

Manuscript Number: IJPara15\_144R2

Title: Two recently sequenced vertebrate genomes are contaminated with apicomplexan species of the Sarcocystidae family

Article Type: Full Length Article

Keywords: Apicomplexa, apicoplast, Sarcocystis, apicortin, host, phylogenetic tree, bat, Nephroisospora

Corresponding Author: Prof. Ferenc Orosz, PhD

Corresponding Author's Institution:

First Author: Ferenc Orosz, PhD

Order of Authors: Ferenc Orosz, PhD

Manuscript Region of Origin: HUNGARY

**Abstract:** The paper enlightens a general problem, namely, that host genomes can easily be contaminated with parasite ones thus careful isolation of genetic material and careful bioinformatics analysis are needed in all cases. I show for two recently published genomes that they are contaminated with apicomplexan parasites which belong to the Sarcocystidae family. Sequences of the characteristic apicomplexan organelle, apicoplast, were used as queries in BLASTN search against nucleotide sequences of various animal groups looking for possible contaminations. I found that the draft genomes of a bird, *Colinus virginianus* (Halley et al., 2014, PLoS ONE 9, e90240), and a bat, *Myotis davidii* (Zhang et al., 2013, Science, 339, 456-460.) contain at least 6 and 17 contigs, respectively, originated from the apicoplast of an apicomplexan species, and other genes specific of this phylum can also be found in the published genomes. Obviously, the sources of the genetic material, the muscle and the kidney of the animals, respectively, contained the parasitic cysts. Phylogenetic analyses using 18S ribosomal RNA and internal transcribed spacer 1 genes show that the parasite contaminating *C. virginianus* is a species of *Sarcocystis* related to ones known to cycle between avian and mammalian hosts; in the case of *M. davidii* it belongs to the *Nephroisospora* genus, the only member of which, *N. eptesici*, has been recently identified from the kidney of big brown bats.

A. Loukas  
Editor-in-Chief  
International Journal for Parasitology  
Queensland Tropical Health Alliance, Fac. of Medicine, Health & Molecular Sciences,  
James Cook University, McGregor Rd, Smithfield, Cairns, Qld 4878, Queensland, Australia

Dear Professor Loukas,

I would like to resubmit a research article entitled "**Two recently sequenced vertebrate genomes are contaminated with apicomplexan species of the Sarcocystidae family**" to your journal for considering of publication. It is the revised version of my manuscript, IJPara15\_144.

Budapest, July 9, 2015

Sincerely yours

Ferenc Orosz, PhD  
Institute of Enzymology  
Research Centre for Natural Sciences,  
Hungarian Academy of Sciences

Jan Slapeta  
Deputy Editor  
International Journal for Parasitology

Dear Editor,

Thank you for your favourable criticism. Here are my answers to you and to the Reviewers.

**Answers to the Editor:**

"Authors must remove any reference to a formal new name."

I removed it from the abstract and the main text, too.

"Reviewer #3 provided an annotated manuscript that should be used to prepare the further revised version."

I used it and accepted his suggestions without exception.

**Answers to Reviewer 2**

"I think merging the manuscripts was a good idea, and think the result is acceptable."

Thank you for your favourable criticism.

**Answers to Reviewer 4**

"I think that this paper is an interesting one, and should be published in the International Journal for Parasitology."

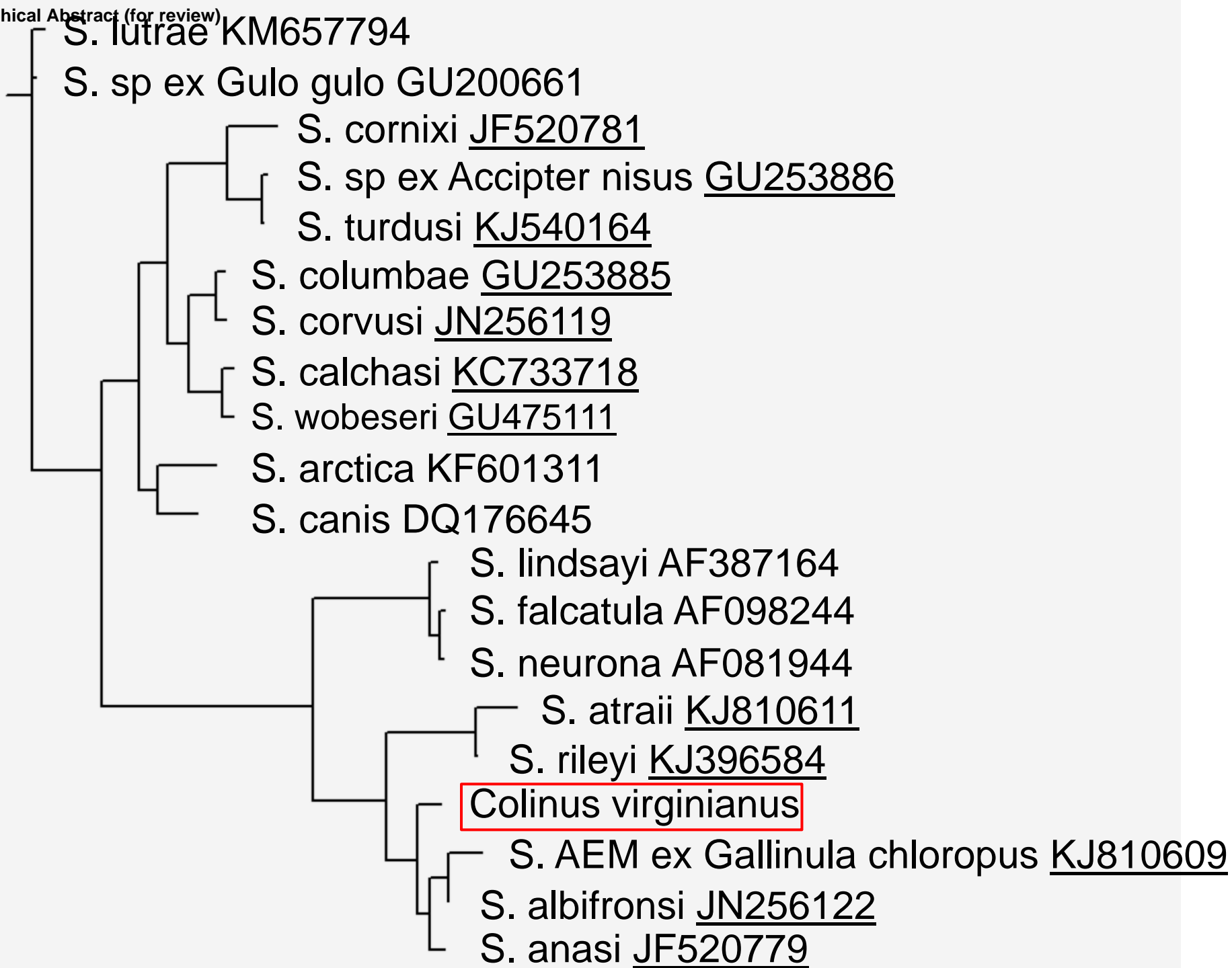
Thank you for your favourable criticism.

"I have only one substantial suggestion. At the end of the manuscript the author says: "which can be named tentatively as *Nephroisospora myotisi*." I think that the author should delete this phrase, which also occurs in the Abstract. This could be interpreted as formally naming the species, which should not be done without considerably more information than is provided here."

I removed it from the abstract and the main text, too.

I have made several suggestions for improving the clarity of the English expression.

Special thanks for the corrections of the grammatical and style shortcomings. Of course, all of them were accepted.



## Highlights

The draft genome of *Colinus virginianus* is contaminated by a *Sarcocystis* species

The parasite is related to ones known to cycle between avian and mammalian hosts

The draft genome of *Myotis davidii*, a bat, is contaminated by a Sarcocystidae species

The parasite is the second (tentative) member of the *Nephroisospora* genus

Several contigs of the hosts originated from the apicomplexan organelle, apicoplast

1

2

3

4 **Two recently sequenced vertebrate genomes are contaminated with apicomplexan species of the**  
5 **Sarcocystidae family**

6

7

8

Ferenc Orosz

9

Institute of Enzymology, Research Centre for Natural Sciences, Hungarian Academy of Sciences

10

11

12

13

14 Address:

15

Ferenc Orosz, Ph.D., D.Sc.

16

Institute of Enzymology,

17

Research Centre for Natural Sciences,

18

Hungarian Academy of Sciences

19

Magyar tudósok körútja 2.

20

H-1117 Budapest, Hungary

21

Phone: (36)-1-3826714

22

E-mail: [orosz.ferenc@ttk.mta.hu](mailto:orosz.ferenc@ttk.mta.hu)

23

24

25

26

27

Note: Supplementary Data are associated with this article'

28

29 **Abstract**

30

31 The paper highlights a general problem, namely, that host genome sequences can easily be contaminated  
32 with parasite ones thus careful isolation of genetic material and careful bioinformatics analysis are needed in  
33 all cases. I show for two recently published genomes that they are contaminated with sequences of  
34 apicomplexan parasites which belong to the Sarcocystidae family. Sequences of the characteristic  
35 apicomplexan organelle, the apicoplast, were used as queries in BLASTN search against nucleotide  
36 sequences of various animal groups looking for possible contaminations. I found that the draft genomes of a  
37 bird, *Colinus virginianus* (Halley et al., 2014, PLoS ONE 9, e90240), and a bat, *Myotis davidii* (Zhang et al.,  
38 2013, Science, 339, 456-460.) contain at least 6 and 17 contigs, respectively, originating from the apicoplast  
39 of an apicomplexan species, and other genes specific to this phylum can also be found in the published  
40 genomes. Obviously, the sources of the genetic material, the muscle and the kidney of the animals,  
41 respectively, contained the parasitic cysts. Phylogenetic analyses using *18S ribosomal RNA* and *internal*  
42 *transcribed spacer 1* genes show that the parasite contaminating *C. virginianus* is a species of *Sarcocystis*  
43 related to ones known to cycle between avian and mammalian hosts. In the case of *M. davidii* it belongs to  
44 the *Nephroisospora* genus, the only member of which, *N. eptesici*, has been recently identified from the  
45 kidney of big brown bats.

