

The BeeBiome Data Portal provides easy access to bee microbiome information

Valentine Rech de Laval^{1,2,3} (ORCID 0000-0002-3020-1490), Benjamin Dainat³ (ORCID 0000-0002-1740-7136), Philippe Engel⁴ (ORCID 0000-0002-4678-6200), Marc Robinson-Rechavi^{1,2,*} (ORCID 0000-0002-3437-3329)

1: Department of Ecology and Evolution, University of Lausanne, 1015 Lausanne, Switzerland.

2: SIB Swiss Institute of Bioinformatics, 1015 Lausanne, Switzerland.

3: Agroscope, Swiss Bee Research Centre, Schwarzenburgstrasse 161, 3003, Bern, Switzerland.

4: Department of Fundamental Microbiology, University of Lausanne, 1015 Lausanne, Switzerland.

* To whom correspondence should be addressed: marc.robinson-rechavi@unil.ch

Abstract

Bees can be colonized by a large diversity of microbes, including beneficial gut symbionts and detrimental pathogens, with implications for bee health. Over the last few years, researchers around the world have collected a huge amount of genomic and transcriptomic data about the composition, genomic content, and gene expression of bee-associated microbial communities. While each of these datasets by itself has provided important insights, the integration of such datasets provides an unprecedented opportunity to obtain a global picture of the microbes associated with bees and their link to bee health. The challenge of such an approach is that datasets are difficult to find within large generalist repositories and are often not readily accessible, which hinders integrative analyses.

Here we present a publicly-available online resource, the BeeBiome data portal (<https://www.beebiome.org>), which provides an overview of and easy access to currently available metagenomic datasets involving bee-associated microbes. Currently the data portal contains 33,678 Sequence Read Archive (SRA) experiments for 278 Apoidae hosts. We present the content and functionalities of this portal.

Importance

By providing access to all bee microbiomes in a single place, with easy filtering on relevant criteria, Beebiome will allow faster progress of applied and fundamental research on bee biology and health. It should be a useful tool for researchers, academics, funding agencies, and governments, with beneficial impacts for stakeholders.

36 Introduction

37 A multitude of factors contribute to bee declines worldwide, but microbes have been identified to
38 play a major role. From pathogens, like viruses or fungi causing severe diseases, to beneficial
39 gut bacteria important for protection against pathogens or digestion of nutrients, microbes play a
40 key role in bee fitness and survival. Sequencing efforts around the world have contributed to a
41 better understanding of the genetic and functional diversity of microbes associated with bees
42 (1). These datasets are publicly available from dedicated databases, and notably include
43 amplicon sequences, isolated genomes, shotgun metagenomes, and transcriptomes. The rapid
44 accumulation of such datasets offers new opportunities for data integration and cross-study
45 analyses (2). However, such approaches are hampered by the lack of standardization in dataset
46 annotation which makes it difficult to systematically (and automatically) search for all sequence
47 resources of one data type.

48
49 There is agreement in the community that a centralized bioinformatics tool which would
50 systematically catalog and provide access to sequence datasets from bee-associated microbes
51 would be a useful resource for both fundamental and applied research (1). One way how this
52 can be achieved is via a data portal. A data portal is an online platform which provides access to
53 data, not by storing the data directly, but by systematically cataloging and linking datasets
54 deposited elsewhere. Many large-scale, data-driven projects in biology have dedicated portals,
55 for example the TARA Oceans data portal (3). There are also portals which provide access to
56 data unified by a common theme, even though they were generated by multiple projects with no
57 prior coordination. For example HumanMetagenomeDB (4) provides access to public human
58 metagenomes, by organizing the relevant metadata. As raw data is deposited in the open
59 databases of NCBI, EBI, and DDBJ (5–7), access to metadata linked to that raw data can be
60 sufficient to empower users.

61 Here, we present the BeeBiome data portal, which automatically collects and systematically
62 stores metadata of publicly available DNA and RNA sequence datasets of bee microbiome
63 projects and hence makes them readily accessible to the growing research community working
64 on bee-associated microbes. This portal will facilitate data integration and cross-study analyses
65 with the ultimate goal to understand the ecology and evolution of bee-associated microbes and
66 viruses, and advance our understanding of their impact on bees and bee health, from managed
67 honey bees to solitary wild bees.

