

1 **PPSDB : A Linked Open Data knowledge base for protist-**
2 **prokaryote symbiotic interactions**

3

4 Brandon K. B. Seah (<http://orcid.org/0000-0002-1878-4363>)

5 Thünen Institute for Biodiversity, Bundesallee 65, 38116 Braunschweig, Germany.

6 Email: brandon.seah@thuenen.de kb.seah@gmail.com

7

8 **Abstract**

9 As the ecological and evolutionary importance of symbiotic interactions between protists
10 (microbial eukaryotes) and prokaryotes (bacteria and archaea) is better appreciated, keeping
11 an overview of their diversity and the literature becomes a growing and ongoing challenge.
12 Here we present the Protist-Prokaryote Symbiosis Database (PPSDB), comprising 789
13 manually curated interaction statements sourced from 410 publications, where biological
14 taxonomy, anatomical localization, and analytical methods applied have been annotated and
15 mapped to external databases and ontologies, such as Wikidata, NCBI Taxonomy and Gene
16 Ontology. We describe how our data model deals practically with challenges such as
17 incomplete information and inconsistent taxon concepts, which will be applicable to similar
18 projects. Both the model and underlying Wikibase software platform are highly extensible, so
19 new items and properties can easily be added. Unlike a static table or list of citations, PPSDB
20 is a structured knowledge base that enables programmatic access and powerful, integrated
21 semantic queries. The database is available at <https://ppsdb.wikibase.cloud/>.

22

23 **Keywords:** linked open data, knowledge graph, FAIR, SPARQL, semantic web,
24 Wikibase.cloud

25 **Introduction**

26 The study of protists (microbial eukaryotes) has revealed a fascinating diversity of
27 interactions with prokaryotes, including symbionts that defend their hosts, help them move,
28 and orient them in the environment (Petroni et al. 2000; Hongoh et al. 2007; Monteil et al.
29 2019). New symbioses are regularly discovered, e.g. by mining protist genome data for
30 prokaryote sequences (Davison, Hurst, and Siozios 2023), while previously described ones
31 can be profitably revisited with modern methods. To contextualize new discoveries, and spot
32 larger trends and knowledge gaps, an accurate overview of the current body of knowledge is
33 indispensable.

34 The existing entry points to knowledge on the diversity of protist-prokaryote
35 interactions are research articles and reviews, e.g. (Ball 1969; Bjorbækmo et al. 2020; Husnik
36 et al. 2021; Kostygov et al. 2021; Shi et al. 2021; Fokin and Serra 2022). These are static,
37 usually not interlinked with primary data, and formatted for human readers rather than
38 programmatic queries, even if some information is presented in tabular form. It is therefore
39 difficult to get a reliable answer to a question such as “what protists have
40 Alphaproteobacteria symbionts localized in the host nucleus” without a deep dive into the
41 literature for oneself.

42 The study of biotic interactions also suffers from poor discoverability of relevant
43 information (Poelen, Simons, and Mungall 2014; Mihara et al. 2016). Ideally, databases
44 should be more than just a list of taxa or citations, but also capture other facets such as
45 phylogenetic affiliation, interaction type, and environmental context. To achieve this, they
46 should be built on a semantic data model that is extensible enough to accommodate new
47 information, concepts, and terminology as they arise.

48 Specific challenges are posed by the ever-evolving biological taxonomy and methods
49 used to identify and describe organisms. Better taxon sampling and methods constantly drive
50 updates to the higher taxonomy and nomenclature of both eukaryotes (Adl et al. 2019) and
51 prokaryotes (Parks et al. 2020). For taxa originally identified or described on the basis of
52 morphology or phenotype, their placement within modern taxonomies may be unclear,
53 although some symbioses have been revisited with sequencing to clarify their phylogenetic
54 identity (Boscaro et al. 2013; Schrallhammer, Castelli, and Petroni 2018). In contrast, many
55 recent studies assign names solely from sequence data. Many organisms, particularly
56 environmental microbes, never receive a formal scientific name and remain known under an
57 informal or provisional name even if otherwise well characterized. Therefore, in addition to

58 names and taxonomy, the analytical methods and evidence base behind each described
59 symbiotic interaction should also be documented.

60 For the uses envisioned above, we argue that the information is best managed in the
61 form of a knowledge graph. Knowledge graphs are data structures that represent
62 concepts/entities and the relationships between them abstractly as the nodes and edges of a
63 directed graph (network) (Chaudhri et al. 2022). A common way to specify a graph's
64 structure is as a collection of linkages, each comprising two nodes and the edge that connects
65 them. The meanings assigned to nodes and edges depend on the domain-specific application
66 of the knowledge graph: for symbioses, nodes can represent biological taxa, and edges their
67 interactions. Complex multi-way or nested relationships as well as incomplete information
68 can thus be represented more naturally and efficiently than in a tabular format or relational
69 database.

70 Database items can be assigned standardized identifiers (uniform resource identifier,
71 URI), and be further linked ("mapped") to equivalent entities in other databases, e.g. for taxa
72 or publications. URIs can be a web address that returns useful information about the entity.
73 Data published following such principles are known as Linked Open Data (LOD) (Bauer and
74 Kaltenböck 2012). Taxon names can be seen as an analogy to URIs within biology, as they
75 also aim to be (ideally) unique and language-agnostic identifiers for real-world entities.
76 Taxon names remain the primary vehicles through which biologists convey and retrieve
77 information about organisms, albeit intended for human use and recall (Patterson et al. 2010).
78 Machine-readable URIs fulfill a similar role but can be processed programmatically. Linking
79 equivalent concepts or entities between datasets with URIs lets us build on other databases
80 and avoid duplicated effort; for example, we do not need to curate a full biological taxonomy
81 in our own database if we map taxa to an existing, programmatically accessible taxonomic
82 database. Other concepts/entities, such as anatomical and environmental terms, can be
83 mapped to ontologies, such as the Gene Ontology (Gene Ontology Consortium et al. 2023)
84 and Environment Ontology (Buttigieg et al. 2016). This not only ensures that terminology is
85 used consistently with the wider community, but also allows sophisticated queries that take
86 advantage of the semantic relationships encoded in those ontologies (Pacheco et al. 2022).

