samples were de-salted on the

202    trap column for 6 minutes using 0.1% formic acid (aq) at 60 μl/min. The trap column was then

203    placed in-line with the analytical nano HPLC column, a 150 mm x 100 μm 300SBC18, 3.5 μm

204    (Agilent Technologies, Australia) for mass spectrometry analysis. For peptide elution and analysis

205    the nano-HPLC pump was initially held at 2% solvent B for 6 minutes followed by a linear gradient

206    of 2-40% solvent B over 80 minutes at 500 nl/minute flow rate and then a steeper gradient from

207    40% to 80% solvent B in 10 minutes was applied. Solvent B was held at 80% for 5 minutes for

208    washing the column and returned to 2% solvent B for equilibration prior to the next sample

209    injection. Solvent A consisted of 0.1% formic acid (aq) and solvent B contained 90/10 acetonitrile/

210    0.1% formic acid (aq). The ionspray voltage was set to 2200V, declustering potential (DP) 100V,

211    curtain gas flow 25, nebulizer gas 1 (GS1) 12 and interface heater at 150°C. The mass spectrometer

212    acquired 250 ms full scan TOF-MS data followed by 20 by 250 ms full scan product ion data in

213    an Information Dependent Acquisition (IDA) mode. Full scan TOF-MS data was acquired over

214    the mass range 300-1600 and for product ion ms/ms 80-1600. Ions observed in the TOF-MS scan

215    exceeding a threshold of 150 counts and a charge state of +2 to +5 were set to trigger the acquisition

216    of product ion, MS/MS spectra of the resultant 20 most intense ions. The data was acquired using

217    Analyst TF 1.6.1 (SCIEX, Canada).

218

219    *Proteogenomics*

220    The mass spectrometry raw data was searched against the *N. americanus* protein database using

221    SequestHT algorithm in Proteome Discoverer 2.1 (Thermo Scientific, Bremen, Germany). Trypsin

222    was used as the protease, allowing a maximum of two missed cleavages. Carbamidomethylation

223    of cysteine was specified as a fixed modification, and oxidation of methionine was included as

224    variable modifications. The minimum peptide length was specified as 6 amino acids. The mass

225    error of parent ions was set to 10 ppm, and for fragment ions it was set to 0.05 Da. Protein inference

226    was based on the rule of parsimony and required one or more unique peptides. False Discovery

227    Rate (FDR) of 1% at peptide and protein levels was applied. The unassigned spectra from the

228    protein database search were further searched against the six-frame translated genome database of

229    *N. americanus* (11). The genome sequences were downloaded from WormBase ParaSite release 9

230    (http://parasite.wormbase.org/Necator_americanus_prjna72135/Info/Index) in FASTA format and

231    translated in all six frames using in-house python scripts. All the sequences greater than 10 amino

232    acids between any two stop codons were added to the database. The search was again performed

233    using SequestHT with precursor mass tolerance of 50 ppm and fragment ion tolerance of 1 Da.

234    Carbamidomethylation of cysteine was specified as a fixed modification, and oxidation of

235    methionine was included as variable modifications. Results were obtained at 1% FDR for both

236    protein and peptide level.

237    The identified peptides in the search against the six-frame translated genome were mapped back

238    to the protein database using standalone BLAST. Any sequences that mapped 100% to the protein

239    database were discarded. The filtered peptides were mapped to the *N. americanus* genome using

240    the standalone tblastn program (11). Peptide identifications that unambiguously mapped to a single

241    region in the genome, also known as Genome Search Specific Peptides (GSSPs), were considered

242    to perform proteogenomics-based annotation of novel coding regions. The known gene annotation

243    GFF3 was downloaded from WormBase ParaSite for *N. americanus* (11). Using in-house scripts,

244    we categorized the GSSPs into various categories: intergenic, intronic, exon-extension, alternative

245    frame, N-terminal extensions or repeat regions with respect to the known regions from GFF3.

246    Additionally, for all the intergenic peptides, we determined if there was a potential open reading

247    frame (ORF) within the stretch of amino acids. Each of the spectra was further manually validated.

248

249    *Genome annotation*

250    The genome annotation from *N. americanus* (11) was updated using the MAKER pipeline v2.31.8

251    (32). The genome assembly (GenBank assembly accession: GCA_000507365.1) was softmasked

252    for repetitive elements with RepeatMasker v4.0.6 using a species-specific repeat library created

253    by RepeatModeler v1.0.8, RepBase repeat libraries (33) and a list of known transposable elements

254    provided by MAKER (32). From the NCBI Sequence Read Archive (SRA), *N. americanus* RNA-

255    Seq data (11) (Adult: SRR609895, SRR831085, SRR892200; L3: SRR609894, SRR831091,

256    SRR89220) were obtained. After adapter and quality trimming using Trimmomatic v0.36 (34),

257    RNA-Seq reads were aligned to the genome using HISAT2 v2.0.5 (35) with the --dta option and

11

258   subsequently assembled using StringTie v1.2.4 (36). The resulting alignment information and

259   transcript assembly were used by BRAKER (37) and MAKER pipelines, respectively, as extrinsic

260   evidence data. In addition, protein sequences from SwissProt UniRef100 (38) and WormBase

261   ParaSite WS258 (39) (*Ancylostoma ceylanicum* PRJNA231479, *Brugia malayi* PRJNA10729,

262   *Caenorhabditis elegans* PRJNA13758, *Onchocerca volvulus* PRJEB513, *Pristionchus pacificus*

263   PRJNA12644, *Trichinella spiralis* PRJNA12603 and *Strongyloides ratti* PRJEB125) were

264   provided to MAKER as protein homology evidence. Following the developer's recommendation

265   (40), the protein-coding gene models of Tang *et al.* (11) were passed to MAKER as pred_gff to

266   update the models by adding new 3' and 5' exons, additional UTRs, and merging split models.

267   This method, however, cannot change internal exons nor create new annotations where evidence

268   suggests a gene but no corresponding model is previously present. To address this shortcoming,

269   additional *ab initio* gene predictions were generated using BRAKER v2.0.1 (37) and passed to

270   MAKER so that the intron-exon model that best matched the evidence could be included in the

271   final annotation set. Within the BRAKER pipeline, the gene prediction tools GeneMark (41) and

272   AUGUSTUS (42) were trained utilizing the *N. americanus* RNA-Seq alignment and protein

273   homology information from *C. elegans*. The GFFPs identified in the present study were used to

274   confirm the validity of proposed annotation changes and resolve competing gene predictions

275   through manual curation. Gene models with no evidence support were not included in the final

276   annotation build to reduce false positives in the existing annotations. However, *ab initio* gene

277   predictions that encoded Pfam domains as detected by InterProScan v5.19 (43) were rescued to

278   enhance overall accuracy by balancing sensitivity and specificity (32, 44). The completeness of

279   the annotated gene set was assessed using BUSCO v3.0 with Eukaryota-specific single copy

280   orthologs (OrthoDB v9) (45).

281

282 *Protein Identification*

283 Mascot version 2.5 (Matrix Science), X!Tandem (The Global Proteome Machine Organisation)

284 version Jackhammer and Comet v2014.02 rev.2 were used to analyze data from the mass

285 spectrometer. Searches were carried out against a database comprised of either the updated genome

286 annotation provided with this study, or the *N. americanus* genome (11), both appended to the

287 common repository of adventitious proteins (cRAP; http://www.thegpm.org/crap/) database (to

288 detect potential contamination). The following parameters were used: enzyme, trypsin; variable

289 modifications, oxidation of methionine, carbamidomethylation of cysteine, deamidation of

290 asparagine and glutamine; maximum missed cleavages, 2; precursor ion mass tolerance, 50 ppm;

291 fragment ion tolerance 0.1 Da; charge states, 2+, 3+ 4+. A FDR of 0.1% was applied, and a filter

292 of greater than 2 significant unique sequences was used to further improve the robustness of data.

293 The mass spectrometry data have been deposited in the ProteomeXchange Consortium via the

294 PRIDE partner repository with the dataset identifier PXD010669.

295

296 *Bioinformatic Analysis of Proteomic Sequence Data*

297 Gene ontology (GO) annotations were assigned using the program Blast2GO and Pfam analysis

298 was performed using HMMER (46). Pfam domains were detected at the P<0.001 threshold for the

299 HMMER software. Putative signal peptides were predicted with SignalP v4.1 and transmembrane

300 domains with TMHMM v2.0 (47, 48). REViGO, an online tool, was used to summarise and plot

301 GO terms (49). UpSetR was used to group proteins based on whether they had a GO term with

302 one, two or all of the GO categories (biological process, molecular function or cellular process)

303 (50).

304

305     *Similarity analysis*

306     A similarity analysis was carried out based on the Parkinson and Blaxter method using an in-house

307     script (26, 51). Given the difficulty of working with *N. americanus* specifically (in terms of

308     accessibility to samples and establishment of life cycle in other hosts), other hookworms are

309     frequently used to model this human parasite. Data from the secreted proteome of adult *N.*

310     *americanus* (described herein) was compared to the published secreted proteomes from related

311     adult nematode species including *H. polygyrus*, *A. caninum* and *Nippostrongylus brasiliensis* (15,

312     19, 52).

313

314     *Protein family similarity visualization*

315     *H. polygyrus*, *A. caninum* and *N. brasiliensis* SCP/TAPS and protease protein sequences were

316     obtained from their respective published secreted proteomes (15, 19, 52). These sequences were

317     aligned with adult *N. americanus* homologous protein sequences using BLAST. Significant

318     sequence alignments were visualized using Circos (53).

319

320     *Phylogenetic Analyses*

321     SCP/TAPS proteins were identified from published secretomes of 6 species of parasitic helminths

322     including *A. caninum, Ascaris suum, H. polygyrus, N. brasiliensis, Trichuris muris* and *Toxocara*

323     *canis.* SCP/TAPS proteins were identified from *N. americanus* ES products using the proteome

324     data generated in this study. Sequences were sorted into single and double SCP/TAPS domain

325     proteins for individual phylogenetic analysis due to distinct differences described previously (54).