68 Results

69 Content

70 The Beebiome portal integrates metadata from all NCBI(7) available genomic or transcriptomic
71 microbiome data for which the host is identified as Apoidea. Thanks to INSDC data sharing, this
72 also includes all data submitted to ENA(8) or DDBJ(5). As of 5 January 2025, Beebiome
73 contains 30427 BioSamples (unique entries), encompassing 453 Bioprojects and 33678 SRA
74 experiments (**Table 1**). This represents data from 278 Apoidea host species, which includes

75 honey bees, bumble bees, stingless bees, sweat bees, and carpenter bees (among others). All
76 data is automatically updated every month (see Data and methods). This offers the advantage
77 that newly deposited datasets are integrated into the Beebiome portal on a regular basis.
78 However, datasets with ambiguous taxonomic annotations can be missing. For example, NCBI
79 Biosamples which are only annotated as “Bombus” without species identification are not
80 included in the Beebiome portal, because “Bombus” refers both to a genus and to a subgenus
81 and hence matches two NCBI TaxIDs: 28641 (genus) and 144708 (subgenus). Such
82 ambiguities need to be avoided to begin with, or corrected at the level of the database to have
83 these samples integrated into the Beebiome portal.
84

85 **Table 1:** High level Beebiome content
86

GSC MixS or NCBI package name	BioSamples
Metagenome or environmental; version 1.0	12016
MIMARKS: survey, host-associated; version 6.0	7947
MIMS: metagenome/environmental, host-associated; version 6.0	4181
MIMARKS: survey, air; version 6.0	2351
Microbe; version 1.0	1171
Virus; version 1.0	867
MIMARKS: specimen, host-associated; version 6.0	745
MIMARKS: survey, plant-associated; version 6.0	281
MIMS: metagenome/environmental, air; version 6.0	182
MIGS: cultured bacteria/archaea, host-associated; version 6.0	179
Invertebrate; version 1.0	166
Pathogen: clinical or host-associated; version 1.0	92
MIMS: metagenome/environmental, plant-associated; version 6.0	53
MIMS: metagenome/environmental, agriculture; version 6.0	45
MIMARKS: survey, microbial; version 6.0	33
MIGS: cultured bacteria/archaea; version 6.0	32
MIUVIG: uncultivated virus genome, host-associated; version 6.0	29

MIMARKS: specimen; version 6.0	15
Pathogen: environmental/food/other; version 1.0	14
Generic	13
MIGS: eukaryote, host-associated; version 6.0	11
MIGS: cultured bacteria/archaea, agriculture; version 6.0	1
MIGS: cultured bacteria/archaea, human-associated; version 6.0	1
MIGS: cultured bacteria/archaea, miscellaneous; version 6.0	1
MIMS: metagenome/environmental, miscellaneous; version 6.0	1

87
88 For each sample, BeeBiome collects metadata which facilitates searching for relevant datasets
89 by key words in different categories such as ‘Organism’ (e.g., *Snodgrassella alvi*) ‘Host’ (e.g.
90 *Apis mellifera*), ‘Library strategies’ (e.g. amplicon or WGS), ‘Library sources’ (e.g.
91 metagenomic), or ‘Collection locality’ and ‘Collection date’. For example, to identify all amplicon
92 sequence datasets, a user would search for Library source ‘Genomic’ and Library strategy
93 ‘Amplicon’. Of note, it is not yet possible to automatically filter the gene amplified, for example to
94 search only 16S rRNA gene amplicons; this information is rarely available. Filtered and sorted
95 data can then be recovered from primary databases through BioProject, BioSample, SRA
96 experiment, or NCBI Nucleotide identifiers, which are all linked back to the source databases.
97 We also store assay type, center name and instrument used (e.g. “Illumina HiSeq 2000”), to
98 allow filtering when relevant.

99 Access

100 The primary access to BeeBiome data is through our webpage, at beebiome.org. The
101 homepage provides direct access to a ‘basic search’, as well as menus to navigate towards an
102 ‘advanced search’, a map, and a wiki. The map simply shows the geographical location of
103 collection for all samples for which this information is available, while ‘Browse table’ allows to
104 see the complete table of all data in BeeBiome (**Figure 1**).
105

BeeBiome browse table

This browser interface allows to discover BeeBiome data. A basic search (in all fields) is available in the right of the table. An **advanced search** is available to do a search on each field; more details on each field are available in our [FAQ page](#). Results are ordered by BeeSample acc. (from 'BeeSample acc.'). The order could be changed by clicking on one column. Rows are sorted with and click on another column. Clicking on the 'i' icon shows the full information for each sample.