87 Here, we describe a knowledge base for protist-prokaryote symbiotic interactions, and
88 showcase how Linked Open Data principles enable powerful, integrated searches across
89 multiple resources. The design objectives were to: (i) Represent curated information from the
90 scientific literature, with citations for each statement. (ii) Focus on named symbiotic

91 interaction partners from low-diversity systems, rather than microbiome studies dealing with
92 higher level taxa or OTUs. (iii) Link records to sequence databases. (iv) Enable multiple
93 entry points for queries, including biological taxonomy, anatomical localization of symbionts,
94 and analytical methods used to identify organisms. (iv) Map concepts and entities in the
95 database to external taxonomies, ontologies, and identifiers, to ensure that they are described
96 consistently and interoperable with other resources and knowledge representations.

97 **Methods**

98 *Software platform and tools*

99 The database was built on a Wikibase instance hosted by Wikibase.cloud, a service provided
100 by Wikimedia Deutschland. The database was edited through the web interface, through
101 batch edits using the QuickStatements tool
102 (<https://www.wikidata.org/wiki/Help:QuickStatements>), and programmatically with Python
103 scripts (<https://github.com/kbseah/ppsdb-utils/>) using the WikibaseIntegrator library v0.12.5
104 (<https://github.com/LeMyst/WikibaseIntegrator>). The data dump to XML was performed with
105 mediawiki-dump-generator ([https://github.com/mediawiki-client-tools/mediawiki-dump-](https://github.com/mediawiki-client-tools/mediawiki-dump-generator)
106 [generator](https://github.com/mediawiki-client-tools/mediawiki-dump-generator)). Periodic exports for indexing by Global Biotic Interactions (GloBI) are hosted on
107 GitHub (<https://github.com/kbseah/ppsdb-globi-export>).

108 *Data model and terminology*

109 In the Wikibase platform, nodes are called “items”, and the edges connecting them are
110 assigned specific meanings, or “properties”. Each connection (two items linked by a
111 property) makes up a “statement”, with one item as the subject and the other as the object
112 (Figure 1). Statements themselves can be treated like items and be the subject of further
113 “qualifier” statements that provide additional information. Statements can also be annotated
114 with references, which are a special type of qualifier statement.

115 Items may belong to one of two types, “classes” and “instances”, following the usage
116 in Wikidata (https://www.wikidata.org/wiki/Help:Basic_membership_properties). A class is a
117 set of items that have common properties; the members of a class are known as instances.
118 Classes may be further subdivided into subclasses. For example, “*Pelomyxa palustris*” is an
119 instance of the class “formally named taxon”, which is a subclass of “taxon”. Other classes in
120 PPSDB represent references, organismal body parts, analytical methods, environmental
121 terms, and interaction types; all items are ultimately descended from the root class “entity”.

122 We modeled each biotic interaction as a statement linking two taxon items with an “interacts
123 with” property (Figure 1, Figure 2). Each statement was further qualified with (i) where the
124 symbiont is localized in the host organism/cell, (ii) the analytical methods used to identify
125 (taxonomically and phylogenetically) host and symbiont, and (iii) the nature of the biotic
126 interaction (e.g. transfer of fixed organic carbon, pathogenic), if known. Further statements
127 on taxon items mapped them to external taxonomy databases and representative sequence
128 records, and described the environmental context from which organisms were isolated or
129 sampled (Figure 1, Figure 2).

130 *Data collation and mappings to external identifiers*

131 Reported symbiotic interactions between protists and prokaryotes were gathered from the
132 published literature through ad hoc keyword searches and relevant review articles. These
133 included studies that specifically focused on symbiosis, as well as morphological or
134 taxonomic studies that incidentally described associated microbes.

135 Relevant information was extracted from original research publications where
136 possible, and mapped to external identifiers if a suitable exact match existed (Table 1). Taxon
137 items were created to represent the interacting organisms; these are understood to be at
138 species rank or below, even if they are only identified to a higher ranked taxon, similar to the
139 concept of a “submittable taxon” in ENA ([https://ena-
140 docs.readthedocs.io/en/latest/faq/taxonomy.html](https://ena-docs.readthedocs.io/en/latest/faq/taxonomy.html)). This was to avoid potential conflation of
141 multiple taxa under the same identifier. If sequence data were available, the taxon was linked
142 to a representative sequence named in the publication describing it, to ensure that the name is
143 associated with empirical data should disambiguation be necessary in the future. Equivalent
144 identifiers in the NCBI Taxonomy (Schoch et al. 2020) were linked if available. Formally
145 described or *Candidatus* (for prokaryotes) taxa were linked to Wikidata and the List of
146 Prokaryotic Names with Standing in Nomenclature (LPSN) (Parte et al. 2020) on the basis of
147 taxon name.

148 The localization of symbionts in the host was mapped to cellular anatomy terms in the
149 Gene Ontology (GO) (Ashburner et al. 2000; Gene Ontology Consortium et al. 2023) or
150 metazoan anatomical terms in Uberon (for protists that are also symbionts of animals)
151 (Mungall et al. 2012). The relationships between anatomical terms are represented in
152 ontologies, which can be exploited when performing queries.

153 Cited publications were linked to digital object identifiers (DOIs) and to Wikidata,
154 where bibliographic data are maintained by a community project, Wikicite. Citations missing
155 from Wikidata are easily imported with the Scholia tool (Nielsen, Mietchen, and Willighagen
156 2017). Formatted citations were obtained from CrossRef from their DOIs, otherwise added
157 manually.

158 Three properties were used to describe the environmental context at different scales—
159 broad scale environmental context, local environmental context, and environmental material—
160 using Environment Ontology (EnvO) terms (Buttigieg et al. 2013, 2016), following the MIXS
161 guidelines (<https://github.com/EnvironmentOntology/envo/wiki/Using-ENVO-with-MIXS>)
162 (Yilmaz et al. 2011). If equivalent terms existed, analytical techniques were mapped to the
163 Ontology for Biomedical Investigations (OBI) (Bandrowski et al. 2016), and interaction types
164 were mapped to the OBO Relations Ontology (<https://github.com/oborel/obo-relations>).

165 *Challenges for data mapping and modeling*

166 Most challenges related to taxonomy, because informal or provisional names are often used
167 for microorganisms, and disparate methods and levels of detail have been used to characterize
168 them. Ideally, each organism would be described in a scientific publication under a formal
169 taxon name, accompanied by published sequence accessions and an equivalent NCBI taxon
170 item. A new Wikidata item for the taxon was created if it did not already exist. The PPSDB
171 item was then linked to the corresponding Wikidata and NCBI Taxonomy items. Here we
172 describe how we dealt with other cases that were not so neatly organized.

