1

2    Review Article

3

4

5

6

7    Proteomes and transcriptomes of Apicomplexa - where's the

8                              message?

9

10    JM Wastling*, D Xia*, A Sohal†, M Chaussepied‡, A Pain†, G Langsley‡

11

12

13

14

15

16    Addresses: * Department of Pre-clinical Veterinary Science, Faculty of Veterinary

17    Science, University of Liverpool, Liverpool, L69 7ZJ, UK. †Pathogen Genomics,

18    Wellcome Trust Sanger Institute, Hinxton, Cambridge CB10 1SA. ‡ Laboratoire de

19    Biologie Cellulaire Comparative des Apicomplexes, Département de Maladies

20    Infectieuses, Institut Cochin, Inserm, U567, CNRS, UMR 8104, Faculté de Médecine

21    Paris V – Hôpital Cochin, 27, rue du Faubourg Saint-Jacques, 75014 Paris, France.

22

23

24 **Abstract**

25 The Apicomplexa now have some of the most comprehensive and integrated

26 proteome datasets of all pathogenic organisms. Proteomic coverage is now at a level

27 where it can be used to predict the potential biological involvement of proteins in

28 these parasites, without having to defer to measurement of mRNA levels.

29 Transcriptomic data for the Apicomplexa (microarrays, EST collections and MPSS)

30 are also abundant, enabling us to investigate the extent to which global mRNA levels

31 correlate with proteomic data. Here, we present a proteomic and transcriptomic

32 perspective of gene expression in key apicomplexan parasites including *Plasmodium*

33 *spp.*, *Toxoplasma gondii*, *Cryptosporidium parvum*, *Neospora caninum* and *Theileria*

34 *spp*. and discuss the alternative views of gene expression that they provide. Although

35 proteomic studies do not detect every gene for which transcripts are seen, many

36 examples of readily detected proteins, whose corresponding genes display little or no

37 detectable transcription are seen across the Apicomplexa. These examples are not

38 easily explained by the "guilt by association", or "stock and go hypotheses" of gene

39 transcription. With the advent of ultra-high-throughput sequencing technologies there

40 will be a quantum shift in transcriptional analysis which, combined with improving

41 quantitative proteome data sets, will provide a core component of a systems-wide

42 approach to studying the Apicomplexa.

43

44

45

46

47

48

## Introduction

The last five years has seen proteomics become established as an integral component of the functional genomics repertoire. This growth, which has resulted from fundamental technical advances in mass spectrometry and bioinformatics, has been accompanied by the emergence of numerous large-scale proteomic experiments with substantial amounts of protein expression data being deposited into increasingly sophisticated on-line proteome resources. Protozoan parasites have not been left-behind in this rush for a proteomic perspective on gene expression; on the contrary, the Apicomplexa, for example, now have some of the most comprehensive and integrated proteomic datasets of all pathogenic organisms. This continuing appetite for proteomic data follows the recognition that examining the proteome has the potential to reveal far more about putative function than can be accounted for by transcriptional data alone. Furthermore, there has been little slow-down in the pace of technological advances in both mass spectrometry and the increasing sophistication of the bioinformatic resources that underpinned the emergence of proteomics little over a decade ago. Importantly, these advances have resulted in a significant increase in the depth and breadth of proteomics coverage that is realistically achievable in an experiment. Whereas a few years ago whole-cell (or so-called "global") proteome surveys could do little more than sample just a small top-slice of the most abundant proteins, deep-mining of the proteome is now becoming increasingly feasible and with it the ability to monitor simultaneously the expression of thousands of proteins in a biological system.

Studies on apicomplexan parasites have been especially prominent promoting a proteomic understanding of gene expression in lower eukaryotes with large-scale

74    proteomic surveys of *Plasmodium falicparum* (for example (Florens, L. et al., 2002;

75    Lasonder, E. et al., 2002), *Cryptosporidium parvum* (Sanderson, S. J. et al., 2008;

76    Snelling, W. J. et al., 2007) and *Toxoplasma gondii* (Xia, D. et al., 2008) being

77    undertaken.  Apicomplexan proteomics has also benefited from a range of advances

78    such as improved sub-fractionation of complex protein mixtures prior to analysis

79    (Nirmalan, N. et al., 2007), separation and analysis of apicomplexan sub-proteomes

80    (Bradley, P. J. et al., 2005; Hu, K. et al., 2006; Zhou, X. W. et al., 2005) and a strong

81    genome bioinformatic resource populated with increasingly accurate gene models

82    (Bahl, A. et al., 2003; Gajria, B. et al., 2008; Heiges, M. et al., 2006).  Proteomic

83    studies have not only provided valuable corroborative evidence for predicted gene

84    models by verifying the existence of thousands of hitherto hypothetical proteins, but

85    have provided sufficient depth of coverage to begin to query the relationship between

86    data acquired from transcriptional surveys, such as those from EST and microarray

87    analysis, and actual protein expression.  Such comparative surveys combining datasets

88    from ESTs, microarray expression and proteomics have already raised fascinating

89    questions pertaining to the link between transcription and translation in the

90    Apicomplexa.

91

92    Despite magnitude advances in the accuracy and sensitivity of mass spectrometry,

93    proteomics still suffers from the disadvantage that, unlike DNA, proteins cannot be

94    amplified to increase the sensitivity of detection.  The debate therefore remains on

95    whether current proteomic technologies can provide sufficient depth and breadth of

96    coverage to describe fully global gene expression.   However, at a time when

97    technological gaps in proteomics seems to be rapidly closing, questions over the

98    relative biological meaning of proteomic and transcriptomic datasets are timely and

99     especially pertinent to apicomplexan biology. In this paper we review advances in

100    proteomic and transcriptional studies in the Apicomplexa, which have enabled us for

101    the first time to examine the relationship between transcription and translation across

102    this important group of parasites and that highlight some fascinating, if not yet fully

103    understood, discrepancies between these types of data.

104

105    Although still imperfect, proteomics does after all provide first hand data on the

106    functional products of gene expression – proteins and hence their putative function.

107    Some argue that we should even look routinely to proteomics, rather than

108    transcriptional patterns, to give us a more meaningful picture of the biological

109    functions of genes. It is perhaps a sign of the breathtaking speed of advance in

110    genomic analysis in the post-genomic era that transcriptional analysis is now seen by

111    some as an "old technology" compared to its younger cousin, proteomics. Here, we

112    argue that a combination of proteomics and transcriptional analysis provides the better

113    perspective on gene expression, but these technologies are still in their infancy and we

114    still have much to learn about the intimate and complex relationship between the two

115    in the Apicomplexa.

116

117    **A global proteomic perspective of the Apicomplexa**

118    Recent global proteomic studies of apicomplexan parasites have massively increased

119    the amount of protein expression data available for these parasites. In order to

120    maximise the depth of coverage obtained in these analyses a combination of

121    specialised separation and mass spectrometry approaches have been adopted. Thus, a

122    typical experiment may involve gel-based analysis of parasite protein (one- or two-

123    dimensional gel electrophoresis) followed by mass spectrometry of trypsin digested

124    bands or spots.  In addition, the parasite will also be analysed by whole shotgun

125    proteome analysis, commonly known as "MudPIT". Whereas gel-based analysis

126    reveals potentially more detailed protein data in the form of semi-quantitation and

127    some post-translational information, shotgun analysis involves the separation of

128    digested peptides in liquid phase, thus avoiding some of the common problems

129    associated with gel separation of hydrophobic proteins, or proteins with extreme

130    mass/pI.  These approaches have enabled up to nearly 50% of the predicted proteome

131    to be resolved on a proteomics platform.  A summary of some of the whole-proteome

132    projects in the Apicomplexa is presented in Table 1 and include those for *P.*

133    *falciparum* (Florens, L. et al., 2002; Lasonder, E. et al., 2002) in which four different

134    life cycle stages were identified using MudPIT and 1-DE gel LC-MS/MS.