326     A list of the proteins and their respective sequences used for this analysis can be found in

327    Supplemental Data 1. A multiple sequence alignment was carried out using the alignment program

328    MUSCLE. Outliers with poor alignment (long unaligned regions) were detected and filtered out

329    using ODseek. PhyML, a phylogeny software, was used for a maximum-likelihood (ML)

330    phylogenetic analyses of SCP/TAPS amino acid sequences. The tree was visualized with The

331    Interactive Tree of Life (iTOF) online phylogeny tool (https://itol.embl.de/) (55).

332

333    **Results**

334    *Proteogenomic analysis and genome annotation*

335    Different proteomes from two life stages (infective L3 somatic extract and adult somatic extracts

336    and ES products) of *N. americanus* were analysed by mass spectrometry in a 5600 ABSciex Triple

337    Tof to perform a proteogenomic analysis of the hookworm genome. After excluding peptides that

338    mapped accurately to the existing protein database entry, a total of 218 novel peptides were

339    identified that did not match any protein sequences of *N. americanus* contained in the annotations

340    of Tang *et al* (11). Of these 218 novel peptides, 83 were found exclusively in adult somatic extracts,

341    50 exclusively in L3, and 67 exclusively in adult ES, while 10 peptides were found in both the

342    adult somatic extract and ES products, and 8 were found in both adult and larval somatic extracts.

343    No common peptides were identified from adult ES products and larval somatic extract.

344    Newly identified peptides could be grouped into the following six categories: (1) peptides mapping

345    to intergenic regions; (2) peptides mapping to introns; (3) peptides mapping to alternative reading

346    frames; (4) peptides extending gene boundaries (exon extensions); (5) peptides mapping to N-

347    terminal extensions; and (6) peptides mapping to repeat regions (not gene) regions (Figure 1). Of

348    the identified peptides, more than half belonged to group 1, 18% to group 2, 15% to group 3 and

349    a small number to groups 4, 5 and 6 (Figure 1). Group 1 peptides were further analyzed for an

350   ORF including whether or not they contained a methionine; see Supplemental data 2. From this

351   data, one-third of group 1 peptides were found to overlap with a viable ORF, with the shortest

352   having just 21 amino acids and the longest 381 amino acids. The translated protein sequences with

353   their respective peptides, six-frame translated genome (SFG) ID and ORF length are provided in

354   Supplemental data 3.

355   These results highlighted the need for improving the previously published gene models, so we

356   updated the genome annotation of *N. americanus* using a new RNA-Seq based gene calling

357   pipeline as outlined in the Materials and Methods (Table 1). The total number of predicted *N.*

358   *americanus* genes decreased by 3,425, with substantial increases in both the number of exons and

359   introns. The total coding sequence (CDS) length increased by 2.22 Mb with the mean gene length

360   (including introns and UTRs) nearly doubling from 4.3 kb to 8.1 kb. Furthermore, the percentage

361   of the genome covered by genes increased by 18.3%, and the percentage of detected BUSCOs

362   among the predicted genes increased from 95.7% to 97.4% (with 31% reduction in the number of

363   fragmented BUSCOs). While these improvements are attributable in most part to the use of more

364   sophisticated genome annotation methods utilizing RNA-Seq data and the inclusion of more

365   extensive, up-to-date homologous protein databases, the peptide sequences generated in this study

366   contributed directly to the refinement of 14 gene models. The newly annotated, improved gene

367   models were used in subsequent proteomic analysis of ES products, and are publicly available on

368   Nematode.net (56, 57).

369

370   *Analysis of the ES products from N. americanus adult worms*

371   A comprehensive analysis of the ES products from adult worms was carried out using in-gel and

372   off-gel fractionation and the tryptic peptides were analyzed using LC-MS/MS. Mascot, X!Tandem

373 and Comet searches were carried out against a database including the predicted proteins from the

374 annotated *N. americanus* genome available in this study and the cRAP sequences available at

375 http://www.thegpm.org/crap/. A total of 186 and 141 proteins were identified using Mascot and

376 X!Tandem/Comet, respectively. All ES proteins were identified with at least two unique peptides,

377 at a 99.0% probability and FRD of 0.1%. The two search methods combined found a total of 198

378 proteins, Supplemental data 4). These 198 proteins were obtained using the updated genome

379 annotation generated as part of this study. In comparison, using the first version of the annotated

380 genome sequence we identified 203 proteins using the same analytical software (11). Using the

381 exponential modified protein abundance index (emPAI) and the newly annotated genome, the most

382 abundant proteins in the ES products of *N. americanus* adult worms were ranked, and the top 30

383 are shown in Table 2. One emPAI value was listed for each of the digestion methods used (Ingel

384 and Offgel) and their combined value was used to rank the proteins.

385 The conserved Pfam domains of the 198 ES proteins identified were analyzed. The most abundant

386 protein family in the ES products was the SCP/TAPS family with 54/198 proteins containing a

387 single or double cysteine-rich secretory protein family (CAP) domain (PF00188) (Figure 2B). The

388 top 10 most abundant protein families are displayed in Figure 2B. Many of these proteins were

389 described by Blast2GO as 'SCP-partial' or 'SCP-like', but for a more standardized annotation in

390 our analysis we have grouped them all as 'SCP/TAPS'. The second most frequently represented

391 group (with 42 proteins) were proteins with one or more domains of unknown function (DUF).

392 Other abundant families included Ancylostoma secreted protein related (ASPR) proteins and

393 metalloproteases with 35 and 9 of each identified respectively (Supplemental data 5). Despite some

394 published reports classifying ASPRs as SCP/TAPS proteins, they are a diverse set of secreted

395 cysteine rich proteins based on Pfam annotation and therefore we have grouped them separately

17

396  (28). Of the 198 identified ES proteins, 51% contained a predicted signal peptide. Supporting the

397  accuracy of the new gene model, 96 of the identified adult ES proteins were predicted to contain

398  a signal peptide compared to just 75 using the previous genome annotation.

399  The adult *N. americanus* ES proteins were annotated using Blast2GO (46). In total, 30 GO terms

400  were identified (following the removal of parent child redundancy) belonging to one of the three

401  GO database categories: biological processes, molecular function or cellular component (Figure

402  2A, for raw data see Supplemental data 6). Blast2GO returned biological process GO terms for

403  87/198 proteins, molecular function GO terms for 99/198 proteins and cellular component GO

404  terms for 83/198 proteins (Figure 2A). The most prominent biological process term was proteolysis

405  (Figure 2D), with 15% (29 proteins) of total ES products being involved in a proteolytic process.

406  The most prominent single molecular function term was "hydrolase activity" followed by

407  "peptidase activity" and "metal ion binding" (Figure 2C). Interestingly, 50/198 proteins did not

408  return any GO terms (Figure 2A). Of these, SCP/TAPS proteins made up 64% with no known

409  biological process, molecular function, or cellular component. This highlights a significant

410  knowledge gap surrounding SCP/TAPS produced by helminths.

411

412

413  *Similarity analysis of the ES products from different gastrointestinal nematode species*

414  A similarity analysis of ES proteomic data from *N. americanus* and three of the most commonly

415  used animal models for human hookworms, *A. caninum*, *H. polygyrus* and *N. brasiliensis* was

416  carried out (Figure 3). A total of 15, 10 and 1 *N. americanus* ES proteins had unique homology to

417  ES proteins from *A. caninum*, *H. polygyrus* and *N. brasiliensis* respectively. This included one

418  SCP/TAPS protein (NAME_13724) which was similar only to an *H. polygyrus* protein, while 3 of

419    the proteins were similar only to *A. caninum* ES proteins with domains of unknown function

420    (NAME_07734, NAME_09181, NAME_09182). Seven proteins shared homology with only *A.*

421    *caninum* and *H. polygyrus* proteins, 11 shared homology with only *A. caninum* and *N. brasiliensis*

422    proteins and 7 shared homology with only *H. polygyrus* and *N. brasiliensis* proteins. One hundred

423    and twenty-two *N. americanus* proteins shared different degrees of homology with proteins from

424    *A. caninum*, *H. polygyrus* and *N. brasiliensis*.

425    From the 54 SCP/TAPS proteins found in *N. americanus* ES products, one (NAME_13724) shared

426    homology with a protein found in only *H. polygyrus*, while another SCP/TAPS protein

427    (NAME_15177) shared homology with proteins from both *A. caninum* and *N. brasiliensis*. Fifty-

428    one (51) of 54 SCP/TAPS proteins were similar to all compared species, leaving a single

429    SCP/TAPS protein (NAME_11218) which did not have homology to any SCP/TAPS protein from

430    the compared species. Twenty-five *N. americanus* ES proteins did not have homology to any

431    proteins in the secretome of *A. caninum*, *H. polygyrus* or *N. brasiliensis*. Notable proteins among

432    these 25 were NAME_01848 (aspartyl protease) and NAME_05081 (zinc metalloprotease).

433    Of the 26 proteases in the *N. americanus* ES products, 22 had homologs in all compared species,

434    while one serine protease (NAME_06735) only had homologs in *H. polygyrus* and *A. caninum* and

435    one metalloprotease (NAME_00535, peptidase family M1) only had homologs in the ES products

436    of *A. caninum*, and *N. brasiliensis*. One zinc metalloprotease (NAME_05081) and one serine

437    protease (NAME_01250) were only found in the ES products of *N. americanus* and therefore did

438    not have any similarity to ES products from the other species.

439

440    *Homology analysis of SCP/TAPS and proteases in the ES products of N. americanus*

441     Since SCP/TAPS proteins and proteases from *N. americanus* numerically dominate the ES protein

442     dataset and likely play key roles in infection, migration and parasite establishment, we performed

443     an in-depth analysis of these families of proteins between the human and three model

444     gastrointestinal nematode species. Adult *N. americanus* ES SCP/TAPS protein sequences were

445     aligned with homologs from *H. polygyrus*, *A. caninum* and *N. brasiliensis* ES SCP/TAPS protein

446     sequences using BLAST. Sequences which aligned with maximum scores >36 were visualized

447     using Circos (Figure 4A). SCP/TAPS from *N. americanus* are more similar to *A. caninum* than *H.*

448     *polygyrus* or *N. brasiliensis*, as denoted by thicker, darker ribbons (Figure 4A). The sequences,

449     their homologs and the corresponding blast scores are detailed in Supplemental data 7. As with the

450     similarity analysis, NAME_11218 had no significant alignment to any of the compared species.

451     A similar analysis was performed for the proteases from each of the aforementioned species

452     (Figure 4B). These proteases have been grouped together into their respective mechanistic classes:

453     aspartyl (ASP), cysteine (CYS), metallo (MET) or serine (SER) proteases. In general, all three

454     comparator species had high protease sequence homology to metalloproteases from *N.*

455     *americanus*. Aspartyl protease sequences were more similar between *H. polygyrus*, *N. brasiliensis*

456     and *N. americanus* than sequences from *A. caninum*. It is interesting to note that *N. americanus*

457     contained more aspartyl proteases than the other nematode species analyzed (*N. americanus* – 8;

458     *A. caninum* – 4; *H. polygyrus* – 6; *N. brasiliensis* – 6). Serine proteases had the lowest sequence

459     homology across all of the compared species and protease subclasses.