46

47

48

49

50

51 Keynotes: Apicomplexa, apicoplast, *Sarcocystis*, apicortin, host, phylogenetic tree, bat, *Nephroisospora*

52

53



54 **1. Introduction**

55

56 The raw data from a genome sequencing project sometimes contains DNA from contaminating  
57 organisms, which may be introduced during sample collection or sequence preparation. In some instances,  
58 these contaminants remain in the sequence even after assembly and deposition of the genome into public  
59 databases. As a consequence, searches of these databases may yield erroneous and confusing results  
60 (Merchant et al., 2014). Human DNA is a common contaminant, from the scientists who handle the samples  
61 during the process of extraction through sequencing (Longo et al., 2011). Computational filters applied to the  
62 raw sequencing reads are usually effective at removal of human DNA and other common laboratory  
63 contaminants such as *E. coli*, but other contaminants may be more difficult to identify. In the present paper I  
64 highlight an additional problem, namely, that host genomes can easily be contaminated with parasite ones,  
65 and thus careful isolation of genetic material and careful bioinformatics analysis are needed in all cases.

66 Apicomplexan parasites cause serious illnesses in humans and domestic animals. Most members of  
67 the phylum Apicomplexa are obligate parasites, with some members being causative agents for diseases in  
68 vertebrates. Species in the genus *Plasmodium* cause malaria, from which over 1 million people die each year.  
69 Members of the apicomplexan families Babesiidae, Theileriidae, Eimeriidae, Sarcocistidae, and  
70 Cryptosporiidae are responsible for numerous infectious diseases in wild and domesticated animals, such as  
71 coccidiosis and babesiosis, resulting in significant economic burden for animal husbandry.

72 One of the apicomplexan families, belonging to the class Coccidia, is the Sarcocystidae. The  
73 members of the genera *Besnoitia*, *Hammondia*, and *Sarcocystis* have obligatory two-host predator-prey life  
74 cycle: asexual stages (sarcocysts) develop in the muscles of the intermediate hosts (prey); ingestion of  
75 muscle sarcocysts through predation or scavenging by the definitive host (predator) propagates the life cycle,  
76 and sexual multiplication takes place in its small intestine that results in sporocyst shedding in feces (Dubey  
77 et al., 1988). Other families of the Sarcocystidae, such as *Toxoplasma*, *Neospora*, and *Cystoisospora* can  
78 complete their life cycle using only one host (Wünschmann et al., 2010).

79 A new member of the Sarcocystidae family has been recently identified, namely, *Nephroisospora*  
80 *eptesici*, from the big brown bat (*Eptesicus fuscus*) (Wünschmann et al., 2010). It is the only known member

81 of the *Nephroisospora* genus, which is similar to *Besnoitia*, *Toxoplasma* and *Hammondia* species. Bats are  
82 known to host almost every kind of apicomplexan parasite (except gregarines which are known to parasitize  
83 only invertebrates): Sarcocystidae (Cabral et al., 2013; Dodd et al., 2014), Eimeriidae (McAllister et al.,  
84 2011; Afonso et al., 2014), Cryptosporidae (Wang et al., 2013), Haemosporidia (Duval et al., 2012; Schaer et  
85 al., 2013).

86 In searches of genomes and sequence data that have become recently available, I found that a  
87 sequence similar to a characteristic apicomplexan protein, apicortin (Orosz, 2009, 2011), is present in the  
88 whole genome shotgun (WGS) sequence of a bird, *Colinus virginianus* (bobwhite), the draft sequence of  
89 which has recently been published (Halley et al., 2014). This is a very surprising finding since in higher level  
90 (Eumetazoa) animals no apicortins have been found to date. It can be either the result of horizontal gene  
91 transfer, or it is, more probably, a contamination. Thus I decided to systematically investigate this problem:  
92 sequences of the characteristic apicomplexan organelle, the apicoplast, were used as queries in BLASTN  
93 search against nucleotide sequences of various animal groups, looking for possible contaminations. I found  
94 that, indeed, the latter case is valid; moreover, further vertebrate genomes are contaminated with  
95 apicomplexan sequences. Additionally, based on the contamination of a recently published bat genome  
96 (Zhang et al., 2013), I suggest the existence of a second member of the *Nephroisospora* genus hosted also by  
97 a bat, *Myotis davidii*.

98

## 99 **2. Methods**

### 100 ***2.1. Database similarity search and phylogenetic analysis***

101

102 Accession numbers of protein and nucleotide sequences refer to the NCBI GenBank database. The  
103 database search was started with an NCBI blast search using the sequences of known apicortin proteins as  
104 queries. BLASTP or TBLASTN analyses (Altschul et al., 1997) were performed on protein or nucleotide  
105 sequences available at the NCBI website. Then the whole nucleotide sequences of known apicomplexan  
106 apicoplasts were used as queries. BLASTN analysis (Altschul et al., 1997) was performed on nucleotide  
107 sequences, including expressed sequenced tags, TSAs (Transcriptome Shotgun Assembly) and WGSs,

108 available at the NCBI website. In further analyses, the *18S ribosomal RNA* and the *internal transcribed*  
109 *spacer 1 (ITS-1)* genes of Sarcocystidae were used as queries against the *C. virginianus* and *M. davidii*  
110 nucleotide sequences using BLASTN.

111 Multiple alignments of protein and nucleotide sequences were carried out by the Clustal Omega  
112 program (Sievers et al., 2011) and were manually refined. Multiple sequence alignments used for  
113 constructing phylogenetic trees are provided in Supplementary Data S1. The alignments were subjected to  
114 Bayesian phylogenetic analysis with the software MrBayes v.3.1.2 (Ronquist and Huelsenbeck, 2003).  
115 Default priors and the GTR model (Tavaré, 1986) including a proportion of invariant sites and a gamma-  
116 shaped distribution of variable sites with four rate categories (GTR+ $\Gamma_{(4)}$ +I) were used. Four chains were run  
117 up to  $2.4 \times 10^6$  generations, with a sampling frequency of 0.01, and the first 25% of the generations were  
118 discarded as burn-in. The tree was drawn using the program Drawgram.

119 The Phylip (Phylogeny Inference Package, version 3.696) program package (Felsenstein, 2008) was  
120 used to build a Maximum Likelihood (ML) tree with bootstrap values. One thousand datasets were generated  
121 using the program Seqboot from the original data, i.e. the multiple alignments done by Clustal Omega. This  
122 was followed by running the program Dnaml (DNA Maximum Likelihood) on each of the datasets in the  
123 group, using the same rate heterogeneity model as above ( $\Gamma_{(4)}$ +I). The values for the gamma distribution  
124 were taken from the Bayesian analysis. A consensus tree (from all the 1000 trees) was generated using the  
125 program Consense. The trees were drawn using the program Drawgram.

126

### 127 **3. Results and Discussion**

128

#### 129 ***3.1. Apicortin is present in the C. virginianus WGS sequence***

130

131 *C. virginianus* putative apicortin is very similar to the apicortins of the Sarcocystidae, *T. gondii*, *N.*  
132 *caninum* and *Hammondia hammondii* (Fig. 1). However, it shows the highest identity with a WGS sequence  
133 of *Sarcocystis neurona* (Fig. 1), the draft (not fully annotated) genome of which has been reported very  
134 recently (Blazejewski et al., 2015). *S. neurona*, an apicomplexan pathogen that cycles in nature between its

135 definitive host, Virginia opossum (*Didelphis virginiana*), and a broad range of mammals and birds (orders  
136 Passeriforme and Psittaciformes) as intermediate hosts, causes equine protozoal myeloencephalitis, a  
137 neurologic disease of horses (Dame et al., 1995). Although reports showing that *C. virginianus* is a host for  
138 *S. neurona* are not yet known, the geographical identity of its habitats with that of *D. virginiana*, the  
139 definitive host for *S. neurona*, makes it reasonable that similarly to other birds, *C. virginianus* is also an  
140 intermediate host of the parasite. (*C. virginianus* belongs to the order Galliformes, which are also known to  
141 be *Sarcocystis* hosts (Odening, 1998).) However, although the similarity between the *C. virginianus* and *S.*  
142 *neurona* potential apicortin-coding sequences is very high (the identity is above 90%), it is not high enough  
143 to suggest that the sequence found in *C. virginianus* genome is a contamination originated from *S. neurona*.  
144 Rather, it is probably a contamination from another species of the *Sarcocystis* genus, the genome of which  
145 has not been sequenced yet.

146

### 147 **3.2. Genes of apicoplast origin are present in the *C. virginianus* WGS sequence**

148

149 Most apicomplexan parasites possess an apicoplast, a plastid with no photosynthetic ability, which is  
150 essential for cell survival (Arisue and Hashimoto, 2015). The apicoplast has its own genome that mainly  
151 encodes the transcription and translation related genes necessary for plastid gene expression. There are six  
152 independent entries at the NCBI web page for the query “*Sarcocystis* + apicoplast” as nucleotide sequences,  
153 five of them for RNA polymerase beta subunit-like (RPOb) gene from various *Sarcocystis* species, and one  
154 of them for small subunit ribosomal RNA gene from *Sarcocystis muris*. Using these sequences as BLASTN  
155 queries against avian nucleotide sequences of the NCBI GenBank database, there are no hits but those of the  
156 *C. virginianus* WGS sequence. Any hit would be very improbable; however, the identities are extremely  
157 high, 94-98% (Table 1). Examples of hits are shown in Supplementary Data S2. All the hits are in duplicate,  
158 since Halley et al. (2014) produced a simple *de novo* (i.e. no scaffolding) and a scaffolded *de novo* assembly,  
159 and both of them were deposited in GenBank.

160 The complete genemap of several apicoplast genomes are available as listed in Arisue and Hashimoto  
161 (2015). These genomes commonly encode rRNAs, tRNAs, ribosomal proteins, bacterial-type RNA

162 polymerase subunits, EF-Tu and ClpC protein. The closest relative of *Sarcocystis* species is *T. gondii*,  
163 belonging to the same family, Sarcocystidae. Thus I carried out a BLASTN search in the NCBI GenBank  
164 databases against avian nucleotide sequences using the complete 35 kb sequence of *T. gondii* apicoplast  
165 (GenBank U87145). Again, there were no hits but those of the *C. virginianus* WGS sequence; moreover,  
166 with high identity values (Table 2, Supplementary Data S3).

167 The apicoplast genomes of *P. falciparum*, *E. tenella* and *T. gondii* have an inverted repeat region,  
168 which contains duplicated SSU and LSU rRNA genes as well as duplicated tRNA genes (Arisue and  
169 Hashimoto, 2015). Genes are bi-directionally encoded in the genome, with half of the circle being mainly  
170 transcribed in a clockwise direction and the other half generally transcribed counter clockwise. The presence  
171 of the inverted repeat region is the reason that the “Total score” is about twice the “Max. score” in the first  
172 two lines of Table 2: both the clockwise and counter clockwise *T. gondii* sequences give hits on the *C.*  
173 *virginianus* contigs in Plus/Minus and Plus/Plus directions (cf. AWGU01003450.1 and AWGT01002297.1  
174 in Supplementary Data S3).