135    Comprehensive proteomic approaches have also been used to analyze the proteome of

136    *P. berghei* and *P. yoelii* (Hall, N. et al., 2005; Khan, S. M. et al., 2005; Tarun, A. S. et

137    al., 2008). Thus, proteomic analysis of *Plasmodium* has resulted in one of the most

138    comprehensive datasets for any micro-organism, with data detailed proteomic

139    coverage of up to 5 stages of the complex life cycle of *Plasmodium* species.  These

140    studies have been aimed at addressing important biological questions such as

141    determining the functional characterisation of previously unknown cellular pathways

142    (e.g. kinase pathways that regulate sex-specific functions in *Plasmodium* described by

143    Khan et al. (Khan, S. M. et al., 2005). Doolan et al. studied combined genome and

144    proteome data to identify a large number of sporozoite antigens that are expressed

145    highly in sporozoites and showed high interferon-gamma response in the PBMCs of

146    human volunteers, thus providing a list of novel candidates that could be tested as

147    vaccine candidates (Doolan, D. L. et al., 2003).  In a study which combined the

148    transcriptome and proteome of *P. berghei*, evidence was obtained to demonstrate the

149    developmental stage-specific translational control of mRNA transcripts and gave rise

150    to the "stock and go" hypothesis (Hall, N. et al., 2005). Patra and co-workers

151    undertook a study on the ookinete/zygote proteome of *P. gallinaceum*, the results of

152    which represent a detailed proteomic view of *Plasmodium*-mosquito midgut

153    interactions, fundamental to the development of a novel transmission blocking

154    vaccine in malaria (Patra, K. P. et al., 2008).

155

156    Large scale protein expression profiling projects have also been carried out on the

157    tachyzoite stage of *T. gondii* (Cohen, A. M. et al., 2002; Xia, D. et al., 2008) and

158    similar approaches have been applied in investigating proteome of *C. parvum*

159    sporozoites (Sanderson, S. J. et al., 2008; Snelling, W. J. et al., 2007). These studies

160    have identified between approximately 30-40% of the "total" proteome.  Further

161    unpublished proteome data for *Toxoplasma* and *Cryptosporidium* are available via

162    ApiDB, most notably a substantial additional data set for the *Toxoplasma* and

163    *Cryptosporidium* proteome (unpublished, http://toro.aecom.yu.edu/biodefense/).

164    More recently, proteome profiling of *N. caninum* has also been carried out (Wastling,

165    unpublished data) and peptide evidence has so far been obtained for 660 of the

166    predicted gene models in the current gene prediction set (www.genenedb.org),

167    although this number is anticipated to increase substantially in the near future.

168

169    **Sub-proteomes of the Apicomplexa**

170    Apicomplexan sub-proteomes have been investigated in some detail, with analysis of

171    the apical invasive organelles leading the field. Bradley and co-workers pioneered the

172    proteomic investigation of apicomplexan rhoptry organelles, identifying many novel

173    components of the rhoptry and rhoptry neck of *T. gondii*  (Bradley, P. J. et al., 2005),

174     whilst other key proteins released during host-cell invasion by tachyzoites have also

175     been characterised using 2-DE and MudPIT (Zhou, X. W. et al., 2005; Zhou, X. W.

176     et al., 2004). Rhoptry-enriched fractions have also been investigated in *Plasmodium*

177     merozoites (Sam-Yellowe, T. Y. et al., 2004). The fractionated surface protein of

178     parasite-infected erythrocytes of *P. falciparum* (Florens, L. et al., 2004) and the

179     enriched cytoskeleton components of *T. gondii* (Hu, K. et al., 2006) and the

180     cytoskeletal and membrane fractions of both *T. gondii* and *C. parvum* have also been

181     examined (unpublished, http://toro.aecom.yu.edu/biodefense/).

182

183     **Gene finding and curation in apicomplexans - a proteomic perspective:**

184     Except in a very small number of cases where protein sequence is generated by *de*

185     *novo* protein sequencing, the quality of proteomics identifications is entirely

186     dependent on the sophistication of the gene models against which mass spectrometry

187     data is searched against. Without accurately predicted gene models, proteomics

188     experiments produce only a partial view of the proteome with considerable

189     uncertainty surrounding the nature and number of proteins that may have been

190     identified in any one experiment. Conversely, MS-generated peptide sequence data

191     can be used in reverse logic as a powerful tool not only to provide confirmation, or

192     correction of predicted protein-coding genes, but also to elucidate splicing patterns

193     and as a key input to train gene finding algorithms (Choudhary, J. S. et al., 2001;

194     Fermin, D. et al., 2006; Foissac, S. and Schiex, T., 2005; Tanner, S. et al., 2007;

195     Sanderson, S. J. et al., 2008; Xia, D. et al., 2008). Many of the large-scale proteomics

196     surveys of the Apicomplexa have focussed on genomes that have relatively accurate

197     gene models such as *Cryptosporidium*, *Toxoplasma* and of course *Plasmodium*. In

198     each of these examples proteomics has proved to be a powerful tool in corroborating

199 thousands of hypothetical gene models. Moreover, in cases where several conflicting

200 gene models exist for a particular region of DNA, MS-generated peptide data has

201 been able to identify the most probable interpretation of gene structure and in some

202 cases suggested completely alternative gene models. In one of the first genome scale

203 proteomic survey studies in *P. falciparum,* a large number of good quality 'orphan'

204 peptides (i.e. peptides that did not match to any existing predicted gene in *P.*

205 *falciparum* during the time of publication in 2002) were used to curate manually gene

206 boundaries and also to add missing exons in a number of genes (Florens, L. et al.,

207 2002).

208

209 **Apicomplexan proteomic database resources**

210 Most apicomplexan proteomic datasets are now fully integrated into their respective

211 publicly accessible online genome repositories (see Table 1). In a model developed

212 first for *Cryptosporidium* (www.cryptodb.org) and *Toxoplasma* (www.toxodb.org),

213 mass spectrometry data are now deposited in a standardised way in ApiDB. Thus,

214 MS data can be interrogated in a variety of ways; for example by individual

215 experiment; by sub-proteome; or by "alternative gene model" if variant gene

216 annotations are suspected. One of the most informative ways of visualising

217 proteomics data is a Genome Browser mode (GBrowse), where MS/MS peptide data

218 are shown aligned against predicted gene structure as shown, for example, for the

219 putative nicotinate phosphoribosyltransferasein gene (25.m01815) of *T. gondii*

220 (Figure 1a). In this and other examples, peptide data can be seen alongside other

221 forms of expression data such as EST analysis. A brief examination of a number of

222 genes for which multiple forms of expression data are displayed (peptide and mRNA

223 transcript) shows that whilst there is often broad agreement between gene expression

224    indicators, discrepancies are also common. For example, Figure 1b illustrates the

225    GBrowse view for a putative *Toxoplasma* oxidoreductase (37.m00770) which shows

226    clearly that whilst substantial peptide evidence exists for this gene covering all four of

227    the predicted exons, no corresponding EST data is present.  Interestingly, this gene

228    also shows microarray transcript expression levels below the 25 percentile, indicating

229    little or no transcript could be detected by microarray. Any potential biological role

230    performed by these proteins would escape the "guilt by association" criteria that is

231    based on inferring potential biological function from mRNA levels (Le Roch, K. G. et

232    al., 2003).

233

234    Since all forms of expression data, including proteomics are now integrated into the

235    same database, it is possible systematically to examine such correlations on a genome-

236    wide scale in a way that would have been impossible in the past. The remainder of

237    this review builds on these resources to examine some fundamental questions

238    regarding the nature of proteomic and transcriptomic data in the Apicomplexa.

