1 **Annotating the 'hypothetical' in hypothetical proteins: *in-silico* analysis of**

2 **uncharacterised proteins for the Apicomplexan parasite, *Neospora***

3 ***caninum***

4

5 Larissa Calarco[1*], John Ellis[1]

6

7 [1] School of Life Sciences, University of Technology Sydney,

8 PO Box 123, Broadway,

9 NSW 2007, Australia

10 Larissa.M.Calarco@student.uts.edu.au

11 John.Ellis@uts.edu.au

12

13 *****Corresponding Author**: Larissa Calarco, University of Technology Sydney, School of Life

14 Sciences, Sydney, Australia, P: +61 2 9514 4161, F: +61 2 9514 8206,

15 E: Larissa.M.Calarco@student.uts.edu.au

16

17

18

19

20

21

22

23

24

25

## Abstract

*Neospora caninum* is a parasite of veterinary and economic importance, affecting beef and dairy cattle industries globally. While this species has been recognised as a serious cause of disease in cattle and dogs for over 30 years, treatment and control options are still not available. Furthermore, whilst vaccination was identified as the most economic control strategy, vaccine discovery programs require new leads to investigate as vaccines.

The current lack of gene annotation available for *N. caninum*, especially compared to the closely related model organism, *Toxoplasma gondii*, considerably hinders vaccine related research. Moreover, due to the high degree of similarity between the two organisms, a significant amount of gene annotation available for *N. caninum* stems from sequence homology between the species. However, there is a plethora of literature identifying conserved virulence factors between members of the Apicomplexa, which suggests that key players are contributing to successful parasite invasion, motility, and host cell attachment.

In this study, bioinformatic approaches classified 125 uncharacterised proteins within the *N. caninum* genome, as transmembrane proteins with signal peptide sequences. Functional annotation assigned enriched gene ontologies for cell-adhesion, ATP binding, protein serine/threonine phosphatase complex, immune system process, antigen binding, and proteolysis. Additionally, 32 of these proteins were also identified as adhesins, or having adhesin-like properties, which were further characterised through the discovery of domains and and gene ontology, to reveal their potential functional significance as virulence factors for *N. caninum*. This study identifies a new, small subset of proteins within *N. caninum*, that may be involved in host-cell interaction, parasite adhesion, and invasion, thereby implicating them as potential targets to exploit in the development of control options against the disease.

49  **Key words**: hypothetical protein, *in*-silico, annotation, virulence factors, adhesin,

50  transmembrane protein

51

52  ## Introduction

53  The Apicomplexa represent a phylum of diverse, ubiquitous, and successful parasites that are

54  responsible for a range of medical, economical, and veterinary diseases. The increasing

55  significance and relevance of this group of parasites has sparked a plethora of research

56  elucidating parasite biology, host cell interaction, and diversity between and within species.

57  Shared amongst Apicomplexans is the presence of specialised secretory organelles that

58  form part of the unique apical complex (Carruthers and Sibley, 1997; Gubbels and Duraisingh,

59  2012; English et al., 2015). The release of effector molecules from these secretory organelles

60  provides a catalyst for the execution of crucial processes, which promote parasite motility, host

61  cell attachment, and subsequent invasion (Carruthers and Sibley, 1997; Sibley, 2004; English

62  et al., 2015). Invasion begins with attachment to the host cell via the apical complex, resulting

63  in the organised secretion of proteins from rhoptries and adhesive micronemes (Sam-Yellowe,

64  1996; Carruthers and Sibley, 1997; Carruthers et al., 1999; Sibley, 2004). This is followed by

65  creation of the protective parasitophorous vacuole (PV), where subsequently the parasite is

66  able to grow, replicate, and disrupt host cell signalling and defence mechanisms (Sibley, 2004;

67  Plattner and Soldati-Favre, 2008; Luder et al., 2009; Pelle et al., 2015; Clough and Frickel,

68  2017).

69  A protein's structure determines its function, and membrane proteins are vital to a

70  plethora of cellular processes, including cellular attachment, invasion, molecule transport and

71  signalling, thereby representing a category of biologically significance proteins (Reynolds et

72  al., 2008). Conversely, proteins that are transported to secretory organelles generally contain

73  an N-terminal signal sequence (Chen et al., 2008), where the mechanisms for coordinated

74  parasite egress and invasion, rely on signal transduction (Gubbels and Duraisingh, 2012).

75  Effector molecules that function to facilitate parasite invasion and direct modulation of host

76  cell signalling in apicomplexans are constantly being identified, many of which contain such

77  important structural features.

78      For example, microneme (MICs), rhoptry (ROPs) and dense granule proteins (GRAs)

79  are classified as excretory/secretory antigens (ESA), representing a group of proteins

80  instrumental in parasite invasion, intracellular survival, and successful replication (Decoster et

81  al., 1988; Cesbron-Delauw and Capron, 1993; Cesbron-Delauw et al., 1996; Hoppe et al., 2000;

82  Nam, 2009; Sheiner et al., 2010). Many of these secreted proteins commonly possess a signal

83  peptide, and/or transmembrane domains, conducive to their function (Ngo et al., 2004; Nam,

84  2009; Sheiner et al., 2010; Cabrera et al., 2012; Huynh et al., 2014). The MIC family of proteins

85  can also be organised based on their adhesive motifs, which are predicted to mediate parasite

86  motility, invasion, and attachment (Sibley et al., 1998; Tomley and Soldati, 2001). These

87  commonly include epidermal growth factor (EGF), von Willebrand Factor A (vWF), and

88  thrombospondin type 1 (TSP-1) (Lawler and Hynes, 1986; Bork and Rohde, 1991; Tordai et

89  al., 1999; Tomley and Soldati, 2001; Chen et al., 2008).

90      *N. caninum* is a cyst forming protozoan parasite of veterinary and economic

91  importance, that affects beef and dairy cattle industries globally (Dubey, 1999, 2003). While

92  neosporosis as a disease has been recognised for over 30 years, the development of treatment

93  and control options is severely lacking, but becoming increasingly vital (Reichel and Ellis,

94  2002). The current extent of genome annotation for *N. caninum* however, presents a hindrance

95  to the crucial identification of key contributors to pathogenicity. Many proteins are termed

96  'hypothetical' or 'unnamed' due to either their unknown function, or lack of sequence

97  homology to recognised proteins (Galperin and Koonin, 2004). Furthermore, recent studies

98  focusing on improving and expanding the available gene structure and annotations for *N.*

99    *caninum* are yet to be integrated into popular online databases, such as NCBI or ToxoDB

100   (Gajria et al., 2008) reference resources.

101        While it is logical to assume that essential protein-coding genes implicated in parasite

102   virulence have been identified, the sheer number of unclassified or hypothetical regions cannot

103   be neglected or deemed unimportant, prompting this study. Current vaccine candidates for

104   parasites within this phylum involve either surface or secreted antigens that appear to be

105   fundamental to parasite invasion, the mechanisms and contributors of which appear to be

106   mostly conserved (Kim and Weiss, 2004; Hemphill, 2015). Consequently, the identification

107   and investigation of uncharacterised proteins through sequence homology, structure, and

108   known hallmarks of parasite virulence, has the power to perpetuate the discovery of targets for

109   vaccine development. This study aimed to exploit bioinformatic techniques to identify

110   previously uncharacterised proteins of biological and functional significance, based on protein

111   sequence topology, structure, and discerning features.

112

113

114

115

116

117

118

119

120

121

122

123

124

125

126

127

128

129

## Methods

A range of tools were used for the functional annotation of hypothetical proteins in this study,

are detailed in Table 1.

**PLACE TABLE 1 HERE**

**Sequence retrieval**

The Entrez GeneIDs of proteins classified as 'hypothetical' or 'unnamed' (collectively referred

to as 'uncharacterised' from here on in), were extracted from NCBI (GenBank assembly

accession                                                                 #GCA_000208865.2:

https://www.ncbi.nlm.nih.gov/genome/proteins/248?genome_assembly_id=28617),        and

uploaded to UniProtKB. The sequences for all uncharacterised proteins that also had no gene

names or annotations in UniProtKB were then extracted in FASTA format.