173 **Informal or provisional taxon names with NCBI Taxonomy equivalent**

174 Many studies described organisms without assigning a formal taxon name, but their
175 phylogenetic affiliation is nonetheless known, sequence data are available, and the
176 corresponding informally named taxon item in the NCBI Taxonomy appears to be equivalent
177 to the taxon concept in the study. If so, the item was mapped to that NCBI taxon ID, and
178 labeled with the informal name used in the cited publication, as well as known aliases from
179 other publications and databases.

180 **Taxon concept not in NCBI Taxonomy, but sequence data available**

181 Taxon items in the NCBI Taxonomy and taxonomic annotations of sequence records may not
182 be up to date, or may differ from the published literature. For example, the ciliate species
183 *Eufolliculina methanicola* (<https://ppbdb.wikibase.cloud/entity/Q52>) was formally described
184 in a scientific publication (Pasulka et al. 2017), but sequences from that study were published

185 in Genbank under a placeholder taxon "Folluculinidae sp." (NCBI:txid1934002), which is
186 used in the NCBI Taxonomy for records which were only identified to the family level and so
187 may represent a mixture of different species. The PPSDB item was therefore not mapped to
188 the NCBI Taxonomy, because it can lead to incorrect results if the identifier is used to
189 retrieve sequence data. For this example, a formal taxon name was published, so a Wikidata
190 item was created for it and mapped to PPSDB. To allow us to track the identity, should the
191 NCBI Taxonomy be updated in the future, a representative SSU rRNA sequence record for
192 this taxon that was cited in the original publication (KX012915) was linked to this item with
193 the property "representative SSU rRNA sequence record".

194 **Taxon name/concept with no sequence data available**

195 An organism may have been identified by morphology alone, without using sequencing
196 methods, or sequence data produced in a study cannot be found. For example, a species of
197 *Arcobacter* (<https://ppbdb.wikibase.cloud/entity/Q410>) was identified as a symbiont of
198 *Bihospites bacati* (<https://ppbdb.wikibase.cloud/entity/Q409>), but although sequencing of a
199 marker gene was reported, the sequence was not published. Alternatively, the organism may
200 have been identified to a higher taxonomic group by its morphology or with methods such as
201 group-specific molecular probes. For example, the ciliate *Frontonia leucas*
202 (<https://ppbdb.wikibase.cloud/entity/Q1782>) is associated with an unclassified
203 Alphaproteobacteria (<https://ppbdb.wikibase.cloud/entity/Q158>) that was identified with
204 group-specific molecular probes. As no direct sequence data were available in these cases,
205 there were no sequence records to anchor the taxon concept empirically (see above).
206 PPSDB items for such taxa were labeled with descriptive names based on what was reported
207 in the publication, e.g. "unclassified Alphaproteobacteria", but were not mapped to external
208 identifiers. Similarly named *incertae sedis* items may exist in the NCBI Taxonomy (e.g.
209 "unclassified Alphaproteobacteria", NCBI:txid33807), but these were deliberately not
210 mapped from PPSDB because they will pull in incorrect results if used for programmatic
211 queries.

212 **Consistent mapping of items to higher taxonomy**

213 Not all organisms described in the literature could be mapped to external taxonomies, nor
214 were they always identified to the species level. Nonetheless, we linked all taxon items to the
215 next-highest-ranking, formally named parent taxon that was represented in both Wikidata and
216 NCBI Taxonomy. This enabled consistent searches by taxonomy, even if the species

217 themselves were not mapped to an external taxonomy. The parent taxon items are instances
218 of a class “higher taxon” that is not used in interaction statements.

219 **Experimentally induced interactions**

220 A number of symbiotic microbes were first identified in one host species but maintained in
221 the laboratory in a different host because the original host was not suitable for experiments.
222 For example, *Acanthamoeba castellanii* has been used as a lab host for various intracellular
223 bacteria isolated from other amoebae (Schulz et al. 2016). These experimentally induced
224 interactions were represented with a different property, “interacts experimentally with”, to
225 distinguish them from naturally occurring interactions.

226 **Results and Discussion**

227 PPSDB is hosted by Wikibase.cloud and is browsable through the web interface at
228 <https://pppsdb.wikibase.cloud/>. A SPARQL endpoint is available for programmatic queries,
229 with examples to help users get started: <https://pppsdb.wikibase.cloud/query/>. The structured
230 data (in the Item: and Property: namespaces) are released under a CC0 1.0 public domain
231 dedication (<https://creativecommons.org/publicdomain/zero/1.0/>).

232 *Database statistics*

233 The database currently (16 Aug 2024) documents 789 biotic interactions between 498 host
234 taxa and 514 symbiont taxa, with 410 references cited. The number of citations is incidentally
235 similar to the 328 works cited by Gordon H. Ball in his 1969 review, “Organisms living on
236 and in Protozoa” (Ball 1969). However, given our focus on phylogenetic identity of the
237 symbiotic partners and linking them to sequence data, there is a bias towards more recent
238 publications in PPSDB and the overlap between the two sets of citations is minimal.

239 The most commonly represented host protist phyla are Ciliophora (136 taxon items),
240 Metamonada (78), and Amoebozoa (55), while the most commonly represented symbiont
241 prokaryotic phyla are *Pseudomonadota* (207 items), *Bacteroidota* (53), and
242 *Methanobacteriota* (36). This undoubtedly reflects the activity of researchers rather than the
243 abundance or ecological significance of these organisms. Some non-protist, non-prokaryote
244 taxa are represented, e.g. termite hosts of metamonad flagellates that themselves have
245 bacterial symbionts, as the host species helps to identify the flagellate. Multipartite
246 interactions, or highly nested ones, are easily modeled, e.g. the bacterial epibionts of
247 spirochaete ectosymbionts of flagellates from termite guts (Utami et al. 2019).

248 *Usage examples*

249 The SPARQL query engine bundled with Wikibase and the mappings to other databases,
250 particularly Wikidata, allow users to query PPSDB in ways that are not easily implemented
251 with other databases of similar scope. Further example SPARQL queries are listed at
252 <https://pppsdb.wikibase.cloud/wiki/Project:SPARQL/examples>.

253 1. Symbionts localized in the host nucleus and their class (<https://tinyurl.com/25fxfn9>)

254 This search showcases how a query can use semantic content in the database and
255 integrate an external database in the search.

256 The localization of intranuclear symbionts may be reported simply as “nucleus”, or
257 more specifically as “nuclear envelope lumen”, or “macronucleus” (in ciliates, which have
258 two developmentally distinct types of nuclei). The relationships between these terms are
259 modeled in the database, e.g. nuclear envelope lumen is a part of, and macronucleus is a
260 subclass of “nucleus”. Users can be more or less specific as required.