Show 10 - 20 entries

Showing 10 of 55483 entries

BeeSample acc.	BeeSample acc.	SRA experiment entries	NCBI Nucleotide entries	Host	Organism	Submission Date	Gen. location name	Library strategy(1)	Library strategy(2)	Library strategy(3)	Institution(1)	Center name	BeeBiome project acc.
PRJNA23484	SRR1271002	1	493	Apis mellifera	honeybee genome	2008		GENOMIC	PAIRED	WGS	Brunel Vilas 2008	University of Chile	WGS as host associated 1.2
PRJNA15857	SRR1271002	1	493	Apis mellifera	honeybee genome	2008		GENOMIC	PAIRED	WGS	Brunel Vilas 2008	National Center for Biotechnology Information	WGS as host associated 1.4
PRJNA70585	SRR1888244	1	5471	Streptococcus pneumoniae	Streptococcus pneumoniae	1992-07	Canada, Lehigh, PA, USA	GENOMIC	PAIRED	WGS	Brunel Vilas 2008	University of California, Riverside	WGS as host associated 1.1
PRJNA70582	SRR1888242	1	5277	Streptococcus pneumoniae	Streptococcus pneumoniae	1979-08	USA, Logan, Utah	GENOMIC	PAIRED	WGS	Brunel Vilas 2008	University of California, Riverside	WGS as host associated 1.3
PRJNA23176	SRR1271002	1	574	Apis mellifera	honeybee genome	2008-08	Canada, Saskatchewan	GENOMIC	PAIRED	WGS	Brunel Vilas 2008	NCBI	Microbi.1.2
PRJNA23176	SRR1271002	1	574	Apis mellifera	honeybee genome	2008-08	Canada, Saskatchewan	GENOMIC	PAIRED	WGS	Brunel Vilas 2008	University of Saskatchewan	Microbi.1.1
PRJNA23176	SRR1271002	1	582	Apis mellifera	honeybee genome	2008-08	Canada, Saskatchewan	GENOMIC	PAIRED	WGS	Brunel Vilas 2008	NCBI	Microbi.1.3
PRJNA23176	SRR1271002	1	586	Apis mellifera	honeybee genome	2008-08	Canada, Saskatchewan	GENOMIC	PAIRED	WGS	Brunel Vilas 2008	University of Saskatchewan	Microbi.1.2
PRJNA23176	SRR1271002	1	547	Apis mellifera	honeybee genome	2008-08	Canada, Saskatchewan	GENOMIC	PAIRED	WGS	Brunel Vilas 2008	NCBI	Microbi.1.2
PRJNA23176	SRR1271002	1	547	Apis mellifera	honeybee genome	2008-08	Canada, Saskatchewan	GENOMIC	PAIRED	WGS	Brunel Vilas 2008	University of Saskatchewan	Microbi.1.2

Page 1 of 1

This work is published under the [Creative Commons Attribution 4.0 International License](#). If you intend to use data from BeeBiome, it would be nice to see our [FAQ page](#).



Figure 1: Main views of BeeBiome data. (A) “Browse Table” view; (B) Map view.

106
107
108
109
110
111
112
113
114
115
116
117
118
119
120
121
122
123
124
125
126
127
128