**Protein topology prediction and annotation**

The final list of uncharacterised protein sequences was submitted to the Philius Prediction

Server     for     individual     classification     by     protein     type

(http://www.yeastrc.org/philius/runPhilius.do) (Reynolds et al., 2008). This software

categorises proteins as globular (G), globular with signal peptides (G+SP), transmembrane

149 (TM), or transmembrane with signal peptides (TM+SP). Sequences were subsequently

150 obtained for proteins classified as TM+SP in FASTA format, and submitted to Blast2GO

151 (version 5), for functional annotation (Conesa et al., 2005). The gene ontology (GO) annotation

152 workflow available in Blast2GO incorporates BLAST analysis, GO, and InterProScan

153 (https://www.ebi.ac.uk/interpro/) (Quevillon et al., 2005; Finn et al., 2017).

154 Integrating the InterProScan database allowed identification of homologous

155 superfamilies, domains, and repeats present within each query protein sequence. It also

156 incorporates the transmembrane topology predictor Phobius (Kall et al., 2004), and signal

157 peptide predictor SignalP (Petersen et al., 2011). This allowed confirmation of protein

158 sequence classification by Philius to be corroborated by these tools. Proteins that were not

159 identified as TM+SP by at least two of these tools were discarded.

160 In an attempt to further assign biological function to the remaining unannotated proteins

161 in this list, the protein sequences were uploaded to the SECLAF webserver

162 (https://pitgroup.org/seclaf/), to identify enriched or over represented gene ontologies in this

163 protein callset. This server uses deep neural networks for the hierarchical classification of

164 biological sequences (Szalkai and Grolmusz, 2018b, a).

165

166 **Identification and annotation of adhesion-like proteins**

167 All TM+SP uncharacterised protein sequences were analysed by MAAP, a malarial adhesins

168 proteins predictor (http://maap.igib.res.in/) (Ansari et al., 2008). This predictor is based on

169 Support Vector Machines, where a default threshold of $P_{maap} = 0$ was used, characterising any

170 protein sequences above this threshold as adhesin or adhesin-like. The identified adhesin

171 proteins were cross-referenced with their predicted Philius protein classification, resulting in a

172 final list of uncharacterised proteins, identified as adhesin or adhesin-like transmembrane

173    proteins, containing signal peptides. The original Blast2GO results were retrieved for proteins

174    in this callset for further analysis. The bioinformatics workflow is summarised in Figure 1.

175

176    **PLACE FIGURE 1 HERE**

177

178         The list of TM+SP proteins were also uploaded to the ExPASy PROSITE database of

179    protein domains, families, and functional sites (de Castro et al., 2006; Sigrist et al., 2013). This

180    involved identifying sequence patterns, sites, and profiles, and also calculating the amino acid

181    composition of each sequence. The protein browser available in ToxoDB, PBrowse, was

182    subsequently used to identify any orthologous sites across each protein, through BLASTP.

183    Lastly, the 32 proteins were searched against the Database of Essential Genes (DEG;

184    http://www.essentialgene.org/) using default BLAST parameters, which consolidates currently

185    available genomic elements considered essential and indispensable for the survival of an

186    organism.

187

188    **Evidence for expression of proteins using RNA-seq data**

189    To obtain experimental evidence supporting the expression of the proteins in the final callset,

190    the *de novo* transcriptome assembled using RNA-seq data generated from NC-Liverpool

191    tachyzoites as per a previous study (Calarco et al., 2018) was exploited. Each protein sequence

192    was subjected to a BLAST analysis against the NC-Liverpool transcriptome, using the

193    command-line NCBI BLAST tool (version 2.7.1), where the most confident transcriptome

194    contig hits (low e-value and high bit score) for each protein were retained. For any proteins not

195    returning a result, data integrated into ToxoDB from Reid *et al.* (2012), generated from the

196    transcriptomes of days three and four NC-Liverpool tachyzoites, was used to determine mRNA

197    expression levels.

198

**Sequence variation within the final protein callset**

199

200 Calarco *et al.* (2018) compared RNA-seq data from tachyzoites of the NC-Liverpool and NC-
201 Nowra isolates. By employing a variant analysis pipeline, sequence variants located within
202 functionally significant genes or regions that differed between the two isolates were identified
203 and reported. The 32 adhesin-like transmembrane proteins with signal peptides presented in
204 this study were investigated for SNPs and whether they were located in a genome hotspot.

205

206

207

208 # Results

209 **Topology prediction and annotation for all uncharacterised proteins**

210 There were 4008 "hypothetical" proteins, and 256 unnamed" proteins extracted from NCBI,
211 from a total of 6936 *N. caninum in-silico* predicted proteins. These proteins are listed in
212 Supplementary File S1, detailing their chromosome location, GeneIDs, locus tags, and lengths.
213 Once the GeneIDs were uploaded to UniProt for sequence retrieval, 981 proteins were removed
214 as they were assigned predicted protein descriptions and annotations based on the data
215 available in UniProt (Supplementary File S2). This includes annotations assigned based on
216 sequence similarity, or experimental evidence at the protein and transcript level. The details of
217 the remaining proteins, whose sequences were retrieved from UniProt in FASTA format, are
218 provided in Supplementary File S3. This process is summarised in Supplementary File S4.

219 Philius identified more than half of the uncharacterised proteins as globular, with no
220 transmembrane domains or signal peptide sequences (Figure 2). There were however 147
221 proteins predicted to be TM+SP proteins by Philius.

222

223    **PLACE FIGURE 2 HERE**

224

225    Of the 147 TM+SP proteins identified by Philius, the topologies of 20 were not

226    corroborated by either Phobius or SignalP following Blast2GO analysis. There were also an

227    additional two proteins with no topology features predicted by Phobius or SignalP (F0V8W3

228    and F0VLJ1). All of these proteins, which also had relatively low confidence scores in Philius,

229    were discarded from functional annotation. The Blast2GO results for all 147 TM+SP proteins

230    are provided in Supplementary File S5.

231    The Blast2GO results identified enriched gene ontologies, featured protein families,

232    and domains occurring across the remaining 125 TM+SP proteins under investigation. This

233    workflow also aimed to assign gene descriptions to each protein based on sequence similarity

234    to closely related organisms. Of the 125 TM+SP proteins, 43 of these were assigned sequence

235    descriptions (Supplementary File S6). The remaining proteins however, were still only

236    classified as 'putative transmembrane protein' or 'hypothetical protein'.

237    The homologous protein superfamilies represented multiple times in this callset of

238    TM+SP proteins, included growth factor receptor cysteine-rich domain superfamily

239    (IPR009030), which includes proteins involved in signal transduction by receptor tyrosine

240    kinases (Ward et al., 1995; Garrett et al., 1998; Cho and Leahy, 2002), major facilitator

241    transporter superfamily (IPR036259), consisting of membrane transport proteins (Pao et al.,

242    1998; Walmsley et al., 1998), and vWF A-like domain superfamily (IPR036465), where such

243    proteins participate in cell adhesion, signal transduction, membrane transport, and immune

244    defence mechanisms (Colombatti et al., 1993).

245    The main GOs related to molecular function included nucleic acid binding

246    (GO:0003676), DNA binding (GO:0003677), ATP binding (GO:0005524), and serine-type

247    endopeptidase activity (GO:0004252). Conversely, the most represented GOs pertaining to

248     biological function were proteolysis (GO:0006508) and regulation of apoptotic process

249     (GO:0042981). As expected, most of the protein sequences were assigned cellular component

250     GOs for 'integral component of membrane' (GO:0016021) and 'membrane' (GO:0016020).

251     While the Blast2GO analysis returned only minimal GOs for all 125 proteins, the

252     SECLAF webserver provided a more thorough and extensive list of enriched gene ontology

253     protein function prediction. The most represented GOs that were associated with almost all

254     proteins in this callset, included cell junction (GO:0030054), protein serine/threonine

255     phosphatase complex (GO:0008287), cell tip of elongated cells (GO:0051286), and binding

256     (GO:000584). Other GOs of functional interest associated with many of these TM+SP proteins

257     included signal transduction (GO:0007165), immune system process (GO:0002376),

258     anchoring junction (GO:0070161), adhesion of symbiont to host (GO:0044406), and

259     interaction with symbiont (GO:0051702).

260

261     **Prediction of adhesin-like proteins and their classification**

262     Of the 3283 uncharacterised proteins investigated, 654 (20%) were identified as having adhesin

263     properties by MAAP (Supplementary File S7). Supplementary File S8 contains a small subset

264     of these proteins that were investigated through InterProScan sequence analysis, to justify the

265     applicability and efficacy of this malarial adhesins predictor for *N. caninum*.

266     Figure 3 is a pie chart presenting the percentage of adhesin proteins, and their predicted

267     protein classification according to Philius. A total of 32 uncharacterised proteins (~1%) were

268     identified as adhesin-like transmembrane proteins, with signal peptides.