239

240    **Merging transcriptional and proteomic data in the Apicomplexa**

241    Extensive stage-specific transcriptional data have been acquired for apicomplexan

242    parasites with the implicit assumption that transcriptional changes will reflect protein

243    changes and that this will in turn enable key functions of proteins to be determined,

244    for example those that play a role in stage-specific adaptations; this concept underlies

245    the "guilt by association" hypothesis.  dbEST (Boguski, M. S. et al., 1993) and

246    ApiDB (Aurrecoechea, C. et al., 2007) host the largest collection of expressed

247    sequence tags (EST) for the Apicomplexa. Serial analysis of gene expression (SAGE)

248    projects have also been carried out for both *P. falciparum* and *T. gondii* (Gunasekera,

249  A. M. et al., 2003; Gunasekera, A. M. et al., 2007; Gunasekera, A. M. et al., 2004;

250  Patankar, S. et al., 2001; Radke, J. R. et al., 2005). Microarray expression data are

251  also available for *P. falciparum, P. berghei* and *T. gondii* (Ben, Mamoun C. et al.,

252  2001; Bozdech, Z. et al., 2003; Hall, N. et al., 2005; Kidgell, C. et al., 2006; LaCount,

253  D. J. et al., 2005).

254

255  At this time microarray data are missing for *Theileria* parasites, so to gain insights

256  into parasite gene expression profiles a collection of ESTs from different *T. annulata*

257  life cycle stages were sequenced, the majority of which (circa 10k) came from

258  parasite infected macrophages (Pain, A. et al., 2005) and in the case of *T. parva-*

259  infected lymphocytes an alternative powerful technique was used called Massively

260  Parallel Signature Sequencing (MPSS) (Bishop, R. et al., 2005). MPSS is a PCR-

261  based technique that gives sort (20bp) sequence tags of very high coverage and

262  generates both sense and anti-sense data for a given gene. Importantly, since more

263  than a million *T. parva* transcripts were sequenced, the number of times a transcript

264  from the same gene was sequenced it generated a score (or a signature) that is an

265  indication of the level of transcription of that gene. For *T. parva* MPSS scores ranged

266  from 4 to 52 thousand per million, indicating a wide-range in gene transcription and

267  more surprisingly, signatures could be detected for greater than 80% of genes

268  (Bishop, R. et al., 2005). This suggests that at a given life cycle stage (schizont

269  infected lymphocytes) the vast majority of *Theileria* genes are being transcribed,

270  albeit at variable levels. Unfortunately, this wealth of MPSS data for *Theileria* is not

271  backed up by proteomic data. Nonetheless, Bishop and colleagues noted that for 7

272  known schizont antigens the MPSS scores for the corresponding genes varied 1000-

273  fold again underlining that protein and mRNA levels do not necessarily correlate.

274  Clearly, proteomic data for *Theileria* and its comparison with the MPSS data set

275  would allow one to see how often abundant message translates into abundant protein.

276

277  There have been a small number of studies designed to obtain a simultaneous system-

278  wide view of transcript and protein expression capable of testing the relationship

279  between transcription and the proteome in *Plasmodium* (Hall, N. et al., 2005; Tarun,

280  A. S. et al., 2008).  Overall these studies have revealed a relatively weak correlation

281  between mRNA and protein expression, with many genes being uniquely detected

282  either by transcriptome or proteome.  Similar discrepancies have been noted in a

283  recent proteomic study of *Toxoplasma* (Xia, D. et al., 2008).  In this study 2252

284  proteins were identified from the tachyzoite stage of the parasite using a

285  multiplatform proteome approach. When these data are compared to genes that have

286  transcriptional evidence from the same life-stage, 626 genes are detected solely by

287  EST evidence and 1131 solely by microarray expression evidence (despite the 68%

288  genome coverage by ESTs and nearly 99% microarray coverage).  Significantly,

289  peptide evidence for 72 tachyzoite genes was obtained from proteomics for which no

290  transcripts were observed either by EST, or by microarray (Figure 2).  This latter

291  observation is particularly fascinating which argues against the common

292  misconception that proteomics is relatively insensitive compared with transcriptional

293  analysis. The presence of proteome evidence in the absence of detectable mRNA

294  transcripts has also been noted in mammalian examples, where large numbers of

295  proteins without transcriptional evidence were detected by proteomics in Hela cells

296  (Cox, J. and Mann, M., 2007).

297

298    Given the abundance of good quality transcriptional and translational data across the

299    Apicomplexa we decided to test systematically two related hypotheses concerning the

300    relationship between proteins and their mRNA message: (1) that discrepancies

301    between proteomic and transcriptional datasets occur frequently across the

302    Apicomplexa (2) that orthologs of proteins that show conflicting transcriptional and

303    proteomics profiles behave in the same way across the Apicomplexa i.e. we hoped to

304    identify apicomplexan-wide groups of proteins which behaved aberrantly with respect

305    to gene transcription and translation. To do this, EST and microarray data (where

306    available) were first compared to their respective proteomics datasets for four species

307    of Apicomplexa including *T. gondii* tachyzoites, *C. parvum* sporozoites, *P.*

308    *falciparum* (all life-stages) and *N. caninum* tachyzoites in order to identify sub-sets of

309    proteins for which transcriptional evidence was apparently missing (Figure 3). All the

310    genes identified by major proteome projects listed in ApiDB were included in the

311    analysis and comparative EST libraries and microarray expression data were used (no

312    microarray data were available for *Neospora* or *Cryptosporidium*). Each column

313    represents the total number of proteins identified by proteomics, with the red portion

314    indicating proteins without any EST evidence and the green proportion showing

315    proteins without either EST, or microarray data (where suitable microarray data are

316    available). These data show clearly that a significant number of genes could be

317    detected by proteomics for which neither EST, nor microarray evidence existed (103

318    for *Plasmodium* and 72 for *Toxoplasma*).

319

320    We reasoned that if the discrepancy between proteome and transcriptome is caused by

321    a biological phenomenon that is conserved across apicomplexan parasites, the

322    orthologs of "proteome only" proteins should have a similar expression pattern in the

323    closely related species, i.e. have proteome evidence, but no transcript evidence. To

324    test this we examined proteome and transcriptome expression signatures for *P.*

325    *falciparum*, *T. gondii* and *N. caninum* (we did not include *C. parvum* because of its

326    relatively poor EST coverage). First, the identities were obtained for every gene for

327    which any form of proteome, EST or microarray expression data were available (in

328    the case of *Plasmodium,* data were included from all life-stages). The criteria for

329    inclusion were any gene that has (i) peptide evidence (ii) an EST hit (iii) $\geq$25%

330    microarray expression.  Next, proteins were sorted into the following categories (a)

331    transcript present but no protein detected (b) protein detected but no EST evidence

332    *and* no transcript detected by microarray $\geq$25% threshold (c) protein detected but no

333    EST evidence. We then determined which proteins from each species were shared

334    between each category using an orthology table derived from a one:many OrthoMCL

335    analysis.  Figure 4(a) shows that of the genes which lacked proteome data, but for

336    which transcripts were present, significant numbers had orthologs in other species,

337    with 313 being common between all three species. This is perhaps an unsurprising

338    result, since it is known that certain types of proteins may be under-represented in

339    proteomic analysis due to their physiochemical composition, low levels of expression

340    or high rates of turn-over and degradation. Further analysis of these orthologous genes

341    would be merited to determine why their corresponding peptide evidence is

342    apparently missing.

343

344    Performing the same analysis in reverse reveals that out of the genes for which protein

345    evidence occurs in the absence of detectable EST and microarray transcripts (356

346    across all species), only a handful are shared as orthologs (Figure 4b), although when

347    the analysis is performed with EST data alone (Figure 4c) a larger number of proteins

348    are shared, including two orthologs seen across all three species. In general however,

349    these data appear to disprove our second hypothesis that a shared biological

350    phenomenon might account for these apparently contradictory expression patterns

351    across the phylum.

352

353    From the analysis performed above, there is no apparent underlying rule that

354    dominates the discrepancy between proteome and transcriptome across apicomplexan

355    parasites, except perhaps for a very small number of genes. There are some interesting

356    candidates in the comparison (59.m00090, coatomer protein gamma 2-subunit) which

357    consistently produces convincing peptide evidence (e.g. 37 peptides and 53 spectra in

358    *T. gondii*), but is without transcript evidence at the EST level in *T. gondii*, *N. caninum*

359    and *C. parvum,* with only a single EST seen in a *P. falicparum* blood-stage EST

360    library. The ortholog of this gene in *T. parva* also appears in the lower than 25

361    percentile MPSS expression analysis (Bishop, R. et al., 2005) and interestingly an

362    orthology search in *Saccharomyces cerevisiae* (YNL287W) also reveals a gene for

363    which no EST evidence has been found, although it is detected by proteomics (The

364    Global Proteome Machine Database) (Craig, R. et al., 2004). It is not known why the

365    coatomer protein, a Golgi-coat associated protein, appears so reluctant to reveal itself

366    at the transcript level across not just the Apicomplexa, but other eukaryotes.