261 We have chosen not to maintain a full biological taxonomy within PPSDB, but
262 instead to map taxa to Wikidata and the NCBI Taxonomy (see “Challenges for data mapping
263 and modeling” above). The search is therefore executed as a federated SPARQL query across
264 both PPSDB and Wikidata. Most described intranuclear symbionts are *Alphaproteobacteria*
265 from ciliates, but there are diverse hosts where the symbiont’s phylogenetic position is
266 unknown.

267 2. Interactions connected to *Ca. Megaira polyxenophila* (<https://tinyurl.com/2948dund>)

268 *Ca. Megaira polyxenophila* is known to associate with diverse host eukaryotes
269 (Schrallhammer et al. 2013), and many of its known hosts in turn have more than one
270 prokaryotic symbiont. These complex linkages can be visualized as a graph, which shows
271 that some of these host species are also linked by other symbionts that have multiple host
272 species: *Polynucleobacter necessarius* and *Caedimonas varicaedens* (Figure 3).

273 3. Symbionts identified by fluorescence in situ hybridization but not sequencing
274 (<https://tinyurl.com/282s455g>)

275 A number of symbionts have been described in publications that employed group-
276 specific molecular probes that could identify them to e.g. class level, but which did not
277 sequence a phylogenetic marker gene, so a more precise classification was not possible.
278 These may be interesting to revisit with modern sequencing methods.

279 4. Symbioses described in publications by a specific author (<https://tinyurl.com/29gnxhw3>)

280 We can query bibliographic metadata of the publications referenced in PPSDB such
281 as authors and publication venues, via reference items mapped to Wikidata. Like the search
282 by biological taxonomy, this is a federated SPARQL query. The search exploits the growing
283 representation of publication and person data in Wikidata, which can be used for
284 scientometric studies, such as investigating coauthorship networks.

285 *Choice of software platform*

286 We chose Wikibase as the platform for this database because it has both a web browser-based
287 interactive interface and an API for programmatic access, and is available as a cloud service.
288 Wikibase was originally developed as the backend for Wikidata, the largest open knowledge
289 graph. As such, its design caters to the Wikidata model, but this potential limitation was
290 outweighed by its ease of use, active user community, ongoing support and development, and
291 integration of a SPARQL engine and other tools. Existing tools and libraries to work with
292 Wikibase can be applied instead of reinventing the wheel. The ease of federated searches
293 with Wikidata was also an advantage. Many current Wikibase users come from the cultural
294 heritage field (Diefenbach, Wilde, and Alipio 2021; Huaman, Huaman, and Huaman 2023;
295 Shimizu et al. 2023), and include institutions like the European Union and German National
296 Library. PPSDB shows that an application in the natural sciences is straightforward.

297 Technical requirements and know-how remain a hurdle to the adoption of knowledge
298 graphs. Wikibase.cloud is a good compromise for smaller projects and prototypes driven by
299 subject-matter experts who may not have a deep background in semantic web technologies.
300 No programming experience is required to get started, as data entry and editing can be
301 performed through the web interface, with users learning additional tools (e.g.
302 QuickStatements for tabular data entry, SPARQL for queries) as they go along. User
303 management, project planning, and discussion pages can be maintained on the same wiki as
304 the database itself, making it self-contained. Other projects with similar aims have built
305 bespoke software, e.g. AQUASYMBIO <http://www.aquasymbio.fr/>, Viral Host Range DB
306 <https://viralhostrangedb.pasteur.cloud/> (Lamy-Besnier et al. 2021), and Virus-Host DB
307 <https://www.genome.jp/virushostdb/> (Mihara et al. 2016). Such software is harder to maintain
308 in the long term, and requires more effort to integrate with other linked data sets. Similar
309 considerations have been cited by Wikibase users who have migrated from other platforms
310 (Koho et al. 2023).

311 *Virtuous cycles of data curation*

312 During data curation, we sometimes discovered outdated records or errors while mapping
313 items to NCBI Taxonomy and Wikidata. We edited Wikidata directly, while the NCBI
314 Taxonomy team was contacted by email with corrections. Commonly encountered issues
315 included NCBI Taxonomy records that still used a provisional name although a formal taxon
316 name or Candidatus name has been published, and taxon names or publications that were not
317 yet represented on Wikidata. Linked open data naturally fosters collaboration and a
318 mutualistic relationship between the linked resources, such that the curation and data cleaning
319 of one benefits the others too (Seah 2023; von Mering et al. 2023).

320 *Data sharing and archiving*

321 Even if data are linked and open, interested users may not be able to find them easily. The
322 core interaction data in PPSDB were therefore exported as a table for indexing by the Global
323 Biotic Interactions (GloBI) database, an aggregator for species interaction data that is
324 searchable from its web site and through an R package (Poelen, Simons, and Mungall 2014).
325 This increases the visibility of protists, which are underrepresented in ecological studies, and
326 of the original publications, which are cited in full. To secure the long term availability of the
327 database, periodic XML dumps are archived on Internet Archive. The export for GloBI is
328 also archived separately on Zenodo.

329 *Remaining challenges for data modeling/mapping*

330 Some types of statements in the literature remain difficult to represent formally in our data
331 model. Broad statements about higher taxa, e.g., “all *Kentrophoros* species are associated
332 with bacteria from genus *Candidatus Kentron*”, are represented in our data model by creating
333 individual items and links for each known species within those taxa. Such a statement implies
334 that as-yet unstudied or undescribed host species will also be found to interact with
335 corresponding symbiont species. However, we cannot create items for unknown species, so
336 we conservatively do not add such implicit statements.

337 The modeling of facets other than biological taxonomy is relatively basic and can be
338 further developed. For symbiont localization, we currently do not distinguish between
339 different types of topological relationships. For example, methanogenic endosymbionts are
340 typically located in the host cytoplasm close to hydrogenosomes, but this detail is not
341 captured by the single “subject body part” property. For interaction types, most terms have
342 not been mapped to the OBO Relations Ontology (RO), because many terms that are

343 meaningful to microbial ecologists, e.g. “syntrophy” and “auxotrophy”, do not appear in RO.
344 The outcome (e.g. mutualistic vs. parasitic) and function of many microbial interactions is
345 also unclear or only inferred. Nonetheless, more elaborate modeling of these aspects may be
346 overly complex for most users, so a simpler representation may be more useful.

347 Microbiome survey studies were excluded from the scope of PPSDB. However, some
348 larger protists are associated with diverse prokaryotes; some may be stable partnerships while
349 others are facultative. Sequencing surveys may reveal dozens of such interactions per host
350 species—should they all be included in the database?

