# Short reads from honey bee (*Apis* sp.) sequencing projects reveal microbial associate diversity

**Michael Gerth** [Corresp., 1] , **Gregory DD Hurst** [1]

[1] Institute of Integrative Biology, University of Liverpool, Liverpool, United Kingdom

Corresponding Author: Michael Gerth
Email address: gerth@liverpool.ac.uk

High throughput (or 'next generation') sequencing has transformed most areas of biological research and is now a standard method that underpins empirical study of organismal biology, and (through comparison of genomes), reveals patterns of evolution. For projects focused on animals, these sequencing methods do not discriminate between the primary target of sequencing (the animal genome) and 'contaminating' material, such as associated microbes. A common first step is to filter out these contaminants to allow better assembly of the animal genome. Here, we aimed to assess if these 'contaminations' provide information with regard to biologically important microorganisms associated with the individual as part of the 'hologenome'. To achieve this, we examined whether the short read data from *Apis* retrieved elements of its well established microbiome. To this end, we screened almost 1,000 short read libraries of honey bee (*Apis* sp.) sequencing project for the presence of microbial sequences, and find sequences from known honey bee microbial associates in at least 9% of them. Further to this, we used the data to reconstruct draft genomes of three *Apis* associated bacteria *de novo*. We conclude that 'contamination' in short read sequencing libraries can provide useful genomic information on microbial taxa known to be associated with the target organisms, and may even lead to the discovery of novel associations. However, we also find that sequences deriving from microbes outside of the natural microbiome may present a challenge to our approach.

1 **Short reads from honey bee (*Apis* sp.) sequencing projects reveal microbial associate**

2 **diversity**

3 Michael Gerth* & Gregory D. D. Hurst

4 University of Liverpool, Institute of Integrative Biology, Biosciences Building, Crown Street

5 Liverpool L69 7ZB, United Kingdom

6 *Correspondence & material requests

7 gerth@liv.ac.uk

8  **Abstract**

9  High throughput (or 'next generation') sequencing has transformed most areas of biological

10  research and is now a standard method that underpins empirical study of organismal biology, and

11  (through comparison of genomes), reveals patterns of evolution. For projects focused on animals,

12  these sequencing methods do not discriminate between the primary target of sequencing (the

13  animal genome) and 'contaminating' material, such as associated microbes. A common first step

14  is to filter out these contaminants to allow better assembly of the animal genome. Here, we aimed

15  to assess if these 'contaminations' provide information with regard to biologically important

16  microorganisms associated with the individual as part of the 'hologenome'. To achieve this, we

17  examined whether the short read data from *Apis* retrieved elements of its well established

18  microbiome. To this end, we screened almost 1,000 short read libraries of honey bee (*Apis* sp.)

19  sequencing project for the presence of microbial sequences, and find sequences from known

20  honey bee microbial associates in at least 9% of them. Further to this, we used the data to

21  reconstruct draft genomes of three *Apis* associated bacteria *de novo*. We conclude that

22  'contamination' in short read sequencing libraries can provide useful genomic information on

23  microbial taxa known to be associated with the target organisms, and may even lead to the

24  discovery of novel associations. However, we also find that sequences deriving from microbes

25  outside of the natural microbiome may present a challenge to our approach.

## Introduction

Novel DNA sequencing methods have revolutionized biological and medical research in the last two decades (Goodwin et al. 2016). High throughput sequencing (or 'massively parallelized sequencing', 'next generation sequencing','NGS') facilitated the creation of enormous amounts of data for a fraction of the costs associated with traditional Sanger sequencing (Kircher & Kelso 2010; Sboner et al. 2011). This 'genomics revolution', has not only enhanced our understanding of molecular and genome evolution (Wolfe & Li 2003), but also contributed to the recognition that eukaryotes are commonly associated with a plethora of microbial taxa.

In eukaryote genome sequencing projects, sequences deriving from these microbes may obstruct genome assembly efforts, and measures directed at removing microbial associates are routinely performed. This is achieved either by antibiotic treatment of the target organism prior to sequencing (Colbourne et al. 2011), or by removing microbial sequences bioinformatically after sequencing (Schmieder & Edwards 2011). While eliminating microbes may facilitate eukaryotic genome reconstruction, it neglects the recently emerging appreciation of microbes as a biologically important component of all multicellular organisms. Numerous examples illustrate the impact of microbes on animal and plant biology, including physiology, behavior, and evolution (McFall-Ngai et al. 2013). These findings have led to a concept that defines an individual eukaryote with all its associated microbes (microbiome) as an entity (holobiont-hologenome) (Bordenstein & Theis 2015). Although this concept is contentious (Moran & Sloan 2015; Douglas & Werren 2016), it is undisputed that some aspects of organismal biology can only be understood by deciphering interactions with microbial symbionts.

To characterize microbiome composition, three approaches are commonly used. First, microbes may be isolated from the host and cultured axenically. Their properties can then be determined through traditional microbiological methods or by sequencing (Browne et al. 2016).

50    This approach has the benefit of providing both biological and genomic information, but limits

51    discovery to culturable taxa. Second, microbiome taxa may be identified by amplicon

52    sequencing. Specific primers are used to amplify a short informative region from all bacterial

53    taxa in a sample (usually a part of the 16S rRNA gene), and then sequenced (today typically via

54    NGS methods) (Caporaso et al. 2012). This mechanism discovers broad patterns of community

55    diversity, but at a coarse scale, and with weaker functional information. Finally, microbiome

56    composition can be determined via metagenomics, i.e., collective genome sequencing of all

57    bacteria present in a sample (Riesenfeld et al. 2004). This is unbiased, fine scaled, and provides

58    an assessment of biological potential at a community scale, but resolution of genome sequences is

59    more complex

60        In this study, we examined if the data generated in eukaryotic sequencing projects can be

61    used to identify microbiome taxa, and thus to inform about the composition of the wider

62    'holobiont'. Previously, this approach was used to recover genomes of heritable microbes that

63    occur in high densities in many arthropod species, and are therefore prone to be retrieved in

64    arthropod sequencing projects. For example, the genomes of multiple *Wolbachia* strains were

65    discovered in *Drosophila* sequencing data, revealing novel *Wolbachia* diversity and patterns of

66    *Wolbachia* evolution (Salzberg et al. 2005; Richardson et al. 2012).