269

270     **PLACE FIGURE 3 HERE**

271

272     **Annotation of adhesin TM+SP proteins**

273     The Blast2GO analysis assigned gene descriptions for 20 proteins, based on sequence

274     similarity. This included proteins identified as MIC2 (F0VIM1), subtilisin SUB2 (F0VNN6),

275     septin (F0VML7), and *T. gondii* family A protein. There were however 12 proteins that

276     remained described as 'hypothetical' or simply 'putative transmembrane protein', due to a lack

277     of sequence homology with related species. Additionally, two hypothetical proteins that were

278     not assigned descriptions, had between 34-38% identity with *T. gondii* GRA11 (F0V9X3 and

279     F0V9Z2). The featured domains identified within this final protein callset included vWF type

280     A domain, (IPR002035), CARD or caspase recruitment domain (IPR001315), subtilisin SUB1-

281     like catalytic domain (IPR034204), and peptidase S8 domain (IPR036852).

282     The only represented GOs from Bast2GO included 'serine-type endopeptidase activity'

283     (GO:0004252) and 'protein binding' (GO:005515) for molecular function, and 'proteolysis'

284     (GO:0006508) and 'regulation of apoptotic process' (GO:0042981) for biological process.

285     Again however, the SECLAF webserver assigned further functionally relevant GOs to the 32

286     proteins, where those over-represented included locomotion (GO:0040011), cell adhesion

287     (GO:0098602), antigen binding (GO:0003823), cofactor transmembrane transporter activity

288     (GO:0051184), and structural molecule activity (GO:0005198).

289     The 32 adhesin-like transmembrane proteins with signal peptides identified in this

290     study are listed in Table 2, along with their Blast2GO descriptions, gene ontologies, and

291     InterProScan features.

292

293     **PLACE TABLE 2 HERE**

294

295     Represented across the 32 adhesin proteins with transmembrane domains and signal

296     peptides, were serine-rich (PS50324) and alanine-rich regions (PS50310), as determined by

297     PROSITE. Further to this, after calculating the amino acid composition for each protein, serine

298  was the most abundant amino acid in 14 of the 32 protein sequences, followed by alanine in 12

299  protein sequences. In relation to sequence similarity, four proteins contained regions with

300  BLASTP hits to gel-forming secreted mucin-19 from mice. An additional four proteins had

301  sequence similarity to serine-rich adhesin for platelets segments, with BLAST hits to various

302  *Staphyloccocus* and *Streptococcus* species. Other notable hits included adhesin-like cell wall

303  proteins from *Candida albicans* (23% PID), and endochitenase from *Aspergillus fumigatus*

304  (24-32% PID). The amino acid composition and BLASTP results for each of the 32 proteins

305  are presented in Supplementary File S9.

306      Three proteins in the final callset returned hits to genes in the Database of Essential

307  Genes. Protein F0VIM1 (MIC2) returned BLAST hits to thrombospondin, integrin subunit

308  alpha  1,    collagen  alpha  1,  and  ADAM  (disintegrin  and  metalloproteinase)

309  metalloendopeptidase genes in both humans and mice. Genes returned for protein sequence

310  F0VNN6 (SUB2) included membrane-bound transcription factor peptidase and proprotein

311  covertase subtilisin/kexin genes from human and mice, and protein F0VM28 (vWF type A

312  domain containing protein), aligned to huge dynein-related AAA-type ATPase (midasin) from

313  *Saccharomyces cerevisiae*, mediating ATP-dependent remodelling of 60S subunits and

314  subsequent export from nucleoplasm to cytoplasm.

315

316  **Evidence of protein expression provided by RNA-seq data**

317  All but three of the final 32 proteins (F0VIG7, F0VEH5, and F0VK21) had high confidence

318  BLAST hits to contigs in the NC-Liverpool transcriptome published in Calarco *et al.* (2018),

319  with percentage identities (PID) > 80%. Additionally, some of the proteins had BLAST hits to

320  the same NC-Liverpool transcriptome contig, indicating that they may be paralogous genes

321  within *N. caninum*. Of the three remaining aforementioned proteins, transcript expression was

322  recorded for each of these based on RNA-seq data generated from either day 3 and 4

323    tachyzoites, published by Reid *et al*. (2012). Supplementary File S10 contains a list of the final

324    32 proteins, along with the recorded FPKM (Fragments Per Kilobase Million) and percentiles

325    from Reid *et al.* (2012), based on the transcriptomes of days three and four tachyzoites.

326

327    **Sequence variation within 32 adhesin-like TM+SP proteins**

328    Calarco *et al.* (2018) identified subtilisin SUB2 protease as present in a SNP hotspot, based on

329    the number of sequence variants identified within the gene sequence, when comparing the NC-

330    Liverpool and NC-Nowra isolates. Protein F0VNN6 in this study was annotated as SUB2, and

331    predicted to contain adhesin-like properties, a signal peptide, and transmembrane domains. The

332    only other protein in this final callset which was previously found to contain SNPs, was

333    F0VNG1.

334

335

336

337

338    **Discussion**

339    Identifying and characterising key players of the parasite invasion process, and elucidating how

340    they could represent treatment, control, and vaccine targets is an important step for any vaccine

341    discovery program. Host-modulating effectors currently of interest include parasite surface

342    antigens, and proteins secreted from the unique Apicomplexan secretory organelles: the

343    rhoptries, micronemes, and dense granules (Carruthers and Sibley, 1997; Sibley, 2004;

344    Gubbels and Duraisingh, 2012). To exploit the current understanding of parasite virulence in

345    the context of *N. caninum*, this study employed various bioinformatic tools to identify a small

346    subset of biologically important proteins, potentially associated with parasite adhesion,

347    invasion, and host cell interactions. The reasoning behind this approach stems from the

348 disturbing lack of genomic annotation available for *N. caninum*, and the sizeable existence of

349 unidentified, theoretically important proteins that await characterisation.

350 With a focus on non-model organisms lacking sequence annotation for many predicted

351 proteins, there is currently a plethora of research dedicated to the description of hypothetical

352 proteins, through *in-silico* analysis of available sequencing data. This can subsequently result

353 in the identification of functionally significant proteins involved in essential processes,

354 pertinent to the organism under investigation. For example, the *in-silico* analysis of

355 hypothetical proteins in the *Plasmodium falciparum* proteome by Oladele *et al.* (2011), resulted

356 in the classification of several sequences as potential biomarkers of malaria. Another study

357 exploited bioinformatics tools to identify potential new drug targets, from a set of hypothetical

358 proteins classified in a previous immunoproteomics study for *Leishmania* spp. (Chavez-

359 Fumagalli et al., 2017). The *in-silico* workflow elucidated the cellular localisation, biological

360 function, and structure of these proteins, thereby presenting a method for the functional

361 annotation and elucidation of potential drug candidates against Leishmaniasis. In *Trypanosoma*

362 *cruzi*, all proteins with predicted transmembrane regions were computationally analysed for

363 potential biological function (Silber and Pereira, 2012). A total of 54 proteins were found to be

364 involved in signal-transduction processes through sequence annotation, which again could

365 represent putative drug targets. Lastly, a comprehensive bioinformatics study on the

366 hypothetical protein dataset for *Leishmania donovani*, assigned putative functions, GO terms,

367 and protein domains to a previously uncharacterised set of proteins (Ravooru et al., 2014). The

368 association of these proteins to specific biological pathways and classification as essential

369 genes, demonstrated the advantage of robust computational strategies for the identification of

370 molecules as potential therapeutic targets against such diseases.

371 The fact that 4264 of 6936 genes in the published *N. caninum* genome are

372 uncharacterised and described as 'hypothetical' proteins, presents a major and concerning

373  hindrance to the study of potential virulence factors. Compounding this problem is the lack of

374  consistency and consensus between popular online databases containing genomic data and

375  gene annotation. For example during sequence retrieval, a total of 981 of these proteins had

376  annotation information in UniProt that was not present in NCBI or integrated into ToxoDB.