367

368    Despite their discrepancies, it is clear that both transcriptomes and proteomes

369    continue to provide experimental evidence for gene expression following the central

370    dogma of Gene-Transcription-Translation. Apparent contradictions between the

371    datasets for a specific set of genes may still be accounted for by genuine biological

372    phenomena such as post-transcriptional control mechanisms as those described by

373    Hall and colleagues (Hall, N. et al., 2005), who combined genome-scale transcriptome

374    and proteome data for several life cycle stages of *P. berghei* and observed evidence

375    for post-transcriptional gene silencing through translational repression of messenger

376    RNA during sexual development of the parasite. A further explanation may be the

377    "stock and go hypothesis" in *Plasmodium* (Mair, G. R. et al., 2006), where

378    translational repression of messenger RNAs (mRNAs) may play an important role in

379    sexual differentiation and gametogenesis.

380

381    **Proteomics and transcriptomics at the host-cell interface**

382    It would be remiss to end a review on gene expression in the Apicomplexa without

383    acknowledging the intimate relationship between parasite and host-cell gene

384    expression. A considerable number of studies have been undertaken to describe global

385    host-cell gene expression changes associated with the infection of Apicomplexa and

386    other intracellular protozoa, but these are dominated by transcriptional rather than

387    proteomic experiments (summarised in Table 2). It is immediately clear that even

388    comparisons between various microarray studies are difficult, because of the

389    considerable experimental variables introduced into each study, including infection

390    time-course (Blader, I. J. et al., 2001; Jensen, K. et al., 2008; Knight, B. C. et al.,

391    2006; Okomo-Adhiambo, M. et al., 2006; Vaena de, Avalos S. et al., 2002), parasite

392    strain (Knight, B. C. et al., 2006), and host cell type (Chaussabel, D. et al., 2003;

393    Jensen, K. et al., 2008). Notably, the importance of the experimental system chosen

394    and especially the host cell type is critical.  For example, infection of macrophages

395    and dendritic cells with various pathogens will elicit quite distinct transcriptional

396    responses (Chaussabel, D. et al., 2003) illustrating not only a pathogen-specific

397    response, but also a cell-type specific response. For technical reasons, the microarrays

398    are often not made from the host cell type that is naturally infected and this

399    complicates further interpretations regarding disease. When comparing different

400    analyses the precise genetic background of the relevant natural host cell type also has

401    to be taken into consideration, as *T. annulata*-infected macrophages from two

402    different breeds of cow (resistant and susceptible to disease) show changes in their

403    expression profiles when infected with the same genetically cloned parasite (Jensen,

404    K. et al., 2008).

405

406    The modulation of the host-cell proteome by *T. gondii* has been examined in depth by

407    quantitative two-dimensional electrophoresis (Nelson, M. M. et al., 2008) providing

408    an opportunity to compare directly proteomic data with transcriptional data from an

409    identically designed experiment (Blader, I. J. et al., 2001). In this analysis only a weak

410    relationship was observed between host-cell transcriptional data and host proteome

411    data at the individual gene level (Nelson, M. M. et al., 2008).  Significantly however,

412    despite differences in detail, both transcriptomic and proteomic analyses came to

413    similar overall conclusions regarding the modulation of key host-cell pathways by

414    *Toxoplasma*. This perhaps illustrates an important overriding principle when dealing

415    with transcript and protein expression data: that they are complementary data which,

416    although linked intimately, are capable of providing a different, rather than conflicting

417    perspective on the same problem.

418

419    **Conclusions and outlook**

420    It is important to acknowledge that both proteomics and transcriptomics are still

421    relatively young technologies, representing some of the first generation of genome-

422    wide data to follow the apicomplexan genome sequencing projects. Until recently we

423    have been in an exploratory phase, systematically cataloguing what is expressed by

424    apicomplexan parasites, when expression occurs (stage-specific expression) and

425    where expression occurs (organelle proteomic).  Whilst these studies have indeed

426    been pioneering, the focus of proteomics is about to be rapidly altered and extended to

427    the proteomics of protein modifications, drug-parasite and host-parasite interactions.

428    In particular the emphasis will shift to more sensitive and accurate proteomic

429    measurements, with quantitative proteomics enabling us to undertake more

430    meaningful comparisons between transcript abundance and protein abundance.

431    Advances in the context of transcriptional analysis are also anticipated such as the

432    application of MPSS to other Apicomplexa over and above *Theileria*.  With the

433    advent of ultra-high-throughput sequencing technologies [e.g.  Roche (454),

434    Illumina(Solexa); ABI-SoliD], there will be a quantum shift in our ability to fine-map

435    the transcript boundaries of the genes by directly sequencing the transcripts to a high

436    coverage (Graveley, B. R., 2008). Recent studies using these state-of-art techniques

437    have provided unprecedented insight into the transcription states (including alternative

438    splice variants and a large number of previously unrecognised transcripts) in the

439    fission yeast *S. pombe* and human at a single nucleotide resolution (Sultan, M. et al.,

440    2008; Wilhelm, B. T. et al., 2008). Similar transcript sequencing studies are now also

441    underway in apicomplexan parasites and thus the accuracy of gene predictions is

442    expected to get significantly higher in the near future that in turn, will prove highly

443    beneficial to the proteomics. As demonstrated for *T. parva*, the depth of transcript

444    sequencing will also allow us to determine the dynamic range (i.e. signature) of a

445    given transcript. The development of these advanced technologies and their

446    application to other Apicomplexa are likely to reveal even more complexity in the

447    relationship between protein and its message. They will also provide an ever more

448    powerful tool to determine the extent of non-coding RNAs (anti-sense, micro and

449    macro) and their eventual contribution to the success Apicomplexa have demonstrated

450    in parasitizing such a wide range of host cells.

451

452    **Acknowledgements**

453    The authors gratefully acknowledge support form the COST 857 action

454    "Apicomplexan biology in the post-genomic era", which provided an invaluable

455    forum for much of the discussion contained in this manuscript.

456

457

458    Reference List
459

460    Aurrecoechea, C., Heiges, M., Wang, H., Wang, Z., Fischer, S., Rhodes, P., Miller, J.,
461          Kraemer, E., Stoeckert, C. J., Jr., Roos, D. S., and Kissinger, J. C., 2007.
462          ApiDB: integrated resources for the apicomplexan bioinformatics resource
463          center. Nucleic Acids Res. 35, D427-D430.

464    Bahl, A., Brunk, B., Crabtree, J., Fraunholz, M. J., Gajria, B., Grant, G. R., Ginsburg,
465          H., Gupta, D., Kissinger, J. C., Labo, P., Li, L., Mailman, M. D., Milgram, A.
466          J., Pearson, D. S., Roos, D. S., Schug, J., Stoeckert, C. J., Jr., and Whetzel, P.,
467          2003. PlasmoDB: the Plasmodium genome resource. A database integrating
468          experimental and computational data. Nucleic Acids Res. 31, 212-215.

469    Ben, Mamoun C., Gluzman, I. Y., Hott, C., MacMillan, S. K., Amarakone, A. S.,
470          Anderson, D. L., Carlton, J. M., Dame, J. B., Chakrabarti, D., Martin, R. K.,
471          Brownstein, B. H., and Goldberg, D. E., 2001. Co-ordinated programme of
472          gene expression during asexual intraerythrocytic development of the human
473          malaria parasite Plasmodium falciparum revealed by microarray analysis. Mol.
474          Microbiol. 39, 26-36.

475    Bishop, R., Shah, T., Pelle, R., Hoyle, D., Pearson, T., Haines, L., Brass, A., Hulme,
476          H., Graham, S. P., Taracha, E. L., Kanga, S., Lu, C., Hass, B., Wortman, J.,
477          White, O., Gardner, M. J., Nene, V., and de Villiers, E. P., 2005. Analysis of
478          the transcriptome of the protozoan Theileria parva using MPSS reveals that the
479          majority of genes are transcriptionally active in the schizont stage. Nucleic
480          Acids Res. 33, 5503-5511.