67        Here, we examine short reads of honey bee (*Apis* sp.) sequencing projects to investigate

68    whether this archived data can be used to retrieve a wider set of microbial associates, including

69    pathogens and gut symbionts. We focus on honey bees because 1) there is a large number of short

70    read sequencing projects targeting *Apis*; 2) the components of healthy and unhealthy *Apis*

71    microbiomes are well established (Evans & Schwarz 2011; Kwong & Moran 2016); 3) managed

72    populations of the economically important honey bees have been in decline worldwide (Neumann

73    & Carreck 2010), and it was hypothesized that certain bacteria and viruses are key players in this

74 decline (Cox-Foster et al. 2007). Thus, any novel genomic data on honey bee symbionts may

75 directly contribute to our understanding of bee disease.

76       To identify 'contaminants', we here use short signature 'bait' sequences of symbionts and

77 pathogens to screen a large number of short read libraries from *Apis* sequencing projects. We

78 demonstrate that the libraries contain non-target sequences from many sources, some of which

79 reflect the natural honey bee microbiome. We further show that highly covered, and possibly

80 novel symbiont genomes can be retrieved from this contamination. Our study highlights the value

81 of database sequences for exploratory symbiont screens and argues against neglecting the filtered

82 'contaminants' in sequencing projects.

83 **Materials & methods**

84       Reference sequences of 18 common *Apis* associated symbionts and pathogens were

85 compiled to be used as baits to detect presence of the microbe (Table S1). In order to reduce the

86 computational expense of all following steps, only short signature sequences were used instead of

87 complete genomes; where possible these were of slowly evolved housekeeping genes to allow a

88 range of diversity to be recovered through sequence similarity to the bait. For previously

89 identified bacterial symbionts for example, we included a 16S rRNA sequence for each known

90 associate. Next, we searched for honey bee sequencing projects in NCBI's short read archive,

91 using the search term *'Apis'*, and excluding transcriptome and microbiome (e.g., metagenome or

92 amplicon sequencing) projects. At the time of the search, 306 experiments matched these criteria,

93 including 32 using museum specimens. We downloaded all short read libraries associated with

94 these experiments (993 in total, Table S2) and mapped all reads of each of the libraries to the

95 reference sequences using NextGenMap version 0.4.12 (Sedlazeck et al. 2013). If at least 1,000

96 reads of a library were aligned to one or more sequence baits, we extracted the matching reads

97 and assembled them using SPAdes version 3.7 (Bankevich et al. 2012). Contigs resulting from

98  this assembly were then subject to taxonomic annotation via BLAST+ (Camacho et al. 2009)

99  searches against a local copy of the NCBI 'nt' database, and the Blobtools package  (Kumar et al.

100  2013). Detailed description of all steps outlined above can be found under

101  https://github.com/gerthmicha/symbiont-sra.

102      Since this approach yielded a high number of hits to various *Lactobacillus* species, we

103  repeated the entire procedure using 620 16S bait sequences from *Lactobacillus* only. These

104  sequences were taken from a previously compiled dataset of Lactobacilli associated with *Apis*,

105  other Hymenoptera, and other *Lactobacillus* sequences retrieved from public databases

106  (McFrederick et al. 2013). All hits short than 250bp were discarded, and remaining contigs were

107  combined with the reference sequences. We used SSU-ALIGN version 0.1 (Nawrocki 2009) to

108  align and mask this dataset based on conserved secondary structure. Original and masked

109  alignments are available from https://github.com/gerthmicha/symbiont-sra. A maximum

110  likelihood phylogeny was reconstructed from the complete 16S alignment (740 sequences in

111  total) using IQTREE version 1.3.10 (Nguyen et al. 2015) with automated model selection and

112  1,000 ultrafast bootstraps (Minh et al. 2013) to assess node support. The resulting tree was

113  visualized using the online tool Evolview (He et al. 2016). Furthermore, as an approximate

114  measure for the number of *Lactobacillus* OTUs recovered with our approach, we used the

115  average neighbor clustering algorithm as implemented in mothur version 1.34.4 (Schloss et al.

116  2009).

117      Although our aim was not to recover all, but only the highly covered symbiont data from

118  honey bee short reads, we wanted to test if our screening approach yields comparable results to

119  more commonly used metagenomic approaches. To this end, we screened the reads of a

120  metagenomic dataset created from the pooled DNA of 150 honeybee worker hindguts (Engel et

121  al. 2012; ~43M 150bp paired-end reads, SRA accession: SRR5237156) for *Lactobacillus* in the

122  same way as described above. We found 6 different *Lactobacillus* 16S sequences, all within the

  
123    Firm-4 and Firm-5 *Lactobacillus* groups (Fig. S1). This was in agreement to the results obtained

124    from taxonomic profiling approaches performed by Engel et al. (2012) and thus confirmed the

125    general effectiveness of our approach (Fig. S1).

126        Next, we aimed to validate that whole symbiont genomes can in principle be recovered

127    from *Apis* sequencing projects. To this end, we chose one sequencing library (SRR1046114,

128    ~85.5M 100bp paired-end reads) that contained 'contamination' from two *Lactobacillus* strains

129    (*Lactobacillus kunkeei* & *Fructobacillus* sp.). We performed a *de novo* assembly using all reads

130    with MEGAHIT version 1.0.4-beta (Li et al. 2015). All resulting contigs of this assembly were

131    taxonomically assigned to either *L. kunkeei*, *Fructobacillus* sp. or 'other' based on BLAST

132    searches, GC distributions, and read coverage. Reads matching to contigs from either

133    *Lactobacillus* strain were then separately re-assembled using SPAdes, and all contigs smaller than

134    500bp discarded. Completeness and contamination of the novel draft genomes were assessed

135    using CheckM version 1.0.6 (Parks et al. 2015), and annotation performed with PROKKA

136    version 1.12 (Seemann 2014). The annotated draft genomes are available under under

137    https://github.com/gerthmicha/symbiont-sra and via NCBI accession numbers XXXX00000000

138    (*L. kunkeei*) and YYYY00000000 (*Fructobacillus* sp.). To evaluate the evolutionary relationships

139    of newly assembled genomes in a broader taxonomic context, we assessed their phylogenetic

140    placement. Whole-genome datasets were compiled for both strains (13 *L. kunkeei* genomes, 9

141    *Fructobacillus* & *Leuconostoc* genomes altogether, Table S3). For each of the datasets, single

142    copy orthologs were identified using OrthoFinder version 0.2.8 (Emms & Kelly 2015).