377  These included functionally important proteins such as apical membrane antigen AMA1

378  (NCLIV_065490), rhoptry proteins including ROP5-ROP8, rhoptry neck protein RON2

379  (NCLIV_064620), MIC8 (NCLIV_062770), and multiple GRA proteins including GRA6,

380  GRA7, GRA10, and GRA14.

381      Based on sequence similarity, Blast2GO assigned descriptions for 43 of the 125

382  predicted TM+SP proteins (Supplementary Files S5 and S6). Protein F0VQ63 was described

383  as rhomboid-like protease ROM6, belonging to a large family of intramembrane-cleaving

384  serine proteases that are ubiquitous in almost all organisms (Urban and Dickey, 2011). In

385  Apicomplexans, rhomboid proteases are involved in the shedding of adhesins from the cell

386  surface during parasite motility and host-cell invasion, and hence play an important role in

387  host-parasite interactions (Santos et al., 2012; Sibley, 2013). Another TM+SP protein was

388  described as a cytoadherence-linked asexual protein (Clag; F0V7G1), which is thought to be

389  essential for the adhesion and survival of *P. falciparum* in vivo and is regarded as a major

390  determinant of the parasite's virulence (Ocampo et al., 2005). Additionally, protein F0VPV9

391  was annotated as lectin C-type domain protein, which are integral membrane proteins that have

392  been shown to play a role in the recognition of glycosylated parasite antigens (Vazquez-

393  Mendoza et al., 2013). These proteins have been implicated in processes such as cell adhesion,

394  platelet activation, and pathogen recognition in various pathogenic organisms (Weis et al.,

395  1998; Kilpatrick, 2002; Kerrigan and Brown, 2009).

396      The processes of parasite invasion are facilitated by organised, sequential protein

397  secretion from specialised apical organelles, to release adhesins for cell attachment and protein

398    transport to the PV membrane (Carruthers and Sibley, 1997; Bradley and Sibley, 2007). Studies

399    have implicated various apicomplexan surface proteins in host cell recognition, where such

400    proteins can be identified by the presence of conserved domains found across a wide range of

401    organisms (Templeton et al., 2004). This class of proteins usually contains adhesion domains,

402    the structural patterns of which can be exploited by an adhesin predictor such as MAAP. An

403    assessment of MAAP indicated it was applicable to the dataset from *N. caninum*, based on the

404    sequence annotation of proteins classified as adhesins in this study (Supplementary File S8).

405    Many of the proteins contained adhesion domains, implicating them in cell adhesion and host

406    cell recognition. Additionally, enriched GOs assigned to the adhesin-like proteins

407    characterising the final callset, such as 'cell adhesion', 'cell-cell adhesion', and 'antigen

408    binding', provided further reassurance and confidence in the use of this tool for *N. caninum*

409    proteins. This was also supported by the BLASTP results, where some of these protein

410    sequences contained 'serine-rich adhesin for platelets' segments present in bacterial species,

411    or 'adhesin-like cell wall protein' segments found in some fungi. Further investigation of the

412    654 adhesin-like proteins identified in this study, may reveal further key players involved in

413    crucial parasite adhesion and invasion mechanisms conducive to their success.

414         While the annotation of many *N. caninum* proteins remains incomplete or insufficient,

415    it is expected that important information can be gained through sequence homology searches

416    with closely related species, especially those part of the Apicomplexa phylum. Arguably, one

417    of the most significant proteins identified and described in the final callset through sequence

418    similarity, was MIC2. Huynh and Carruthers (2006) demonstrated that reduced MIC2

419    expression led to ineffective host-cell attachment and parasite invasion, as well as reduced

420    gliding motility. This implicated the MIC2 complex as a major determinant of virulence in

421    *Toxoplasma* infection, and identified the potential for MIC2-deficient parasites as an effective

422    live attenuated vaccine against the disease. However, although MIC2 was previously described

423    for *N. caninum* by Lovett *et al.* (2000), it still remains annotated as a hypothetical protein in

424    the NCBI, UniProt, and ToxoDB reference databases.

425          Another protein described in the final callset was SUB2 (F0VNN6). The success of

426    host cell invasion by Apicomplexan species is contingent on the secretion of proteins from

427    specialised apical organelles (Carruthers et al., 1999). Much of the research concerning these

428    important secretory organelles and their protein contents implicates proteolytic processing as

429    central to the maturation of these crucial proteins (Sam-Yellowe, 1996; Miller et al., 2003).

430    Studies have shown that serine proteinase inhibitors obstruct host cell invasion, implicating

431    subtilisin-like serine proteinases as biologically important in Apicomplexans (Conseil et al.,

432    1999; Blackman, 2000; Miller et al., 2003). The MIC2 and SUB2 proteins identified here also

433    returned BLAST hits to protein-coding genes within the eukaryotic Database of Essential

434    Genes, further cementing their functional significance within the *N. caninum* proteome. The

435    annotation of SUB2 in this study, as well as the previous identification of the SUB2 gene as a

436    SNP hotspot (Calarco et al., 2018), suggests that this protein could represent a potential

437    virulence factor of *N. caninum* that warrants future investigation.

438          What was still surprising despite the efforts and measures taken in this study, was the

439    lack of descriptions and sequence features present for 12 of the 32 final proteins, such as

440    domains, repeats, and gene ontology. However, the annotation workflow implemented in this

441    study identified functionally significant proteins such as MIC2 and SUB2, within the *N.*

442    *caninum* proteome. This therefore suggests that the remaining proteins identified represent

443    previously uncharacterised, but biologically important proteins, based on their sequence

444    topology and predicted adhesin-like properties. Such proteins may potentially be involved in

445    adhesion, invasion, and secretion processes that are responsible for the success of such parasite

446    species, despite the lack of sequence features assigned.

447    Most of the adhesin TM+SP proteins were rich in serine, alanine, and threonine

448    (Supplementary File S9), which likely reflects the training dataset incorporated into MAAP.

449    Furthermore, based on the integrated protein browser in ToxoDB (PBrowse), many of these

450    protein sequences were found to contain mucin-like segments, identified through BLAST

451    analysis (Supplementary File S9). In *Cryptosporidium parvum* for example, there are more

452    than 30 unique mucin-like surface proteins, which are also characterised by serine- and

453    threonine- rich repeats in their extracellular regions, and hence proposed to facilitate adhesion

454    between the parasite and host cell surface (Barnes et al., 1998; Ward and Cevallos, 1998;

455    Cevallos et al., 2000a; Cevallos et al., 2000b; Winter et al., 2000; Templeton et al., 2004).

456    These *C. parvum* mucins, which include gp900 and gp40/gp15, are described as highly

457    immunogenic, and hence potentially important vaccine candidates (Barnes et al., 1998;

458    Cevallos et al., 2000b; O'Connor et al., 2007; O'Connor et al., 2009; Chatterjee et al., 2010).

459    This class of mucin-like proteins is also shared with *T. gondii*, however these proteins are not

460    present in *Plasmodium* and *Theileria* species, and therefore may represent adaptations of

461    coccidians to harsh environments, and immune system evasion mechanisms (Templeton et al.,

462    2004). As the antigens shown thus far to be crucial for the attachment and invasion of

463    *Cryptosporidium* species into host cells, are all mucin-like glycoproteins (O'Connor et al.,

464    2009), the identification of similar proteins in this study hence bolsters their potential

465    significance as part of the *N. caninum* proteome.

466

467

468    **Conclusion**

469    This study characterised previously unannotated proteins of the *N. caninum* proteome, using

470    *in-silico* tools. The workflow implemented resulted in the identification of 125 proteins with

471    predicted transmembrane domains and signal peptide sequences. Further analysis of these

472 proteins classified 32 as having adhesin features, which suggests they may be part of crucial

473 parasite mechanisms of cell invasion, adhesion, and motility. Such processes are conserved in

474 the Apicomplexan phylum, the key contributors of which may represent virulence factors to

475 target in the development of therapeutic drugs or vaccines against the disease.

476       The relevance and value of the bioinformatics tools exploited in this study, was

477 supported by the biologically significant annotations collated. Enriched gene ontologies for the

478 prioritised proteins included proteolysis, cell adhesion, protein serine/threonine phosphatase

479 complex, and integral component of membrane. The *in-silico* approach described is especially

480 useful for non-model organisms or those in the early stages of genomic and proteomic

481 exploration, which may lack sufficient or robust characterisation of functionally significant

482 proteins.

483

## Acknowledgments

486

## Competing interests

488 The authors declare they have no competing interests.

489

## Funding

## References

494 Ansari, F.A., Kumar, N., Bala Subramanyam, M., Gnanamani, M., Ramachandran, S., 2008.
495     MAAP: malarial adhesins and adhesin-like proteins predictor. Proteins 70, 659-666.