460

461     *Phylogenetic analysis of SCP/TAPS proteins*

462     SCP/TAPS sequences of 6 parasitic nematodes were obtained from published secretomes and

463     compared to *N. americanus* SCP/TAPS identified in this study. Sequences were grouped by

464  whether they had a single or double domain and then aligned using MUSCLE and PhyML. From

465  the 7 total species, 232 SCP/TAPS were reported with 134 single-domain proteins and 98 double-

466  domain. For the single SCP/TAPS-domain proteins an unrooted tree was generated. The analysis

467  identified seven main clades. Three of these clades consisted entirely of sequences from *N.*

468  *americanus* and *A. caninum*, while sub-clades in 3 of the other main clades followed a similar

469  trend. A majority of *T. muris* single-domain sequences formed a sub-clade with *A. suum* and *T.*

470  *canis*, grouping together the non-clade V helminths. An unrooted tree was also generated for

471  double domain sequences. The analysis identified five main clades. *N. americanus* SCP/TAPS

472  clustered almost exclusively with sequences from *A. caninum* again, representing one main clade

473  and several sub-clades. *H. polygyrus* and *N. brasiliensis* also formed a number of distinct sub-

474  clades. *T. canis*, the only non-clade V helminth, was reported to only produce one double-domain

475  SCP/TAPS protein in its ES products which was not closely related to any of the compared

476  sequences. Interestingly, no double-domain SCP/TAPS were reported for *A. suum* or *T. muris*.

477  Another trend across both single and double domain trees was the diversity of evolution between

478  SCP/TAPS within a single species. For example, the single-domain tree included a sub-clade of 7

479  *N. americanus*-only sequences while other *N. americanus* sequences had more similarity to mouse

480  hookworm sequences.

481

**Discussion**

483  *N. americanus* affects more than 400 million people worldwide and is the most important soil

484  transmitted helminth in terms of morbidity (2). The genome of *N. americanus* was sequenced in

485  2014, providing an important dataset to facilitate efforts to combat hookworm disease; however,

486  the tools available at that time for annotating genes from parasitic helminths were limited. For

487    instance, the number of proteins with a top hit to a 'hypothetical protein' present in the original

488    genome annotation was 3,043, corresponding to 15.8% of the total predicted proteins (11). In

489    addition, inferring gene and protein functions for parasitic nematodes is a major challenge as most

490    species are genetically intractable and databases and algorithms are biased towards (free-living)

491    model nematodes (29). Proteogenomics is a relatively new approach in which proteomic data is

492    used to improve genome annotation (31, 58), and although it had never been applied to parasitic

493    helminths (until now), its potential utility in this area has been suggested (59). High-throughput

494    sequencing and gene prediction tools are prone to false-negative and false-positive predictions

495    which can lead to missed genes, false exons or exon boundaries and/or incorrect translational

496    start/stop sites, so knowing the sequences of the proteins expressed by an organism will help to

497    improve gene predictions.

498    The proteogenomic analysis carried out in this study addresses some of these challenges by

499    improving the characterization of predicted proteins from the annotated *N. americanus* genome.

500    Overall, we identified 121 peptides that map to intergenic regions in the first draft genome

501    sequence for *N. americanus*. Peptides that map to intergenic regions are highly significant as they

502    can lead to identification of novel protein-coding genes or corrections of pre-existing models (31).

503    To investigate whether these peptides are likely to be new genes we checked for any potential

504    ORFs where the peptides map. Thirty-two (32) of the peptides identified mapped to alternative

505    ORFs than those described in the current gene model. While these peptides map to known coding

506    regions, they highlight out-of-frame ORFs which is likely to, once correctly annotated, result in an

507    entirely different protein. Of the total newly identified peptides, 39 of them mapped to introns.

508    Peptides mapping to introns can lead to identification of novel splice isoforms or amendments in

509    gene structure. Peptides mapping to exon extensions and N-terminal extensions were less

510    abundant, with 10 and 4 peptides respectively. These groups of peptides suggest a possible

511    correction in reading frame or an incorrect start site annotation.

512    The decrease in gene number seen in this study is in line with other genome annotations. For

513    example, the *Schistosoma mansoni* genome was originally thought to encode 11,809 genes,

514    however further annotation has reduced this number to 10,772 (60). Given that this is the first re-

515    annotation of the *N. americanus* genome since the original draft was published, a substantial

516    decrease in gene number was to be expected, and 15,728 genes is comparable with the predicted

517    gene numbers of other nematode genomes (61). Despite being a parasitic nematode, *N. americanus*

518    has    almost    5,000    fewer    genes    than    its    free-living    relative    *C.    elegans*

519    (https://parasite.wormbase.org/).

520    Of particular importance in the characterization of parasite-specific genes is the presence of a

521    signal peptide. Of the secreted ES proteins from *N. americanus*, 51% were predicted to have a

522    signal peptide; this is also in agreement with ES proteomes from other parasitic nematodes (15,

523    19, 52). The presence of extracellular proteins without predicted signal peptides in the ES products

524    of *N. americanus* could be due to one of three reasons: (a) the protein is secreted via an alternative

525    pathway, including release of parasite exosomes (62-65) or non-classical secretory signals; (b) the

526    lack of full-length RNA transcript sequence to confirm gene model accuracy resulted in an error

527    in the predicted sequence (i.e. truncation or ORF shift); (c) the pre-set D-cutoff threshold of 0.33

528    in SignalP results in false-negative predictions.

529    Helminth secretomes represent the molecular host-parasite interface (10), and provide useful

530    insights into the biological strategies employed by these parasites to ensure longevity inside their

531    respective hosts (10). At this interface, ES products have been implicated in numerous roles from

532    initial penetration/invasion and feeding to host immune regulation (10). Obtaining sufficient ES

533    products to generate the proteome described in this study was time consuming due to the difficulty

534    in culturing sufficient quantities of parasite material in hamsters. For this reason, human infection

535    with *N. americanus* is frequently modelled using other hookworms and related nematodes that

536    survive in rodents or larger animals, including *Ancylostoma sp., H. polygyrus* and *N. brasiliensis.*

537    The protein family analysis of *N. americanus* ES products revealed a diverse number of known

538    and unknown domains (391 domains total), which attests to the many biological functions of ES

539    proteins as well as to the lack of information on these proteomes. The ES products of the three

540    comparator species used here had similar protein family profiles with 458 (*A. caninum*), 434 (*H.*

541    *polygyrus*) and 628 (*N. brasiliensis*) unique domains present. To assess the overall usefulness of

542    these models, a similarity analysis was carried out on their ES products. This analysis relates

543    relative protein sequence similarities in a plot where each of the identified *N. americanus* ES

544    proteins is compared to the ES proteome of the comparator species (19, 51). The majority of the

545    *N. americanus* ES proteins (173/198), including 51/54 SCP/TAPS proteins, had homologs in the

546    ES products of all three comparators, highlighting the relevance and usefulness of all three models.

547    A total of 25 proteins did not have similarity to ES proteins from the other 3 nematodes analyzed.

548    Since *A. caninum*, *H. polygyrus* and *N. brasiliensis* are animal parasites, these 25 proteins might

549    have evolved to such an extent that they target human-specific pathways. Two unique proteins of

550    interest are the aspartyl protease NAME_01848 (PF00026) and the zinc metalloprotease

551    NAME_05081 (PF01546). Metalloproteases and aspartyl proteases play crucial roles in host tissue

552    penetration and parasite feeding, and as such, proteolytic enzymes from helminths may be of

553    particular interest as vaccine and/or drug targets (66, 67).

554    The most abundantly represented protein family (54/198; 27%) in *N. americanus* ES products was

555    proteins containing a single or double SCP/TAPS domain (PF00188). This finding aligns with

24

556 previous work that highlighted the abundance of this protein family in nematode ES products in

557 particular (15, 68). For instance, a total of 45, 90 and 25 SCP/TAPS proteins were found in the ES

558 products from *N. brasiliensis* (45/313; 14%), *A. caninum* (90/315; 29%) and *H. polygyrus*

559 respectively (25/374; 12%). Interestingly, this family of proteins is also abundant in the

560 extracellular vesicles secreted by different nematodes (64, 65). Given that SCP/TAPS proteins

561 from *N. brasiliensis* are almost exclusively secreted by the adult developmental stage (19), these

562 proteins are likely coordinating specific roles in the gastrointestinal tract of the host. In fact, it has

563 also been shown that SCP/TAPS are overrepresented at the transcript level in *N. americanus* adult

564 worms (11). While relatively little functional information is available for SCP/TAPS proteins,

565 neutrophil inhibitory factor (NIF), an SCP/TAPS protein in the ES of *A. caninum*, was reported to

566 abrogate neutrophil adhesion to the endothelium (69). However, a *N. americanus* homolog of this

567 protein was not detected in the current study, despite the presence of a NIF-encoding gene in the

568 draft genome (11). SCP/TAPS proteins are thought to play numerous and diverse roles at the host-

569 parasite interface, from defense mechanisms, normal body formation and lifespan (54). The

570 diverse nature of *N. americanus* SCP/TAPS sequences is evidenced by their phylogenetic

571 relationships (Figure 5A and B). Blast analyses presented in the Circos plot (Figure 4A) reveal

572 varying degrees of sequence homology to SCP/TAPS proteins from *A. caninum*, *H. polygyrus* or

573 *N. brasiliensis*. These SCP/TAPS proteins should be further explored to understand their roles in

574 *N. americanus*-human host interactions. Given the limited availability of information regarding

575 the function of helminth SCP/TAPS proteins in general, it was unsurprising that GO analyses

576 revealed no known molecular function or biological process for 32/54 of the SCP/TAPS from *N.*

577 *americanus* ES products, and more studies should be performed to characterize the properties of

578 this intriguing family of proteins.