The ‘basic search’ is a query on all metadata. Thus a query for ‘Lactobacillus’ will return 511 entries (as of 5 January 2025). The ‘advanced search’ allows querying by specific metadata fields, joined by ‘and’ (**Figure 2**). There is also an option to ‘browse’ all metadata. All primary views are tables which can be filtered by terms; filtering the browse view reproduces the same results as a basic search. All tables can also be sorted by clicking on column names. Which metadata is shown adapts dynamically to the window size, while all other metadata remains accessible by unfolding each row. Whatever is shown, all metadata is used for filtering by terms. The order of metadata columns was established following a poll of the bee microbiome community, to make sure that the relevant information is visible even when window size is limited. As of writing, the first columns are thus: BioProject accession, BioSample accession, SRA experiment entries, NCBI Nucleotide entries, Host, and Organism. It should be noted that for the microbiome of e.g. the honey bee, *Apis mellifera* is the Host, while the Organism is the microbe or pest (e.g. *Snodgrassella alvi*) or the type of microbial community (e.g. insect gut metagenome). All search or filtering results, as well as the complete contents of Beebiome, can be downloaded in TSV or copied to clipboard for easy re-use. For example, the results of a query can be downloaded as a TSV and imported to a spreadsheet software such as MS Excel, from which the Biosample IDs can be copied then used to query NCBI simply by pasting them into the NCBI search. Then they can be simply batch downloaded from e.g. NCBI SRA. The results of advanced search can also be shown on the map (Figure 2C).

A

Advanced search

This advanced search interface allows to do a text-based search on each field and combine them using a "AND" boolean logic. More details on each field are available in our [help page](#).

Bioproject accession Ex: PRJNA200000	BioSample accession Ex: SAMN0171000	Host Ex: <i>Apis mellifera</i>
Organism Ex: <i>Lasioglossum</i>	Collection date Ex: 2014-05-14	Dec. location name Ex: <i>Sentinel</i>
Library sources Ex: SRA	Library layouts Ex: SRA	Library strategies Ex: SRA
Instruments Ex: Illumina	Center name Ex: University of Leuven	Sample package accession Ex: SRA
Ex: Name	Ex: University of Leuven	Ex: SRA

Submit Clear form

B

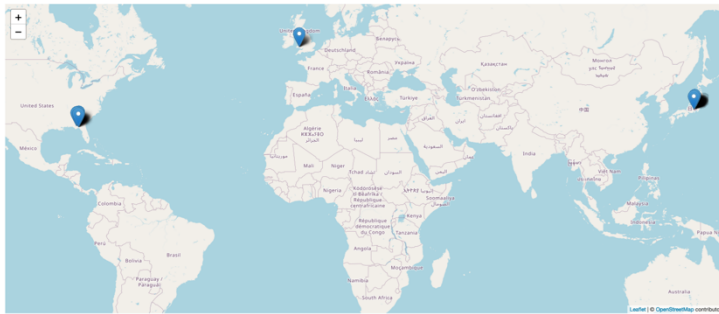
Results are ordered by "Bioproject acc.", then "Accession date". The order could be changed by clicking on any column. Then press "Refresh" and click on another column.

Showing 10 of 100 entries

Bioproject acc.	Accession date	SRA	NCBI	Host	Organization	Accession date	Seq. method name	Library technology	Library layout	Library strategy	Host taxon
PRJNA200000	SRA00047000	1	0	<i>Lasioglossum</i>	public	2013-07-10	Japan Satsuma Chikoku	METAGENOMIC	PAIRED	AMPLICON	<i>Apis mellifera</i>
PRJNA200000	SRA00047001	1	0	<i>Lasioglossum</i>	public	2013-07-10	Japan Satsuma Chikoku	METAGENOMIC	PAIRED	AMPLICON	<i>Apis mellifera</i>
PRJNA200000	SRA00047004	1	0	<i>Lasioglossum</i>	public	2013-07-10	Japan Satsuma Chikoku	METAGENOMIC	PAIRED	AMPLICON	<i>Apis mellifera</i>
PRJNA200000	SRA00047005	1	0	<i>Lasioglossum</i>	public	2013-07-10	Japan Satsuma Chikoku	METAGENOMIC	PAIRED	AMPLICON	<i>Apis mellifera</i>
PRJNA200000	SRA00047007	1	0	<i>Lasioglossum</i>	public	2013-07-10	Japan Satsuma Chikoku	METAGENOMIC	PAIRED	AMPLICON	<i>Apis mellifera</i>
PRJNA200000	SRA00047003	1	0	<i>Lasioglossum</i>	public	2013-07-21	Japan Satsuma Chikoku	METAGENOMIC	PAIRED	AMPLICON	<i>Apis mellifera</i>
PRJNA200000	SRA00047002	1	0	<i>Lasioglossum</i>	public	2013-07-21	Japan Satsuma Chikoku	METAGENOMIC	PAIRED	AMPLICON	<i>Apis mellifera</i>
PRJNA200000	SRA00047006	1	0	<i>Lasioglossum</i>	public	2013-07-20	Japan Satsuma Chikoku	METAGENOMIC	PAIRED	AMPLICON	<i>Apis mellifera</i>
PRJNA200000	SRA00047007	1	0	<i>Lasioglossum</i>	public	2013-07-18	Japan Satsuma Chikoku	METAGENOMIC	PAIRED	AMPLICON	<i>Apis mellifera</i>
PRJNA200000	SRA00047008	1	0	<i>Lasioglossum</i>	public	2013-07-18	Japan Satsuma Chikoku	METAGENOMIC	PAIRED	AMPLICON	<i>Apis mellifera</i>