175 Although there are no available complete apicoplast sequences for *Sarcocystis* spp. in the database yet,  
176 a significant part of the *S. neurona* apicoplast sequence can be identified by using the *T. gondii* apicoplast  
177 sequence as a query. The JAQE01002351.1 (24 001 bp) of the *S. neurona* WGS sequence covers about two-  
178 thirds of the *T. gondii* apicoplast sequence, while the JAQE01002350.1 (5698 bp) corresponds to the  
179 inverted repeat region. Using these sequences as queries against avian nucleotide sequences, only the *C.*  
180 *virginianus* WGS sequence produced hits, with even higher cover and identity values than the *T. gondii*  
181 query (Table 3, Supplementary Data S4, 5). Beside the contigs identified in Tables 1 and 2 (DeNovo\_contigs  
182 9595, 3450, 45101, 73205), further hits were found; the full sequences of the six contigs, where the score is  
183 higher than 1000, seem to originate from *Sarcocystis* contamination (cf. Supplementary Data S4, 5).

184

### 185 **3.3. A general method for fast identification of genome contaminations of apicomplexan origin**

186

187 The available apicoplast genomes include *T. gondii* (GenBank U87145; KE138841), *E. tenella* (Cai  
188 et al., 2003), *Theileria parva* (Gardner et al., 2005), *B. bovis* (Brayton et al., 2007), *Leucocytozoon caulleryi*

189 (GenBank NC\_022667), *Plasmodium falciparum* (Wilson et al., 1996) and several other *Plasmodium*  
190 species. Thus, I carried out a BLASTN search (Altschul et al., 1997) in the NCBI GenBank databases against  
191 nucleotide sequences (including ESTs, TSAs and WGSs) of various animal groups (mammals, birds, insects  
192 etc.) using the complete 29-35 kb sequences of the above mentioned apicoplasts as queries, looking for  
193 possible contaminations. As expected, there were very few hits; they included, beside *C. virginianus*, several  
194 contigs of the *M. davidii* and *Ovis aries musimon* WGS sequences; moreover, with high identity values. The  
195 highest identities were found when the *T. gondii* apicoplast was used as a query. In the case of the unfinished  
196 genome of *O. aries musimon* (contig 124943, 5098 bp, GenBank CBYI010124943.1), it was obvious that the  
197 source of the contamination is *T. gondii* since the sequence identity is practically 100% (Supplementary Data  
198 S6). For *M. davidii*, the 17 contigs with 299–4986 bp length are shown in Table 4 and Supplementary Data  
199 S7. Since *T. gondii* is a member of the Sarcocystidae family of coccidian parasites, it is a reasonable  
200 suggestion that the source of contamination found in *M. davidii* also belongs to this family.

201

#### 202 **3.4. Which apicomplexan species may cause the contaminations?**

203

204 The question arises whether the Sarcocystidae species causing the contamination of the bird and bat  
205 genomes can be identified. Unfortunately, no complete *Sarcocystis* genomes except that of *S. neurona* have  
206 been sequenced; however, the sequences of two characteristic genes, often used for phylogenetic analysis,  
207 are known in many cases. These are the *18S ribosomal RNA* and the *internal transcribed spacer 1 (ITS-1)*  
208 genes. I used these genes of *Sarcocystis* species as queries in BLASTN search against the *C. virginianus*  
209 WSG sequences, and identified the best hits. The *18S ribosomal RNA* gene is rather conservative; e.g. other  
210 avian sequences, including *C. virginianus* itself, show higher than 70% identity with the *Sarcocystis* genes.  
211 However, there are no hits above 80% identity, and thus the difference is high enough to differentiate  
212 between contamination and endogenous sequences. On the contrary, the *ITS-1* gene differs much more  
213 among the various species.

214 In the other case (*M. davidii*), not the genus but only the family can be identified on the basis of the  
215 BLAST search using apicoplast sequences as queries. Thus the highly divergent *ITS-1* sequences cannot be  
216 used for search, and only the *18S ribosomal RNA* genes are suitable for the analysis.

217 Table 5 and Supplementary Table S1 show the identity of the *18S ribosomal RNA* gene of various  
218 *Sarcocystis* species with the best hit in the *C. virginianus* WGS sequence. (The hits can be found on the  
219 AWGU01001200 and AWGT01000758 contigs.) The identity is always higher than 90%, and in some cases  
220 higher than 99%. The highest identity values are produced, not surprisingly, by *Sarcocystis* species of avian  
221 hosts (Table 5). The highest identities (99.8%) were found in the case of *S. albifronsi* (GenBank:  
222 EU502868.2) and *S. anasi* (GenBank: EU553477.2). The difference is only 3 nucleotides in both cases. The  
223 third most similar sequence is that of the *S. rileyi* (GenBank: KJ396583.1; 99.7%), where the difference is 6  
224 nucleotides. It is worth noting that the sources of the isolation, the intermediate hosts of the parasites, belong  
225 to the order Anseriformes (e.g. ducks and geese) in these three cases. Anseriformes and Galliformes form the  
226 clade Galloanserae. *C. virginianus* belongs to Odontophoridae, which is one of the families of Galliformes.  
227 Although there are examples where a species belonging to Galliformes hosts a *Sarcocystis* species (Wenzel  
228 et al., 1982; Odening, 1998; Chen et al., 2012), there are no sequence data available in these cases. Thus  
229 Anseriformes are the closest relatives of *C. virginianus* for which these data can be used.

230 Much fewer hits were found when the *ITS-1* genes were used as queries. Besides *S. albifronsi*  
231 (GenBank: JN256122.1; 87%) and *S. anasi* (GenBank: JF520779.2; 88%) only a sequence of *Sarcocystis*  
232 AEM-2014a (GenBank: KJ810609.1; 80%), hosted by *Gallinula chloropus*, was recognised in the BLASTN  
233 search as similar to the *C. virginianus* sequence, the 4300-3367 nucleotides of AWG1001200.1  
234 SimpleDeNovo\_contig\_1200. When I used this corresponding *C. virginianus* sequence as query in BLASTN  
235 search (in general, not only against Sarcocystidae nucleotide sequences), only the above mentioned three  
236 *Sarcocystis* nucleotide sequences were obtained as hits.

237 A phylogenetic tree of *Sarcocystis* spp was constructed by Bayesian analysis using the available *ITS-1*  
238 nucleotide sequences, with *S. lutrae* sequence as an outgroup (Fig. 2). The corresponding nucleotide  
239 sequence of SimpleDeNovo\_contig\_1200 of *C. virginianus* was also involved. The tree fits to the known  
240 phylogeny (Kutkieni  et al., 2012; Prakas et al., 2013); e.g., *S. wobeseri* is sister to *S. calchasi*, they are sister  
241 to *S. columbae* and all of them are sister to *S. cornixi* + *S. sp* ex *Accipter nisus*. Similarly, it was also found

242 that *S. anasi* and *S. albifronsi* are sister to *S. rileyi*, and they are sister to (*S. neurona* +*S. falcatula*) + *S.*  
243 *lindsayi* whose definitive hosts are opossums, while a wide range of birds are known as intermediate hosts  
244 (cf. Fig. 2). (It was suggested by Dame et al (1995) that *S. neurona* is identical with *S. falcatula*. The  
245 sequence identity of the *ITS-1* region is 97.3%, one of the highest values of the known sequences.) The  
246 position of the sequence contaminating *C. virginianus* is clearly within the clade represented by the above  
247 mentioned Anseriformes-hosted *Sarcocystis* species (*S. anasi*, *S. albifronsi*, *S. rileyi*). Both the other two  
248 species in this clade (*S. atraii* and *S. AEM-14a*) were isolated from hosts belonging to the order Gruiformes.

249         There are at least 30 *Sarcocystis* species that are known to be hosted by birds (Odening, 1998;  
250 Kutkiené et al., 2012). Granivorous, insectivorous, and omnivorous birds serve as intermediate hosts, and  
251 carnivorous birds are usually definitive hosts. Most of the avian orders are known to be infected by  
252 *Sarcocystis* species. Unfortunately, the majority of investigations into bird *Sarcocystis* lack molecular data or  
253 are characterized to the genus only (“*Sarcocystis* sp.”) (Kutkiené et al., 2012). As Fig. 2 shows, in  
254 accordance with the results of others (Kutkiené et al., 2012; Prakas et al., 2013), *Sarcocystis* species from  
255 birds form two groups. In one of them, the “upper group” on Fig. 2, predatory birds are thought to be  
256 definitive hosts (Olias et al., 2011; Prakas et al., 2013). In the “lower group” with *S. anasi* and others,  
257 mammalian definitive hosts are known: Arctic fox (*Alopex lagopus*) for *S. albifronsi* (Kutkiené et al., 2012);  
258 and striped skunk (*Mephitis mephitis*), red fox (*Vulpes vulpes*) and raccoon dogs (*Nyctereutes procyonoides*)  
259 for *S. rileyi* (Odening, 1998; Prakas et al., 2015). Thus, the *Sarcocystis* species hosted by *C. virginianus* as  
260 intermediate host probably has a mammalian definitive host as well. It is known that some *Sarcocystis*  
261 species are not host specific and can parasitize a wide range of hosts. For example, *S. neurona* is hosted not  
262 only by horses and various birds as originally thought (Dame et al. 1995), but its host range expands to  
263 raccoons, cats, skunks, a variety of mustelids, pinnipeds, cetaceans and more recently sea otters, harbour  
264 seals, and harbour porpoises (Blazejewski et al., 2015). In the case of bird-hosted parasites, *S. wobeseri* is  
265 hosted by mallard duck (*Anas platyrhynchos*) and white-fronted goose (*Anser albifrons*) from the order  
266 Anseriformes as well as by herring gull (*Larus argentatus*) from the order Charadriiformes (Kutkiené et al.,  
267 2010; Prakas et al., 2011). Thus it could be possible that a known *Sarcocystis* species is shared by *C.*  
268 *virginianus* and another bird of the order Galliformes or even by a bird of Anseriformes, Gruiformes or  
269 Charadriiformes. However, we can state at this moment only that there is no *Sarcocystis* species with a



270 known genome or partial sequence which can be identified as being responsible for the contamination of the  
271 *C. virginianus* genome. The contamination could be due to the fact that for isolation of the genomic DNA,  
272 skeletal muscle derived from the legs of a bobwhite was utilized (Halley et al., 2014), and sarcocysts develop  
273 in the muscles of the intermediate hosts.