25

579   Despite the lack of functional information on the SCP/TAPS proteins in parasitic helminths, of the

580   species we compared in this study, *N. americanus* proteins were generally most similar to those

581   from *A. caninum*, which could simply be a reflection of the phylogenetic similarity between the

582   two species. (Figure 5A and 5B). The trees generated in this study highlight strong clade-specific

583   similarities between SCP/TAPS in the ES products of the compared species. In support of this, the

584   vast majority of the SCP/TAPS proteins came from the four clade V species, while *A. suum*, *T.*

585   *muris* and *T. canis* had only 2, 5 and 2 SCP/TAPS proteins respectively. The clustering of *A.*

586   *caninum* with *N. americanus* and *H. polygyrus* with *N. brasiliensis* supports the notion of host-

587   specific roles for SCP/TAPS. Another trend across both single and double domain trees was the

588   diversity of evolution between SCP/TAPS within a single species. For example, the single-domain

589   tree included a sub-clade of 7 *N. americanus*-only sequences indicating an important human-

590   specific role for this evolutionary cluster. This compares well with a number of other sub-clades

591   which included proteins from the four predominant species. These SCP/TAPS proteins are more

592   likely to share a common function, potentially in host-infection or parasite development.

593   The phylogenetic analysis strongly attests to the preferred use of *A. caninum* for studying

594   hookworm molecular biology and ES products in general (70). This finding is reinforced by the

595   Circos plot of SCP/TAPS (Figure 4A). Not only was there a greater number of SCP/TAPS

596   homologs in the ES products of *A. caninum* but these proteins also had relatively higher blast

597   scores (denoted by link ribbon thickness) and higher percent identity scores (denoted by ribbon

598   darkness). This type of analysis can prove useful since it also reveals which species to consider for

599   investigating a specific *Necator* SCP/TAPS protein. For example, NAME_13850 has significant

600   sequence homology to a SCP/TAPS protein from each of the compared species; however, the

601   species with the highest sequence homology is *A. caninun* (ANCCAN_19759), making this protein

602    the most relevant to study as a model for NAME_13850. Prior to updating the genome, all three

603    species of nematode used in this comparison had longer average gene sequence lengths than the

604    *N. americanus* SCP/TAPS sequences. In the previous annotation, the average *N. americanus*

605    SCP/TAPS sequence length was 244 predicted amino acids, compared with 355 residues in the

606    updated genome. This was likely due to some of the sequences being truncated, yielding similar

607    results to the published *H. polygyrus* proteome (52).

608    The Circos plot representing *N. americanus* proteases and homologous proteins from the three

609    comparator species (Figure 4B) provides insight into the high degree of sequence similarity.

610    Unlike the SCP/TAPS Circos analysis, all the *N. americanus* ES proteases had homologs with high

611    similarity in the compared species (average 50% identity between *N. americanus* and comparator

612    species protease). This finding supports the notion of using any one of these three parasites to

613    study *N. americanus* proteases in general. We identified 8 aspartyl, 6 cysteine, 9 metallo and 3

614    serine proteases in the ES products of *N. americanus*. As the adult stage of *N. americanus* feeds

615    on blood, the high abundance of aspartyl proteases was expected (25, 71). Yet when compared to

616    other species, the human hookworm ES products included more of these proteases. Aspartyl

617    proteases play a fundamental role in the digestion of host hemoglobin and have also been

618    implicated in skin penetration, feeding, and host tissue degradation (72, 73). The finding that *N.*

619    *americanus* has more of this mechanistic class of proteases than the other nematodes assessed here

620    is likely due to split gene models and/or may be a true gene family expansion. Due to the vital role

621    that aspartyl proteases play in parasite feeding, *Na*-APR-1 an *N. americanus* aspartyl protease, was

622    targeted as a vaccine candidate (74). While we were unable to detect *Na*-APR-1 in the current ES

623    proteome - probably because it is anchored to the gut epithelium (75) - 9 other aspartyl proteases

624    were detected which could potentially be targeted as novel vaccine candidates.

625 Cysteine proteases, particularly the group belonging to the papain superfamily, are common in

626 nematodes (76). They have been specifically described for their proteolytic activity against

627 hemoglobin, antibodies and fibrinogen in the *N. americanus* lifecycle (77). Similarly, another

628 study highlighted the importance of four cysteine proteases that were upregulated in the *N.*

629 *americanus* transition from free-living larvae to blood-feeding adult worm, indicating that these

630 proteins are likely to be important for nutrient acquisition (78). Metalloproteases - particularly the

631 astacins - identified in this study are most likely important for larval and adult migration and

632 invasion through human host tissue (79). In support of this, astacin metalloproteases were found

633 to be upregulated in *N. brasiliensis* larvae when compared with adult stage parasites (19).

634 Interestingly, *N. americanus* metalloproteases were reported to inhibit eosinophil recruitment

635 through the cleavage of eotaxin, a potent eosinophil chemoattractant (23). The least abundant

636 family of proteases in the adult *N. americanus* ES products were the serine proteases. Relatively

637 little is known about these proteases from *N. americanus* specifically; however, a serine protease

638 from the whipworm *T. muris* is involved in degradation of the mucus barrier to facilitate feeding

639 (80). Due to the importance of these various proteases in parasite feeding, infection, migration and

640 defense, they represent potential targets for chemotherapies and vaccines to limit infection (81).

641 In the current study, we have provided the first proteogenomic analysis of a helminth parasite,

642 resulting in a more accurate genome annotation. While the above annotations are valuable

643 additions to the current genome, further improvement requires access to substantially more

644 parasitic material including adult stage parasites which are difficult to obtain. Furthermore, we

645 have carried out the first proteomic analysis of the ES products of the human hookworm *N.*

646 *americanus*. The results presented herein offer significant insight into the validity of these model

647 species while also highlighting differences between these important parasites.

Proteomic analysis of *Necator americanus* secreted proteins

648

649

650  **Acknowledgements**

658

659  **Data availability**

660  The mass spectrometry proteomics data have been deposited to the ProteomeXchange Consortium

661  via the PRIDE partner repository with the dataset identifier PXD010669.

662

663  To access the data, please use the below login details.

664  **Username:** reviewer49665@ebi.ac.uk

665  **Password:** fybW0mgJ

666

667  The newly annotated, improved gene models were used in subsequent proteomic analysis of ES

668  products, and are publicly available on Nematode.net.

669

670

671

672 **References**

673 1.      Ngui, R., Lim, Y. A., Traub, R., Mahmud, R., and Mistam, M. S. (2012) Epidemiological

674 and genetic data supporting the transmission of *Ancylostoma ceylanicum* among human and

675 domestic animals. *PLoS Negl Trop Dis* 6, e1522

676 2.      Loukas, A., Hotez, P. J., Diemert, D., Yazdanbakhsh, M., McCarthy, J. S., Correa-Oliveira,

677 R., Croese, J., and Bethony, J. M. (2016) Hookworm infection. *Nat Rev Dis Primers* 2, 16088

678 3.      Hotez, P., Brooker, S., Bethony, J., Bottazzi, M. E., Loukas, A., and Xiao, S. (2004)

679 Hookworm Infection. *N Engl J Med* 351, 799-807

680 4.      Brooker, S., Bethony, J., and Hotez, P. J. (2004) Human Hookworm Infection in the 21st

681 Century. *Adv Parasitol* 58, 197-288

682 5.      Christian, P., Khatry, S. K., and West, K. P. (2004) Antenatal anthelmintic treatment,

683 birthweight, and infant survival in rural Nepal. *Lancet* 364, 981-983

684 6.      Brooker, S., Hotez, P. J., and Bundy, D. A. (2008) Hookworm-related anaemia among

685 pregnant women: a systematic review. *PLoS Negl Trop Dis* 2, e291

686 7.      Murray, C. J. L., Vos, T., Lozano, R., Naghavi, M., Flaxman, A. D., Michaud, C., Ezzati,