Showing 10 of 100 entries

C



129

130

131

Figure 2: Example of advanced search in BeeBiome. (A) Search for pollen metagenomes in host *Lasioglossum*; (B) table of results; (C) map of results.

132

133

134

135

136

The Beebiome wiki currently contains a comprehensive catalog of Apis and non-Apis diseases and microbes, providing an overview of most of the microorganisms known to date (1): Apis bee diseases, including known hosts and known effects on hosts; non-Apis bee diseases, including known hosts and known effects on hosts; Apis bee microbes; and non-Apis bee microorganisms.

137

Discussion

138

139

140

141

142

143

144

145

146

147

148

149

150

151

While generalist microbiome databases can be very powerful (9), the volume and diversity of data can be daunting, and make it difficult for small teams and researchers from diverse backgrounds to find what they need (2). Thus we also need dedicated database portals to organize and access relevant metadata, and allow researchers to easily find datasets of interest. For example, HumanMetagenomeDB provides access to standardized metadata for human metagenomes (4), and TerrestrialMetagenomeDB to terrestrial metagenomes (10). Unlike Beebiome, these databases also include manual curation. Manual curation poses problems of sustainability for a small community such as bee microbiome researchers, while our automated filtering according to criteria defined by the community allows to keep Beebiome updated continuously. Importantly, these criteria can easily be adapted or updated according to community needs and feedback. For example, following the discovery by users that our criteria could include bee associated beetle pests, we updated these criteria to exclude all of Metazoa and Viridiplantae, thus restricting to bacterial and eukaryotic microbes. To avoid manual curation, ensure metadata standardization, and enable the Beebiome portal to correctly detect

152 as much data as possible, an important future goal of the community should be to establish
 153 guidelines for how to deposit beebiome datasets into public repositories (e.g. ENA EBI
 154 checklist). This would also allow more fine-grained filtering of the datasets. For example,
 155 datasets coming from different life stages or body sites of the same host species should ideally
 156 be distinguishable. Also, while the possibility to filter datasets by library source or library
 157 strategy allows to download datasets of different types, amplicon sequence data e.g. will include
 158 datasets of amplicons coming from different genes (.e.g rpoB or 16S rRNA gene) or different
 159 regions of a given gene.

160
 161 Another possible area to explore in the future is to provide access to processed datasets, and
 162 make analysis tools or pipelines available via the portal. This will facilitate data usage and help
 163 in standardizing analysis pipelines as much as possible. We hope that the community will find
 164 the current tool already helpful and help us to develop the portal further into the directions
 165 discussed above.

166 Data and methods

167 Data

168 Beebiome stores metadata on bee microbiomes which come from NCBI Biosample, Bioproject,
 169 and SRA entries (7). Data is retrieved using a Perl script generated by Ebot (11), modified to
 170 retrieve relevant metadata. We do not restrict metadata to those which follow a specific
 171 standard, such as GSC MlxS (12, 13) or FAANG metadata standards (14). However, NCBI
 172 recommends submitting data according to GSC MlxS packages. These packages include
 173 attributes defined by the GSC to formally describe and standardize sample metadata. NCBI
 174 submission asks the use either of the GSC MlxS packages or of NCBI packages, forcing the
 175 submitter to give a minimum of information. Entries in Beebiome are represented by the fields
 176 detailed in **Table 2**.

177 **Table 2:** Beebiome entries and corresponding standards.