481    Blader, I. J., Manger, I. D., and Boothroyd, J. C., 2001. Microarray analysis reveals
482          previously unknown changes in Toxoplasma gondii-infected human cells
483          1. J. Biol. Chem. 276, 24223-24231.

484    Boguski, M. S., Lowe, T. M., and Tolstoshev, C. M., 1993. dbEST--database for
485          "expressed sequence tags"
486          1. Nat. Genet. 4, 332-333.

487    Bozdech, Z., Llinas, M., Pulliam, B. L., Wong, E. D., Zhu, J., and DeRisi, J. L., 2003.
488          The transcriptome of the intraerythrocytic developmental cycle of Plasmodium
489          falciparum
490          1. PLoS. Biol. 1, E5-

491    Bradley, P. J., Ward, C., Cheng, S. J., Alexander, D. L., Coller, S., Coombs, G. H.,
492          Dunn, J. D., Ferguson, D. J., Sanderson, S. J., Wastling, J. M., and Boothroyd,
493          J. C., 2005. Proteomic analysis of rhoptry organelles reveals many novel
494          constituents for host-parasite interactions in *Toxoplasma gondii*. J. Biol.
495          Chem. 280, 34245-34258.

496    Chaussabel, D., Semnani, R. T., McDowell, M. A., Sacks, D., Sher, A., and Nutman,
497          T. B., 2003. Unique gene expression profiles of human macrophages and
498          dendritic cells to phylogenetically distinct parasites
499          1. Blood. 102, 672-681.

500 Choudhary, J. S., Blackstock, W. P., Creasy, D. M., and Cottrell, J. S., 2001.
501      Matching peptide mass spectra to EST and genomic DNA databases. Trends
502      Biotechnol. 19, S17-S22.

503 Cohen, A. M., Rumpel, K., Coombs, G. H., and Wastling, J. M., 2002.
504      Characterisation of global protein expression by two-dimensional
505      electrophoresis and mass spectrometry: proteomics of Toxoplasma gondii1.
506      Int. J. Parasitol. 32, 39-51.

507 Cox, J. and Mann, M., 2007. Is proteomics the new genomics? Cell. 130, 395-398.

508 Craig, R., Cortens, J. P., and Beavis, R. C., 2004. Open source system for analyzing,
509      validating, and storing protein identification data. J. Proteome. Res. 3, 1234-
510      1242.

511 Deng, M., Lancto, C. A., and Abrahamsen, M. S., 2004. Cryptosporidium parvum
512      regulation of human epithelial cell gene expression
513      1. Int. J. Parasitol. 34, 73-82.

514 Doolan, D. L., Southwood, S., Freilich, D. A., Sidney, J., Graber, N. L., Shatney, L.,
515      Bebris, L., Florens, L., Dobano, C., Witney, A. A., Appella, E., Hoffman, S.
516      L., Yates, J. R., III, Carucci, D. J., and Sette, A., 2003. Identification of
517      Plasmodium falciparum antigens by antigenic analysis of genomic and
518      proteomic data. Proc. Natl. Acad. Sci. U. S. A. 100, 9952-9957.

519 Fermin, D., Allen, B. B., Blackwell, T. W., Menon, R., Adamski, M., Xu, Y., Ulintz,
520      P., Omenn, G. S., and States, D. J., 2006. Novel gene and gene model
521      detection using a whole genome open reading frame analysis in proteomics.
522      Genome Biol. 7, R35-

523 Florens, L., Liu, X., Wang, Y., Yang, S., Schwartz, O., Peglar, M., Carucci, D. J.,
524      Yates, J. R., III, and Wub, Y., 2004. Proteomics approach reveals novel
525      proteins on the surface of malaria-infected erythrocytes. Mol. Biochem.
526      Parasitol. 135, 1-11.

527 Florens, L., Washburn, M. P., Raine, J. D., Anthony, R. M., Grainger, M., Haynes, J.
528      D., Moch, J. K., Muster, N., Sacci, J. B., Tabb, D. L., Witney, A. A., Wolters,
529      D., Wu, Y., Gardner, M. J., Holder, A. A., Sinden, R. E., Yates, J. R., and
530      Carucci, D. J., 2002. A proteomic view of the Plasmodium falciparum life
531      cycle. Nature. 419, 520-526.

532 Foissac, S. and Schiex, T., 2005. Integrating alternative splicing detection into gene
533      prediction. BMC. Bioinformatics. 6, 25-34.

534 Gail, M., Gross, U., and Bohne, W., 2001. Transcriptional profile of Toxoplasma
535      gondii-infected human fibroblasts as revealed by gene-array hybridization
536      1. Mol. Genet. Genomics. 265, 905-912.

537 Gajria, B., Bahl, A., Brestelli, J., Dommer, J., Fischer, S., Gao, X., Heiges, M., Iodice,
538      J., Kissinger, J. C., Mackey, A. J., Pinney, D. F., Roos, D. S., Stoeckert, C. J.,
539      Jr., Wang, H., and Brunk, B. P., 2008. ToxoDB: an integrated Toxoplasma
540      gondii database resource. Nucleic Acids Res. 36, D553-D556.

541    Graveley, B. R., 2008. Molecular biology: power sequencing. Nature. 453, 1197-
542        1198.

543    Gunasekera, A. M., Myrick, A., Le, Roch K., Winzeler, E., and Wirth, D. F., 2007.
544        Plasmodium falciparum: genome wide perturbations in transcript profiles
545        among mixed stage cultures after chloroquine treatment
546        1. Exp. Parasitol. 117, 87-92.

547    Gunasekera, A. M., Patankar, S., Schug, J., Eisen, G., Kissinger, J., Roos, D., and
548        Wirth, D. F., 2004. Widespread distribution of antisense transcripts in the
549        Plasmodium falciparum genome
550        1. Mol. Biochem. Parasitol. 136, 35-42.

551    Gunasekera, A. M., Patankar, S., Schug, J., Eisen, G., and Wirth, D. F., 2003. Drug-
552        induced alterations in gene expression of the asexual blood forms of
553        Plasmodium falciparum
554        1. Mol. Microbiol. 50, 1229-1239.

555    Hall, N., Karras, M., Raine, J. D., Carlton, J. M., Kooij, T. W., Berriman, M., Florens,
556        L., Janssen, C. S., Pain, A., Christophides, G. K., James, K., Rutherford, K.,
557        Harris, B., Harris, D., Churcher, C., Quail, M. A., Ormond, D., Doggett, J.,
558        Trueman, H. E., Mendoza, J., Bidwell, S. L., Rajandream, M. A., Carucci, D.
559        J., Yates, J. R., III, Kafatos, F. C., Janse, C. J., Barrell, B., Turner, C. M.,
560        Waters, A. P., and Sinden, R. E., 2005. A comprehensive survey of the
561        Plasmodium life cycle by genomic, transcriptomic, and proteomic analyses1.
562        Science. 307, 82-86.

563    Heiges, M., Wang, H., Robinson, E., Aurrecoechea, C., Gao, X., Kaluskar, N.,
564        Rhodes, P., Wang, S., He, C. Z., Su, Y., Miller, J., Kraemer, E., and Kissinger,
565        J. C., 2006. CryptoDB: a Cryptosporidium bioinformatics resource update.
566        Nucleic Acids Res. 34, D419-D422.

567    Hu, K., Johnson, J., Florens, L., Fraunholz, M., Suravajjala, S., DiLullo, C., Yates, J.,
568        Roos, D. S., and Murray, J. M., 2006. Cytoskeletal components of an invasion
569        machine--the apical complex of Toxoplasma gondii. PLoS. Pathog. 2, e13-

570    Jensen, K., Paxton, E., Waddington, D., Talbot, R., Darghouth, M. A., and Glass, E.
571        J., 2008. Differences in the transcriptional responses induced by Theileria
572        annulata infection in bovine monocytes derived from resistant and susceptible
573        cattle breeds
574        1. Int. J. Parasitol. 38, 313-325.

575    Khan, S. M., Franke-Fayard, B., Mair, G. R., Lasonder, E., Janse, C. J., Mann, M.,
576        and Waters, A. P., 2005. Proteome analysis of separated male and female
577        gametocytes reveals novel sex-specific Plasmodium biology. Cell. 121, 675-
578        687.