687 M., Shibuya, K., Salomon, J. A., Abdalla, S., Aboyans, V., Abraham, J., Ackerman, I., Aggarwal,

688 R., Ahn, S. Y., Ali, M. K., AlMazroa, M. A., Alvarado, M., Anderson, H. R., Anderson, L. M.,

689 Andrews, K. G., Atkinson, C., Baddour, L. M., Bahalim, A. N., Barker-Collo, S., Barrero, L. H.,

690 Bartels, D. H., Basáñez, M.-G., Baxter, A., Bell, M. L., Benjamin, E. J., Bennett, D., Bernabé, E.,

691 Bhalla, K., Bhandari, B., Bikbov, B., Abdulhak, A. B., Birbeck, G., Black, J. A., Blencowe, H.,

692 Blore, J. D., Blyth, F., Bolliger, I., Bonaventure, A., Boufous, S., Bourne, R., Boussinesq, M.,

693 Braithwaite, T., Brayne, C., Bridgett, L., Brooker, S., Brooks, P., Brugha, T. S., Bryan-Hancock,

694 C., Bucello, C., Buchbinder, R., Buckle, G., Budke, C. M., Burch, M., Burney, P., Burstein, R.,

695    Calabria, B., Campbell, B., Canter, C. E., Carabin, H., Carapetis, J., Carmona, L., Cella, C.,

696    Charlson, F., Chen, H., Cheng, A. T.-A., Chou, D., Chugh, S. S., Coffeng, L. E., Colan, S. D.,

697    Colquhoun, S., Colson, K. E., Condon, J., Connor, M. D., Cooper, L. T., Corriere, M., Cortinovis,

698    M., de Vaccaro, K. C., Couser, W., Cowie, B. C., Criqui, M. H., Cross, M., Dabhadkar, K. C.,

699    Dahiya, M., Dahodwala, N., Damsere-Derry, J., Danaei, G., Davis, A., Leo, D. D., Degenhardt,

700    L., Dellavalle, R., Delossantos, A., Denenberg, J., Derrett, S., Des Jarlais, D. C., Dharmaratne, S.

701    D., Dherani, M., Diaz-Torne, C., Dolk, H., Dorsey, E. R., Driscoll, T., Duber, H., Ebel, B.,

702    Edmond, K., Elbaz, A., Ali, S. E., Erskine, H., Erwin, P. J., Espindola, P., Ewoigbokhan, S. E.,

703    Farzadfar, F., Feigin, V., Felson, D. T., Ferrari, A., Ferri, C. P., Fèvre, E. M., Finucane, M. M.,

704    Flaxman, S., Flood, L., Foreman, K., Forouzanfar, M. H., Fowkes, F. G. R., Fransen, M., Freeman,

705    M. K., Gabbe, B. J., Gabriel, S. E., Gakidou, E., Ganatra, H. A., Garcia, B., Gaspari, F., Gillum,

706    R. F., Gmel, G., Gonzalez-Medina, D., Gosselin, R., Grainger, R., Grant, B., Groeger, J.,

707    Guillemin, F., Gunnell, D., Gupta, R., Haagsma, J., Hagan, H., Halasa, Y. A., Hall, W., Haring,

708    D., Haro, J. M., Harrison, J. E., Havmoeller, R., Hay, R. J., Higashi, H., Hill, C., Hoen, B.,

709    Hoffman, H., Hotez, P. J., Hoy, D., Huang, J. J., Ibeanusi, S. E., Jacobsen, K. H., James, S. L.,

710    Jarvis, D., Jasrasaria, R., Jayaraman, S., Johns, N., Jonas, J. B., Karthikeyan, G., Kassebaum, N.,

711    Kawakami, N., Keren, A., Khoo, J.-P., King, C. H., Knowlton, L. M., Kobusingye, O., Koranteng,

712    A., Krishnamurthi, R., Laden, F., Lalloo, R., Laslett, L. L., Lathlean, T., Leasher, J. L., Lee, Y. Y.,

713    Leigh, J., Levinson, D., Lim, S. S., Limb, E., Lin, J. K., Lipnick, M., Lipshultz, S. E., Liu, W.,

714    Loane, M., Ohno, S. L., Lyons, R., Mabweijano, J., MacIntyre, M. F., Malekzadeh, R., Mallinger,

715    L., Manivannan, S., Marcenes, W., March, L., Margolis, D. J., Marks, G. B., Marks, R.,

716    Matsumori, A., Matzopoulos, R., Mayosi, B. M., McAnulty, J. H., McDermott, M. M., McGill,

717    N., McGrath, J., Medina-Mora, M. E., Meltzer, M., Memish, Z. A., Mensah, G. A., Merriman, T.

Proteomic analysis of *Necator americanus* secreted proteins

718    R., Meyer, A.-C., Miglioli, V., Miller, M., Miller, T. R., Mitchell, P. B., Mock, C., Mocumbi, A.

719    O., Moffitt, T. E., Mokdad, A. A., Monasta, L., Montico, M., Moradi-Lakeh, M., Moran, A.,

720    Morawska, L., Mori, R., Murdoch, M. E., Mwaniki, M. K., Naidoo, K., Nair, M. N., Naldi, L.,

721    Narayan, K. M. V., Nelson, P. K., Nelson, R. G., Nevitt, M. C., Newton, C. R., Nolte, S., Norman,

722    P., Norman, R., O'Donnell, M., O'Hanlon, S., Olives, C., Omer, S. B., Ortblad, K., Osborne, R.,

723    Ozgediz, D., Page, A., Pahari, B., Pandian, J. D., Rivero, A. P., Patten, S. B., Pearce, N., Padilla,

724    R. P., Perez-Ruiz, F., Perico, N., Pesudovs, K., Phillips, D., Phillips, M. R., Pierce, K., Pion, S.,

725    Polanczyk, G. V., Polinder, S., Pope Iii, C. A., Popova, S., Porrini, E., Pourmalek, F., Prince, M.,

726    Pullan, R. L., Ramaiah, K. D., Ranganathan, D., Razavi, H., Regan, M., Rehm, J. T., Rein, D. B.,

727    Remuzzi, G., Richardson, K., Rivara, F. P., Roberts, T., Robinson, C., De Leòn, F. R., Ronfani,

728    L., Room, R., Rosenfeld, L. C., Rushton, L., Sacco, R. L., Saha, S., Sampson, U., Sanchez-Riera,

729    L., Sanman, E., Schwebel, D. C., Scott, J. G., Segui-Gomez, M., Shahraz, S., Shepard, D. S., Shin,

730    H., Shivakoti, R., Silberberg, D., Singh, D., Singh, G. M., Singh, J. A., Singleton, J., Sleet, D. A.,

731    Sliwa, K., Smith, E., Smith, J. L., Stapelberg, N. J. C., Steer, A., Steiner, T., Stolk, W. A., Stovner,

732    L. J., Sudfeld, C., Syed, S., Tamburlini, G., Tavakkoli, M., Taylor, H. R., Taylor, J. A., Taylor, W.

733    J., Thomas, B., Thomson, W. M., Thurston, G. D., Tleyjeh, I. M., Tonelli, M., Towbin, J. A.,

734    Truelsen, T., Tsilimbaris, M. K., Ubeda, C., Undurraga, E. A., van der Werf, M. J., van Os, J.,

735    Vavilala, M. S., Venketasubramanian, N., Wang, M., Wang, W., Watt, K., Weatherall, D. J.,

736    Weinstock, M. A., Weintraub, R., Weisskopf, M. G., Weissman, M. M., White, R. A., Whiteford,

737    H., Wiebe, N., Wiersma, S. T., Wilkinson, J. D., Williams, H. C., Williams, S. R. M., Witt, E.,

738    Wolfe, F., Woolf, A. D., Wulf, S., Yeh, P.-H., Zaidi, A. K. M., Zheng, Z.-J., Zonies, D., and Lopez,

739    A. D. (2012) Disability-adjusted life years (DALYs) for 291 diseases and injuries in 21 regions,

740    1990–2010: a systematic analysis for the Global Burden of Disease Study 2010. *Lancet* 380, 2197-

741    2223

742    8.      Hotez, P. (2008) Hookworm and poverty. *Ann N Y Acad Sci* 1136, 38-44

743    9.      Hotez, P. J., Bethony, J., Bottazzi, M. E., Brooker, S., and Buss, P. (2005) Hookworm: "the

744    great infection of mankind". *PLoS Med* 2, e67

745    10.     Hewitson, J. P., Grainger, J. R., and Maizels, R. M. (2009) Helminth immunoregulation:

746    The role of parasite secreted proteins in modulating host immunity. *Mol Biochem Parasitol* 167,

747    1-11

748    11.     Tang, Y. T., Gao, X., Rosa, B. A., Abubucker, S., Hallsworth-Pepin, K., Martin, J., Tyagi,

749    R., Heizer, E., Zhang, X., Bhonagiri-Palsikar, V., Minx, P., Warren, W. C., Wang, Q., Zhan, B.,

750    Hotez, P. J., Sternberg, P. W., Dougall, A., Gaze, S. T., Mulvenna, J., Sotillo, J., Ranganathan, S.,

751    Rabelo, E. M., Wilson, R. K., Felgner, P. L., Bethony, J., Hawdon, J. M., Gasser, R. B., Loukas,

752    A., and Mitreva, M. (2014) Genome of the human hookworm *Necator americanus*. *Nature Genet*

753    46, 261-269

754    12.     Navarro, S., Pickering, D. A., Ferreira, I. B., Jones, L., Ryan, S., Troy, S., Leech, A., Hotez,

755    P. J., Zhan, B., Laha, T., Prentice, R., Sparwasser, T., Croese, J., Engwerda, C. R., Upham, J. W.,

756    Julia, V., Giacomin, P. R., and Loukas, A. (2016) Hookworm recombinant protein promotes

757    regulatory T cell responses that suppress experimental asthma. *Sci Transl Med* 8, 362ra143

758    13.     Hotez, P. J., Bethony, J. M., Diemert, D. J., Pearson, M., and Loukas, A. (2010) Developing

759    vaccines to combat hookworm infection and intestinal schistosomiasis. *Nat Rev Microbiol* 8, 814-

760    826

761    14.     Cancado, G. G., Fiuza, J. A., de Paiva, N. C., Lemos Lde, C., Ricci, N. D., Gazzinelli-

762    Guimaraes, P. H., Martins, V. G., Bartholomeu, D. C., Negrao-Correa, D. A., Carneiro, C. M., and

763  Fujiwara, R. T. (2011) Hookworm products ameliorate dextran sodium sulfate-induced colitis in

764  BALB/c mice. *Inflamm Bowel Dis* 17, 2275-2286

765  15.   Mulvenna, J., Hamilton, B., Nagaraj, S. H., Smyth, D., Loukas, A., and Gorman, J. J.

766  (2009) Proteomics analysis of the excretory/secretory component of the blood-feeding stage of the

767  hookworm, *Ancylostoma caninum*. *Mol Cell Proteomics* 8, 109-121

768  16.   Hotez, P. J., Beaumier, C. M., Gillespie, P. M., Strych, U., Hayward, T., and Bottazzi, M.

769  E. (2016) Advancing a vaccine to prevent hookworm disease and anemia. *Vaccine* 34, 3001-3005

770  17.   Diemert, D. J., Freire, J., Valente, V., Fraga, C. G., Talles, F., Grahek, S., Campbell, D.,

771  Jariwala, A., Periago, M. V., Enk, M., Gazzinelli, M. F., Bottazzi, M. E., Hamilton, R., Brelsford,

772  J., Yakovleva, A., Li, G., Peng, J., Correa-Oliveira, R., Hotez, P., and Bethony, J. (2017) Safety

773  and immunogenicity of the Na-GST-1 hookworm vaccine in Brazilian and American adults. *PLoS*

774  *Negl Trop Dis* 11, e0005574

775  18.   Hunt, V. L., Tsai, I. J., Coghlan, A., Reid, A. J., Holroyd, N., Foth, B. J., Tracey, A., Cotton,

776  J. A., Stanley, E. J., Beasley, H., Bennett, H. M., Brooks, K., Harsha, B., Kajitani, R., Kulkarni,

777  A., Harbecke, D., Nagayasu, E., Nichol, S., Ogura, Y., Quail, M. A., Randle, N., Xia, D., Brattig,

778  N. W., Soblik, H., Ribeiro, D. M., Sanchez-Flores, A., Hayashi, T., Itoh, T., Denver, D. R., Grant,

779  W., Stoltzfus, J. D., Lok, J. B., Murayama, H., Wastling, J., Streit, A., Kikuchi, T., Viney, M., and

780  Berriman, M. (2016) The genomic basis of parasitism in the Strongyloides clade of nematodes.

781  *Nature Genet* 48, 299-307

782  19.   Sotillo, J., Sanchez-Flores, A., Cantacessi, C., Harcus, Y., Pickering, D., Bouchery, T.,

783  Camberis, M., Tang, S., Giacomin, P., Mulvenna, J., Mitreva, M., Berriman, M., LeGros, G.,

784  Maizels, R., and Loukas, A. (2014) Secreted proteomes of different developmental stages of the

785  gastrointestinal nematode *Nippostrongylus brasiliensis*. *Mol Cell Proteomics* 13, 2736-2751

786    20.    Morante, T., Shepherd, C., Constantinoiu, C., Loukas, A., and Sotillo, J. (2017) Revisiting

787    the *Ancylostoma caninum* secretome provides new information on hookworm-host interactions.