351 Finally there are the practical hurdles we experienced when extracting relevant
352 information from publications. Taxon names and sequence accessions may be scattered in
353 different parts of a manuscript and its supplements or even across multiple publications,
354 different names may be used for the same organism at different times, methods may be
355 incompletely reported, and in a few (thankfully rare) cases, which symbiont belongs to which
356 host was not reported at all even though both are separately characterized. We suggest that
357 authors summarize biotic interaction results in tabular format where any sequence accessions
358 or identifiers are also directly listed.

359 *Future directions*

360 We have described how we model biotic interactions in a knowledge graph, and our solutions
361 to challenges such as taxonomic inconsistency, uncertainty, and the proliferation of names
362 and identifiers. The model can easily be adapted to represent other types and facets of biotic
363 interactions by adding new classes and properties. For example, a statement representing an
364 allelopathic interaction between two plant species could have a qualifier that links to the
365 phytochemical responsible, the latter represented as an item of class “chemical” and mapped
366 to databases like ChEBI. PPSDB itself can be extended to encompass other taxonomic
367 groups; viruses are particularly relevant as some giant viruses of protists were initially
368 thought to be bacterial symbionts (La Scola et al. 2003). As mentioned above, a limited
369 number of non-protist, non-prokaryote taxa are already represented in the database.

370 The current bottleneck in data curation is in discovering and parsing the relevant
371 scientific literature. Other projects have data-mined molecular sequence metadata, e.g. the
372 “host” tag in Genbank records, to produce large-scale species interaction datasets (Wardeh et
373 al. 2015; Albrycht et al. 2022). Such pipelines are most suitable for taxa, such as viruses, that
374 are routinely described with sequencing, where relevant metadata are generally accessible in

375 standardized form. We are however also interested in the historical literature, and in
376 symbioses described with other methods, such as microscopy, where data deposition and
377 metadata reporting are not yet as standardized as for sequencing. For the foreseeable future,
378 we still need humans to verify that the data are reliable, but this could be supplemented by
379 provisional interaction claims derived from data mining. Natural language processing could
380 also help screen for relevant publications in full text databases such as Europe PMC.

381 **Acknowledgements**

382 Many thanks are due to the Wikibase.cloud user community and the Wikibase.cloud team at
383 Wikimedia Deutschland (WMDE) for answering questions and support in troubleshooting; to
384 Jorrit Poelen for help with GloBI integration and discussions on name alignment; and to
385 WMDE for making Wikibase.cloud available in the open beta phase.

386 **Data and code availability**

387 Main URL for PPSDB: <https://ppsdb.wikibase.cloud/>
388 Export to tabular format for indexing in GloBI: [https://github.com/kbseah/ppsdb-globi-](https://github.com/kbseah/ppsdb-globi-export)
389 [export](https://github.com/kbseah/ppsdb-globi-export), archived on Zenodo: <https://doi.org/10.5281/zenodo.12687626>
390 Scripts for database maintenance: <https://github.com/kbseah/ppsdb-utils>, archived on Zenodo:
391 <https://doi.org/10.5281/zenodo.12805883>
392 XML export of the entire PPSDB database on Internet Archive (16 Aug 2024):
393 https://archive.org/details/wiki-ppsdbwikibasecloud_w

Mixotricha paradoxa (Q296) *Q-number (identifier)*

species of flagellates edit

▼ In more languages edit

Configure *multilingual labels, descriptions, aliases*

Language	Label	Description	Also known as
English	Mixotricha paradoxa	species of flagellates	

Statements

instance of edit

property **formally named taxon** *statement object*

▼ 0 references

+ add reference
+ add value

interacts with edit

property **Treponema sp. mp1** *statement object*

subject body part	cell surface	<i>statement qualifiers</i>
method used to identify object	electron microscopy	
	fluorescence in situ hybridization	
	light microscopy	
	phylogenetic marker sequencing	
interaction type	motility facilitated by	
method used to identify subject	light microscopy	
▼ 2 references		
reference DOI	10.1098/rspb.1964.0025	<i>statement references</i>
stated in	The fine structure of the flagellate <i>Mixotricha paradoxa</i> and its associated micro-organisms	
reference DOI	10.1078/0932-4739-00893	
stated in	Identification of the ectosymbiotic bacteria of <i>Mixotricha paradoxa</i> involved in movement symbiosis	
+ add reference		

Treponema sp. mp3 edit

subject body part	cell surface	
method used to identify object	electron microscopy	
	fluorescence in situ hybridization	
	light microscopy	
	phylogenetic marker sequencing	
interaction type	motility facilitated by	
method used to identify subject	light microscopy	
▼ 2 references		