274

### 275 **3.5. *In silico* identification of the second member of the *Nephroisospora* genus**

276

277 In the case of *M. davidii*, using the *18S ribosomal RNA* genes of Sarcocystidae as queries against the  
278 *M. davidii* WGS sequence, the hits can be found on contig111512 (GenBank: ALWT01111512.1). The  
279 identity is always higher than 95%, and in one case it is even higher than 99%. The highest identity (99.3%;  
280 14 bp difference) was found in the case of *N. eptesici* (GenBank: EU334134.1).

281 The identities of the *18S ribosomal RNA* genes of various Sarcocystidae species with the 211-2019  
282 nucleotides of contig111512 of the *M. davidii* WGS sequence are shown in Table S6. The genes listed there  
283 were used for the construction of a phylogenetic tree of *18S ribosomal RNA* genes of Sarcocystidae family  
284 by Bayesian and maximum likelihood analysis involving some Eimeriidae genes and using *Goussia*  
285 *balatonica* as outgroup (Fig. 3). The corresponding nucleotide sequence of contig111512 of *M. davidii* was  
286 also involved. (I used two phylogenetic methods since some branches of the tree were weakly supported.) In  
287 accordance with previous phylogenies (Morrison et al., 2004; Wünschmann et al., 2010; Matsubayashi et al.,  
288 2011), the Sarcocystidae genes are well-separated from Eimeriidae (*Eimeria*, *Lankesterella*, *Goussia*) genes;  
289 and within the Sarcocystidae clade *Sarcocystis* (including *Frenkelia*) is sister to all other genera in both trees.  
290 The *Cystoisospora/Isospora* (Barta et al., 2005) and the *Besnoitia* genera as well as the  
291 *Toxoplasma/Hammondia/Neospora* group form evidently three different clades, but their mutual  
292 relationships are not clear since this part of the tree is weakly supported in the cases of both methods.  
293 (According to the above mentioned three analyses, *Cystoisospora/Isospora* forms the sister to *Besnoitia* +  
294 *Toxoplasma/Hammondia/Neospora*.) However, most importantly, the sequence from *M. davidii* forms a  
295 clade with the *N. eptesici* gene in both trees with maximal support. *N. eptesici* is the only known member of  
296 the *Nephroisospora* genus, from the most common bat, *E. fuscus* (Wünschmann et al., 2010). It was found in

297 the kidney of the big brown bat and it was concluded that it is closely related to *Besnoitia*, *Hammondia*,  
298 *Neospora* and *Toxoplasma*. Similarly to *Toxoplasma* and *Neospora*, it can complete its life cycle using only  
299 one host. Moreover, the entire cycle is completed in the kidney of the single host (Wünschmann et al., 2010).  
300 The source of the genetic material used for genome sequencing of *M. davidii* was not given by Zhang et al.  
301 (2013). However, in the GenBank web-page  
302 (<http://www.ncbi.nlm.nih.gov/Traces/wgs/?val=ALWT&display=contigs&search=ALWT01000000>), where  
303 the nucleotide sequences were deposited by the authors of the paper, it is stated that the kidney (beside  
304 spleen and small intestine) of the animal was used. The only nucleotide sequence established for *N. eptesici*  
305 is the *18S ribosomal RNA* gene, and thus no more data are available for comparison. However, the  
306 phylogenetic position of the tentative species, the bat host and the renal localization of the parasite make it  
307 reasonable to suggest that the contamination found in *M. davidii* originates from a new species of the  
308 *Nephroisospora* genus.

309 Finally, I hope that this paper highlights the general problem, namely, that hosts genomes can be  
310 easily contaminated with parasite ones thus careful isolation of genetic material and careful bioinformatic  
311 analysis are needed in all cases.

312

313 **Acknowledgments:** The author thanks Dr. Judit Oláh for the careful reading of the manuscript as well as the  
314 unknown reviewers of the paper for their suggestions.

315

316 **References**

317

- 318 Afonso, E., Baurand, P.E., Tournant, P., Capelli, N., 2014. First amplification of *Eimeria hessei* DNA from  
319 the lesser horseshoe bat (*Rhinolophus hipposideros*) and its phylogenetic relationships with *Eimeria*  
320 species from other bats and rodents. *Exp. Parasitol.* 139, 58-62. doi: 10.1016/j.exppara.2014.02.013.
- 321 Altschul, S.F., Madden, T.L., Schäffer, A.A., Zhang, J., Zhang, Z., Miller, W., Lipman, D.J., 1997. Gapped  
322 BLAST and PSI-BLAST: a new generation of protein database search programs. *Nucleic Acids Res.*  
323 25, 3389–3402.
- 324 Arisue, N., Hashimoto, T., 2015. Phylogeny and evolution of apicoplasts and apicomplexan parasites.  
325 *Parasitol. Int.* 64, 254-259. doi: 10.1016/j.parint.2014.10.005.
- 326 Barta, J.R., Schrenzel, M.D., Carreno, R., Rideout, B.A., 2005. The genus *Atoxoplasma* (Garnham 1950) as a  
327 junior objective synonym of the genus *Isospora* (Schneider 1881) species infecting birds and  
328 resurrection of *Cystoisospora* (Frenkel 1977) as the correct genus for *Isospora* species infecting  
329 mammals. *J. Parasitol.* 91, 726–727.
- 330 Blazejewski, T., Nursimulu, N., Pszenny, V., Dangoudoubiyam, S., Namasivayam, S., Chiasson, M.A.,  
331 Chessman, K., Tonkin, M., Swapna, L.S., Hung, S.S., Bridgers, J., Ricklefs, S.M., Boulanger, M.J.,  
332 Dubey, J.P., Porcella, S.F., Kissinger, J.C., Howe, D.K., Grigg, M.E., Parkinson, J., 2015. Systems-  
333 based analysis of the *Sarcocystis neurona* genome identifies pathways that contribute to a  
334 heteroxenous life cycle. *MBio.* 6, e02445-14. doi: 10.1128/mBio.02445-14.
- 335 Brayton, K.A., Lau, A.O., Herndon, D.R., Hannick, L., Kappmeyer, L.S., Berens, S.J., Bidwell, S.L., Brown,  
336 W.C., Crabtree, J., Fadrosch, D., Feldblum, T., Forberger, H.A., Haas, B.J., Howell, J.M., Khouri, H.,  
337 Koo, H., Mann, D.J., Norimine, J., Paulsen, I.T., Radune, D., Ren, Q., Smith, R.K. Jr., Suarez, C.E.,  
338 White, O., Wortman, J.R., Knowles, D.P. Jr., McElwain, T.F., Nene, V.M., 2007. Genome sequence  
339 of *Babesia bovis* and comparative analysis of apicomplexan hemoprotozoa. *PLoS Pathog.* 3, 1401–  
340 1413. doi:10.1371/journal.ppat.0030148
- 341 Cabral, A.D., Gama, A.R., Sodré, M.M., Savani, E.S., Galvão-Dias, M.A., Jordão, L.R., Maeda, M.M., Yai,  
342 L.E., Gennari, S.M., Pena, H.F., 2013. First isolation and genotyping of *Toxoplasma gondii* from bats  
343 (Mammalia: Chiroptera). *Vet. Parasitol.* 193, 100-104. doi: 10.1016/j.vetpar.2012.11.015.

344 Cai, X., Fuller, A.L., McDougald, L.R., Zhu, G., 2003. Apicoplast genome of the coccidian *Eimeria tenella*.  
345 Gene 321, 39–46.

346 Chen, X., He, Y., Liu, Y., Olias, P., Rosenthal, B.M., Cui, L., Zuo, Y., Yang, Z., 2012. Infections with  
347 *Sarcocystis wenzeli* are prevalent in the chickens of Yunnan Province, China, but not in the flocks of  
348 domesticated pigeons or ducks. *Exp. Parasitol.* 131, 31–34.

349 Dame, J.B., MacKay, R.J., Yowell, C.A., Cutler, T.J., Marsht, A., Greiner, E.C., 1995. *Sarcocystis falcatula*  
350 from passerine and psittacine birds: synonymy with *Sarcocystis neurona* agent of equine protozoal  
351 myeloencephalitis. *J. Parasitol.* 81, 930-935.

352 Dodd, N.S., Lord, J.S., Jehle, R., Parker, S., Parker, F., Brooks, D.R., Hide, G., 2014. *Toxoplasma gondii*:  
353 prevalence in species and genotypes of British bats (*Pipistrellus pipistrellus* and *P. pygmaeus*). *Exp.*  
354 *Parasitol.* 139, 6-11. doi: 10.1016/j.exppara.2014.02.007.

355 Dubey, J.P., Speer, C.A., Fayer, R., 1988. Sarcocystosis of animals and man. CRC Press, Boca Raton, FL

356 Duval, L., Mejean, C., Maganga, G.D., Makanga, B.K., Mangama Koumba, L.B., Peirce, M.A., Arie, F.,  
357 Bourgarel, M., 2012. The chiropteran haemosporidian *Polychromophilus melanipherus*: a worldwide  
358 species complex restricted to the family Miniopteridae. *Infect. Genet. Evol.* 12, 1558-1566. doi:  
359 10.1016/j.meegid.2012.06.006.

360 Felsenstein, J., 2008. PHYLIP (Phylogeny Inference Package) version 3.696. Department of Genome  
361 Sciences and Department of Biology University of Washington, Seattle, WA.  
362 <http://evolution.genetics.washington.edu/phylip.html>

363 Gardner, M.J., Bishop, R., Shah, T., de Villiers, E.P., Carlton, J.M., Hall, N., Ren, Q., Paulsen, I.T., Pain, A.,  
364 Berriman, M., Wilson, R.J., Sato, S., Ralph, S.A., Mann, DJ., Xiong, Z., Shallom, S.J., Weidman, J.,  
365 Jiang, L., Lynn, J., Weaver, B., Shoaibi, A., Domingo, A.R., Wasawo, D., Crabtree, J., Wortman, J.R.,  
366 Haas, B., Angiuoli, S.V., Creasy, T.H., Lu, C., Suh, B., Silva, J.C., Utterback, T.R., Feldblyum, T.V.,  
367 Pertea, M., Allen, J., Nierman, W.C., Taracha, E.L., Salzberg, S.L., White, O.R., Fitzhugh, H.A.,  
368 Morzaria, S., Venter, J.C., Fraser, C.M., Nene, V., 2005. Genome sequence of *Theileria parva*, a  
369 bovine pathogen that transforms lymphocytes. *Science* 309, 134–137.