Barnes, D.A., Bonnin, A., Huang, J.X., Gousset, L., Wu, J., Gut, J., Doyle, P., Dubremetz, J.F., Ward, H., Petersen, C., 1998. A novel multi-domain mucin-like glycoprotein of Cryptosporidium parvum mediates invasion. Mol Biochem Parasitol 96, 93-110.

Blackman, M.J., 2000. Proteases involved in erythrocyte invasion by the malaria parasite: function and potential as chemotherapeutic targets. Curr Drug Targets 1, 59-83.

Bork, P., Rohde, K., 1991. More von Willebrand factor type A domains? Sequence similarities with malaria thrombospondin-related anonymous protein, dihydropyridine-sensitive calcium channel and inter-alpha-trypsin inhibitor. Biochem J 279 ( Pt 3), 908-910.

Bradley, P.J., Sibley, L.D., 2007. Rhoptries: an arsenal of secreted virulence factors. Current opinion in microbiology 10, 582-587.

Cabrera, A., Herrmann, S., Warszta, D., Santos, J.M., John Peter, A.T., Kono, M., Debrouver, S., Jacobs, T., Spielmann, T., Ungermann, C., Soldati-Favre, D., Gilberger, T.W., 2012. Dissection of minimal sequence requirements for rhoptry membrane targeting in the malaria parasite. Traffic 13, 1335-1350.

Calarco, L., Barratt, J., Ellis, J., 2018. Genome Wide Identification of Mutational Hotspots in the Apicomplexan Parasite Neospora caninum and the Implications for Virulence. Genome Biol Evol.

Carruthers, V.B., Giddings, O.K., Sibley, L.D., 1999. Secretion of micronemal proteins is associated with toxoplasma invasion of host cells. Cell Microbiol 1, 225-235.

Carruthers, V.B., Sibley, L.D., 1997. Sequential protein secretion from three distinct organelles of Toxoplasma gondii accompanies invasion of human fibroblasts. Eur J Cell Biol 73, 114-123.

Cesbron-Delauw, M.F., Capron, A., 1993. Excreted/secreted antigens of Toxoplasma gondii-- their origin and role in the host-parasite interaction. Res Immunol 144, 41-44.

Cesbron-Delauw, M.F., Lecordier, L., Mercier, C., 1996. Role of secretory dense granule organelles in the pathogenesis of toxoplasmosis. Curr Top Microbiol Immunol 219, 59-65.

Cevallos, A.M., Bhat, N., Verdon, R., Hamer, D.H., Stein, B., Tzipori, S., Pereira, M.E., Keusch, G.T., Ward, H.D., 2000a. Mediation of Cryptosporidium parvum infection in vitro by mucin-like glycoproteins defined by a neutralizing monoclonal antibody. Infection and immunity 68, 5167-5175.

Cevallos, A.M., Zhang, X., Waldor, M.K., Jaison, S., Zhou, X., Tzipori, S., Neutra, M.R., Ward, H.D., 2000b. Molecular cloning and expression of a gene encoding Cryptosporidium parvum glycoproteins gp40 and gp15. Infection and immunity 68, 4108-4116.

Chatterjee, A., Banerjee, S., Steffen, M., O'Connor, R.M., Ward, H.D., Robbins, P.W., Samuelson, J., 2010. Evidence for mucin-like glycoproteins that tether sporozoites of Cryptosporidium parvum to the inner surface of the oocyst wall. Eukaryot Cell 9, 84-96.

Chavez-Fumagalli, M.A., Schneider, M.S., Lage, D.P., Machado-de-Avila, R.A., Coelho, E.A., 2017. An in silico functional annotation and screening of potential drug targets derived from Leishmania spp. hypothetical proteins identified by immunoproteomics. Exp Parasitol 176, 66-74.

Chen, Z., Harb, O.S., Roos, D.S., 2008. In silico identification of specialized secretory-organelle proteins in apicomplexan parasites and in vivo validation in Toxoplasma gondii. PLoS One 3, e3611.

Cho, H.S., Leahy, D.J., 2002. Structure of the extracellular region of HER3 reveals an interdomain tether. Science 297, 1330-1333.

Clough, B., Frickel, E.M., 2017. The Toxoplasma Parasitophorous Vacuole: An Evolving Host-Parasite Frontier. Trends Parasitol 33, 473-488.

Colombatti, A., Bonaldo, P., Doliana, R., 1993. Type A modules: interacting domains found in several non-fibrillar collagens and in other extracellular matrix proteins. Matrix 13, 297-306.

Conesa, A., Gotz, S., Garcia-Gomez, J.M., Terol, J., Talon, M., Robles, M., 2005. Blast2GO: a universal tool for annotation, visualization and analysis in functional genomics research. Bioinformatics 21, 3674-3676.

Conseil, V., Soete, M., Dubremetz, J.F., 1999. Serine protease inhibitors block invasion of host cells by Toxoplasma gondii. Antimicrob Agents Chemother 43, 1358-1361.

de Castro, E., Sigrist, C.J., Gattiker, A., Bulliard, V., Langendijk-Genevaux, P.S., Gasteiger, E., Bairoch, A., Hulo, N., 2006. ScanProsite: detection of PROSITE signature matches and ProRule-associated functional and structural residues in proteins. Nucleic Acids Res 34, W362-365.

Decoster, A., Darcy, F., Capron, A., 1988. Recognition of Toxoplasma gondii excreted and secreted antigens by human sera from acquired and congenital toxoplasmosis: identification of markers of acute and chronic infection. Clin Exp Immunol 73, 376-382.

Dubey, J.P., 1999. Neosporosis in cattle: biology and economic impact. Journal of the American Veterinary Medical Association 214, 1160-1163.

Dubey, J.P., 2003. Review of Neospora caninum and neosporosis in animals. Korean J Parasitol 41, 1-16.

English, E.D., Adomako-Ankomah, Y., Boyle, J.P., 2015. Secreted effectors in Toxoplasma gondii and related species: determinants of host range and pathogenesis? Parasite Immunol 37, 127-140.

Finn, R.D., Attwood, T.K., Babbitt, P.C., Bateman, A., Bork, P., Bridge, A.J., Chang, H.Y., Dosztanyi, Z., El-Gebali, S., Fraser, M., Gough, J., Haft, D., Holliday, G.L., Huang, H., Huang, X., Letunic, I., Lopez, R., Lu, S., Marchler-Bauer, A., Mi, H., Mistry, J., Natale, D.A., Necci, M., Nuka, G., Orengo, C.A., Park, Y., Pesseat, S., Piovesan, D., Potter, S.C., Rawlings, N.D., Redaschi, N., Richardson, L., Rivoire, C., Sangrador-Vegas, A., Sigrist, C., Sillitoe, I., Smithers, B., Squizzato, S., Sutton, G., Thanki, N., Thomas, P.D., Tosatto, S.C., Wu, C.H., Xenarios, I., Yeh, L.S., Young, S.Y., Mitchell, A.L., 2017. InterPro in 2017-beyond protein family and domain annotations. Nucleic Acids Res 45, D190-D199.

Gajria, B., Bahl, A., Brestelli, J., Dommer, J., Fischer, S., Gao, X., Heiges, M., Iodice, J., Kissinger, J.C., Mackey, A.J., Pinney, D.F., Roos, D.S., Stoeckert, C.J., Jr., Wang, H., Brunk, B.P., 2008. ToxoDB: an integrated Toxoplasma gondii database resource. Nucleic Acids Res 36, D553-556.

Galperin, M.Y., Koonin, E.V., 2004. 'Conserved hypothetical' proteins: prioritization of targets for experimental study. Nucleic Acids Res 32, 5452-5463.

Garrett, T.P., McKern, N.M., Lou, M., Frenkel, M.J., Bentley, J.D., Lovrecz, G.O., Elleman, T.C., Cosgrove, L.J., Ward, C.W., 1998. Crystal structure of the first three domains of the type-1 insulin-like growth factor receptor. Nature 394, 395-399.

Gubbels, M.J., Duraisingh, M.T., 2012. Evolution of apicomplexan secretory organelles. Int J Parasitol 42, 1071-1081.

Hemphill, A., 2015. Vaccines and drugs against Neospora caninum, an important apicomplexan causing abortion in cattle and other farm animals. Reports in Parasitology 2015:4, 31-41.

Hoppe, H.C., Ngo, H.M., Yang, M., Joiner, K.A., 2000. Targeting to rhoptry organelles of Toxoplasma gondii involves evolutionarily conserved mechanisms. Nat Cell Biol 2, 449-456.