143    Recombining loci were identified by using the pairwise homoplasy index test (Bruen et al. 2006),

144    and removed from subsequent analyses (window size = 20 amino acid positions, significance

145    cutoff at 0.05). Using IQTREE, we performed maximum likelihood analysis of two final

146    supermatrices (947 loci and 290,774 aa for the *L. kunkeii* dataset, 435 loci and 145,069 positions

147 for the *Fructobacillus*/*Leuconostoc* dataset). Prior to this, best-fitting partitioning schemes and

148 models were selected using the 'greedy' scheme implemented in IQTREE (Lanfear et al. 2012).

149 Using the same approach, we assembled and annotated a *Spiroplasma melliferum* genome

150 (NCBI accession ZZZZ00000000) from library SRR957082, (~224.5M 50bp single end reads).

151 Phylogenetic analysis was performed based on a dataset of 206 concatenated single copy genes

152 (58,950 amino acid positions) shared among 17 *Spiroplasma* strains (Table S3). Furthermore, to

153 assess synteny, the newly assembled draft genome was ordered against and aligned with other

154 *Spiroplasma melliferum* genomes (one genome each of strains IPMB4A and KC3) using the

155 progressiveMauve algorithm of Mauve development snapshot version 2015-02-13 (Darling et al.

156 2010).

157 **Results**

158 Using bait sequences of 18 common *Apis*- associated microbes, we found non-target

159 symbiont reads in 89 of the 993 investigated libraries (~9%). Taxonomic annotation revealed that

160 the detected sequences belong to one of three categories (Fig. 1): 1) *Apis*- associated symbionts

161 that were targeted with our bait sequences, 2) *Apis*- associated taxa that we did not target with our

162 approach, 3) microbial sequences from other sources for which there is no current evidence of

163 *Apis* association. Category 1 included sequences from 3 of the 18 targeted *Apis*- associated taxa

164 (*Crithidia*, *Nosema*, and *Spiroplasma*, Fig. 1, see also Table S4). The second category included

165 mostly honey bee gut bacteria, such as *Lactobacillus*, *Gilliamella* and *Bartonella* (Fig. 1, Table

166 S4). The third category included sequences from fungi (Ascomycota), plants, and the bacterium

167 *Thermus*, that were likely not part of the native microbiome of the sequenced samples. All of

168 these contaminations were crossed-checked via manual online BLAST searches and were

169 confirmed to represent 'true' hits with high and continuous identities with the respective database

170 sequences.

171    Because the majority of hits in this first screening process were Lactobacilli, we repeated

172    the screening, this time using only *Lactobacillus* 16S sequences as baits. We found 121

173    *Lactobacillus* sequences in 40 of the 993 investigated libraries, corresponding to 25 OTUs

174    (estimated with mothur using a 5% cutoff). In our phylogenetic analysis based on 16S rRNA

175    sequences, most of the detected strains clustered within *Lactobacillus* groups known to be

176    associated with honey bees (Fig. 2a). Of the recovered sequences not clustering within these

177    lineages, three were found to group with other *Apis*- associated Lactobacilli as sister group to the

178    *Lactobacillus coryniformis* group (Fig. 2a). Online BLAST searches revealed *Fructobacillus*

179    species as closest matches based on 16S rRNA sequence.

180    Next, we aimed at recovering draft genome sequences of bee-associated Lactobacilli. We

181    chose a sequencing library from which 16S sequences of both *L. kunkeei* and *Fructobacillus*

182    isolates were detected in our screen. The contigs of a meta-assembly were taxonomically

183    annotated, and reads matching to the respective target taxa were then assembled and annotated

184    separately. For each assembly, we performed a phylogenetic analysis based on all single copy

185    orthologs shared with related genomes (Fig. 2b, c), thus confirming the identity of the strains as

186    *L. kunkeei* (Fig. 2b) and *Fructobacillus* (Fig. 2c). Both genomes were highly covered and mostly

187    complete based on the presence of conserved markers (Fig. 2d). Finally, we recovered the

188    genome of a *Spiroplasma melliferum* strain from another *Apis* sequencing library (Fig. 3). In the

189    meta-assembly, *Spiroplasma* and *Apis* contigs could be clearly separated by coverage and

190    taxonomic annotations (Fig. 3b). The refined assembly resulted in a highly covered draft genome

191    of *Spiroplasma melliferum*, which is very similar to the two previously sequenced *Spiroplasma*

192    *melliferum* strains (Alexeev et al. 2012; Lo et al. 2013), based on shared ortholog clusters,

193    genome organisation, and phylogeny (Fig. 3a, c, d).

194    **Discussion**

195    We used two screens to determine if microbial symbiont data can be retrieved from

196    sequencing projects targeting *Apis* (honey bees). First, by using bait sequences of *Apis* symbionts

197    and pathogens, we found evidence for the presence of these taxa in 9% of 993 *Apis* short read

198    libraries. This measure of non-target 'contamination' can be considered as conservative, since our

199    approach only reports relatively high levels of contamination (at least 1000 reads per bait

200    sequence). Three common honey bee pathogens were detected with this approach: *Nosema*,

201    *Crithidia*, and *Spiroplasma. Nosema* are microsporidian gut parasites of various honey bee

202    species, and while the sampling of our screen is not representative, this finding corroborates the

203    recognition of *Nosema* as widespread pathogen of honey bee colonies worldwide (Nixon 1982;

204    Klee et al. 2007). *Crithidia* (Trypanosomatidae), another gut pathogen of *Apis* and related bee

205    species (Schwarz et al. 2015) was detected at an even higher frequency (Fig. 1, Table S4).