370 Halley, Y.A., Dowd, S.E., Decker, J.E., Seabury, P.M., Bhattarai, E., Johnson, C.D., Rollins, D., Tizard, I.R.,  
371 Brightsmith, D.J., Peterson, M.J., Taylor, J.F., Seabury, C.M., 2014. A draft *de novo* genome assembly  
372 for the northern bobwhite (*Colinus virginianus*) reveals evidence for a rapid decline in effective

373 population size beginning in the late Pleistocene. PLoS One 9, e90240.  
374 doi:10.1371/journal.pone.0090240

375 Kutkienė, L., Prakas, P., Sruoga, A., Butkauskas, D. 2010. The mallard duck (*Anas platyrhynchos*) as  
376 intermediate host for *Sarcocystis wobeseri* sp. nov. from the barnacle goose (*Branta leucopsis*)  
377 Parasitol. Res. 107, 879-888.

378 Kutkienė, L., Prakas, P., Sruoga, A., Butkauskas, D., 2012. Description of *Sarcocystis anasi* sp. nov. and  
379 *Sarcocystis albifronsi* sp. nov. in birds of the order Anseriformes. Parasitol. Res. 110, 1043–1046.

380 Longo, M.S., O’Neill, M.J., O’Neill, R.J., 2011. Abundant human DNA contamination identified in non-  
381 primate genome databases. PLoS ONE 6, e16410; doi: 10.1371/journal.pone.0016410.

382 Matsubayashi, M., Carreno, R.A., Tani, H., Yoshiuchi, R., Kanai, T., Kimata, I., Uni, S., Furuya, M., Sasai,  
383 K., 2011. Phylogenetic identification of *Cystoisospora* spp. from dogs, cats, and raccoon dogs in  
384 Japan. Vet. Parasitol. 176, 270–274.

385 McAllister, C.T., Burt, S., Seville, R.S., Robison, H.W., 2011. A new species of *Eimeria* (Apicomplexa:  
386 Eimeriidae) from the eastern pipistrelle, *Perimyotis subflavus* (Chiroptera: Vespertilionidae), in  
387 Arkansas. J. Parasitol. 97, 896-898. doi: 10.1645/GE-2761.1.

388 Merchant, S., Wood, D.E., Salzberg, S.L., 2014. Unexpected cross-species contamination in genome  
389 sequencing projects. PeerJ 2, e675; doi: 10.7717/peerj.675

390 Morrison, D.A., Bornstein, S., Thebo, P., Wernery, U., Kinne, J., Mattsson, J. G., 2004. The current status of  
391 the small subunit rRNA phylogeny of the coccidia (Sporozoa). Int. J. Parasitol. 34, 501-514.

392 Odening, K. 1998. The present state of species-systematics in *Sarcocystis* Lankester, 1882 (Protista,  
393 Sporozoa, Coccidia). Syst. Parasitol. 41, 209–233.

394 Olias, P., Olias, L., Krücken, J., Lierz, M., Gruber, A.D., 2011. High prevalence of *Sarcocystis calchasi*  
395 sporocysts in European Accipiter hawks. Vet. Parasitol. 175, 230–236.

396 Orosz, F., 2009. Apicortin a unique protein with a putative cytoskeletal role shared only by apicomplexan  
397 parasites and the placozoan *Trichoplax adhaerens*. Infect. Genet. Evol. 9, 1275–1286.

398 Orosz, F., 2011. Apicomplexan apicortins possess a long disordered N-terminal extension. Infect. Genet.  
399 Evol. 11, 1037–1044.

- 400 Prakas, P., Kutkienė, L., Butkauskas, D., Sruoga, A., Zalakevičius, M., 2013. Molecular and morphological  
401 investigations of *Sarcocystis corvusi* sp. nov. from the jackdaw (*Corvus monedula*). Parasitol Res.  
402 112, 1163-1167. doi: 10.1007/s00436-012-3247-5.
- 403 Prakas, P., Kutkienė, L., Sruoga, A., Butkauskas, D., 2011. *Sarcocystis* sp. from the herring gull (*Larus*  
404 *argentatus*) identity to *Sarcocystis wobeseri* based on cyst morphology and DNA results. Parasitol.  
405 Res. 109, 1603-1608. doi: 10.1007/s00436-011-2421-5.
- 406 Prakas, P., Liaugaudaitė, S., Kutkienė, L., Sruoga, A., Švažas, S., 2015. Molecular identification of  
407 *Sarcocystis rileyi* sporocysts in red foxes (*Vulpes vulpes*) and raccoon dogs (*Nyctereutes*  
408 *procyonoides*) in Lithuania. Parasitol. Res. 114, 1671-1676. doi: 10.1007/s00436-015-4348-8
- 409 Ronquist, F., Huelsenbeck, J.P., 2003. MrBayes 3: Bayesian phylogenetic inference under mixture models.  
410 Bioinformatics 19, 1572–1574.
- 411 Schaer, J., Perkins, S.L., Decher, J., Leendertz, F.H., Fahr, J., Weber, N., Matuschewski, K., 2013. High  
412 diversity of West African bat malaria parasites and a tight link with rodent *Plasmodium* taxa. Proc.  
413 Natl. Acad. Sci. U.S.A. 110, 17415-17419. doi: 10.1073/pnas.1311016110.
- 414 Sievers, F., Wilm, A., Dineen, D., Gibson, T.J., Karplus, K., Li, W., Lopez, R., McWilliam, H., Remmert,  
415 M., Söding, J., Thompson, J.D., Higgins, D.G., 2011. Fast scalable generation of high-quality protein  
416 multiple sequence alignments using Clustal Omega. Mol. Syst. Biol. 7, 539. doi: 10.1038/msb.2011.75
- 417 Tavare, S. 1986. Some probabilistic and statistical problems on the analysis of DNA sequences. Lect. Math.  
418 Life Sci. 17, 57-86.
- 419 Wang, W., Cao, L., He, B., Li, J., Hu, T., Zhang, F., Fan, Q., Tu, C., Liu, Q., 2013. Molecular  
420 characterization of *Cryptosporidium* in bats from Yunnan province, southwestern China. J. Parasitol.  
421 99, 1148-1150. doi: 10.1645/13-322.1.
- 422 Wenzel, R., Erber, M., Boch, J., Schellner, H.P., 1982. Sarcosporidia infections in domestic fowl, pheasant  
423 and coot. Berl. Munch. Tierarztl. Wochenschr. 95, 188–193.
- 424 Wilson, R.J., Denny, P.W., Preiser, P.R., Rangachari, K., Roberts, K., Roy, A., Whyte, A., Strath, M.,  
425 Moore, D.J., Moore, P.W., Williamson, D.H., 1996. Complete gene map of the plastid-like DNA of  
426 the malaria parasite *Plasmodium falciparum*. J. Mol. Biol. 261, 155–172.
- 427 Wünschmann, A., Wellehan, J.F. Jr., Armien, A., Bemrick, W.J., Barnes, D., Averbeck, G.A., Roback, R.,  
428 Schwabenlander, M., D'Almeida, E., Joki, R., Childress, A.L., Cortinas, R., Gardiner, C.H., Greiner,

429 E.C., 2010. Renal infection by a new coccidian genus in big brown bats (*Eptesicus fuscus*). J.  
430 Parasitol. 96, 178-183.

431 Zhang, G., Cowled, C., Shi, Z., Huang, Z., Bishop-Lilly, KA., Fang, X., Wynne, J.W., Xiong, Z., Baker,  
432 M.L., Zhao W., Tachedjian, M., Zhu, Y., Zhou, P., Jiang, X., Ng, J., Yang, L., Wu, L., Xiao, J., Feng,  
433 Y., Chen, Y., Sun, X., Zhang, Y., Marsh, G.A., Cramer, G., Broder, C.C., Frey, K.G., Wang, L.F.,  
434 Wang, J., 2013. Comparative analysis of bat genomes provides insight into the evolution of flight and  
435 immunity. Science, 339, 456-460.

436

437 **Legends to Figures**

438 **Fig. 1. Multiple alignment of hypothesized apicortin sequences of Sarcocystidae and *C. virginianus*.**

439 Identical and similar amino acids in all cases are labelled by asterisks and colons, respectively. Grey  
440 background indicates whether the amino acids are identical only in *Sarcocystis neurona* and *C. virginianus*.  
441 Hammondia: *H. hammondi*, XP\_008888750; Neospora: *N. caninum*, NCLIV\_029060; Toxoplasma: *T.*  
442 *gondii* XP\_002364910; Colinus: *C. virginianus* AWGU01108821; Sarcocystis: *S. neurona* JAQE01000498  
443 (JXWP01000002).

444 **Fig. 2. Bayesian phylogenetic tree of the *Sarcocystis* spp. based on sequences of the ITS-1 region.** The  
445 tree was rooted with *S. lutrae*. Posterior probability support was maximal in all cases. The GenBank  
446 accession numbers of ITS-1 genes are given for each taxon and underlined if it was isolated from bird-hosted  
447 species. For *Colinus virginianus* the 4300-3367 nucleotides of AWG1001200.1 SimpleDeNovo\_contig\_1200  
448 were used.

449 **Fig. 3. Phylogenetic trees of the Sarcocystidae 18S rRNA genes.** Three genes of members of the  
450 Eimeriidae were also included, using *G. balatonica* as the outgroup. The GenBank accession numbers are  
451 given in Supplementary Data S9. For *M. davidii* the 211-2019 nucleotides of contig111512 (GenBank:  
452 ALWT01111512.1) were used. The position of *N. eptesici* gene and *M. davidii* contig are indicated in both  
453 cases. (A) Bayesian analysis. The branch lengths indicate the inferred amount of evolutionary change,  
454 according to the scale bar shown. Bayesian posterior probability values are shown for the main branches of  
455 the tree. Black circles represent full support for the node. (B) Maximum Likelihood tree. Confidence of the  
456 tree topology obtained is shown by maximum likelihood bootstrap values calculated from 1000 replicates.  
457



458 **Table 1. Hits in the *C. virginianus* WGS sequence when using as queries the**  
 459 **available nucleotide sequences from *Sarcocystis* apicoplasts**

Species	GenBank	gene	<i>C. virginianus</i>	Identity
<i>S. campestris</i>	GQ851963	RPOb	AWGU01009595	94%
			AWGT01005489	
<i>S. canis</i>	KC191642	RPOb	AWGU01009595	94%
			AWGT01005489	
<i>S. falcatula</i>	GQ851962	RPOb	AWGU01009595	94%
			AWGT01005489	
<i>S. neurona</i>	GQ851961	RPOb	AWGU01009595	95%
			AWGT01005489	
<i>S. sp.</i>	KC191641	RPOb	AWGU01009595	95%
			AWGT01005489	
<i>S. muris</i>	AF255924	small subunit ribosomal RNA	AWGU01003450	98%
			AWGT01002297	

460 RPOb - RNA polymerase beta subunit-like gene. The AWGU and AWGT sequence IDs stand

461 for simple *de novo* (i.e. no scaffolding) and scaffolded *de novo* contigs, respectively.