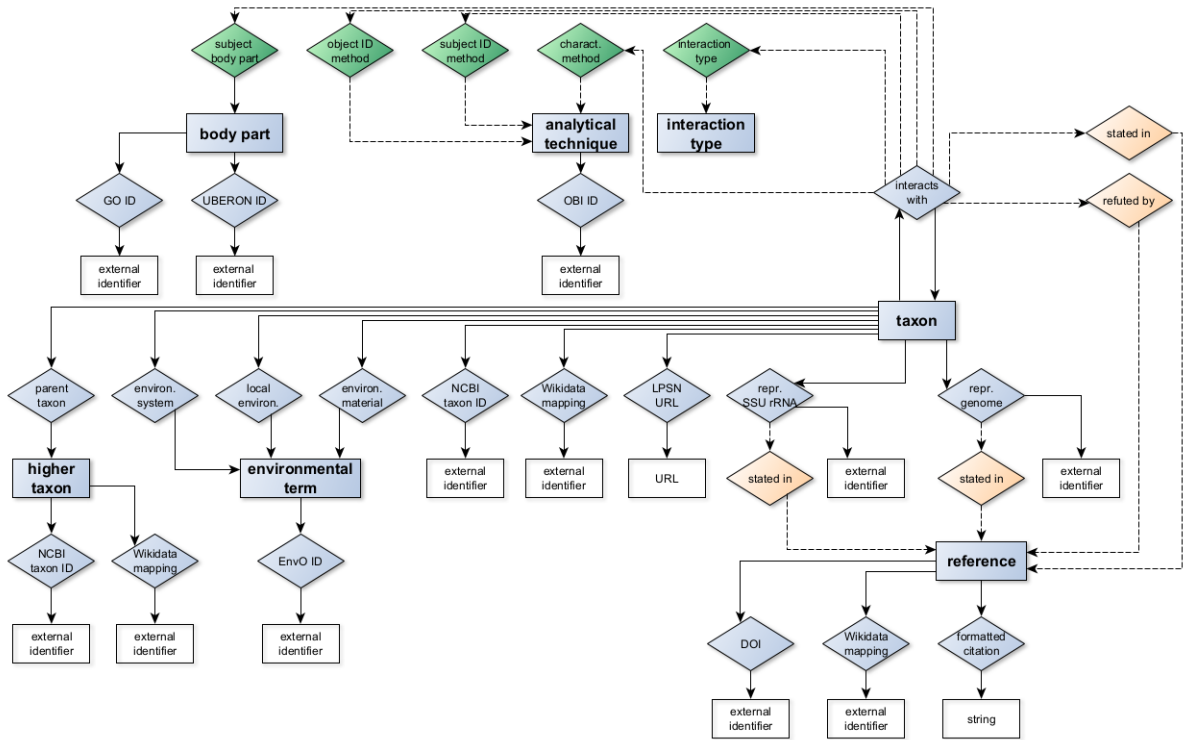
Treponema sp. mpsp15 edit

subject body part	cell surface
-------------------	--------------

multiple statements of the same property are grouped together

395

396 **Figure 1.** Annotated screenshot from the Wikibase item page for a host species, *Mixotricha*
 397 *paradoxa*, which illustrates how symbiotic interactions are represented as statements that link
 398 it to respective symbiont items, and which can be further qualified by additional information
 399 and references.



401

402

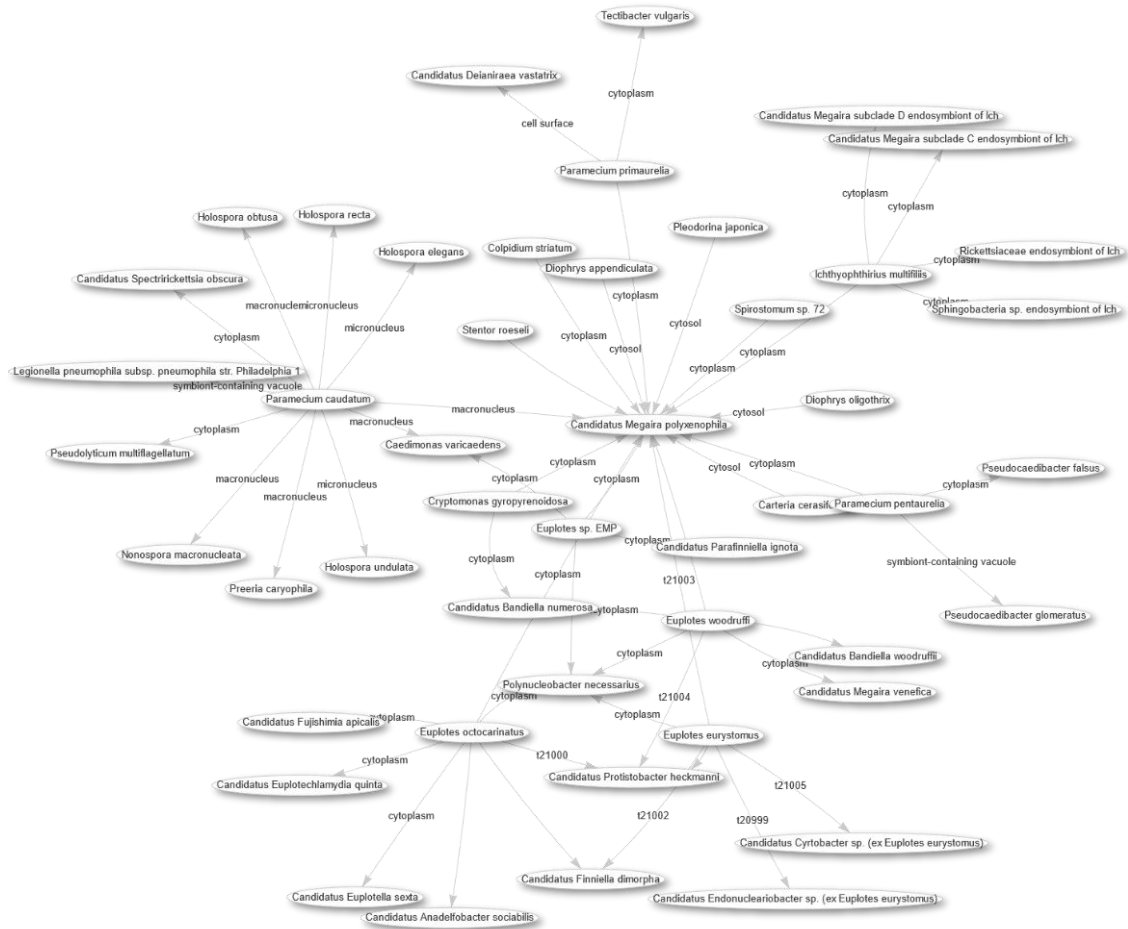
403

404

405

406

Figure 2. Relationships between items (blue rectangles), properties (diamonds), and other data types (white rectangles) in the PPSDB data model. Qualifier and reference relationships are depicted with dashed lines. Properties used in qualifiers or references are colored green and orange respectively. Subclasses of “taxon” and “environmental term” are not shown, for simplicity.



408

409 **Figure 3.** Graphical display of protist host species of the bacterial symbiont *Ca. Megaira*
 410 *polyxenophila*, and their other respective symbionts. The visualization was produced from a
 411 SPARQL query result by the Query Service engine included in Wikibase.cloud
 412 (<https://tinyurl.com/2948dund>).
 413

414 **Table 1.** Concepts or entities represented in the database and relevant external databases,
 415 ontologies, or identifiers that they are linked to, if exact matches are available.