206    Finally, we found *Spiroplasma melliferum* in one of the investigated sequencing libraries.

207    *Spiroplasma* are common symbiotic bacteria of many invertebrates (Duron et al. 2008) and have

208    been connected to pathogenicity in honey bees (Clark 1977). The bait sequences of all of these

209    pathogens showed a high coverage in our screen, suggesting that novel genetic variants can be

210    recovered from already available data, or from data that will become available as by-product of

211    future honey bee sequencing projects. We did not find any viral sequences in our screen (Table

212    S1), probably because most honey bee viruses are RNA viruses (Chen et al. 2004), that are in

213    retrospect unlikely to be picked up with WGS approaches (but could potentially be retrieved from

214    RNAseq data).

215    This first screen also revealed the presence of many reads originating from *Apis* gut

216    microbes. These reads were the most common 'contamination' detected in the libraries, despite

217    these taxa not being specifically targeted. The microbiome of healthy honey bees is dominated by

218    Lactobacilli (Kwong & Moran 2016), and this is also reflected in our results (Fig1, Table S1).

219    Furthermore, a number of taxa that are likely not part of the natural *Apis* microbiome were

220   detected. For example, we detected *Aspergillus* in several sequencing libraries that originated

221   from museum material, which likely represents post mortem saprophytic growth. We also

222   retrieved hits to plant sequences which might originate from co-amplified and sequenced pollen

223   DNA (Fig. 1). We further detected *Thermus*, which is best explained by contaminated laboratory

224   reagents or sequencing kits (Salter et al. 2014). This 'false discovery' illustrates an important

225   caveat in our approach: the differentiation between host-associated microbes and microbes from

226   other sources may not always be possible, and will be particularly difficult for museum

227   specimens. Though not problematic in the examples we present, the situation is likely more

228   complicated in hosts with a less well-investigated microbiome, or for symbionts that are very

229   similar to environmental taxa. In these cases, the approach will establish candidates that will then

230   require direct validation.

231      In the second screen, targeted only at *Lactobacillus*, our protocol detected 25 taxonomically

232   different *Lactobacillus* strains. Our phylogenetic reconstruction of *Lactobacillus* relationships

233   based on 16S rRNA generally reflected the current understanding of this genus' taxonomy (Felis

234   & Dellaglio 2007; Salvetti et al. 2012), and revealed that most Lactobacilli known to be

235   associated with honey bees are also present in *Apis* short read libraries. This includes Firm-4 and

236   Firm-5 Lactobacilli, both of which are honey bee hindgut colonizers, and *L. kunkeei*, which is

237   common in nectar and other hive material, and sometimes found in honey bee crops (Kwong &

238   Moran 2016). Furthermore, we found *Fructobacillus,* which share an ecological niche with *L.*

239   *kunkeei,* i.e., they are found in flowers, nectar, and in honey bee guts (Endo et al. 2009; Endo &

240   Salminen 2013) . Although not classified as such, recent phylogenomic evidence suggests that

241   *Fructobacillus* (and the closely related *Leuconostoc*) are part of the *Lactobacillus* radiation (Sun

242   et al. 2015). Here, we also infer *Fructobacillus* grouping within, rather than outside of

243   *Lactobacillus* (Fig. 2a). These results show that a reasonably accurate understanding of

244     *Lactobacillus* community composition in honey bees can be gained from non-target sequences

245     produced as a by-product of honey bee sequencing projects.

246        Finally, we demonstrate that draft genomes of microbial symbionts can be recovered from

247     *Apis* short reads. For example, inspecting the non-target components of just a single *Apis*

248     sequencing library produced novel, highly covered, and near complete draft genomes of

249     *Lactobacillus kunkeii* and a *Fructobacillus* strain (Fig. 2 b,c,d). Although the 16S sequence of the

250     *Fructobacillus* strain best matched *F. fructosus*, our analysis suggests it belongs to a species so

251     far not represented by genomic sequences in public databases, or even a novel species (Fig. 2c).

252     Conceivably, many additional *Lactobacillus* variants could be retrieved from the libraries

253     investigated here, potentially providing a more complete picture of the *Apis* microbiome

254     composition and function. It should be noted that draft genomes reconstructed this way must be

255     regarded as 'population consensus' genomes, as opposed to genomes sequenced from cultured

256     bacterial clones. While these genomes cannot be linked to a bacterial clone, they still provide

257     information of metabolic capacities within the *Apis* microbiome.

258        Although our study was focused on *Apis*, it is conceivable that the amount of non-target

259     'contamination' is similar for other sequencing projects. As a best practice in any sequencing

260     project, we therefore suggest that all non-target taxa should be identified, and their genomes

261     assembled, annotated, and published alongside the target genome. This requires less effort than it

262     may seem, as de-contamination is already a standard post-processing step. Instead of discarding

263     the contaminated reads, they can be processed with one of many available software solutions that

264     automate the process of identifying and assembling genomes from metagenomes (Oulas et al.

265     2015), thus minimizing the additional workload. Not only would this provide the community

266     with valuable genomic data of microbial symbionts from known host taxa, but it can additionally

267     be argued that this is the most sensible thing to do from a biological point of view. Evidence is

268     mounting that symbiotic microbes influence almost all aspects of their host's biology (Douglas

269    2014; Bordenstein & Theis 2015). Taking into account the total genomic information recovered

270    in sequencing projects may therefore provide a more complete picture of the target organism's

271    biology.


272    **Conclusion**

273       The biological properties of an individual are a composite of the functions encoded in their

274    genome and that of microbial associates, the 'hologenome'. We here revisited published short

275    read data from *Apis* spp. sequencing projects to investigate if these give insight into the wider set

276    of associates that are commonly disgarded as 'contaminants'. We found that a large variety of

277    distinct *Apis*-associated microbial symbionts and pathogens can be detected as 'contamination' in

278    these data. Further, due to the large depths of today's sequencing projects, the genomes of some

279    microbial associates (which are typically much smaller than the target genomes) can often be

280    recovered in high quality. Honey bees have a comparatively simple microbiota (Kwong & Moran

281    2016) and are thus considered suitable models for microbiome-animal interactions and evolution

282    (Engel et al. 2016). Their enormous economic importance (Calderone 2012), has driven the large

283    (and still increasing) number of honey bee sequencing projects. Our examination of the output of

284    these projects suggests that large amounts of genomic information on bee-associated microbes

285    are included in these data. While genomes gained from contaminated bee samples cannot and

286    should not replace focused microbiological and metagenomic investigations, they might still

287    improve our understanding of honey bee microbiome composition and functioning.

## References

294 **References**

295  Alexeev D, Kostrjukova E, Aliper A, Popenko A, Bazaleev N, Tyakht A, Selezneva O, Akopian T,
296      Prichodko E, Kondratov I, et al (2012) Application of *Spiroplasma melliferum*
297      proteogenomic profiling for the discovery of virulence factors and pathogenicity
298      mechanisms in host-associated spiroplasmas. *Journal of Proteome Research* **11**, 224–236.

299  Bankevich A, Nurk S, Antipov D, Gurevich AA, Dvorkin M, Kulikov AS, Lesin VM, Nikolenko
300      SI, Pham S, Prjibelski AD, Pyshkin AV, Sirotkin AV, Vyahhi N, Tesler G, Alekseyev MA,
301      Pevzner PA (2012) SPAdes: a new genome assembly algorithm and its applications to
302      single-cell sequencing. *Journal of Computational Biology* **19**, 455–477.