BeeBiome entry fields	Standards
BioProject acc	PRJD# or PRJEB or PRJNA# (NCBI BioProject accession)
BioSample acc	SAMN# (NCBI BioSample accession)
SRA experiment entries	Integer
NCBI Nucleotide entries	Integer
Assay types	NCBI Strategy enum
Center name	Free text

Library layouts	SINGLE or PAIRED
Library sources	NCBI list, subset to: GENOMIC, TRANSCRIPTOMIC, METAGENOMIC, VIRAL RNA, or OTHER
Organism	NCBI Taxonomy scientific name
Host	NCBI Taxonomy scientific name
Intrument	NCBI Instrument enum
Geo. loc. name	Free text
Collection date	YYYY-MM-DD or YYYY-MM or YYYY

178

179 Figure 3 presents a broad overview of Beebiome generation. The details of Step 1, to retrieve
 180 metadata, are as follows:

181 1. Request NCBI Taxonomy: retrieve species under the taxonomic level 'Apoidea'. We use
 182 scientific names, common names, genbank common names and synonyms for the next
 183 point (called *all_names*).

184 2. Request to NCBI BioSample using request 1 result: retrieve samples having one of the
 185 *all_names* (any fields), having an attribute named host and having an organism that is
 186 not a Metazoa or Viridiplantae.

187 host[Attribute Name] AND (Apis mellifera OR honey bee OR
 188 European honey bee OR Western honey bee OR bee OR honeybee OR ...)
 189 NOT Metazoa[Organism] NOT Viridiplantae[Organism]

190 A query is built on the template 'host[Attribute Name] AND (<names>) NOT
 191 Metazoa[Organism] NOT Viridiplantae[Organism]', where <names> is names
 192 contained in *all_names* separated by 'OR'.

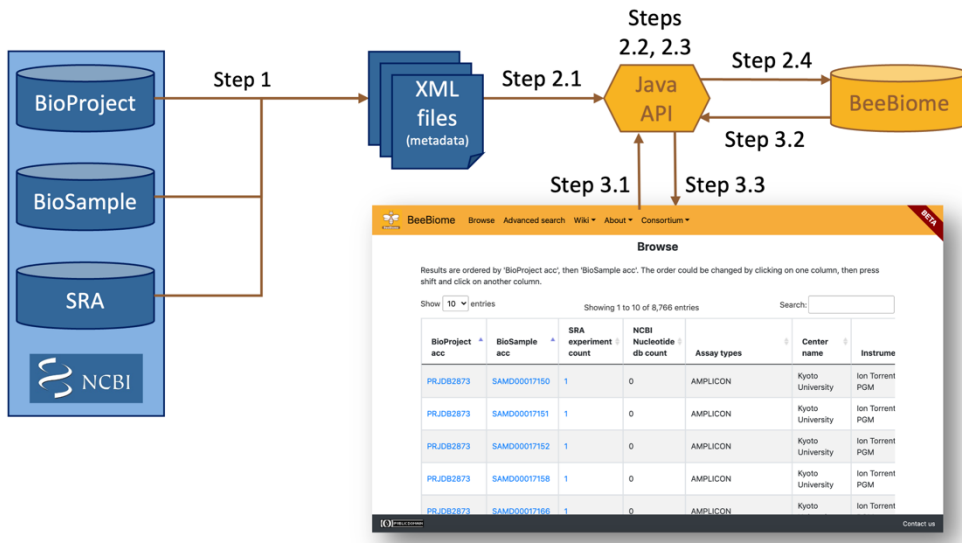
193 To avoid network and technical problems due to large files, we do several requests by
 194 generating <names_separated_by_OR> with a maximum of 300 names. Thus the same
 195 data can be recovered in different files (for instance, a BioSample can be retrieved
 196 several times).

197 3. Request to NCBI BioProject, NCBI SRA and NCBI Nucleotide using request 2 result:
 198 retrieve more metadata and/or links between Biosample and these databases.

199 Beebiome is automatically updated every first Saturday of the month, according to these steps.

200

Pipeline



201

202 **Figure 3:** Overview of the Beebiome database generation.