579    Kidgell, C., Volkman, S. K., Daily, J., Borevitz, J. O., Plouffe, D., Zhou, Y., Johnson,
580        J. R., Le, Roch K., Sarr, O., Ndir, O., Mboup, S., Batalov, S., Wirth, D. F., and
581        Winzeler, E. A., 2006. A systematic map of genetic variation in Plasmodium
582        falciparum. PLoS. Pathog. 2, e57-

583    Knight, B. C., Kissane, S., Falciani, F., Salmon, M., Stanford, M. R., and Wallace, G.
584          R., 2006. Expression analysis of immune response genes of Muller cells
585          infected with Toxoplasma gondii
586          1. J. Neuroimmunol. 179, 126-131.

587    LaCount, D. J., Vignali, M., Chettier, R., Phansalkar, A., Bell, R., Hesselberth, J. R.,
588          Schoenfeld, L. W., Ota, I., Sahasrabudhe, S., Kurschner, C., Fields, S., and
589          Hughes, R. E., 2005. A protein interaction network of the malaria parasite
590          Plasmodium falciparum. Nature. 438, 103-107.

591    Lasonder, E., Ishihama, Y., Andersen, J. S., Vermunt, A. M., Pain, A., Sauerwein, R.
592          W., Eling, W. M., Hall, N., Waters, A. P., Stunnenberg, H. G., and Mann, M.,
593          2002. Analysis of the Plasmodium falciparum proteome by high-accuracy
594          mass spectrometry. Nature. 419, 537-542.

595    Le Roch, K. G., Zhou, Y., Blair, P. L., Grainger, M., Moch, J. K., Haynes, J. D., de,
596          la, V, Holder, A. A., Batalov, S., Carucci, D. J., and Winzeler, E. A., 2003.
597          Discovery of gene function by expression profiling of the malaria parasite life
598          cycle. Science. 301, 1503-1508.

599    Mair, G. R., Braks, J. A., Garver, L. S., Wiegant, J. C., Hall, N., Dirks, R. W., Khan,
600          S. M., Dimopoulos, G., Janse, C. J., and Waters, A. P., 2006. Regulation of
601          sexual development of Plasmodium by translational repression. Science. 313,
602          667-669.

603    Nelson, M. M., Jones, A. R., Carmen, J. C., Sinai, A. P., Burchmore, R., and
604          Wastling, J. M., 2008. Modulation of the host cell proteome by the
605          intracellular apicomplexan parasite Toxoplasma gondii
606          1. Infect. Immun. 76, 828-844.

607    Nirmalan, N., Flett, F., Skinner, T., Hyde, J. E., and Sims, P. F., 2007. Microscale
608          solution isoelectric focusing as an effective strategy enabling containment of
609          hemeoglobin-derived products for high-resolution gel-based analysis of the
610          Plasmodium falciparum proteome. J. Proteome. Res. 6, 3780-3787.

611    Okomo-Adhiambo, M., Beattie, C., and Rink, A., 2006. cDNA microarray analysis of
612          host-pathogen interactions in a porcine in vitro model for Toxoplasma gondii
613          infection
614          1. Infect. Immun. 74, 4254-4265.

615    Pain, A., Renauld, H., Berriman, M., Murphy, L., Yeats, C. A., Weir, W., Kerhornou,
616          A., Aslett, M., Bishop, R., Bouchier, C., Cochet, M., Coulson, R. M., Cronin,
617          A., de Villiers, E. P., Fraser, A., Fosker, N., Gardner, M., Goble, A., Griffiths-
618          Jones, S., Harris, D. E., Katzer, F., Larke, N., Lord, A., Maser, P., McKellar,
619          S., Mooney, P., Morton, F., Nene, V., O'Neil, S., Price, C., Quail, M. A.,
620          Rabbinowitsch, E., Rawlings, N. D., Rutter, S., Saunders, D., Seeger, K.,
621          Shah, T., Squares, R., Squares, S., Tivey, A., Walker, A. R., Woodward, J.,
622          Dobbelaere, D. A., Langsley, G., Rajandream, M. A., McKeever, D., Shiels,
623          B., Tait, A., Barrell, B., and Hall, N., 2005. Genome of the host-cell
624          transforming parasite Theileria annulata compared with T. parva
625          1. Science. 309, 131-133.

626 Patankar, S., Munasinghe, A., Shoaibi, A., Cummings, L. M., and Wirth, D. F., 2001.
627    Serial analysis of gene expression in Plasmodium falciparum reveals the
628    global expression profile of erythrocytic stages and the presence of anti-sense
629    transcripts in the malarial parasite
630    1. Mol. Biol. Cell. 12, 3114-3125.

631 Patra, K. P., Johnson, J. R., Cantin, G. T., Yates, J. R., III, and Vinetz, J. M., 2008.
632    Proteomic analysis of zygote and ookinete stages of the avian malaria parasite
633    Plasmodium gallinaceum delineates the homologous proteomes of the lethal
634    human malaria parasite Plasmodium falciparum. Proteomics. 8, 2492-2499.

635 Radke, J. R., Behnke, M. S., Mackey, A. J., Radke, J. B., Roos, D. S., and White, M.
636    W., 2005. The transcriptome of Toxoplasma gondii. BMC. Biol. 3, 26-

637 Sam-Yellowe, T. Y., Florens, L., Wang, T., Raine, J. D., Carucci, D. J., Sinden, R.,
638    and Yates, J. R., III, 2004. Proteome analysis of rhoptry-enriched fractions
639    isolated from Plasmodium merozoites. J. Proteome. Res. 3, 995-1001.

640 Sanderson, S. J., Xia, D., Prieto, H., Yates, J., Heiges, M., Kissinger, J. C., Bromley,
641    E., Lal, K., Sinden, R. E., Tomley, F., and Wastling, J. M., 2008. Determining
642    the protein repertoire of Cryptosporidium parvum sporozoites. Proteomics. 8,
643    1398-1414.

644 Slater, G. S. and Birney, E., 2005. Automated generation of heuristics for biological
645    sequence comparison
646    1. BMC. Bioinformatics. 6, 31-

647 Snelling, W. J., Lin, Q., Moore, J. E., Millar, B. C., Tosini, F., Pozio, E., Dooley, J.
648    S., and Lowery, C. J., 2007. Proteomics analysis and protein expression during
649    sporozoite excystation of Cryptosporidium parvum (Coccidia, Apicomplexa).
650    Mol. Cell Proteomics. 6, 346-355.

651 Sultan, M., Schulz, M. H., Richard, H., Magen, A., Klingenhoff, A., Scherf, M.,
652    Seifert, M., Borodina, T., Soldatov, A., Parkhomchuk, D., Schmidt, D.,
653    O'Keeffe, S., Haas, S., Vingron, M., Lehrach, H., and Yaspo, M. L., 2008. A
654    global view of gene activity and alternative splicing by deep sequencing of the
655    human transcriptome. Science. 321, 956-960.

656 Tanner, S., Shen, Z., Ng, J., Florea, L., Guigo, R., Briggs, S. P., and Bafna, V., 2007.
657    Improving gene annotation using peptide mass spectrometry. Genome Res. 17,
658    231-239.

659 Tarun, A. S., Peng, X., Dumpit, R. F., Ogata, Y., Silva-Rivera, H., Camargo, N.,
660    Daly, T. M., Bergman, L. W., and Kappe, S. H., 2008. A combined
661    transcriptome and proteome survey of malaria parasite liver stages. Proc. Natl.
662    Acad. Sci. U. S. A. 105, 305-310.

663 Vaena de, Avalos S., Blader, I. J., Fisher, M., Boothroyd, J. C., and Burleigh, B. A.,
664    2002. Immediate/early response to Trypanosoma cruzi infection involves
665    minimal modulation of host cell transcription
666    1. J. Biol. Chem. 277, 639-644.

667 Wilhelm, B. T., Marguerat, S., Watt, S., Schubert, F., Wood, V., Goodhead, I.,
668        Penkett, C. J., Rogers, J., and Bahler, J., 2008. Dynamic repertoire of a
669        eukaryotic transcriptome surveyed at single-nucleotide resolution. Nature.
670        453, 1239-1243.