303  Bordenstein SR, Theis KR (2015) Host biology in light of the microbiome: ten principles of
304      holobionts and hologenomes. *PLoS Biology* **13**, e1002226.

305  Browne HP, Forster SC, Anonye BO, Kumar N, Neville BA, Stares MD, Goulding D, Lawley TD
306      (2016) Culturing of `unculturable' human microbiota reveals novel taxa and extensive
307      sporulation. *Nature* **533**, 543–546.

308  Bruen TC, Philippe H, Bryant D (2006) A simple and robust statistical test for detecting the
309      presence of recombination. *Genetics* **172**, 2665–2681.

310  Calderone NW (2012) Insect pollinated crops, insect pollinators and US agriculture: trend
311      analysis of aggregate data for the period 1992–2009. *PLoS ONE* **7**, e37235.

312  Camacho C, Coulouris G, Avagyan V, Ma N, Papadopoulos J, Bealer K, Madden TL (2009)
313      BLAST+: architecture and applications. *BMC Bioinformatics* **10**, 421.

314  Caporaso JG, Lauber CL, Walters WA, Berg-Lyons D, Huntley J, Fierer N, Owens SM, Betley J,
315      Fraser L, Bauer M, Gormley N, Gilbert JA, Smith G, Knight R (2012) Ultra-high-
316      throughput microbial community analysis on the Illumina HiSeq and MiSeq platforms. *The*
317      *ISME Journal* **6**, 1621–1624.

318  Chen Y, Zhao Y, Hammond J, Hsu H-t, Evans J, Feldlaufer M (2004) Multiple virus infections in
319      the honey bee and genome divergence of honey bee viruses. *Journal of Invertebrate*
320      *Pathology* **87**, 84–93.

321  Clark T (1977) *Spiroplasma* sp, a new pathogen in honey bees. *Journal of Invertebrate Pathology*
322      **113**, 112–113.

323  Colbourne JK, Pfrender ME, Gilbert D, Thomas WK, Tucker A, Oakley TH, Tokishita S, Aerts A,
324      Arnold GJ, Basu MK, Bauer DJ, Caceres CE, Carmel L, Casola C, Choi J-H, Detter JC,

Dong Q, Dusheyko S, Eads BD, Frohlich T, Geiler-Samerotte KA, Gerlach D, Hatcher P, Jogdeo S, Krijgsveld J, Kriventseva EV, Kultz D, Laforsch C, Lindquist E, Lopez J, Manak JR, Muller J, Pangilinan J, Patwardhan RP, Pitluck S, Pritham EJ, Rechtsteiner A, Rho M, Rogozin IB, Sakarya O, Salamov A, Schaack S, Shapiro H, Shiga Y, Skalitzky C, Smith Z, Souvorov A, Sung W, Tang Z, Tsuchiya D, Tu H, Vos H, Wang M, Wolf YI, Yamagata H, Yamada T, Ye Y, Shaw JR, Andrews J, Crease TJ, Tang H, Lucas SM, Robertson HM, Bork P, Koonin EV, Zdobnov EM, Grigoriev IV, Lynch M, Boore JL (2011) The ecoresponsive genome of *Daphnia pulex*. *Science* **331**, 555–561.

Cox-Foster DL, Conlan S, Holmes EC, Palacios G, Evans JD, Moran NA, Quan P-L, Briese T, Hornig M, Geiser DM, Martinson V, vanEngelsdorp D, Kalkstein AL, Drysdale A, Hui J, Zhai J, Cui L, Hutchison SK, Simons JF, Egholm M, Pettis JS, Lipkin WI (2007) A metagenomic survey of microbes in honey bee colony collapse disorder. *Science* **318**, 283–287.

Darling AE, Mau B, Perna NT (2010) Progressivemauve: multiple genome alignment with gene gain, loss and rearrangement. *PLoS ONE* **5**, e11147.

Douglas AE (2014) Multiorganismal insects: diversity and function of resident microorganisms. *Annual Review of Entomology* **60**, 1–18.

Douglas AE, Werren JH (2016) Holes in the hologenome: why host-microbe symbioses are not holobionts. *mBio* **7**, e02099–15.

Duron O, Bouchon D, Boutin S, Bellamy L, Zhou L, Engelstädter J, Hurst GDD (2008) The diversity of reproductive parasites among arthropods: *Wolbachia* do not walk alone. *BMC Biology* **6**, 27.

Emms DM, Kelly S (2015) OrthoFinder: solving fundamental biases in whole genome comparisons dramatically improves orthogroup inference accuracy. *Genome Biology* **16**, 157.

Endo A, Futagawa-Endo Y, Dicks LM (2009) Isolation and characterization of fructophilic lactic acid bacteria from fructose-rich niches. *Systematic and Applied Microbiology* **32**, 593–600.

Endo A, Salminen S (2013) Honeybees and beehives are rich sources for fructophilic lactic acid bacteria. *Systematic and Applied Microbiology* **36**, 444–448.

Engel P, Kwong WK, McFrederick Q, Anderson KE, Barribeau SM, Chandler JA, Cornman RS, Dainat J, de Miranda JR, Doublet V, Emery O, Evans JD, Farinelli L, Flenniken ML, Granberg F, Grasis JA, Gauthier L, Hayer J, Koch H, Kocher S, Martinson VG, Moran N, Munoz-Torres M, Newton I, Paxton RJ, Powell E, Sadd BM, Schmid-Hempel P, Schmid-

358  Hempel R, Song SJ, Schwarz RS, vanEngelsdorp D, Dainat B (2016) The bee microbiome:
359       impact on bee health and model for evolution and ecology of host-microbe interactions.
360       *mBio* **7**, e02164–15.

361  Engel P, Martinson VG, Moran NA (2012) Functional diversity within the simple gut microbiota
362       of the honey bee. *Proceedings of the National Academy of Sciences of the United States of*
363       *America* **109**, 11002–11007.

364  Evans JD, Schwarz RS (2011) Bees brought to their knees: microbes affecting honey bee health.
365       *Trends in Microbiology* **19**, 614–620.

366  Felis GE, Dellaglio F (2007) Taxonomy of lactobacilli and bifidobacteria. *Current Issues in*
367       *Intestinal Microbiology* **8**, 44.