203 Step 1: Perl script; Step 2: Java API; Step 3: Java API and React webapp.

204

205 Database and views

206 BeeBiome data is stored in a PostgreSQL database. The API is in Java and Spring boot, and
207 the webapp in React. The API is used to import data from NCBI NCBI XML files:

- 208 1. Read files to put them into NCBI Java objects (built from NCBI XSD files)
- 209 2. Filter out BioSamples where the host is not one of the *all_names* (value of the attribute
210 host (which is a free text in NCBI submit format) should be an exact match with one of
211 the *all_names*)
- 212 3. Convert NCBI Java objects to BeeBiome Java objects (each BeeBiome Java object is
213 equals to a table into the database)
- 214 4. Save data into database
- 215 5. An SQL view is generated to save time when there is a request. To generate this view,
216 the query filters out biosamples with any SRA experiment.

217

218 The same API which is used to generated these views also allows to retrieve metadata in JSON
219 format with the following URLs:

220 `https://beebiome.org/beebiome/sample/all` for 'Browse' page, the 'basic
221 search' restriction is done by the webapp.

222 `https://beebiome.org/beebiome/sample/{query}` for 'Advanced search' page to
223 retrieve entries with a BioSample accession containing {query}

224 Code and data availability

225 All code is available at <https://github.com/BeeBiome-consortium/beebiome-data-portal> under
226 GPL 3.0 license. All data in Beebiome is distributed under CC0. Other information follows the
227 original licenses, e.g. supplemental data from Engel et al (1) is under CC-BY-NC-SA 3.0.

228 Acknowledgments

229 This work was supported by grant ECTA_20181209_D.E from the Eva Crane Trust. We thank
230 the team of Rodrigo Ortega Polo from Agriculture and Agri-Food Canada, the entire Beebiome
231 consortium as well as Vincent Doublet, and Méline Garcia from the laboratory of Philipp Engel
232 for their valuable feedback and fruitful discussions on the BeeBiome portal.
233

234 References

- 235 1. Engel P, Kwong WK, McFrederick Q, Anderson KE, Barribeau SM, Chandler JA, Cornman
236 RS, Dainat J, Miranda JR de, Doublet V, Emery O, Evans JD, Farinelli L, Flenniken ML,
237 Granberg F, Grasis JA, Gauthier L, Hayer J, Koch H, Kocher S, Martinson VG, Moran N,
238 Munoz-Torres M, Newton I, Paxton RJ, Powell E, Sadd BM, Schmid-Hempel P, Schmid-
239 Hempel R, Song SJ, Schwarz RS, vanEngelsdorp D, Dainat B. 2016. The Bee Microbiome:
240 Impact on Bee Health and Model for Evolution and Ecology of Host-Microbe Interactions.
241 *mBio* 7.
- 242 2. Gkantiragas AG, Gabrielli J. 2021. A Meta-Analysis of the 16S-rRNA Gut Microbiome Data
243 in Honeybees (*Apis Mellifera*). *bioRxiv* <https://doi.org/10.1101/2021.12.18.473299>.
- 244 3. Pesant S, Not F, Picheral M, Kandels-Lewis S, Le Bescot N, Gorsky G, Iudicone D, Karsenti
245 E, Speich S, Troublé R, Dimier C, Searson S. 2015. Open science resources for the
246 discovery and analysis of Tara Oceans data. 1. *Sci Data* 2:150023.
- 247 4. Kasmanas JC, Bartholomäus A, Corrêa FB, Tal T, Jehmlich N, Herberth G, von Bergen M,
248 Stadler PF, Carvalho ACP de LF de, Nunes da Rocha U. 2021. HumanMetagenomeDB: a
249 public repository of curated and standardized metadata for human metagenomes. *Nucleic
250 Acids Res* 49:D743–D750.
- 251 5. Okido T, Kodama Y, Mashima J, Kosuge T, Fujisawa T, Ogasawara O. 2022. DNA Data
252 Bank of Japan (DDBJ) update report 2021. *Nucleic Acids Res* 50:D102–D105.
- 253 6. Cantelli G, Bateman A, Brooksbank C, Petrov AI, Malik-Sheriff RS, Ide-Smith M, Hermjakob
254 H, Flicek P, Apweiler R, Birney E, McEntyre J. 2022. The European Bioinformatics Institute
255 (EMBL-EBI) in 2021. *Nucleic Acids Res* 50:D11–D19.
- 256 7. Sayers EW, Beck J, Bolton EE, Bourexis D, Brister JR, Canese K, Comeau DC, Funk K,
257 Kim S, Klimke W, Marchler-Bauer A, Landrum M, Lathrop S, Lu Z, Madden TL, O’Leary N,
258 Phan L, Rangwala SH, Schneider VA, Skripchenko Y, Wang J, Ye J, Trawick BW, Pruitt KD,
259 Sherry ST. 2021. Database resources of the National Center for Biotechnology Information.
260 *Nucleic Acids Res* 49:D10–D17.
- 261 8. Cummins C, Ahamed A, Aslam R, Burgin J, Devraj R, Edbali O, Gupta D, Harrison PW,
262 Haseeb M, Holt S, Ibrahim T, Ivanov E, Jayathilaka S, Kadirvelu V, Kay S, Kumar M, Lathi
263 A, Leinonen R, Madeira F, Madhusoodanan N, Mansurova M, O’Cathail C, Pearce M,
264 Pesant S, Rahman N, Rajan J, Rinck G, Selvakumar S, Sokolov A, Suman S, Thorne R,
265 Tootoo P, Vijayaraja S, Waheed Z, Zyoud A, Lopez R, Burdett T, Cochrane G. 2022. The