462

463 **Table 2. Hits in the *C. virginianus* WGS sequence when using as a query the *T. gondii* apicoplast**  
 464 **complete genome (GenBank U87145.2)**

Description	Max score	Total score	Query cover	E value	Identity	Accession
SimpleDeNovo_contig_3450	6835	13510	30%	0.0	89%	AWGU01003450.1
ScaffoldedDeNovo_contig_2297	6835	13510	30%	0.0	89%	AWGT01002297.1
SimpleDeNovo_contig_9595	6630	6630	23%	0.0	82%	AWGU01009595.1
ScaffoldedDeNovo_contig_5489	6630	6630	23%	0.0	82%	AWGT01005489.1
SimpleDeNovo_contig_45101	3236	4781	23%	0.0	77%	AWGU01045101.1
ScaffoldedDeNovo_contig_42584	3236	4781	23%	0.0	77%	AWGT01042584.1
SimpleDeNovo_contig_73205	2663	2663	11%	0.0	80%	AWGU01073205.1
ScaffoldedDeNovo_contig_42585	2663	2663	11%	0.0	80%	AWGT01042585.1
SimpleDeNovo_contig_47460	198	397	0%	1e-46	94%	AWGU01047460.1
ScaffoldedDeNovo_contig_11590	198	397	0%	1e-46	94%	AWGT01011590.1
ScaffoldedDeNovo_contig_54661	91.6	183	0%	2e-14	95%	AWGT01054661.1
SimpleDeNovo_contig_66517	71.3	142	0%	3e-08	95%	AWGU01066517.1
ScaffoldedDeNovo_contig_21304	71.3	142	0%	3e-08	95%	AWGT01021304.1
SimpleDeNovo_contig_137357	69.4	138	0%	1e-07	98%	AWGU01137357.1
ScaffoldedDeNovo_contig_33333	69.4	138	0%	1e-07	98%	AWGT01033333.1
SimpleDeNovo_contig_257380	63.9	127	0%	4e-06	95%	AWGU01257380.1

465

466 **Table 3. Hits in the *C. virginianus* WGS sequence when using as queries the *S. neurona* potential**  
 467 **apicoplast nucleotide sequences**

Description	MaxScore	TotalScore	Query cover	E value	Identity	Accession
<i>Query: JAQE01002351 (24001 bp)</i>						
SimpleDeNovo_contig_45101	15005	15005	42%	0.0	93%	AWGU01045101.1
SimpleDeNovo_contig_9595	12759	12759	34%	0.0	95%	AWGU01009595.1
SimpleDeNovo_contig_73205	5775	5775	16%	0.0	93%	AWGU01073205.1
SimpleDeNovo_contig_206802	1341	1341	3%	0.0	93%	AWGU01206802.1
SimpleDeNovo_contig_78021	1194	1194	3%	0.0	95%	AWGU01078021.1
SimpleDeNovo_contig_70976	154	154	0%	2e-33	89%	AWGU01070976.1
SimpleDeNovo_contig_68381	121	121	0%	2e-23	99%	AWGU01068381.1
SimpleDeNovo_contig_282453	110	110	0%	4e-20	97%	AWGU01282453.1
SimpleDeNovo_contig_270464	86.1	86.1	0%	6e-13	95%	AWGU01270464.1
SimpleDeNovo_contig_106377	84.2	84.2	0%	2e-12	92%	AWGU01106377.1
SimpleDeNovo_contig_56616	75.0	75.0	0%	1e-09	91%	AWGU01056616.1
SimpleDeNovo_contig_168981	67.6	67.6	0%	2e-07	95%	AWGU01168981.1
SimpleDeNovo_contig_6309	56.5	56.5	0%	5e-04	86%	AWGU01006309.1
<i>Query: JAQE01002350 (5698 bp)</i>						
SimpleDeNovo_contig_3450	9788	10121	99%	0.0	99%	AWGU01003450.1
SimpleDeNovo_contig_151619	244	443	5%	2e-61	95%	AWGU01151619.1
SimpleDeNovo_contig_135140	137	270	4%	4e-29	88%	AWGU01135140.1
SimpleDeNovo_contig_158980	110	110	1%	9e-21	94%	AWGU01158980.1
SimpleDeNovo_contig_257380	69.4	69.4	0%	2e-08	98%	AWGU01257380.1
SimpleDeNovo_contig_133852	69.4	69.4	0%	2e-08	95%	AWGU01133852.1
SimpleDeNovo_contig_57220	69.4	124	1%	2e-08	98%	AWGU01057220.1
SimpleDeNovo_contig_39406	60.2	60.2	0%	9e-06	100%	AWGU01039406.1
SimpleDeNovo_contig_226848	56.5	56.5	0%	1e-04	97%	AWGU01226848.1
SimpleDeNovo_contig_110794	54.7	54.7	0%	4e-04	94%	AWGU01110794.1

468 ScaffoldedDeNovo\_contig hits are not shown here. See Supplementary Data S4 and S5 for them.

469

470 **Table 4. Hits in the *M. davidii* WGS sequence when using as a query the *T. gondii* apicoplast complete**  
 471 **genome (GenBank U87145.2)**

Contig	Max score	Total score	BLAST E value	Identity	GenBank Accession	Length, bp
314352	3493	11197	0.0	88%	ALWT01314352.1	4986
309146	1395	1395	0.0	87%	ALWT01309146.1	1290
307364	1375	1375	0.0	89%	ALWT01307364.1	1162
304579	924	924	0.0	84%	ALWT01304579.1	1027
295304	922	922	0.0	89%	ALWT01295304.1	749
297403	880	880	0.0	87%	ALWT01297403.1	796
293953	822	822	0.0	87%	ALWT01293953.1	721
297566	817	817	0.0	85%	ALWT01297566.1	801
293275	769	769	0.0	86%	ALWT01293275.1	709
274256	490	490	3e-132	86%	ALWT01274256.1	450
257569	370	370	4e-96	89%	ALWT01257569.1	305
266684	348	348	2e-89	88%	ALWT01266684.1	376
256340	339	339	1e-86	87%	ALWT01256340.1	297
293156	337	1348	4e-86	84%	ALWT01293156.1	706
256707	287	287	4e-71	84%	ALWT01256707.1	299
262086	276	276	8e-68	84%	ALWT01262086.1	338
261058	267	267	5e-65	82%	ALWT01261058.1	331

472

473 **Table 5. *Sarcocystis* 18S ribosomal RNA genes showing highest identity to *C. colinus* contamination**

Species	GenBank	Identity, %	bp	Source of isolation (host)	Order (birds)
<i>S. albifronsi</i>	EU502868	99.8	1792	<i>Anser albifrons</i>	Anseriformes
<i>S. anasi</i>	EU553477	99.8	1792	<i>Anas platyrhynchos</i>	Anseriformes
<i>S. arctica</i>	KF601301	99.0	1803	<i>Vulpes lagopus</i>	
<i>S. arctosi</i>	EF564590	99.3	1484	<i>Ursus arctos</i>	
<i>S. atraii</i>	KJ810606	99.3	1493	<i>Fulica atra</i>	Gruiformes
<i>S. calchasi</i>	GQ245670	99.3	1804	<i>Columba livia</i>	Columbiformes
<i>S. canis</i>	DQ146148	99.5	994	<i>Canis canis</i>	
<i>S. columbae</i>	HM125054	99.3	1765	<i>Columba palumbus</i>	Columbiformes
<i>S. cornixi</i>	EU553478	98.7	1795	<i>Corvus cornix</i>	Passeriformes
<i>S. corvusi</i>	JN256117	99.3	1792	<i>Corvus monedula</i>	Passeriformes
<i>S. dispersa</i>	AF120115	98.6	1610	<i>Tyto alba</i>	Strigiformes
<i>S. lutrae</i>	KM657770	99.3	1804	<i>Lutra lutra</i>	
<i>S. mucosa</i>	AF109679	99.2	1824		
<i>S. muris</i>	SARRR16S	98.0	1809		
<i>S. neurona</i>	U07812	98.9	1803	<i>Bos taurus</i>	
<i>S. neurona</i>	HQ709144	99.1	763	<i>Martes pennant</i>	
<i>S. neurona</i>	JXWP01000699	98.4	1384	<i>Enhydra lutris nereis</i>	
<i>S. rileyi</i>	KJ396583	99.7	1803	<i>Somateria mollissima</i>	Anseriformes
<i>S. rileyi Europa</i>	HM185742	99.7	1792	<i>Anas platyrhynchos</i>	Anseriformes
<i>S. sp</i>	KM362427	98.9	1669	<i>Canis familiaris</i>	
<i>S. sp</i>	KF309699	99.3	1749	<i>Eothenomys miletus pocok</i>	
<i>S. sp</i>	GU253884	99.4	1630	<i>Accipiter nisus</i>	Falconiformes

<i>S. sp.</i>	JQ733511	99.4	1801	<i>Phalacrocorax carbo</i>	Pelecaniformes
<i>S. sp.</i>	JQ733508	99.4	1802	<i>Larus marinus</i>	Charadriiformes
<i>S. sp.</i>	AF513487	98.9	1594	<i>Sorex araneus</i>	
<i>S. sp.</i>	KF278953	99.2	1448	<i>Didelphis virginiana</i>	
<i>S. tupaia</i>	FJ827486	99.1	1311	<i>Tupaia belangeri</i> <i>chinensis</i>	
<i>S. turdusi</i>	JF975681	99.3	1793	<i>Turdus merula</i>	Passeriformes
	HM159419			<i>Larus argentatus</i>	
<i>S. wobeseri</i>	GQ922885	99.4	1792	<i>Anas platyrhynchos</i>	Charadriiformes
	EU502869			<i>Branta leucopsis</i>	Anseriformes
				<i>Anser albifrons</i>	

475 **Table 6. : Sarcocystidae 18S ribosomal RNA genes showing highest identity to *M. davidii***  
 476 **contamination**