368  Goodwin S, McPherson JD, McCombie WR (2016) Coming of age: ten years of next-generation
369       sequencing technologies. *Nature Reviews Genetics* **17**, 333–351.

370  He Z, Zhang H, Gao S, Lercher MJ, Chen W-H, Hu S (2016) Evolview v2: an online
371       visualization and management tool for customized and annotated phylogenetic trees.
372       *Nucleic Acids Research* **44**, W236–W241.

373  Kircher M, Kelso J (2010) High-throughput DNA sequencing - concepts and limitations.
374       *BioEssays* **32**, 524–536.

375  Klee J, Besana AM, Genersch E, Gisder S, Nanetti A, Tam DQ, Chinh TX, Puerta F, Ruz JM,
376       Kryger P, Message D, Hatjina F, Korpela S, Fries I, Paxton RJ (2007) Widespread dispersal
377       of the microsporidian *Nosema ceranae*, an emergent pathogen of the western honey bee,
378       *Apis mellifera*. *Journal of Invertebrate Pathology* **96**, 1–10.

379  Kumar S, Jones M, Koutsovoulos G, Clarke M, Blaxter ML (2013) Blobology: exploring raw
380       genome data for contaminants, symbionts and parasites using taxon-annotated GC-coverage
381       plots. *Frontiers in Genetics* **4**, 237.

382  Kwong WK, Moran NA (2016) Gut microbial communities of social bees. *Nature Reviews*
383       *Microbiology* **14**, 374–384.

384  Lanfear R, Calcott B, Ho SYW, Guindon S (2012) Partitionfinder: combined selection of
385       partitioning schemes and substitution models for phylogenetic analyses. *Molecular Biology*
386       *and Evolution* **29**, 1695–701.

387  Li D, Liu C-M, Luo R, Sadakane K, Lam T-W (2015) MEGAHIT: an ultra-fast single-node
388       solution for large and complex metagenomics assembly via succinct de Bruijn graph.
389       *Bioinformatics* **31**, 1674–1676.

390 Lo W-S, Chen L-L, Chung W-C, Gasparich GE, Kuo C-H (2013) Comparative genome analysis
391       of *Spiroplasma melliferum* IPMB4A, a honeybee-associated bacterium. *BMC Genomics* **14**,
392       22.

393 McFall-Ngai M, Hadfield MG, Bosch TCG, Carey HV, Domazet-Lošo T, Douglas AE, Dubilier
394       N, Eberl G, Fukami T, Gilbert SF, Hentschel U, King N, Kjelleberg S, Knoll AH, Kremer
395       N, Mazmanian SK, Metcalf JL, Nealson K, Pierce NE, Rawls JF, Reid A, Ruby EG,
396       Rumpho M, Sanders JG, Tautz D, Wernegreen JJ (2013) Animals in a bacterial world, a
397       new imperative for the life sciences. *Proceedings of the National Academy of Sciences of*
398       *the United States of America* **110**, 3229–3236.

399 McFrederick QS, Cannone JJ, Gutell RR, Kellner K, Plowes RM, Mueller UG (2013) Specificity
400       between lactobacilli and hymenopteran hosts is the exception rather than the rule. *Applied*
401       *and Environmental Microbiology* **79**, 1803–12.

402 Minh BQ, Nguyen MAT, von Haeseler A (2013) Ultrafast Approximation for Phylogenetic
403       Bootstrap. *Molecular Biology and Evolution* **30**, 1188–1195.

404 Moran NA, Sloan DB (2015) The Hologenome Concept: Helpful or Hollow?. *PLOS Biology* **13**,
405       e1002311.

406 Nawrocki EP (2009) *Structural RNA homology search and alignment using covariance models*.
407       PhD Thesis, Washington University in Saint Louis.

408 Neumann P, Carreck NL (2010) Honey bee colony losses. *Journal of Apicultural Research* **49**, 1–
409       6.

410 Nguyen L-T, Schmidt HA, von Haeseler A, Minh BQ (2015) IQ-TREE: a fast and effective
411       stochastic algorithm for estimating maximum-likelihood phylogenies. *Molecular Biology*
412       *and Evolution* **32**, 268–274.

413 Nixon M (1982) Preliminary world maps of honeybee diseases and parasites. *Bee World* **63**, 23–
414       42.

415 Oulas A, Pavloudi C, Polymenakou P, Pavlopoulos GA, Papanikolaou N, Kotoulas G, Arvanitidis
416       C, Iliopoulos4 I (2015) Metagenomics: tools and insights for analyzing next-generation
417       sequencing data derived from biodiversity studies. *Bioinformatics and Biology Insights* ,
418       75.

419 Parks DH, Imelfort M, Skennerton CT, Hugenholtz P, Tyson GW (2015) CheckM: assessing the
420       quality of microbial genomes recovered from isolates, single cells, and metagenomes.
421       *Genome Research* **25**, 1043–1055.

422 Richardson MF, Weinert LA, Welch JJ, Linheiro RS, Magwire MM, Jiggins FM, Bergman CM
423     (2012) Population genomics of the *Wolbachia* endosymbiont in *Drosophila melanogaster*.
424     *PLoS Genetics* **8**, e1003129.

425 Riesenfeld CS, Schloss PD, Handelsman J (2004) Metagenomics: genomic analysis of microbial
426     communities. *Annual Review of Genetics* **38**, 525–552.

427 Salter SJ, Cox MJ, Turek EM, Calus ST, Cookson WO, Moffatt MF, Turner P, Parkhill J, Loman
428     NJ, Walker AW (2014) Reagent and laboratory contamination can critically impact
429     sequence-based microbiome analyses. *BMC Biology* **12**, .

430 Salvetti E, Torriani S, Felis GE (2012) The genus *Lactobacillus*: a taxonomic update. *Probiotics*
431     *& Antimicrobial Proteins* **4**, 217–226.

432 Salzberg SL, Hotopp J, Delcher A, Pop M, Smith D, Eisen M, Nelson W (2005) Serendipitous
433     discovery of *Wolbachia* genomes in multiple *Drosophila* species. *Genome Biology* **6**, R23.

434 Sboner A, Mu X, Greenbaum D, Auerbach RK, Gerstein MB (2011) The real cost of sequencing:
435     higher than you think!. *Genome Biology* **12**, 125.