Concept/entity	Relevant database, ontology, or identifier
Taxonomy of the interacting organisms	NCBI Taxonomy Wikidata List of Prokaryotic Names with Standing in Nomenclature (LPSN)
Localization of symbionts in the host organism	Gene Ontology UBERON
Nature of the biotic interaction, if known or inferred	OBO Relations Ontology
Analytical methods used to identify organisms or determine the interaction type	OBI Evidence Ontology
Environment where the host organism was collected or isolated	Environment Ontology
Publication describing the symbiosis	DOI Wikidata
Nucleotide sequence records	Genbank accession

416

417 **References**

- 418 Adl S. M., Bass D., Lane C. E., Lukeš J., Schoch C. L., Smirnov A., Agatha S., Berney C.,
419 Brown M. W., Burki F., et al. 2019. Revisions to the classification, nomenclature, and
420 diversity of eukaryotes. *J. Eukaryot. Microbiol.*, **66**:4–119.
- 421 Albrycht K., Rynkiewicz A. A., Harasymczuk M., Barylski J. & Zielezinski A. 2022. Daily
422 Reports on Phage-Host Interactions. *Front. Microbiol.*, **13**:946070.
- 423 Ashburner M., Ball C. A., Blake J. A., Botstein D., Butler H., Cherry J. M., Davis A. P.,
424 Dolinski K., Dwight S. S., Eppig J. T., et al. 2000. Gene Ontology: tool for the
425 unification of biology. *Nat. Genet.*, **25**:25–29.
- 426 Ball G. H. 1969. Organisms living on and in protozoa. *In: Research in Protozoology.*
427 Elsevier. p. 565–718.
- 428 Bandrowski A., Brinkman R., Brochhausen M., Brush M. H., Bug B., Chibucos M. C.,
429 Clancy K., Courtot M., Derom D., Dumontier M., et al. 2016. The ontology for
430 biomedical investigations. *PLoS ONE*, **11**:e0154556.
- 431 Bauer F. & Kaltenböck M. 2012. Linked Open Data: The Essentials.
- 432 Bjorbækmo M. F. M., Evenstad A., Røsæg L. L., Krabberød A. K. & Logares R. 2020. The
433 planktonic protist interactome: where do we stand after a century of research? *ISME J.*,
434 **14**:544–559.
- 435 Boscaro V., Schrollhammer M., Benken K. A., Krenek S., Szokoli F., Berendonk T. U.,
436 Schweikert M., Verni F., Sabaneyeva E. V. & Petroni G. 2013. Rediscovering the genus
437 *Lyticum*, multiflagellated symbionts of the order Rickettsiales. *Sci. Rep.*, **3**:3305.
- 438 Buttigieg P. L., Morrison N., Smith B., Mungall C. J., Lewis S. E. & ENVO Consortium.
439 2013. The environment ontology: contextualising biological and biomedical entities. *J.*
440 *Biomed. Semantics*, **4**:43.
- 441 Buttigieg P. L., Pafilis E., Lewis S. E., Schildhauer M. P., Walls R. L. & Mungall C. J. 2016.
442 The environment ontology in 2016: bridging domains with increased scope, semantic
443 density, and interoperation. *J. Biomed. Semantics*, **7**:57.
- 444 Chaudhri V. K., Baru C., Chittar N., Dong X. L., Genesereth M., Hendler J., Kalyanpur A.,
445 Lenat D. B., Sequeda J., Vrandečić D., et al. 2022. Knowledge graphs: Introduction,
446 history, and perspectives. *AIMag*, **43**:17–29.
- 447 Davison H. R., Hurst G. D. D. & Siozios S. 2023. “Candidatus Megaira” are diverse
448 symbionts of algae and ciliates with the potential for defensive symbiosis. *Microb.*
449 *Genom.*, **9**.

450 Diefenbach D., Wilde M. D. & Alipio S. 2021. Wikibase as an Infrastructure for Knowledge
451 Graphs: The EU Knowledge Graph. *In: Hotho A., Blomqvist E., Dietze S., Fokoue A.,*
452 *Ding Y., Barnaghi P., Haller A., Dragoni M. & Alani H. (eds.), The semantic web –*
453 *ISWC 2021. Vol. 12922. Lecture notes in computer science. Cham, Springer*
454 *International Publishing. p. 631–647.*

455 Fokin S. I. & Serra V. 2022. Bacterial symbiosis in ciliates (Alveolata, Ciliophora): Roads
456 traveled and those still to be taken. *J. Eukaryot. Microbiol.*, **69**:e12886.

457 Gene Ontology Consortium, Aleksander S. A., Balhoff J., Carbon S., Cherry J. M., Drabkin
458 H. J., Ebert D., Feuermann M., Gaudet P., Harris N. L., et al. 2023. The Gene Ontology
459 knowledgebase in 2023. *Genetics*, **224**.

460 Hongoh Y., Sato T., Dolan M. F., Noda S., Ui S., Kudo T. & Ohkuma M. 2007. The motility
461 symbiont of the termite gut flagellate *Caduceia versatilis* is a member of the
462 “Synergistes” group. *Appl. Environ. Microbiol.*, **73**:6270–6276.

463 Huaman E., Huaman J. L. & Huaman W. 2023. Getting Quechua Closer to Final Users
464 Through Knowledge Graphs. *In: Lossio-Ventura J. A., Valverde-Rebaza J., Díaz E. &*
465 *Alatrística-Salas H. (eds.), Information management and big data: 9th annual international*
466 *conference, simbig 2022, lima, peru, november 16–18, 2022, proceedings. Vol. 1837.*
467 *Communications in computer and information science. Cham, Springer Nature*
468 *Switzerland. p. 61–69.*

469 Husnik F., Tashyreva D., Boscaro V., George E. E., Lukeš J. & Keeling P. J. 2021. Bacterial
470 and archaeal symbioses with protists. *Curr. Biol.*, **31**:R862–R877.

471 Koho M., Coladangelo L. P., Ransom L. & Emery D. 2023. Wikibase model for premodern
472 manuscript metadata harmonization, linked data integration, and discovery. *J. Comput.*
473 *Cult. Herit.*, **16**:1–25.

474 Kostygov A. Y., Karnkowska A., Votýpka J., Tashyreva D., Maciszewski K., Yurchenko V.
475 & Lukeš J. 2021. Euglenozoa: taxonomy, diversity and ecology, symbioses and viruses.
476 *Open Biol.*, **11**:200407.

477 La Scola B., Audic S., Robert C., Jungang L., de Lamballerie X., Drancourt M., Birtles R.,
478 Claverie J.-M. & Raoult D. 2003. A giant virus in amoebae. *Science*, **299**:2033.

479 Lamy-Besnier Q., Brancotte B., Ménager H. & Debarbieux L. 2021. Viral Host Range
480 database, an online tool for recording, analyzing and disseminating virus-host
481 interactions. *Bioinformatics*, **37**:2798–2801.

482 von Mering S., Gardiner L. M., Knapp S., Lindon H., Leachman S., Ulloa Ulloa C., Vincent
483 S. & Vorontsova M. S. 2023. Creating a multi-linked dynamic dataset: a case study of
484 plant genera named for women. *Biodivers. Data J.*, **11**:e114408.