788    *Proteomics* 17

789    21.    Hsieh, G. C. F., Loukas, A., Wahl, A. M., Bhatia, M., Wang, Y., Williamson, A. L., Kehn,

790    K. W., Maruyama, H., Hotez, P. J., Leitenberg, D., Bethony, J., and Constant, S. L. (2004) A

791    secreted protein from the human hookworm *Necator americanus* binds selectively to NK cells and

792    induces IFN- production. *J Immunol* 173, 2699-2704

793    22.    Winter, J., Davies, O., Brown, A., Garnett, M., Stolnik, S., and Pritchard, D. (2005) The

794    assessment of hookworm calreticulin as a potential vaccine for necatoriasis. *Parasite Immunol* 27,

795    139–146

796    23.    Culley, F. J., Brown, A., Conroy, D. M., Sabroe, I., Pritchard, D. I., and Williams, T. J.

797    (2000) Eotaxin is specifically cleaved by hookworm metalloproteases preventing its action *in vitro*

798    and *in vivo*. *J Immunol* 165, 6447-6453

799    24.    Kumar, S., and Pritchard, D. I. (1992) Secretion of metalloproteases by living infective

800    larvae of *Necator americanus*. *J Parasitol* 78, 917-919

801    25.    Ranjit, N., Zhan, B., Hamilton, B., Stenzel, D., Lowther, J., Pearson, M., Gorman, J.,

802    Hotez, P., and Loukas, A. (2009) Proteolytic degradation of hemoglobin in the intestine of the

803    human hookworm *Necator americanus*. *J Infect Dis* 199, 904-912

804    26.    Cantacessi, C., Mitreva, M., Jex, A. R., Young, N. D., Campbell, B. E., Hall, R. S., Doyle,

805    M. A., Ralph, S. A., Rabelo, E. M., Ranganathan, S., Sternberg, P. W., Loukas, A., and Gasser, R.

806    B. (2010) Massively parallel sequencing and analysis of the *Necator americanus* transcriptome.

807    *PLoS Negl Trop Dis* 4, e684

808    27.      Wang, Z., Abubucker, S., Martin, J., Wilson, R. K., Hawdon, J., and Mitreva, M. (2010)

809    Characterizing *Ancylostoma caninum* transcriptome and exploring nematode parasitic adaptation.

810    *BMC Genomics* 11, 307-307

811    28.      Schwarz, E. M., Hu, Y., Antoshechkin, I., Miller, M. M., Sternberg, P. W., and Aroian, R.

812    V. (2015) The genome and transcriptome of the zoonotic hookworm *Ancylostoma ceylanicum*

813    identify infection-specific gene families. *Nat Genet* 47, 416-422

814    29.      Palevich, N., Britton, C., Kamenetzky, L., Mitreva, M., de Moraes Mourão, M., Bennuru,

815    S., Quack, T., Scholte, L. L. S., Tyagi, R., and Slatko, B. E. (2018) Tackling hypotheticals in

816    helminth genomes. *Trends Parasitol* 34, 179-183

817    30.      Borchert, N., Dieterich, C., Krug, K., Schutz, W., Jung, S., Nordheim, A., Sommer, R. J.,

818    and Macek, B. (2010) Proteogenomics of *Pristionchus pacificus* reveals distinct proteome

819    structure of nematode models. *Genome Res* 20, 837-846

820    31.      Nesvizhskii, A. I. (2014) Proteogenomics: concepts, applications and computational

821    strategies. *Nat Methods* 11, 1114-1125

822    32.      Holt, C., and Yandell, M. (2011) MAKER2: an annotation pipeline and genome-database

823    management tool for second-generation genome projects. *BMC Bioinformatics* 12, 491

824    33.      Bao, W., Kojima, K. K., and Kohany, O. (2015) Repbase Update, a database of repetitive

825    elements in eukaryotic genomes. *Mob DNA* 6, 11

826    34.      Bolger, A. M., Lohse, M., and Usadel, B. (2014) Trimmomatic: a flexible trimmer for

827    Illumina sequence data. *Bioinformatics* 30, 2114-2120

828    35.      Kim, D., Langmead, B., and Salzberg, S. L. (2015) HISAT: a fast spliced aligner with low

829    memory requirements. *Nat Methods* 12, 357-360

830    36.    Pertea, M., Pertea, G. M., Antonescu, C. M., Chang, T. C., Mendell, J. T., and Salzberg, S.

831    L. (2015) StringTie enables improved reconstruction of a transcriptome from RNA-seq reads. *Nat*

832    *Biotechnol* 33, 290-295

833    37.    Hoff, K. J., Lange, S., Lomsadze, A., Borodovsky, M., and Stanke, M. (2016) BRAKER1:

834    Unsupervised RNA-Seq-Based genome annotation with GeneMark-ET and AUGUSTUS.

835    *Bioinformatics* 32, 767-769

836    38.    The UniProt, C. (2017) UniProt: the universal protein knowledgebase. *Nucleic Acids Res*

837    45, D158-D169

838    39.    Howe, K. L., Bolt, B. J., Shafie, M., Kersey, P., and Berriman, M. (2017) WormBase

839    ParaSite - a comprehensive resource for helminth genomics. *Mol Biochem Parasitol* 215, 2-10

840    40.    Campbell, M. S., Holt, C., Moore, B., and Yandell, M. (2014) Genome annotation and

841    curation ising MAKER and MAKER-P. *Curr Protoc Bioinformatics* 48, 4 11 11-39

842    41.    Lomsadze, A., Burns, P. D., and Borodovsky, M. (2014) Integration of mapped RNA-Seq

843    reads into automatic training of eukaryotic gene finding algorithm. *Nucleic Acids Res* 42, e119

844    42.    Stanke, M., Diekhans, M., Baertsch, R., and Haussler, D. (2008) Using native and

845    syntenically mapped cDNA alignments to improve *de novo* gene finding. *Bioinformatics* 24, 637-

846    644

847    43.    Jones, P., Binns, D., Chang, H. Y., Fraser, M., Li, W., McAnulla, C., McWilliam, H.,

848    Maslen, J., Mitchell, A., Nuka, G., Pesseat, S., Quinn, A. F., Sangrador-Vegas, A., Scheremetjew,

849    M., Yong, S. Y., Lopez, R., and Hunter, S. (2014) InterProScan 5: genome-scale protein function

850    classification. *Bioinformatics* 30, 1236-1240

851    44.    Campbell, M. S., Law, M., Holt, C., Stein, J. C., Moghe, G. D., Hufnagel, D. E., Lei, J.,

852    Achawanantakun, R., Jiao, D., Lawrence, C. J., Ware, D., Shiu, S. H., Childs, K. L., Sun, Y., Jiang,

853    N., and Yandell, M. (2014) MAKER-P: a tool kit for the rapid creation, management, and quality

854    control of plant genome annotations. *Plant Physiol* 164, 513-524

855    45.    Simao, F. A., Waterhouse, R. M., Ioannidis, P., Kriventseva, E. V., and Zdobnov, E. M.

856    (2015) BUSCO: assessing genome assembly and annotation completeness with single-copy

857    orthologs. *Bioinformatics* 31, 3210-3212

858    46.    Conesa, A., Gotz, S., Garcia-Gomez, J. M., Terol, J., Talon, M., and Robles, M. (2005)

859    Blast2GO: a universal tool for annotation, visualization and analysis in functional genomics

860    research. *Bioinformatics* 21, 3674-3676

861    47.    Emanuelsson, O., Brunak, S., von Heijne, G., and Nielsen, H. (2007) Locating proteins in

862    the cell using TargetP, SignalP and related tools. *Nat Protoc* 2, 953-971

863    48.    Krogh, A., Larsson, B., von Heijne, G., and Sonnhammer, E. L. (2001) Predicting

864    transmembrane protein topology with a hidden Markov model: application to complete genomes.

865    *J Mol Biol* 305, 567-580

866    49.    Supek, F., Bosnjak, M., Skunca, N., and Smuc, T. (2011) REVIGO summarizes and

867    visualizes long lists of gene ontology terms. *PLoS One* 6, e21800

868    50.    Conway, J. R., Lex, A., and Gehlenborg, N. (2017) UpSetR: an R package for the

869    visualization of intersecting sets and their properties. *Bioinformatics* 33, 2938-2940

870    51.    Parkinson, J., and Blaxter, M. (2003) SimiTri--visualizing similarity relationships for

871    groups of sequences. *Bioinformatics* 19, 390-395

872    52.    Hewitson, J. P., Harcus, Y., Murray, J., van Agtmaal, M., Filbey, K. J., Grainger, J. R.,

873    Bridgett, S., Blaxter, M. L., Ashton, P. D., Ashford, D. A., Curwen, R. S., Wilson, R. A., Dowle,

874    A. A., and Maizels, R. M. (2011) Proteomic analysis of secretory products from the model

875    gastrointestinal nematode *Heligmosomoides polygyrus* reveals dominance of venom allergen-like

876    (VAL) proteins. *J Proteomics* 74, 1573-1594

877    53.    Krzywinski, M., Schein, J., Birol, I., Connors, J., Gascoyne, R., Horsman, D., Jones, S. J.,

878    and Marra, M. A. (2009) Circos: an information aesthetic for comparative genomics. *Genome Res*

879    19, 1639-1645

880    54.    Cantacessi, C., Campbell, B. E., Visser, A., Geldhof, P., Nolan, M. J., Nisbet, A. J.,

881    Matthews, J. B., Loukas, A., Hofmann, A., Otranto, D., Sternberg, P. W., and Gasser, R. B. (2009)

882    A portrait of the "SCP/TAPS" proteins of eukaryotes--developing a framework for fundamental

883    research and biotechnological outcomes. *Biotechnol Adv* 27, 376-388

884    55.    Letunic, I., and Bork, P. (2007) Interactive Tree Of Life (iTOL): an online tool for

885    phylogenetic tree display and annotation. *Bioinformatics* 23, 127-128

886    56.    Martin, J., Rosa, B. A., Ozersky, P., Hallsworth-Pepin, K., Zhang, X., Bhonagiri-Palsikar,

887    V., Tyagi, R., Wang, Q., Choi, Y.-J., Gao, X., McNulty, S. N., Brindley, P. J., and Mitreva, M.

888    (2015) Helminth.net: expansions to Nematode.net and an introduction to Trematode.net. *Nucleic*