Huynh, M.H., Boulanger, M.J., Carruthers, V.B., 2014. A conserved apicomplexan microneme protein contributes to Toxoplasma gondii invasion and virulence. Infection and immunity 82, 4358-4368.

Huynh, M.H., Carruthers, V.B., 2006. Toxoplasma MIC2 is a major determinant of invasion and virulence. PLoS Pathog 2, e84.

Kall, L., Krogh, A., Sonnhammer, E.L., 2004. A combined transmembrane topology and signal peptide prediction method. J Mol Biol 338, 1027-1036.

Kerrigan, A.M., Brown, G.D., 2009. C-type lectins and phagocytosis. Immunobiology 214, 562-575.

Kilpatrick, D.C., 2002. Animal lectins: a historical introduction and overview. Biochim Biophys Acta 1572, 187-197.

Kim, K., Weiss, L.M., 2004. Toxoplasma gondii: the model apicomplexan. Int J Parasitol 34, 423-432.

Lawler, J., Hynes, R.O., 1986. The structure of human thrombospondin, an adhesive glycoprotein with multiple calcium-binding sites and homologies with several different proteins. J Cell Biol 103, 1635-1648.

Lovett, J.L., Howe, D.K., Sibley, L.D., 2000. Molecular characterization of a thrombospondin-related anonymous protein homologue in Neospora caninum. Mol Biochem Parasitol 107, 33-43.

Luder, C.G., Stanway, R.R., Chaussepied, M., Langsley, G., Heussler, V.T., 2009. Intracellular survival of apicomplexan parasites and host cell modification. Int J Parasitol 39, 163-173.

Miller, S.A., Thathy, V., Ajioka, J.W., Blackman, M.J., Kim, K., 2003. TgSUB2 is a Toxoplasma gondii rhoptry organelle processing proteinase. Mol Microbiol 49, 883-894.

Nam, H.W., 2009. GRA proteins of Toxoplasma gondii: maintenance of host-parasite interactions across the parasitophorous vacuolar membrane. Korean J Parasitol 47 Suppl, S29-37.

Ngo, H.M., Yang, M., Joiner, K.A., 2004. Are rhoptries in Apicomplexan parasites secretory granules or secretory lysosomal granules? Mol Microbiol 52, 1531-1541.

O'Connor, R.M., Burns, P.B., Ha-Ngoc, T., Scarpato, K., Khan, W., Kang, G., Ward, H., 2009. Polymorphic mucin antigens CpMuc4 and CpMuc5 are integral to Cryptosporidium parvum infection in vitro. Eukaryot Cell 8, 461-469.

O'Connor, R.M., Wanyiri, J.W., Cevallos, A.M., Priest, J.W., Ward, H.D., 2007. Cryptosporidium parvum glycoprotein gp40 localizes to the sporozoite surface by association with gp15. Mol Biochem Parasitol 156, 80-83.

Ocampo, M., Rodriguez, L.E., Curtidor, H., Puentes, A., Vera, R., Valbuena, J.J., Lopez, R., Garcia, J.E., Ramirez, L.E., Torres, E., Cortes, J., Tovar, D., Lopez, Y., Patarroyo, M.A., Patarroyo, M.E., 2005. Identifying Plasmodium falciparum cytoadherence-linked asexual protein 3 (CLAG 3) sequences that specifically bind to C32 cells and erythrocytes. Protein Sci 14, 504-513.

Oladele, T.O., Sadiku, J.S., Bewaji, C.O., 2011. *In silico* characterisation of some hypothetical proteins in the proteome of *Plasmodium falciparum*. Centrepoint J. 17, 129-139.

Pao, S.S., Paulsen, I.T., Saier, M.H., Jr., 1998. Major facilitator superfamily. Microbiol Mol Biol Rev 62, 1-34.

642 Pelle, K.G., Jiang, R.H., Mantel, P.Y., Xiao, Y.P., Hjelmqvist, D., Gallego-Lopez, G.M., A,
643         O.T.L., Kang, B.H., Allred, D.R., Marti, M., 2015. Shared elements of host-targeting
644         pathways among apicomplexan parasites of differing lifestyles. Cell Microbiol 17,
645         1618-1639.
646 Petersen, T.N., Brunak, S., von Heijne, G., Nielsen, H., 2011. SignalP 4.0: discriminating
647         signal peptides from transmembrane regions. Nat Methods 8, 785-786.
648 Plattner, F., Soldati-Favre, D., 2008. Hijacking of host cellular functions by the
649         Apicomplexa. Annu Rev Microbiol 62, 471-487.
650 Quevillon, E., Silventoinen, V., Pillai, S., Harte, N., Mulder, N., Apweiler, R., Lopez, R.,
651         2005. InterProScan: protein domains identifier. Nucleic Acids Res 33, W116-120.
652 Ravooru, N., Ganji, S., Sathyanarayanan, N., Nagendra, H.G., 2014. Insilico analysis of
653         hypothetical proteins unveils putative metabolic pathways and essential genes in
654         Leishmania donovani. Front Genet 5, 291.
655 Reichel, M.P., Ellis, J.T., 2002. Control options for Neospora caninum infections in cattle--
656         current state of knowledge. New Zealand veterinary journal 50, 86-92.
657 Reid, A.J., Vermont, S.J., Cotton, J.A., Harris, D., Hill-Cawthorne, G.A., Konen-Waisman,
658         S., Latham, S.M., Mourier, T., Norton, R., Quail, M.A., Sanders, M., Shanmugam, D.,
659         Sohal, A., Wasmuth, J.D., Brunk, B., Grigg, M.E., Howard, J.C., Parkinson, J., Roos,
660         D.S., Trees, A.J., Berriman, M., Pain, A., Wastling, J.M., 2012. Comparative
661         genomics of the apicomplexan parasites Toxoplasma gondii and Neospora caninum:
662         Coccidia differing in host range and transmission strategy. PLoS Pathog 8, e1002567.
663 Reynolds, S.M., Kall, L., Riffle, M.E., Bilmes, J.A., Noble, W.S., 2008. Transmembrane
664         topology and signal peptide prediction using dynamic bayesian networks. PLoS
665         Comput Biol 4, e1000213.
666 Sam-Yellowe, T.Y., 1996. Rhoptry organelles of the apicomplexa: Their role in host cell
667         invasion and intracellular survival. Parasitol Today 12, 308-316.
668 Santos, J.M., Graindorge, A., Soldati-Favre, D., 2012. New insights into parasite rhomboid
669         proteases. Mol Biochem Parasitol 182, 27-36.
670 Sheiner, L., Santos, J.M., Klages, N., Parussini, F., Jemmely, N., Friedrich, N., Ward, G.E.,
671         Soldati-Favre, D., 2010. Toxoplasma gondii transmembrane microneme proteins and
672         their modular design. Mol Microbiol 77, 912-929.
673 Sibley, L.D., 2004. Intracellular parasite invasion strategies. Science 304, 248-253.
674 Sibley, L.D., 2013. The roles of intramembrane proteases in protozoan parasites. Biochim
675         Biophys Acta 1828, 2908-2915.
676 Sibley, L.D., Hakansson, S., Carruthers, V.B., 1998. Gliding motility: an efficient mechanism
677         for cell penetration. Curr Biol 8, R12-14.
678 Sigrist, C.J., de Castro, E., Cerutti, L., Cuche, B.A., Hulo, N., Bridge, A., Bougueleret, L.,
679         Xenarios, I., 2013. New and continuing developments at PROSITE. Nucleic Acids
680         Res 41, D344-347.
681 Silber, A.M., Pereira, C.A., 2012. Assignment of putative functions to membrane
682         "hypothetical proteins" from the Trypanosoma cruzi genome. J Membr Biol 245, 125-
683         129.
684 Szalkai, B., Grolmusz, V., 2018a. Near perfect protein multi-label classification with deep
685         neural networks. Methods 132, 50-56.
686 Szalkai, B., Grolmusz, V., 2018b. SECLAF: a webserver and deep neural network design
687         tool for hierarchical biological sequence classification. Bioinformatics 34, 2487-2489.
688 Templeton, T.J., Iyer, L.M., Anantharaman, V., Enomoto, S., Abrahante, J.E., Subramanian,
689         G.M., Hoffman, S.L., Abrahamsen, M.S., Aravind, L., 2004. Comparative analysis of
690         apicomplexa and genomic diversity in eukaryotes. Genome Res 14, 1686-1695.