485 Mihara T., Nishimura Y., Shimizu Y., Nishiyama H., Yoshikawa G., Uehara H., Hingamp P.,
486 Goto S. & Ogata H. 2016. Linking Virus Genomes with Host Taxonomy. *Viruses*, **8**:66.

487 Monteil C. L., Vallenet D., Menguy N., Benzerara K., Barbe V., Fouteau S., Cruaud C.,
488 Floriani M., Viollier E., Adryanczyk G., et al. 2019. Ectosymbiotic bacteria at the origin
489 of magnetoreception in a marine protist. *Nat. Microbiol.*, **4**:1088–1095.

490 Mungall C. J., Torniai C., Gkoutos G. V., Lewis S. E. & Haendel M. A. 2012. Uberon, an
491 integrative multi-species anatomy ontology. *Genome Biol.*, **13**:R5.

492 Nielsen F. Å., Mietchen D. & Willighagen E. 2017. Scholia, scientometrics and wikidata. *In*:
493 Blomqvist E., Hose K., Paulheim H., Ławrynowicz A., Ciravegna F. & Hartig O. (eds.),
494 The semantic web: ESWC 2017 satellite events. Vol. 10577. Lecture notes in computer
495 science. Cham, Springer International Publishing. p. 237–259.

496 Pacheco A. R., Pauvert C., Kishore D. & Segrè D. 2022. Toward FAIR representations of
497 microbial interactions. *mSystems*, **7**:e0065922.

498 Parks D. H., Chuvochina M., Chaumeil P.-A., Rinke C., Mussig A. J. & Hugenholtz P. 2020.
499 A complete domain-to-species taxonomy for Bacteria and Archaea. *Nat. Biotechnol.*,
500 **38**:1079–1086.

501 Parte A. C., Sardà Carbasse J., Meier-Kolthoff J. P., Reimer L. C. & Göker M. 2020. List of
502 Prokaryotic names with Standing in Nomenclature (LPSN) moves to the DSMZ. *Int. J.*
503 *Syst. Evol. Microbiol.*, **70**:5607–5612.

504 Pasulka A. L., Goffredi S. K., Tavormina P. L., Dawson K. S., Levin L. A., Rouse G. W. &
505 Orphan V. J. 2017. Colonial Tube-Dwelling Ciliates Influence Methane Cycling and
506 Microbial Diversity within Methane Seep Ecosystems. *Front. Mar. Sci.*, **3**.

507 Patterson D. J., Cooper J., Kirk P. M., Pyle R. L. & Remsen D. P. 2010. Names are key to the
508 big new biology. *Trends Ecol. Evol.*, **25**:686–691.

509 Petroni G., Spring S., Schleifer K. H., Verni F. & Rosati G. 2000. Defensive extrusive
510 ectosymbionts of Euplotidium (Ciliophora) that contain microtubule-like structures are
511 bacteria related to Verrucomicrobia. *Proc Natl Acad Sci USA*, **97**:1813–1817.

512 Poelen J. H., Simons J. D. & Mungall C. J. 2014. Global biotic interactions: An open
513 infrastructure to share and analyze species-interaction datasets. *Ecol. Inform.*, **24**:148–
514 159.

515 Schoch C. L., Ciufo S., Domrachev M., Hotton C. L., Kannan S., Khovanskaya R., Leipe D.,
516 Mcveigh R., O'Neill K., Robbertse B., et al. 2020. NCBI Taxonomy: a comprehensive
517 update on curation, resources and tools. *Database (Oxford)*, **2020**.

518 Schrallhammer M., Castelli M. & Petroni G. 2018. Phylogenetic relationships among
519 endosymbiotic R-body producer: Bacteria providing their host the killer trait. *Syst. Appl.*
520 *Microbiol.*, **41**:213–220.

521 Schrallhammer M., Ferrantini F., Vannini C., Galati S., Schweikert M., Görtz H.-D., Verni F.
522 & Petroni G. 2013. “Candidatus Megaira polyxenophila” gen. nov., sp. nov.:
523 considerations on evolutionary history, host range and shift of early divergent rickettsiae.
524 *PLoS ONE*, **8**:e72581.

525 Schulz F., Martijn J., Wascher F., Lagkouvardos I., Kostanjšek R., Ettema T. J. G. & Horn
526 M. 2016. A Rickettsiales symbiont of amoebae with ancient features. *Environ.*
527 *Microbiol.*, **18**:2326–2342.

528 Seah B. K. B. 2023. Paying it forward: Crowdsourcing the harmonisation and linking of
529 taxon names and biodiversity identifiers. *Biodivers. Data J.*, **11**:e114076.

530 Shimizu C., Hitzler P., Gonzalez-Estrecha S., Goeke-Smith J., Rehberger D., Foley C. &
531 Sheill A. 2023. The Wikibase Approach to the Enslaved.Org Hub Knowledge Graph. *In*:
532 Payne T. R., Presutti V., Qi G., Poveda-Villalón M., Stoilos G., Hollink L., Kaoudi Z.,
533 Cheng G. & Li J. (eds.), The semantic web – ISWC 2023: 22nd international semantic
534 web conference, athens, greece, november 6–10, 2023, proceedings, part II. Vol. 14266.
535 Lecture notes in computer science. Cham, Springer Nature Switzerland. p. 419–434.

536 Shi Y., Queller D. C., Tian Y., Zhang S., Yan Q., He Zhili, He Zhenzhen, Wu C., Wang C. &
537 Shu L. 2021. The Ecology and Evolution of Amoeba-Bacterium Interactions. *Appl.*
538 *Environ. Microbiol.*, **87**.

539 Utami Y. D., Kuwahara H., Igai K., Murakami T., Sugaya K., Morikawa T., Nagura Y., Yuki
540 M., Deevong P., Inoue T., et al. 2019. Genome analyses of uncultured TG2/ZB3 bacteria
541 in “Margulisbacteria” specifically attached to ectosymbiotic spirochetes of protists in the
542 termite gut. *ISME J.*, **13**:455–467.

543 Wardeh M., Risley C., McIntyre M. K., Setzkorn C. & Baylis M. 2015. Database of host-
544 pathogen and related species interactions, and their global distribution. *Sci. Data*,
545 **2**:150049.

546 Yilmaz P., Kottmann R., Field D., Knight R., Cole J. R., Amaral-Zettler L., Gilbert J. A.,
547 Karsch-Mizrachi I., Johnston A., Cochrane G., et al. 2011. Minimum information about a

548 marker gene sequence (MIMARKS) and minimum information about any (x) sequence
549 (MIxS) specifications. *Nat. Biotechnol.*, **29**:415–420.