- 266 European Nucleotide Archive in 2021. *Nucleic Acids Res* 50:D106–D110.
- 267 9. Oliveira FS, Brestelli J, Cade S, Zheng J, Iodice J, Fischer S, Aurrecoechea C, Kissinger
268 JC, Brunk BP, Stoeckert CJ Jr, Fernandes GR, Roos DS, Beiting DP. 2018. MicrobiomeDB:
269 a systems biology platform for integrating, mining and analyzing microbiome experiments.
270 *Nucleic Acids Res* 46:D684–D691.
- 271 10. Corrêa FB, Saraiva JP, Stadler PF, da Rocha UN. 2020. TerrestrialMetagenomeDB: a
272 public repository of curated and standardized metadata for terrestrial metagenomes. *Nucleic*
273 *Acids Res* 48:D626–D632.
- 274 11. Sayers E. 2018. E-utilities Quick StartEntrez Programming Utilities Help [Internet]. National
275 Center for Biotechnology Information (US). <https://www.ncbi.nlm.nih.gov/books/NBK25500/>.
276 Retrieved 15 June 2021.
- 277 12. Field D, Garrity G, Gray T, Morrison N, Selengut J, Sterk P, Tatusova T, Thomson N, Allen
278 MJ, Angiuoli SV, Ashburner M, Axelrod N, Baldauf S, Ballard S, Boore J, Cochrane G, Cole
279 J, Dawyndt P, De Vos P, DePamphilis C, Edwards R, Faruque N, Feldman R, Gilbert J,
280 Gilna P, Glöckner FO, Goldstein P, Guralnick R, Haft D, Hancock D, Hermjakob H, Hertz-
281 Fowler C, Hugenholtz P, Joint I, Kagan L, Kane M, Kennedy J, Kowalchuk G, Kottmann R,
282 Kolker E, Kravitz S, Kyrpides N, Leebens-Mack J, Lewis SE, Li K, Lister AL, Lord P, Maltsev
283 N, Markowitz V, Martiny J, Methe B, Mizrahi I, Moxon R, Nelson K, Parkhill J, Proctor L,
284 White O, Sansone S-A, Spiers A, Stevens R, Swift P, Taylor C, Tateno Y, Tett A, Turner S,
285 Ussery D, Vaughan B, Ward N, Whetzel T, San Gil I, Wilson G, Wipat A. 2008. The
286 minimum information about a genome sequence (MIGS) specification. *Nat Biotechnol*
287 26:541–547.
- 288 13. Kottmann R, Gray T, Murphy S, Kagan L, Kravitz S, Lombardot T, Field D, Glöckner FO,
289 Genomic Standards Consortium. 2008. A standard MIGS/MIMS compliant XML Schema:
290 toward the development of the Genomic Contextual Data Markup Language (GCDML).
291 *Omics J Integr Biol* 12:115–121.
- 292 14. Harrison PW, Fan J, Richardson D, Clarke L, Zerbino D, Cochrane G, Archibald AL,
293 Schmidt CJ, Flicek P. 2018. FAANG, establishing metadata standards, validation and best
294 practices for the farmed and companion animal community. *Anim Genet* 49:520–526.
- 295