691 Tomley, F.M., Soldati, D.S., 2001. Mix and match modules: structure and function of
692     microneme proteins in apicomplexan parasites. Trends Parasitol 17, 81-88.
693 Tordai, H., Banyai, L., Patthy, L., 1999. The PAN module: the N-terminal domains of
694     plasminogen and hepatocyte growth factor are homologous with the apple domains of
695     the prekallikrein family and with a novel domain found in numerous nematode
696     proteins. FEBS Lett 461, 63-67.
697 Urban, S., Dickey, S.W., 2011. The rhomboid protease family: a decade of progress on
698     function and mechanism. Genome Biol 12, 231.
699 Vazquez-Mendoza, A., Carrero, J.C., Rodriguez-Sosa, M., 2013. Parasitic infections: a role
700     for C-type lectins receptors. Biomed Res Int 2013, 456352.
701 Walmsley, A.R., Barrett, M.P., Bringaud, F., Gould, G.W., 1998. Sugar transporters from
702     bacteria, parasites and mammals: structure-activity relationships. Trends Biochem Sci
703     23, 476-481.
704 Ward, C.W., Hoyne, P.A., Flegg, R.H., 1995. Insulin and epidermal growth factor receptors
705     contain the cysteine repeat motif found in the tumor necrosis factor receptor. Proteins
706     22, 141-153.
707 Ward, H., Cevallos, A.M., 1998. Cryptosporidium: molecular basis of host-parasite
708     interaction. Adv Parasitol 40, 151-185.
709 Weis, W.I., Taylor, M.E., Drickamer, K., 1998. The C-type lectin superfamily in the immune
710     system. Immunol Rev 163, 19-34.
711 Winter, G., Gooley, A.A., Williams, K.L., Slade, M.B., 2000. Characterization of a major
712     sporozoite surface glycoprotein of Cryptosporidum parvum. Funct Integr Genomics 1,
713     207-217.
714 Zhang, R., Ou, H.Y., Zhang, C.T., 2004. DEG: a database of essential genes. Nucleic Acids
715     Res 32, D271-272.
716

717

718

719

720

721

722

723

724

725

726

727

728

729

730

731

**List of Tables**

**Table 1. Summary of the tools used in the annotation of uncharacterised *N. caninum***

**proteins.**

| Tool/Database | Description | Reference |
|---|---|---|
| Philius | Prediction of transmembrane topology and signal peptides | (Reynolds et al., 2008) |
| Blast2GO 5 PRO | Bioinformatics platform for the functional annotation and analysis of datasets | (Conesa et al., 2005) |
| InterProScan (v68.0) | Scans sequences against InterPro protein signature databases to identify protein families, domains, and repeats | (Finn et al., 2017) |
| MAAP | Prediction of adhesins and adhesin-like proteins | (Ansari et al., 2008) |
| ExPASy PROSITE | Database of protein families, functional sites, and sequence patterns | (de Castro et al., 2006; Sigrist et al., 2013) |
| ToxoDB PBrowse (v2.48) | Interactive and integrated protein browser | (Gajria et al., 2008) |
| SECLAF | Webserver that uses deep neural networks for the hierarchical classification of protein sequences | (Szalkai and Grolmusz, 2018b) |
| Database of Essential Genes (DEG) | Contains records of currently available essential genomic elements among bacteria, archaea, and eukaryotes | (Zhang et al., 2004) |

735

736 **Table 2. Annotations for the final 32 proteins, including their gene descriptions, InterProScan results, and gene ontologies.**

| Protein | Description | GO | InterPro IDs |
|---|---|---|---|
| F0VIM1 | Microneme protein MIC2 | F: GO:005515, protein binding | **IPR002035**, von Willebrand factor, type A domain<br>**IPR036465,** von Willebrand factor A-like domain superfamily<br>**IPR036383**, Thrombospondin type-1 (TSP-1) repeat superfamily<br>**IPR000884,** TSP-1 repeat |
| F0VIG7 | Putative transmembrane protein | C: GO:0016021, integral component of membrane | |
| F0V9X3 | Putative transmembrane protein | C: GO:0016021, integral component of membrane | |
| F0VII0 | Hypothetical protein | C: GO:0016021, integral component of membrane | |
| F0VEH5 | Hypothetical protein | | |
| F0VNG1 | Transmembrane protein | C: GO:0016021, integral component of membrane | |
| F0VPX6 | Transmembrane protein | C: GO:0016021, integral component of membrane | |
| F0VFU2 | Hypothetical protein | C: GO:0016021, integral component of membrane, GO:0016020, membrane | |
| F0V9Z2 | Putative transmembrane protein | C: GO:0016021, integral component of membrane<br>P: GO:0042981, regulation of apoptotic process | **IPR001315**, CARD (caspase recruitment) domain |
| F0VM28 | Hypothetical protein | - | - |
| F0VML7 | Septin | C: GO:0016021, integral component of membrane, GO:0016020, membrane | CATH Superfamily **G3DSA:3.40.50.300**, P-loop containing nucleotide triphosphate hydrolases |

| | | | |
|---|---|---|---|
| F0VK21 | Putative transmembrane protein | C: GO:0016021, integral component of membrane, GO:0016020, membrane | |
| F0VQ97 | Hypothetical protein | C: GO:0016021, integral component of membrane, GO:0016020, membrane | |
| F0VQ28 | Putative transmembrane protein | C: GO:0016021, integral component of membrane, GO:0016020, membrane | |
| F0VE57 | *T. gondii* family A protein | C: GO:0016020, membrane | |
| F0VE64 | *T. gondii* family A protein | C: GO:0016020, membrane | |
| F0VNN6 | Subtilisin SUB2 | F: GO:0004252, serine-type endopeptidase activity | **IPR015500**, Peptidase S8, subtilisin-related protein family |
| | | P: GO:0006508, proteolysis | **IPR036852**, Peptidase S8/S53 domain superfamily |
| | | C: GO:0016021, integral component of membrane | **IPR000209**, Peptidase S8/S53 domain **IPR034204**, Subtilisin SUB1-like catalytic domain |
| F0VRI3 | *T. gondii* family A protein | C: GO:0016021, integral component of membrane | |
| F0VRI6 | *T. gondii* family A protein | C: GO:0016021, integral component of membrane | **PS51257**, prokaryotic membrane lipoprotein lipid attachment site profile |
| F0VRI7 | *T. gondii* family A protein | C: GO:0016021, integral component of membrane | |
| F0VRI8 | *T. gondii* family A protein | C: GO:0016021, integral component of membrane, GO:0016020, membrane | **PS51257**, prokaryotic membrane lipoprotein lipid attachment site profile |
| F0VRI9 | *T. gondii* family A protein | C: GO:0016021, integral component of membrane, GO:0016020, membrane | **PS51257**, prokaryotic membrane lipoprotein lipid attachment site profile |
| F0VRJ2 | T. *gondii* family A protein | C: GO:0016021, integral component of membrane, GO:0016020, membrane | **PS51257**, prokaryotic membrane lipoprotein lipid attachment site profile |
| F0VRJ3 | *T. gondii* family A protein | C: GO:0016021, integral component of membrane, GO:0016020, membrane | **PS51257**, prokaryotic membrane lipoprotein lipid attachment site profile |

| F0VRJ6 | *T. gondii* family A protein | C: GO:0016021, integral component of membrane, GO:0016020, membrane | **PS51257**, prokaryotic membrane lipoprotein lipid attachment site profile |
|--------|------------------------------|----------------------------------------------------------------------|------------------------------------------------------------------------------|
| F0VRJ8 | *T. gondii* family A protein | C: GO:0016020, membrane | **PS51257**, prokaryotic membrane lipoprotein lipid attachment site profile |
| F0VRK0 | *T. gondii* family A protein | C: GO:0016020, membrane | **PS51257**, prokaryotic membrane lipoprotein lipid attachment site profile |
| F0VRL7 | *T. gondii* family A protein | C: GO:0016021, integral component of membrane | |
| F0VRL8 | *T. gondii* family A protein | C: GO:0016021, integral component of membrane | |
| F0VRL9 | *T. gondii* family A protein | C: GO:0016021, integral component of membrane | |
| F0VRM0 | *T. gondii* family A protein | C: GO:0016021, integral component of membrane | |
| F0VRM4 | *T. gondii* family A protein | C: GO:0016021, integral component of membrane | |

737

738 **List of Figures**

739 **Figure 1. A summary of the workflow used to annotate uncharacterised proteins for *N.***

740 ***caninum*.** After retrieving the sequences for uncharacterised proteins from online reference

741 databases (i.e. those described as "hypothetical" or "unnamed"), Philius was used to identify

742 those with transmembrane domains and signal peptides (TM+SP proteins). After further

743 defining this set of proteins by whether they had adhesin-like properties using the MAAP

744 predictor, various tools were used to annotate the sequences and identify features such as

745 domains, repeats, sequence similarity, and gene ontology.