Species	Max score	Query cover	E value	Identity	GenBank
<i>Nephroisospora eptesici</i>	3212	98%	0.0	99%	EU334134
<i>Besnoitia besnoiti</i>	3134	100%	0.0	98%	EU789637
<i>Toxoplasma gondii</i>	3116	100%	0.0	98%	M97703
<i>Neospora caninum</i>	3112	99%	0.0	98%	U16159
<i>Cystoisospora timoni</i>	3081	99%	0.0	98%	AY279205
<i>Besnoitia jellisoni</i>	3077	99%	0.0	98%	AF291426
<i>Isospora belli</i>	3075	99%	0.0	98%	DQ060661
<i>Hammondia hammondi</i>	3035	97%	0.0	98%	AF096498
<i>Hammondia truffittae</i>	3031	97%	0.0	98%	GQ984222
<i>Hammondia heydorni</i>	3031	97%	0.0	98%	JX220986
<i>Cystoisospora belli</i>	3005	97%	0.0	98%	AB268326
<i>Isospora suis</i>	2990	99%	0.0	97%	U97523
<i>Cystoisospora sp.</i>	2968	96%	0.0	98%	AB519674
<i>Sarcocystis sp. ex Phalacrocorax carbo</i>	2944	99%	0.0	96%	JQ733511
<i>Sarcocystis sp. ex Columba livia</i>	2944	100%	0.0	96%	GQ245670
<i>Sarcocystis arctica</i>	2942	99%	0.0	96%	KF601301
<i>Sarcocystis mucosa</i>	2940	99%	0.0	96%	AF109679
<i>Sarcocystis lutrae</i>	2933	99%	0.0	96%	KM657769
<i>Sarcocystis sp. ex Corvus monedula</i>	2933	99%	0.0	96%	JN256117
<i>Isospora felis</i>	2933	97%	0.0	97%	L76471
<i>Sarcocystis rileyi</i>	2931	99%	0.0	96%	KJ396583
<i>Sarcocystis sp. ex Larus marinus</i>	2929	99%	0.0	96%	JQ733508
<i>Cystoisospora ohioensis</i>	2929	96%	0.0	97%	GU292304
<i>Sarcocystis anasi</i>	2922	99%	0.0	96%	EU553477
<i>Sarcocystis albifronsi</i>	2916	99%	0.0	96%	EU502868
<i>Sarcocystis sp. ex Columba palumbus</i>	2872	97%	0.0	96%	HM125054
<i>Besnoitia darlingi</i>	2826	91%	0.0	98%	GU479631
<i>Besnoitia oryctofelisi</i>	2822	91%	0.0	98%	GU479632
<i>Besnoitia bennetti</i>	2804	90%	0.0	98%	AY665399
<i>Frenkelia glareoli</i>	2678	90%	0.0	96%	AF009245
<i>Frenkelia microti</i>	2639	90%	0.0	96%	AF009244
<i>Hyaloklossia lieberkuehni</i>	2580	87%	0.0	96%	AF298623
<i>Isospora ohioensis</i>	2357	77%	0.0	97%	AF029303
<i>Eimeria variabilis</i>	2102	98%	0.0	88%	<u>GU479674</u>
<i>Lankesterella minima</i>	1925	93%	0.0	87%	<u>AF080611</u>
<i>Goussia balatonica</i>	2396	94%	0.0	92%	<u>GU479650</u>

477 Query: *Myotis davidii*, contig111512, 211-2019 nucleotides (GenBank: ALWT01111512)

478 Subject: Sarcocystidae nucleotide sequences

479 The last three species belong to the Eimeridiidae, the nearest relative of Sarcocystidae.

480

481 **Supplementary data**

482

483 **Supplementary Data S1.** Multiple sequence alignments used for constructing phylogenetic trees of Fig. 2  
484 and Fig. 3.

485 **Supplementary Data S2.** BLASTN search used *S. neurona* RNA polymerase beta subunit (RPOb)  
486 (GenBank: GQ851961.1) and *Sarcocystis muris* small subunit ribosomal RNA gene (GenBank: AF255924.1)  
487 genes as queries against *C. virginianus* WGS nucleotide sequences.

488 **Supplementary Data S3.** Hits in the *C. virginianus* WGS sequence when using as a query the *T. gondii*  
489 apicoplast complete genome (U87145.2)

490 **Supplementary Data S4.** Hits in the *C. virginianus* WGS sequence when using as a query a *S. neurona*  
491 potential apicoplast nucleotide sequence (GenBank JAQE01002351).

492 **Supplementary Data S5.** Hits in the *C. virginianus* WGS sequence when using as a query a *S. neurona*  
493 potential apicoplast nucleotide sequence (GenBank JAQE01002350).

494 **Supplementary Data S6.** Hits in the *Ovis aries musimon* WSG sequence when using as a query the *T.*  
495 *gondii* apicoplast complete genome

496 **Supplementary Data S7.** Sequences producing significant alignments by BLASTN search. Query: *T. gondii*  
497 apicoplast, complete genome (GenBank: U87145.2). Subject: *M. davidii*, WGS sequence.

498 **Supplementary Table S1.** Further *Sarcocystis* 18S ribosomal RNA genes

499



Figure

Hammondia -----MACGIPWKL----ARRDELMATRQAERPGEYFPPPYPPCPPTVVMPLRTSA  
Toxoplasma -----MACGIPWKL----ARRDELMATRQAERPGEYFPPPYPPCPPTVVTPLRRTSA  
Neospora -----MACGIPWKL----ARRDELMEARQAERQGDYFPPPYPPCPPTVVMPLRTTA  
Sarcocystis MEIFISLLDLICGCKLDFCVGVTAQKTVVWGTGFEEETGIHFFPPPYPPCPPTVVTPFDVSE  
Colinus LEISISSLPLC\*QRYLCTNVTQKTDWGAVSEETGIRFFPPPYPPCPPTVATPFDVSE  
\* : : : : \* :\*\*\*\*\*. \* : .:

Hammondia YDFPEATFVTRPCLPAKKATGHKNVFERLTDYAYTGSHRERFDEFNGNGRIAGREYLY  
Toxoplasma YDFPEATFVTRPCLPAKKATGHKNVFERLTDYAYTGSHRERFDEFNGNGRIAGREYLY  
Neospora YDLPEATFVTRPCPRPRKATGHKNVFERLTDYAYTGSHRERFDEFNGNGRIAGREYLY  
Sarcocystis YDFPEASVLKAQRLATRPPTTRHRNVFDRLTDSQYYTGTTHRERFDEFGNGRIAGRECVY  
Colinus YDFPDASVLKAQRLVTRPRTRHRNVFDRLTDSQFYTGTHRERFDEFGNGRIAGRECVY  
\*\*.\*:\*.:. : \* \* :\*\*\*:\*\*\*: :\*\*\*:\*\*\*\*\*:\*\*\*\*\* :\*

Hammondia AYDGLTESPSRCHVEYSSVIKRPKPVVTPGTLGVQRFVQIPAPRLMWLYRNGDKHDD  
Toxoplasma AYDGLTESPSRCHVEYSSVIKRPKPVVTPGTLGIQRFVQIPAPRLMWLYRNGDKHDD  
Neospora AYDGLTESPSRCHVEYSSVVKRPRKPVVTPGTLGVQRFVQIPAPRLMWLYRNGDKHDD  
Sarcocystis TVDGFTESPSRSHVEYSSVIKRPKPVVTPGTLGIQRFVQIATPRLMWLYRNGDKHDD  
Colinus TVDGLTESPSRSHVEYSSVIKRPKPVVTPGTLGIQRFVQIATPRLMWLYRNGDKHDD  
: \*\*.\*\*\*\*\*\*.\*\*\*\*\*.\*:\*\*\*\*\*.\*\*\*\*\*.\*\*\*\*\*.\*\*\*\*\*.\*\*\*\*\*

Hammondia GTPFFVRPYIKSMEALYQQITKEITPIAGPVRRIFDQNFVITDLDIVDGAKYLCTSG  
Toxoplasma GTPFFVRPYIKSMESLYQQITKEITPIAGPVRRIFDQNFVITDLDIVDGAKYLCTSG  
Neospora GTPFFVRPYIKTMESLYQQITKEITPIAGPVRRIFDQNFVITDLDIVDGAKYLCTSG  
Sarcocystis GTPFFVRSFIRSMEALYQQISKITPIAGPVRRIFDQNFRLITNLEDIVDGAKYLCTSG  
Colinus GTPFFVRPFIKSMEGLYQQISKITPIAGPVRRIFDQNFRLITNLEDIVDGAKYLCTSG  
\*\*\*\*\* :\*:\*.\*.\*\*\*\*\*.\*:\*\*\*\*\*.\*\*\*\*\*.\*\*\*\*\*.\*\*\*\*\*.\*\*\*\*\*

Hammondia EPPAAYDRLEKFLSEWVIQKSQPKVPSQFFV  
Toxoplasma EPPAAYDRLEKFLSEWVIQKSQTKVPSQFFV  
Neospora EPPAAYDRLEKFLSEWVIQKSQTKVPSQFFV  
Sarcocystis EPPAAYDRLEKFLSEWVQKAYSQKVPSEFFIL  
Colinus EPPAAYDRLEKFLSEWVQKAYSQKVPSEFFIV  
\*\*\* \* :\*\*\*\*\*.\*\*\*: \* :\*\*\*:\*\*\*:\*