671 Xia, D., Sanderson, S. J., Jones, A. R., Prieto, J. H., Yates, J. R., Bromley, E.,
672        Tomley, F. M., Lal, K., Sinden, R. E., Brunk, B. P., Roos, D. S., and Wastling,
673        J. M., 2008. The proteome of Toxoplasma gondii: integration with the genome
674        provides novel insights into gene expression and annotation. Genome Biol. 9,
675        R116-

676 Zhou, X. W., Blackman, M. J., Howell, S. A., and Carruthers, V. B., 2004. Proteomic
677        analysis of cleavage events reveals a dynamic two-step mechanism for
678        proteolysis of a key parasite adhesive complex. Mol. Cell Proteomics. 3, 565-
679        576.

680 Zhou, X. W., Kafsack, B. F., Cole, R. N., Beckett, P., Shen, R. F., and Carruthers, V.
681        B., 2005. The opportunistic pathogen *Toxoplasma gondii* deploys a diverse
682        legion of invasion and survival proteins. J. Biol. Chem. 280, 34233-34244.
683
684
685

686 **Table 1**

687

688 **Summary of global proteomic studies in the Apicomplexa**

689

| Species | Life Cycle Stage | Platform | References | Database resource? | Number of unique proteins identified | Estimated Proportion of Proteome | Transcript Expression Data? |
|---|---|---|---|---|---|---|---|
| *P.falciparum* | Sporozoite, Merozoite, Trophozoite Gametocyte Trophozoite/ Schizont | 1-DE Gel-LC MS/MS MudPIT | (Florens, L. et al., 2002; Florens, L. et al., 2004; Lasonder, E. et al., 2002) | ApiDB | 2427 | ~45% | EST SAGE, Microarray |
| *P.berghei* | Gametocyte, Asexual blood stage, Ookinete | 1-DE Gel LC-MS/MS MudPIT | (Hall, N. et al., 2005; Khan, S. M. et al., 2005) | ApiDB | 2924 | ~24% | EST Microarray |
| *P. yoelii* | Liver Stage Schizont | 1-DE Gel LC-MS/MS | (Tarun, A. S. et al., 2008) | None | 816 | ~10% | Microarray |
| *T.gondii* | Tachyzoite | 1-DE Gel LC-MS/MS, 2-DE Gel LC-MS/MS MudPIT | (Bradley, P. J. et al., 2005; Hu, K. et al., 2006; Xia, D. et al., 2008) | ApiDB | 2457 | ~31% | EST SAGE Microarray |
| *C.parvum* | Oocyst/ sporozoite | 1-DE Gel LC-MS/MS, 2-DE Gel LC-MS/MS MudPIT | (Sanderson, S. J. et al., 2008; Snelling, W. J. et al., 2007) | ApiDB | 1322 | ~30% | EST |
| *N.caninum* | Tachyzoite | MudPIT | Un-published | None | 660 genes | ~15% | EST |

690

691

692 **Table 2**

693

694 **Summary of host-cell transcriptional studies in the apicomplexan infections**

695

| Parasite | Target cells | Species | Time points | Microarray | References |
|---|---|---|---|---|---|
| *Theileria annulata* sporozoïtes | Peripheral-blood monocytes | *Bos taurus* (S) & *B. indicus* (R) | 0, 2, 72hrs | Cattle 5K Immune cDNA (ARK-Genomics) | (Jensen, K. et al., 2008) |
| *Toxoplasma gondii* tachyzoite strain TS-4 | PK13, porcine kidney epithelial cell line | *Sus scrofa* | 0, 1, 2, 4, 6, 24, 48, 72hrs | Porcine custom cDNA | (Okomo-Adhiambo, M. et al., 2006) |
| - *Toxoplasma gondii* tachyzoite strain RH | Peripheral-blood monocytes differentiated to macrophages or dendritic cells | *Homo sapiens* | 0, 16hrs | HU95A (Affymetrix) probe array | (Chaussabel, D. et al., 2003) |
| *Toxoplasma gondii* tachizoites RH strain | Human foreskin fibroblasts (HFF) | *Homo sapiens* | 0, 24hrs | Human cDNA array (Human Atlas Array, Clontech) | (Gail, M. et al., 2001) |
| *Toxoplasma gondii* RH strain tachizoites and Prugniaud strain cysts | Human Müller cell line (MOI-M1) | *Homo sapiens* | 0, 2, 24hrs | Human apoptosis and custom probe arrays (Affymetrix) | (Knight, B. C. et al., 2006) |
| *Toxoplasma gondii* | Human foreskin fibroblasts (HFF) | *Homo sapiens* | 0, 1, 2, 4, 6, 24hrs | Human custom cDNA | (Blader, I. J. et al., 2001) |
| *Cryptosporidium parvum* oocytsts | HCT-8 epithelial cell line | *Homo sapiens* | 0, 24hrs | HG-U95Av2 probe array (Affymetrix) | (Deng, M. et al., 2004) |

696

697

698

699  **Figure 1**

700  **Visualisation of proteomic and transcriptomic expression data in ToxoDB**

701  (a) A screenshot of the annotated *T. gondii* gene 25.m01815 (nicotinate

702  phosphoribosyltransferase, putative) on ToxoDB Genome Browser

703  (www.toxodb.org). Predicted gene structures of gene 25.m01815, where blue boxes

704  represent exons, are shown on the top of the figure. EST and proteome (MS/MS

705  peptide) evidence identified for this gene are aligned underneath the gene sequence.

706  The relationship between proteomic (peptide) and transcriptomic (EST) data can be

707  directly visualised.  Note that peptide evidence confirms several predicted intron-exon

708  boundaries (shown by the joins between peptides). (b) GBrowse view for a putative

709  *Toxoplasma* oxidoreductase (37.m00770) gene which shows clearly that whilst

710  substantial peptide evidence exists for this gene covering all four of the predicted

711  exons, no corresponding EST data is present.  Interestingly, this gene also shows

712  microarrary transcript levels below the 25 percentile, indicating little or no transcript

713  could be detected by microarray.

714

715  **Figure 2**

716  **Genes with proteome and transcriptome evidence in *T. gondii***

717  Diagram illustrating the relationship between proteomics, EST and microarray gene

718  expression data in *T. gondii* (data from (Xia, D. et al., 2008). In total 2252 non-

719  redundant proteins were identified from *T.gondii* tachyzoites (blue circle).  These

720  were compared with genes that have tachyzoite EST evidence (green circle) and

721  microarray expression data (orange circle), where higher than 25 expression percentile

722  is observed. The data show that 626 genes have uniquely EST evidence, 1131 genes

723     have uniquely microarray expression evidence, whilst 72 tachyzoite genes are

724     uniquely identified by peptide data and have no transcript expression evidence.

725

726     **Figure 3**

727
728     **Proteome and transcriptome comparisons across four species of Apicomplexa**

729     The numbers of proteins identified by peptide evidence in *T. gondii* tachyzoites, *C.*

730     *parvum* sporozoites, *P. falciparum* (all life-stages) and *N. caninum* tachyzoites are

731     shown. The red portion indicates proteins without EST evidence and the green portion

732     indicates genes without EST and microarray evidence (less than 25 expression

733     percentile). Note that no microarray data were available for *Neospora* or

734     *Cryptosporidium*. All the genes identified by major proteome projects listed in

735     ApiDB are included and comparative EST libraries and microarray expression data

736     were used in the analysis. For *N.caninum*, ESTs were downloaded from dbEST and

737     were aligned to genes that have proteomic evidence under whole genome scaffold

738     using software Exonerate (Slater, G. S. and Birney, E., 2005).