746

747 **Figure 2. Classification of all 3283 uncharacterised proteins under investigation by**

748 **Philius.** The Philius protein topology predictor classifies protein sequences as either globular,

749 globular with a signal peptide, containing transmembrane domains, or containing both

750 transmembrane domains and a signal peptide sequence. As presented, over half of the protein

751 sequences submitted to Philius were classified simply as globular proteins, where the 4%

752 identified as transmembrane proteins with signal peptides were selected for further annotation.

753

754 **Figure 3. Pie of a pie chart presenting the predicted topology for uncharacterised proteins**

755 **predicted to be adhesins.** A total of 654 proteins were identified as having adhesin-like

756 properties, which were further classified by Philius. The 1% in the second pie chart represents

757 the 32 predicted adhesin-like transmembrane proteins with signal peptides, which were
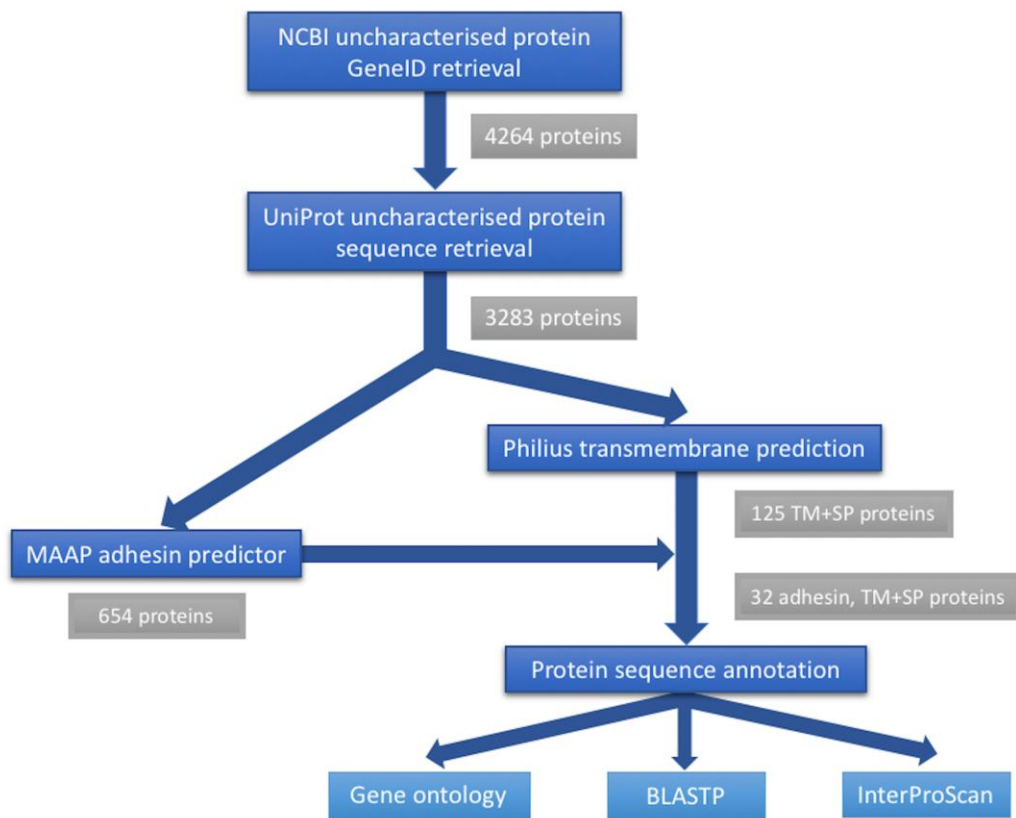
758 subjected to further sequence annotation.

759

760 **\* Black and white to be used for Figures 1-3 in print.**
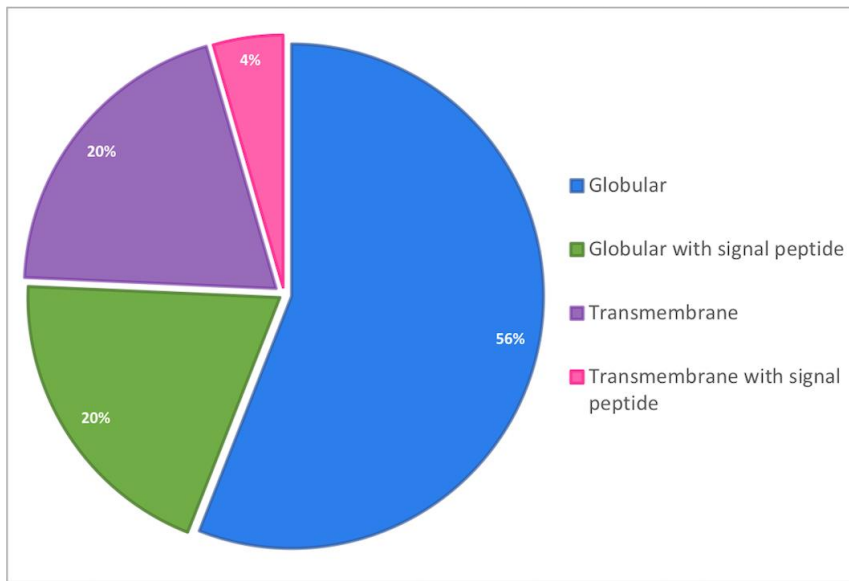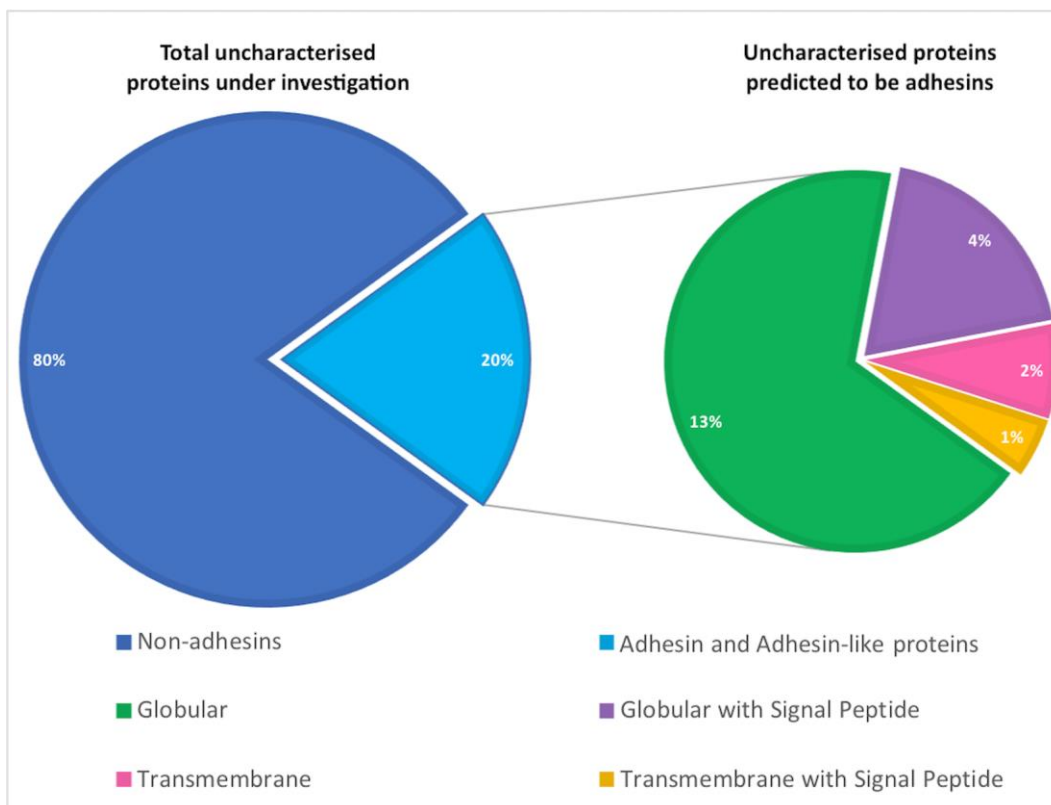
761

762

763  **Figure 1**

764



765

766

767

768    **Figure 2**



769

770

771    **Figure 3**



772