Figure

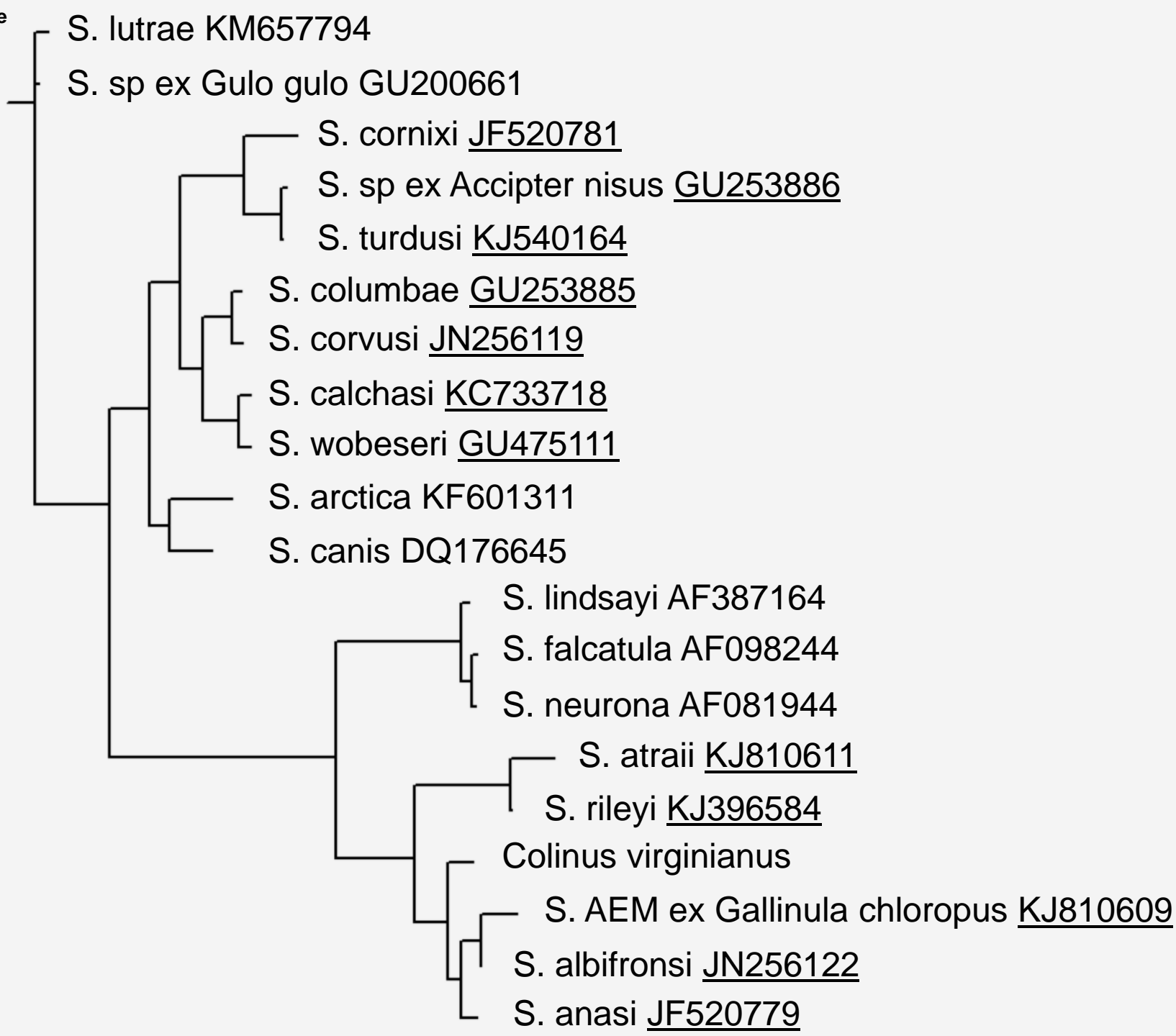


Figure A

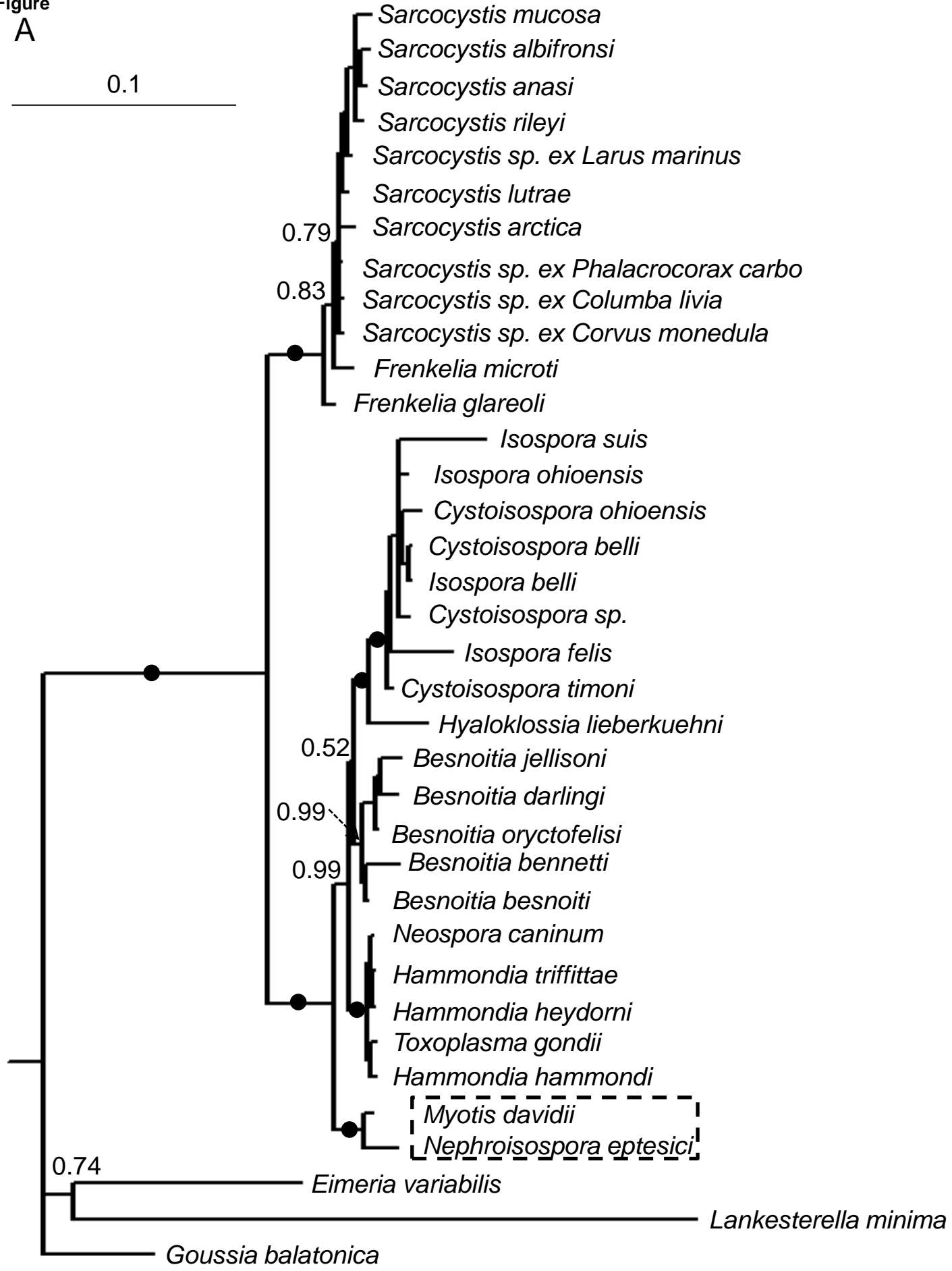
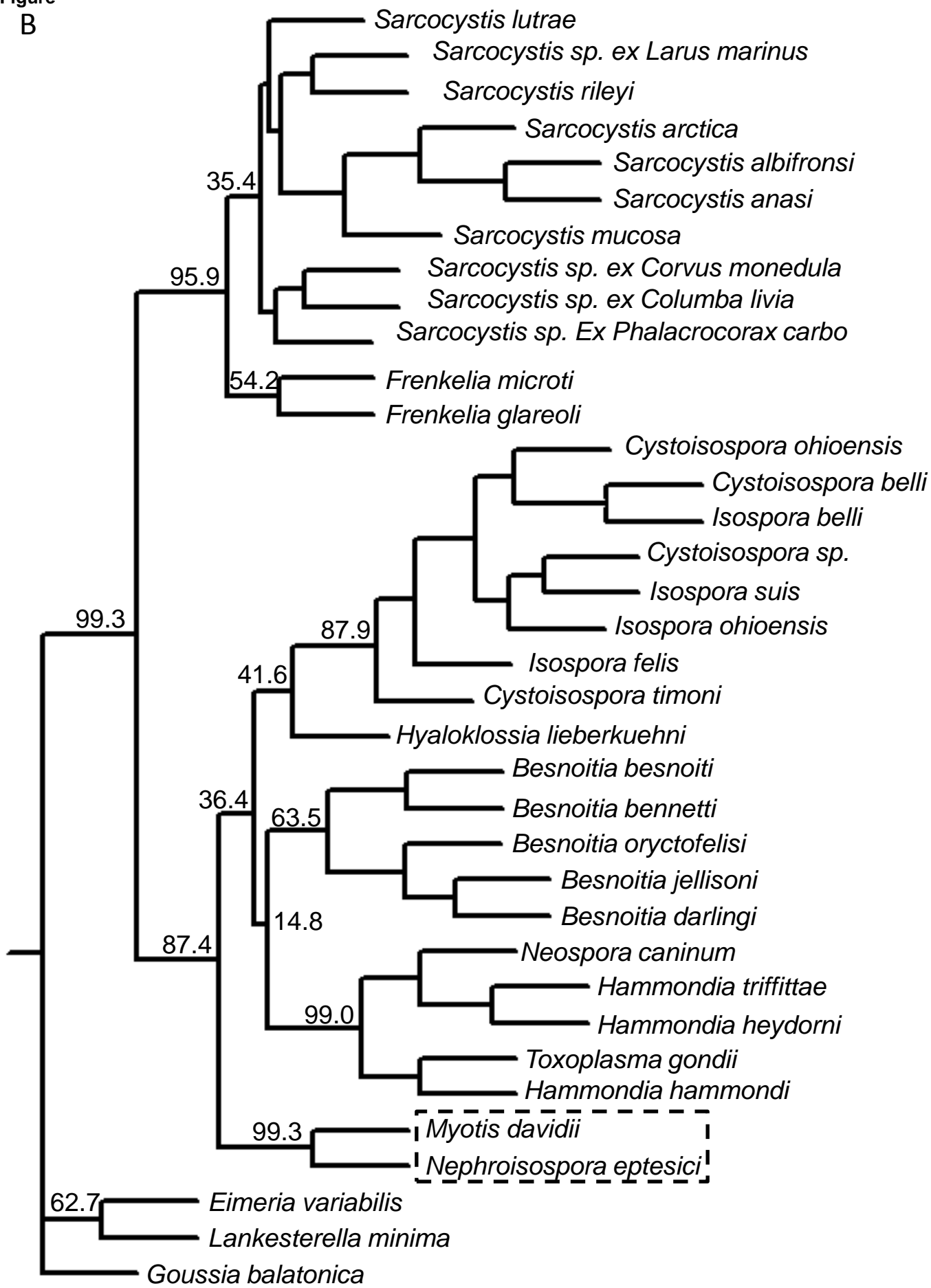


Figure B



**Multi-media supplement**

[Click here to download Multi-media supplement: SupplementaryDataS1.docx](#)

**Multi-media supplement**

[Click here to download Multi-media supplement: SupplementaryDataS2.docx](#)

**Multi-media supplement**

[Click here to download Multi-media supplement: SupplementaryDataS3.docx](#)

**Multi-media supplement**

[Click here to download Multi-media supplement: SupplementaryDataS4.docx](#)



**Multi-media supplement**

[Click here to download Multi-media supplement: SupplementaryDataS5.docx](#)

**Multi-media supplement**

**[Click here to download Multi-media supplement: SupplementaryDataS6.docx](#)**

**Multi-media supplement**

**[Click here to download Multi-media supplement: SupplementarydataS7.docx](#)**

**Multi-media supplement**

[Click here to download Multi-media supplement: SupplementaryTableS1.docx](#)