436 Schloss PD, Westcott SL, Ryabin T, Hall JR, Hartmann M, Hollister EB, Lesniewski RA, Oakley
437     BB, Parks DH, Robinson CJ, Sahl JW, Stres B, Thallinger GG, Horn DJV, Weber CF
438     (2009) Introducing mothur: open-source, platform-independent, community-supported
439     software for describing and comparing microbial communities. *Applied and Environmental*
440     *Microbiology* **75**, 7537–7541.

441 Schmieder R, Edwards R (2011) Fast identification and removal of sequence contamination from
442     genomic and metagenomic datasets. *PLoS ONE* **6**, e17288.

443 Schwarz RS, Bauchan GR, Murphy CA, Ravoet J, Graaf DC, Evans JD (2015) Characterization
444     of two species of Trypanosomatidae from the honey bee *Apis mellifera*: *Crithidia mellificae*
445     Langridge and McGhee, and *Lotmaria passim* n gen, n sp. *Journal of Eukaryotic*
446     *Microbiology* **62**, 567–583.

447 Sedlazeck FJ, Rescheneder P, von Haeseler A (2013) NextGenMap: fast and accurate read
448     mapping in highly polymorphic genomes. *Bioinformatics* **29**, 2790–2791.

449 Seemann T (2014) Prokka: rapid prokaryotic genome annotation. *Bioinformatics* **30**, 2068–2069.

450 Sun Z, Harris HMB, McCann A, Guo C, Argimón S, Zhang W, Yang X, Jeffery IB, Cooney JC,
451     Kagawa TF, Liu W, Song Y, Salvetti E, Wrobel A, Rasinkangas P, Parkhill J, Rea MC,
452     O'Sullivan O, Ritari J, Douillard FP, Ross RP, Yang R, Briner AE, Felis GE, de Vos WM,
453     Barrangou R, Klaenhammer TR, Caufield PW, Cui Y, Zhang H, O'Toole PW (2015)

454  Expanding the biotechnology potential of lactobacilli through comparative genomics of 213

455  strains and associated genera. *Nature Communications* **6**, 8322.

456 Tamarit D, Ellegaard KM, Wikander J, Olofsson T, Vasquez A, Andersson SGE (2015)

457  Functionally structured genomes in *Lactobacillus kunkeei* colonizing the honey crop and

458  food products of honeybees and stingless bees. *Genome Biology and Evolution* **7**, 1455–

459  1473.

460 Wolfe KH, Li W-H (2003) Molecular evolution meets the genomics revolution. *Nature Genetics*

461  **33**, 255–265.

462  **Figure legends**

463  **Figure 1:** Taxonomic annotation of contigs assembled from 'contaminated' *Apis* short read

464  libraries. Bar chart shows the frequency of each taxonomic category assigned by best BLAST

465  matches against NCBI's 'nt' database, as the number of libraries in which that taxon was detected

466  (in the sample of 993 SRA libraries). Bold categories are 'phyla', as defined in

467  https://www.ncbi.nlm.nih.gov/taxonomy, taxa in italics represent typical genera that were

468  recovered within each phylum. See Table S4 for a complete list.

469  **Figure 2:** 'Contamination' from Lactobacilli in *Apis* short read libraries. a) Maximum likelihood

470  tree of 720 16S rRNA sequences from Lactobacilli. Branch colors and the color of the outer

471  annotation circle correspond to *Lactobacillus* species groups according to Felis & Dellaglio

472  (2007). Inner circle demarks taxa found Hymenoptera (grey squares) and in corbiculate apids

473  (honey bees and relatives, black squares). *Lactobacillus* sequences recovered in this study from

474  contaminated *Apis* libraries are labeled with blue triangles. The Lactobacilli typically associated

475  with honey bees (Firm-4, Firm-5, *L. kunkeei*) are further highlighted with a blue background

476  color. Two dotted blue lines denote the taxa of which whole draft genomes were recovered. See

477  text for details. An interactive version of the tree containing all node labels is available under

478  http://www.evolgenius.info/evolview/#shared/wZcKHbwJuT. Abbreviations: al-far- alimentarius-

479  farciminis, bre- brevis, buch- buchneri, cas- casei, cor- coryniformis, del- delbrueckii, fru-

480  fructivorans, per- perolens, plan- plantarum, reu- reuteri, sak- sakei, sal- salivarius, OUT-

481  outgroup. b) Phylogeny of *Lactobacillus kunkeei* strains based on maximum likelihood analyses

482  of 947 concatenated single copy orthologs (290,774 amino acid positions). Tree is rooted with

483  *Lactobacillus apinorum* Fhon13 (taxon not shown). Strain names correspond to the names used

484  in Tamarit et al. (2015; see Table S3). Blue taxon label corresponds to the *L. kunkeei* strain

485  recovered from 'contaminants' in library SRR1046114. Bootstrap values are given on nodes. See

486   Table S3 for sources of genomes. c) Maximum likelihood tree of *Fructobacillus* (F.) and

487   *Leuconostoc* (L.) species based on 435 concatenated single copy orthologs (145,069 amino acid

488   positions). Tree is rooted with *Lactobacillus delbruecki*. Numbers on nodes correspond to

489   bootstrap values. Again, blue taxon label denotes the *Fructobacillus* genome recovered from the

490   'contaminated' library SRR1046114. Note that the phylogenetic distance between *Fructobacillus*

491   *fructosus* and the novel genome is similar to other between-species distances in this tree. See

492   Table S3 for accession numbers of all genomes used for phylogenetic analysis. d) Assembly

493   statistics for the two novel draft genomes recovered from library SRR1046114. Abbreviations:

494   CDS- coding sequences predicted with PROKKA, Comp. & Cont.- completeness and

495   contamination as estimated with CheckM version 1.0.6 (Parks et al. 2015) based on the number

496   of conserved marker loci. Phylogenetic affiliations of the two strains are depicted in Fig. 3b and

497   3c, respectively.

498   **Figure 3:** Characteristics of *Spiroplasma melliferum* isolated from a 'contaminated' *Apis*

499   sequencing library (SRR957082). a) Venn diagram illustrating the number of orthologs shared

500   between the novel strain and its closest sequenced relatives IBMB4A (Lo et al. 2013) and KC3

501   (Alexeev et al. 2012). b) Taxon-annotated GC-coverage plot of SRR951082 metaassembly

502   created with Blobology. *Spiroplasma* and *Apis* contigs can be differentiated by coverage. c)

503   Synteny across *Spiroplasma melliferum* genomes. Contigs from assemblies SRR957082 and

504   IPMB4A were ordered against KC3, the most complete of the three *S. melliferum* genomes. d)

505   Phylogenetic relationships within the genus *Spiroplasma*. Maximum likelihood tree is based on

506   206 concatenated loci (58,950 amino acid positions), numbers on branches correspond to

507   bootstrap values. *Spiroplasma* groups are highlighted with colors. The taxon label of the novel

508   genome is highlighted in bold. Accession numbers for all taxa are listed in Table S4.