889    *Acids Res* 43, D698-D706

890    57.    Martin, J., Tyagi, R., Rosa, B. A., and Mitreva, M. (2018) A multi-omics database for

891    parasitic nematodes and trematodes. *Methods Mol Biol* 1757, 371-397

892    58.    Jaffe, J. D., Berg, H. C., and Church, G. M. (2004) Proteogenomic mapping as a

893    complementary method to perform genome annotation. *Proteomics* 4, 59-77

894    59.    Sotillo, J., Toledo, R., Mulvenna, J., and Loukas, A. (2017) Exploiting helminth-host

895    interactomes through big data. *Trends Parasitol* 33, 875-888

896    60.    Protasio, A. V., Tsai, I. J., Babbage, A., Nichol, S., Hunt, M., Aslett, M. A., De Silva, N.,

897    Velarde, G. S., Anderson, T. J., Clark, R. C., Davidson, C., Dillon, G. P., Holroyd, N. E., LoVerde,

898    P. T., Lloyd, C., McQuillan, J., Oliveira, G., Otto, T. D., Parker-Manuel, S. J., Quail, M. A.,

899    Wilson, R. A., Zerlotini, A., Dunne, D. W., and Berriman, M. (2012) A systematically improved

900    high quality genome and transcriptome of the human blood fluke *Schistosoma mansoni*. *PLoS Negl*

901    *Trop Dis* 6, e1455

902    61.    Zarowiecki, M., and Berriman, M. (2015) What helminth genomes have taught us about

903    parasite evolution. *Parasitology* 142, S85-S97

904    62.    Coakley, G., Buck, A. H., and Maizels, R. M. (2016) Host parasite communications—

905    Messages from helminths for the immune system: Parasite communication and cell-cell

906    interactions. *Mol Biochem Parasitol* 208, 33-40

907    63.    Marcilla, A., Martin-Jaular, L., Trelis, M., de Menezes-Neto, A., Osuna, A., Bernal, D.,

908    Fernandez-Becerra, C., Almeida, I. C., and del Portillo, H. A. (2014) Extracellular vesicles in

909    parasitic diseases. *J Extracell Vesicles* 3, 10.3402/jev.v3403.25040

910    64.    Eichenberger, R. M., Ryan, S., Jones, L., Buitrago, G., Polster, R., Montes de Oca, M.,

911    Zuvelek, J., Giacomin, P. R., Dent, L. A., Engwerda, C. R., Field, M. A., Sotillo, J., and Loukas,

912    A. (2018) Hookworm secreted extracellular vesicles interact with host cells and prevent inducible

913    colitis in mice. *Front Immunol* 9

914    65.    Eichenberger, R. M., Talukder, M. H., Field, M. A., Wangchuk, P., Giacomin, P., Loukas,

915    A., and Sotillo, J. (2018) Characterization of *Trichuris muris* secreted proteins and extracellular

916    vesicles provides new insights into host–parasite communication. *J Extracell Vesicles* 7, 1428004

917    66.    Knox, D. (2011) Proteases in blood-feeding nematodes and their potential as vaccine

918    candidates. *Adv Exp Med Biol* 712, 155-176

919    67.    Abdulla, M. H., Lim, K. C., Sajid, M., McKerrow, J. H., and Caffrey, C. R. (2007)

920    *Schistosomiasis mansoni*: novel chemotherapy using a cysteine protease inhibitor. *PLoS Med* 4,

921    e14

922    68.    Hewitson, J. P., Ivens, A. C., Harcus, Y., Filbey, K. J., McSorley, H. J., Murray, J.,

923    Bridgett, S., Ashford, D., Dowle, A. A., and Maizels, R. M. (2013) Secretion of protective antigens

924    by tissue-stage nematode larvae revealed by proteomic analysis and vaccination-induced sterile

925    immunity. *PLoS Pathog* 9, e1003492

926    69.    Lo, S. K., Rahman, A., Xu, N., Zhou, M. Y., Nagpala, P., Jaffe, H. A., and Malik, A. B.

927    (1999) Neutrophil inhibitory factor abrogates neutrophil adhesion by blockade of CD11a and

928    CD11b beta(2) integrins. *Mol Pharmacol* 56, 926-932

929    70.    Shepherd, C., Wangchuk, P., and Loukas, A. (2018) Of dogs and hookworms: man's best

930    friend and his parasites as a model for translational biomedical research. *Parasit Vectors* 11, 59

931    71.    Ranjit, N., Jones, M. K., Stenzel, D. J., Gasser, R. B., and Loukas, A. (2006) A survey of

932    the intestinal transcriptomes of the hookworms, *Necator americanus* and *Ancylostoma caninum*,

933    using tissues isolated by laser microdissection microscopy. *Int J Parasitol* 36, 701-710

934    72.    Williamson, A. L., Brindley, P. J., Abbenante, G., Prociv, P., Berry, C., Girdwood, K.,

935    Pritchard, D. I., Fairlie, D. P., Hotez, P. J., Dalton, J. P., and Loukas, A. (2002) Cleavage of

936    hemoglobin by hookworm cathepsin D aspartic proteases and its potential contribution to host

937    specificity. *FASEB J* 16, 1458-1460

938    73.    Lin, X., Koelsch, G., Wu, S., Downs, D., Dashti, A., and Tang, J. (2000) Human aspartic

939    protease memapsin 2 cleaves the β-secretase site of β-amyloid precursor protein. *Proc Natl Acad*

940    *Sci USA* 97, 1456-1460

941  74.    Hotez, P. J., Diemert, D., Bacon, K. M., Beaumier, C., Bethony, J. M., Bottazzi, M. E.,

942    Brooker, S., Couto, A. R., da Silva Freire, M., Homma, A., Lee, B. Y., Loukas, A., Loblack, M.,

943    Morel, C. M., Oliveira, R. C., and Russell, P. K. (2013) The Human Hookworm Vaccine. *Vaccine*

944    31, B227-B232

945  75.    Loukas, A., Bethony, J. M., Mendez, S., Fujiwara, R. T., Goud, G. N., Ranjit, N., Zhan, B.,

946    Jones, K., Bottazzi, M. E., and Hotez, P. J. (2005) Vaccination with recombinant aspartic

947    hemoglobinase reduces parasite load and blood loss after hookworm infection in dogs. *PLoS Med*

948    2, e295

949  76.    Ono, Y., and Sorimachi, H. (2012) Calpains — An elaborate proteolytic system. *Biochim.*

950    *Biophys. Acta, Proteins Proteomics* 1824, 224-236

951  77.    Baig, S., Damian, R. T., and Peterson, D. S. (2002) A novel cathepsin B active site motif

952    is shared by helminth bloodfeeders. *Exp Parasitol* 101, 83-89

953  78.    Ranjit, N., Zhan, B., Stenzel, D. J., Mulvenna, J., Fujiwara, R., Hotez, P. J., and Loukas,

954    A. (2008) A family of cathepsin B cysteine proteases expressed in the gut of the human hookworm,

955    *Necator americanus*. *Mol Biochem Parasitol* 160, 90-99

956  79.    Williamson, A. L., Lustigman, S., Oksov, Y., Deumic, V., Plieskatt, J., Mendez, S., Zhan,

957    B., Bottazzi, M. E., Hotez, P. J., and Loukas, A. (2006) *Ancylostoma caninum* MTP-1, an astacin-

958    like metalloprotease secreted by infective hookworm larvae, is involved in tissue migration. *Infect*

959    *Immun* 74, 961-967

960  80.    Hasnain, S. Z., McGuckin, M. A., Grencis, R. K., and Thornton, D. J. (2012) Serine

961    protease(s) secreted by the nematode *Trichuris muris* degrade the mucus barrier. *PLoS Negl Trop*

962    *Dis* 6, e1856

963    81.    Mendez, S., Zhan, B., Goud, G., Ghosh, K., Dobardzic, A., Wu, W., Liu, S., Deumic, V.,

964    Dobardzic, R., Liu, Y., Bethony, J., and Hotez, P. J. (2005) Effect of combining the larval antigens

965    Ancylostoma secreted protein 2 (ASP-2) and metalloprotease 1 (MTP-1) in protecting hamsters

966    against hookworm infection and disease caused by *Ancylostoma ceylanicum*. *Vaccine* 23, 3123-

967    3130

968

969 ## Table 1: Comparison of original and updated genome annotations
970

| | Original annotation | Updated annotation |
|---|---|---|
| Number of genes | 19,153 | 15,728 |
| Number of exons | 122,849 | 148,780 |
| Number of introns | 103,696 | 133,052 |
| Number of CDS | 19,153 | 15,728 |
| Overlapping genes | 395 | 2,424 |
| Contained genes | 2 | 386 |
| Total gene length (bp) | 82,090,364 | 126,651,725 |
| Total exon length (bp) | 15,470,227 | 24,553,753 |
| Total intron length (bp) | 66,827,529 | 102,364,076 |
| Total CDS length (bp) | 15,461,420 | 17,683,227 |
| Mean gene length (bp) | 4,286 | 8,053 |
| Mean exon length (bp) | 126 | 165 |
| Mean intron length (bp) | 644 | 769 |
| Mean CDS length (bp) | 807 | 1,124 |
| % of genome covered by genes | 33.6 | 51.9 |
| % of genome covered by CDS | 6.3 | 7.2 |
| Mean exons per mRNA | 6.4 | 9.5 |
| Mean introns per mRNA | 5.4 | 8.5 |
| Complete BUSCOs | 82.84% | 88.45% |
| Fragmented BUSCOs | 12.87% | 8.91% |
| Missing BUSCOs | 4.29% | 2.64% |