739

740     **Figure 4**

741     **Genes from three Apicomplexa which exhibit discrepancies between**

742     **transcriptional data and proteome data**

743     Each circle represents the number of genes for which a discrepancy was seen between

744     transcriptional data and proteome data for *P. falciparum*, *T. gondii* and *N. caninum*

745     based on (a) transcript present but no protein detected (b) protein detected but no EST

746     evidence *and* no transcript detected by microarray ≥25% threshold (c) protein

747     detected but no EST evidence. The intersections show the numbers of orthologs (as
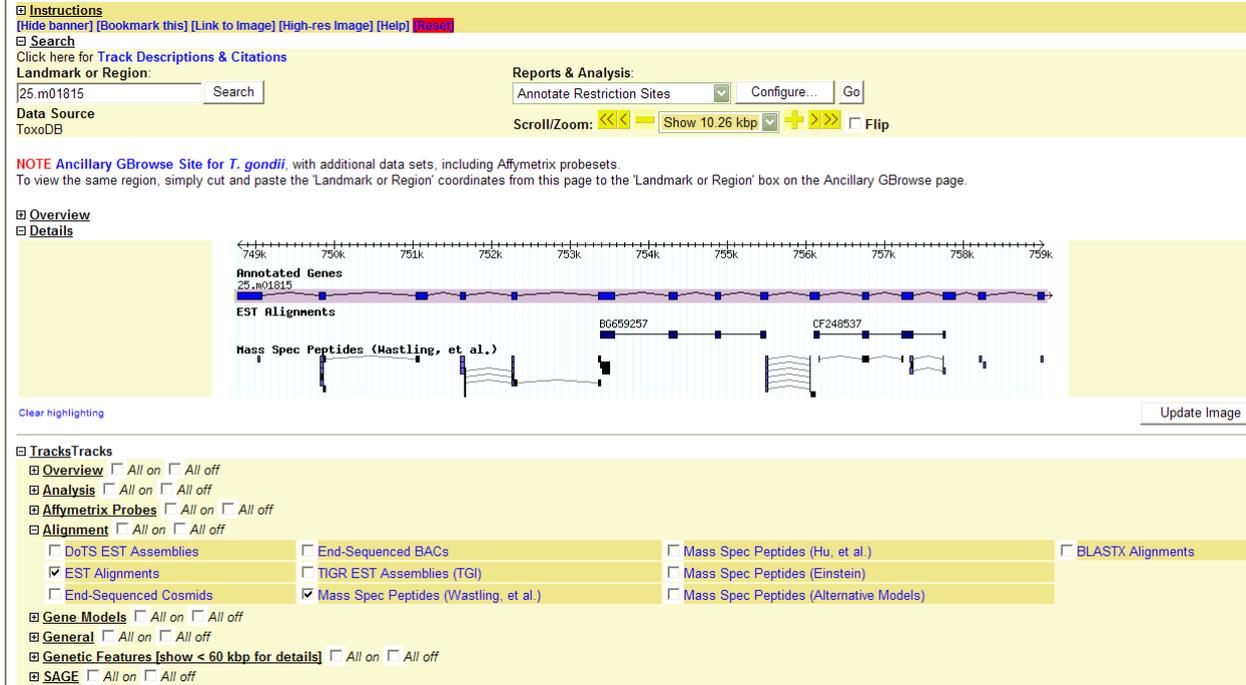
748    determined by OrthoMCL) shared between the species that exhibit contradictory

749    transcriptional and protein expression patterns.
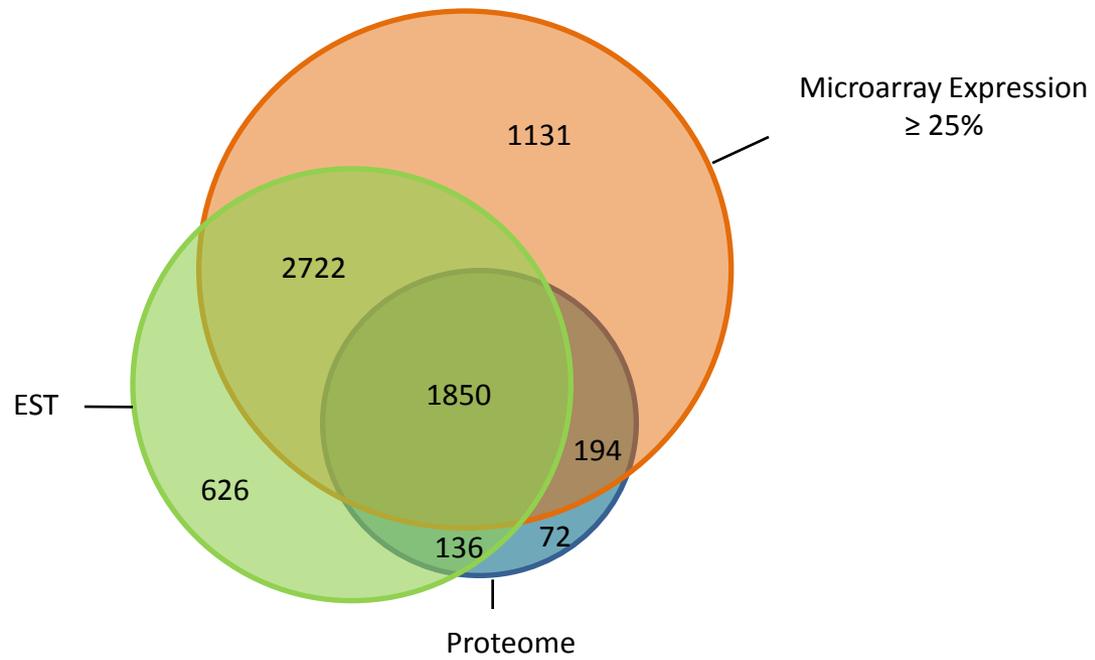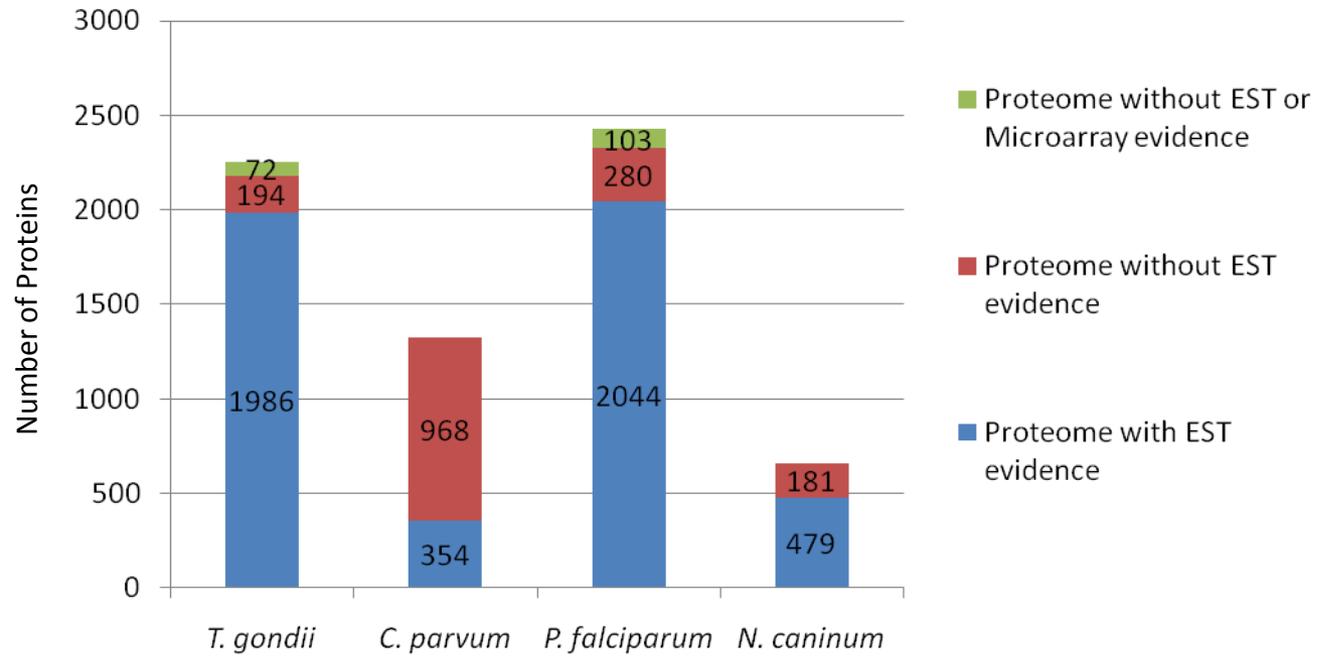
750

# Figure 1

## (a)



## (b)

Figure 2

Figure 3

Figure 4. Proteome vs transcriptome cross apicomplexan parasites