971

972 ## Table 2: Top 30 most abundant proteins in the ES products of *Necator americanus*

| Accession | Blast2GO Description | MASCOT score Ingel | MASCOT score Offgel | iProphet probabilty Ingel | iProphet probabilty Offgel | NP Ingel | NP Offgel | SP | TD | emPAI Ingel | emPAI Offgel | Domain |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| NAME_05596 | Hypothetical protein | 19496 | 15287 | 1 | 1 | 11 | 11 | N | 0 | 867.19 | 308.05 | No conserved domains |
| NAME_07794 | Hypothetical protein | 21210 | 39596 | 1 | 1 | 18 | 22 | Y | 0 | 47.88 | 199.86 | Single SCP |
| NAME_11917 | SCP | 35077 | 43195 | 1 | 1 | 18 | 31 | N | 0 | 11.87 | 68.4 | Double SCP |
| NAME_13724 | Hypothetical protein | 7701 | 23099 | 1 | 1 | 5 | 7 | N | 0 | 7.06 | 64.55 | No conserved domains |
| NAME_15197 | SCP | 9446 | 12866 | 1 | 1 | 6 | 7 | N | 0 | 9.55 | 50.11 | No conserved domains |
| NAME_14329 | SCP | 10911 | 15142 | 1 | 1 | 10 | 15 | N | 0 | 14.34 | 43.25 | Single SCP |
| NAME_11146 | SCP | 21158 | 29296 | 1 | 1 | 24 | 27 | Y | 0 | 17.03 | 22.85 | Double SCP |
| NAME_10941 | Hypothetical protein | 983 | 3870 | 1 | 1 | 3 | 8 | Y | 0 | 2.6 | 33.23 | Single SCP |
| NAME_09098 | SCP | 6908 | 47379 | 1 | 1 | 13 | 25 | Y | 0 | 3.27 | 31.87 | Double SCP |
| NAME_05595 | SCP-like protein, partial | 0 | 23369 | 0 | 1 | 0 | 10 | Y | 0 | 0 | 26.61 | Single SCP |
| NAME_11145 | SCP | 19923 | 25558 | 1 | 1 | 18 | 22 | N | 0 | 6.35 | 16.62 | Double SCP |
| NAME_02907 | Glu Leu Phe Val dimerization | 6645 | 6095 | 1 | 1 | 20 | 22 | N | 0 | 6.82 | 7.63 | Single ELFV_dehydrog_N |
| NAME_10752 | Nematode fatty acid retinoid | 12416 | 3051 | 1 | 1 | 11 | 8 | Y | 0 | 9.76 | 4.65 | Single Gp-FAR-1 |
| NAME_05865 | Globin | 3566 | 5667 | 1 | 1 | 6 | 10 | N | 0 | 2.63 | 9.67 | Single Globin |
| NAME_15231 | SCP | 117 | 1844 | 0 | 1 | 3 | 6 | N | 0 | 1.75 | 9.66 | Single SCP |
| NAME_05592 | Hypothetical protein | 0 | 4150 | 0 | 1 | 0 | 4 | Y | 0 | 0 | 11.02 | Single SCP |
| NAME_13132 | Ferritin | 1808 | 2833 | 1 | 1 | 7 | 8 | N | 0 | 3.39 | 7.31 | Single Ferritin-like |
| NAME_05794 | Hypothetical protein | 2681 | 2271 | 1 | 1 | 7 | 9 | N | 0 | 5.18 | 5.21 | No conserved domains |
| NAME_13809 | SCP | 8287 | 14385 | 1 | 1 | 7 | 9 | Y | 0 | 3.99 | 5.92 | Single SCP |
| NAME_07942 | Copper zinc superoxide dismutase | 4241 | 4469 | 1 | 1 | 4 | 8 | N | 0 | 1.41 | 6.25 | Single Copper/zinc superoxide |
| NAME_09181 | Hypothetical protein | 27774 | 68153 | 1 | 1 | 41 | 50 | N | 0 | 3.17 | 4.43 | No conserved domains |
| NAME_10294 | SCP | 6557 | 11684 | 1 | 1 | 7 | 9 | Y | 0 | 2.58 | 4.8 | Single SCP |
| NAME_02523 | Histidine operon leader | 23085 | 24957 | 1 | 1 | 66 | 63 | Y | 0 | 3.68 | 3.3 | Many NPA |
| NAME_13379 | Flavodoxin | 0 | 1089 | 0.9985 | 0 | 0 | 4 | N | 0 | 0 | 6.89 | No conserved domains |
| NAME_02090 | SCP | 3916 | 12638 | 1 | 1 | 8 | 14 | Y | 0 | 1.35 | 5.47 | Double SCP |

| NAME_01069 | SCP | 3816 | 11201 | 1 | 1 | 10 | 12 | Y | 0 | 1.91 | 4.63 | Double SCP |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| NAME_12806 | SCP | 8094 | 8066 | 1 | 1 | 6 | 8 | Y | 0 | 1.62 | 4.42 | Single SCP |
| NAME_02548 | Globin | 887 | 1473 | 1 | 1 | 4 | 6 | N | 0 | 3.01 | 3.03 | Single Globin |
| NAME_05396 | SCP | 1339 | 5930 | 1 | 1 | 4 | 8 | Y | 0 | 1.48 | 4.13 | Single SCP |
| NAME_13850 | SCP | 3510 | 4006 | 1 | 0.9933 | 7 | 8 | Y | 0 | 1.97 | 3.07 | Single SCP |

973
974
975

976

977

978

979    **Table legends**

980    **Table 1:** Summary of the updated genome annotation of *Necator americanus*. CDS – coding

981    DNA sequence.

982

983    **Table 2:** Top 30 most abundant proteins identified with in-gel and OFF-GEL fractionation and

984    ranked using summed exponential modified protein abundance index (emPAI) score. Blast2GO

985    was used to obtain descriptions of each protein. Abbreviations key: NP – number of significant

986    peptides; SP – signal peptide; TD – transmembrane domain; Mascot scores >28 indicate identity

987    or extensive homology (p < 0.05).

988

989     **Figure legends**

990     **Figure 1**: Depiction of the proteogenomic process as well as the types and numbers of peptide

991     corrections identified.

992

993     **Figure 2: (A)** UpSetR plot displaying the number of each gene ontology (GO) term categories

994     (biological process, molecular function and/or cellular component) available (Blast2GO) for

995     adult *Necator americanus* excretory/secretory (ES) proteins. Proteins with no available GO terms

996     are broken down in a pie chart into 'SCP/TAPS' proteins 'other'. **(B)** Top 10 most abundant

997     protein families in the ES products of adult *N. americanus*. **(C)** Biological processes of adult *N.*

998     *americanus* ES proteins ranked by nodescore (Blast2GO) and plotted using REViGO.

999     Semantically similar GO terms plot close together, increasing heatmap score signifies increasing

1000    nodescore from Blast2GO, while circle size denotes the frequency of the GO term from the

1001    underlying database.  **(D)** Molecular functions of adult *N. americanus* ES proteins ranked by

1002    nodescore (Blast2GO) plotted using REViGO.

1003

1004    **Figure 3**: ES products from *Ancylostoma caninum*, *Heligmosomoides polygyrus, and*

1005    *Nippostrongylus brasiliensis* were compared with the excretory/secretory proteome of *Necator*

1006    *americanus* and displayed in a Simitri plot. SCP/TAPS are represented by an orange square,

1007    proteases by a teal-colored triangle and any other protein by a grey circle. Points on the diagram

1008    triangle represent sequences which only had similarity to the labelled species. Points along the

1009    edges of the triangle are sequences which had similarity to two of the three species (given at the

1010    respective ends of the edge). Any sequence in the middle area of the triangle represents a sequence

1011    with similarity to all three compared species.

1012

1013    **Figure 4:** SCP/TAPS proteins in the excretory/secretory products of *Necator americanus* are

1014    most closely related to SCP/TAPS proteins in the ES of *Ancylostoma caninum*. SCP/TAPS **(A)**

1015    and protease **(B)** protein names are displayed in a circle with *N. americanus* (purple), *A. caninum*

1016    (blue), *Heligmosomoides polygyrus* (pink), and *Nippostrongylus brasiliensis* (green). Ribbon

1017    thickness is relative to the maximum score obtained in the BLAST search while darker ribbons

1018    denote higher sequence percent identity. The corresponding bars provide relative sequence

1019    length of each protein. Respective protease mechanistic classes: aspartic (ASP), cysteine (CYS),

1020    metallo (MET) or serine (SER).

1021

1022    **Figure 5:** Phylogenetic relationships of **(A)** single-domain and **(B)** double-domain SCP/TAPS

1023    proteins determined with MUSCLE alignment software. PhyML was used for a maximum-

1024    likelihood phylogenetic analysis and results were visualized with The Interactive Tree of Life

1025    (iTOF) online phylogeny tool. *Necator americanus* sequences are highlighted in purple with

1026    comparator species each denoted by a different color (see key).

1027

1028

1029

Figure 1

MS/MS spectra → *N. americanus* database → Peptides

| | |
|---|---|
| Existing gene model | 5' UTR / Exon 1 / Exon 2 / Exon 3 / 3' UTR |
| 121 peptides mapped to intergenic regions | Novel protein-coding loci / Exon 1 |
| 39 peptides mapped to intronic regions | |
| 10 peptides mapped to exon-extensions | |
| 32 peptides mapped to alternative frames | |
| 4 peptides mapped to N-terminal extensions | |
| 12 peptides mapped to repeat regions | |

# Figure 2

Figure 3

# Figure 4A



Figure 4A. Circos plot showing relationships among *A. caninum*, *H. polygyrus*, *N. brasiliensis*, and *N. americanus*.

Figure 4B

A. caninum

H. polygyrus

N. brasiliensis

N. americanus

ASP
CYS
MET
SER

Hp_I08139_IG...
Hp_I14662_IG...
Hp_I14314_IG...
Hp_I15488_IG...
Hp_I12336_IG...
Hp_I12444_IG...
Hp_I13080_IG...
Hp_I09306_IG...
Hp_I13357_IG...

SER
ANCCAN_25722
ANCCAN_08123
ANCCAN_24137
ANCCAN_20473
ANCCAN_03692
ANCCAN_00393
ANCCAN_30567
ANCCAN_06649
ANCCAN_06644
ANCCAN_18339
ANCCAN_13546
ANCCAN_04557

MET
CYS
ASP

m_235636
m_276633
m_40513
m_263079
m_21970
m_55169
m_223522
m_265239
m_165518

ASP
CYS
MET
SER

NAME_10772
NAME_06735
NAME_13447
NAME_13146
NAME_02469
NAME_00886
NAME_00885
NAME_00535
NAME_10345
NAME_10342
NAME_05649
NAME_02362
NAME_02076
NAME_01240
NAME_09309
NAME_05268
NAME_05145
NAME_02055
NAME_00332
NAME_00331
NAME_00330
NAME_00320

SER
MET
CYS
ASP

# Figure 5A

Legend:
- *A. caninum*
- *A. suum*
- *H . polygyrus*
- *N. americanus*
- *N. brasiliensis*
- *T. canis*
- *T. muris*

# Figure 5B

*A. caninum*
*H . polygyrus*
*N. americanus*
*N. brasiliensis*
*T. canis*