509 **Supplementary files**

510 **Fig S1:** Verification of screening approach employed here using the dataset of Engel et al. (2012).

511 All short reads from this dataset were mapped against *Lactobacillus* 16S reference sequences as

512 detailed in the materials & methods section. Thus retrieved 16S sequences are highlighted with

513 thick, dark blue lines. All other taxa in this tree are identical to the ones in Fig. 2A, as is the color

514 scheme. Although the topology differs between these two *Lactobacillus* trees, it is evident that

515 the strains recovered from the Engel et al. (2012) dataset cluster within the Firm-4 and Firm-5

516 *Lactobacillus* groups. Engel et al. (2012) essentially find the same ("These distinct clusters reflect

517 the eight dominant species with the two closely related Firmicutes (Firm-4 and Firm-5) [...]"; see

518 also their Fig. 1c) using the programs MetaPhyler (http://metaphyler.cbcb.umd.edu/) and IMG/M

519 (https://img.jgi.doe.gov/) for taxonomic profiling.

520 **Table S1:** Accession numbers for all signature reference sequences used in the initial screen. The

521 sequence for *Arsenophonus* 16S was recovered from honey bee short read data (unpublished).

522 **Table S2:** A list of NCBI accession numbers for all short reads downloaded and screened in this

523 work.

524 **Table S3:** NCBI accession numbers for all genomes employed for comparative/phylogenetic

525 analyses of *Lactobacillus kunkeei*, *Fructobacillus* sp., and *Spiroplasma* sp.

526 **Table S4:** Taxonomic summary of BLAST hits for contigs created in the first round of screening.

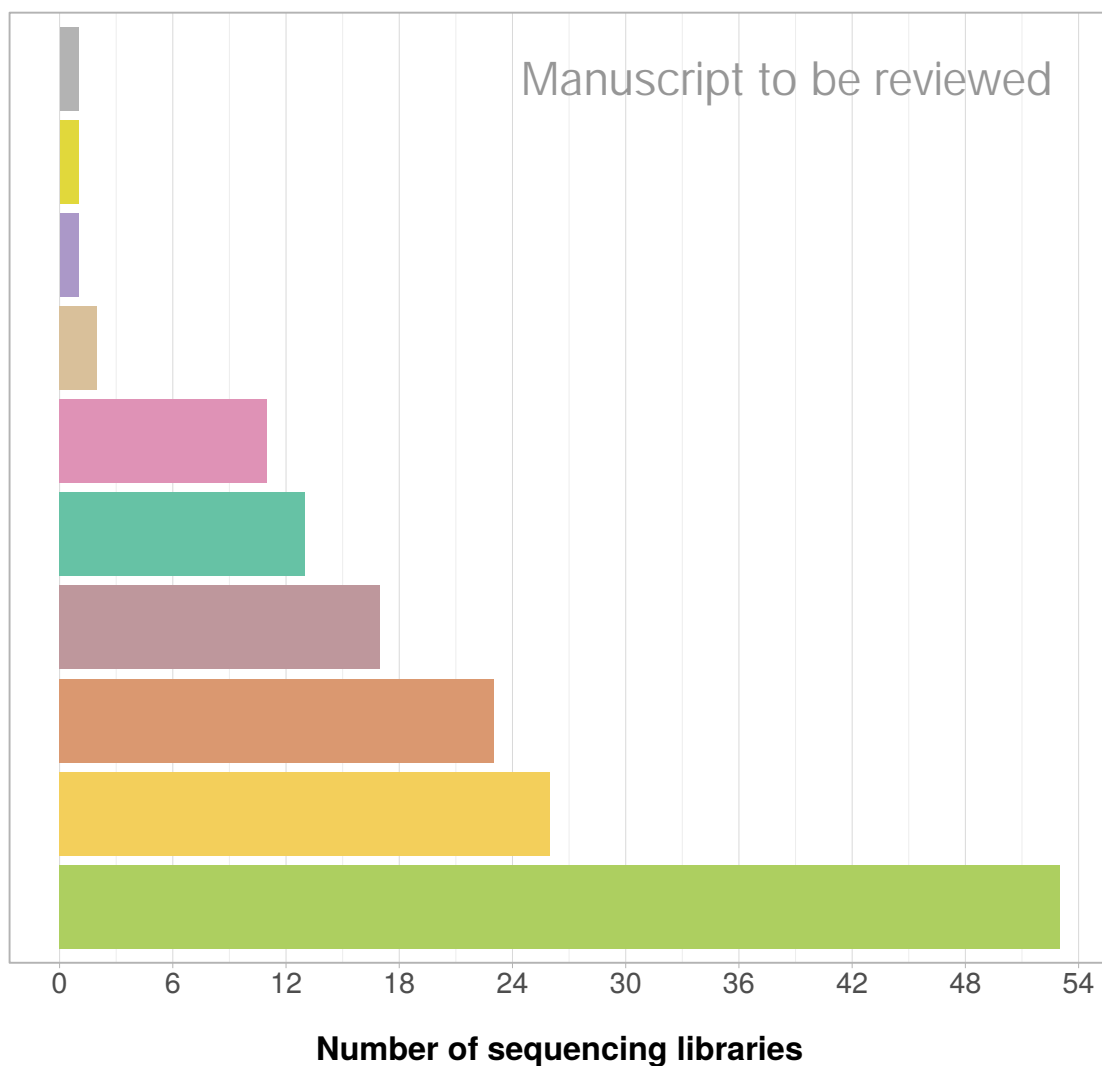527 File created with a Blobtools script (see Materials & methods).

# Figure 1(on next page)

Taxonomic annotation of contigs assembled from 'contaminated' *Apis* short read libraries.

Bar chart shows the frequency of each taxonomic category assigned by best BLAST matches against NCBI's 'nt' database, as the number of libraries in which that taxon was detected (in the sample of 993 SRA libraries). Bold categories are 'phyla', as defined in

https://www.ncbi.nlm.nih.gov/taxonomy , taxa in italics represent typical genera that were recovered within each phylum. See Table S4 for a complete list.
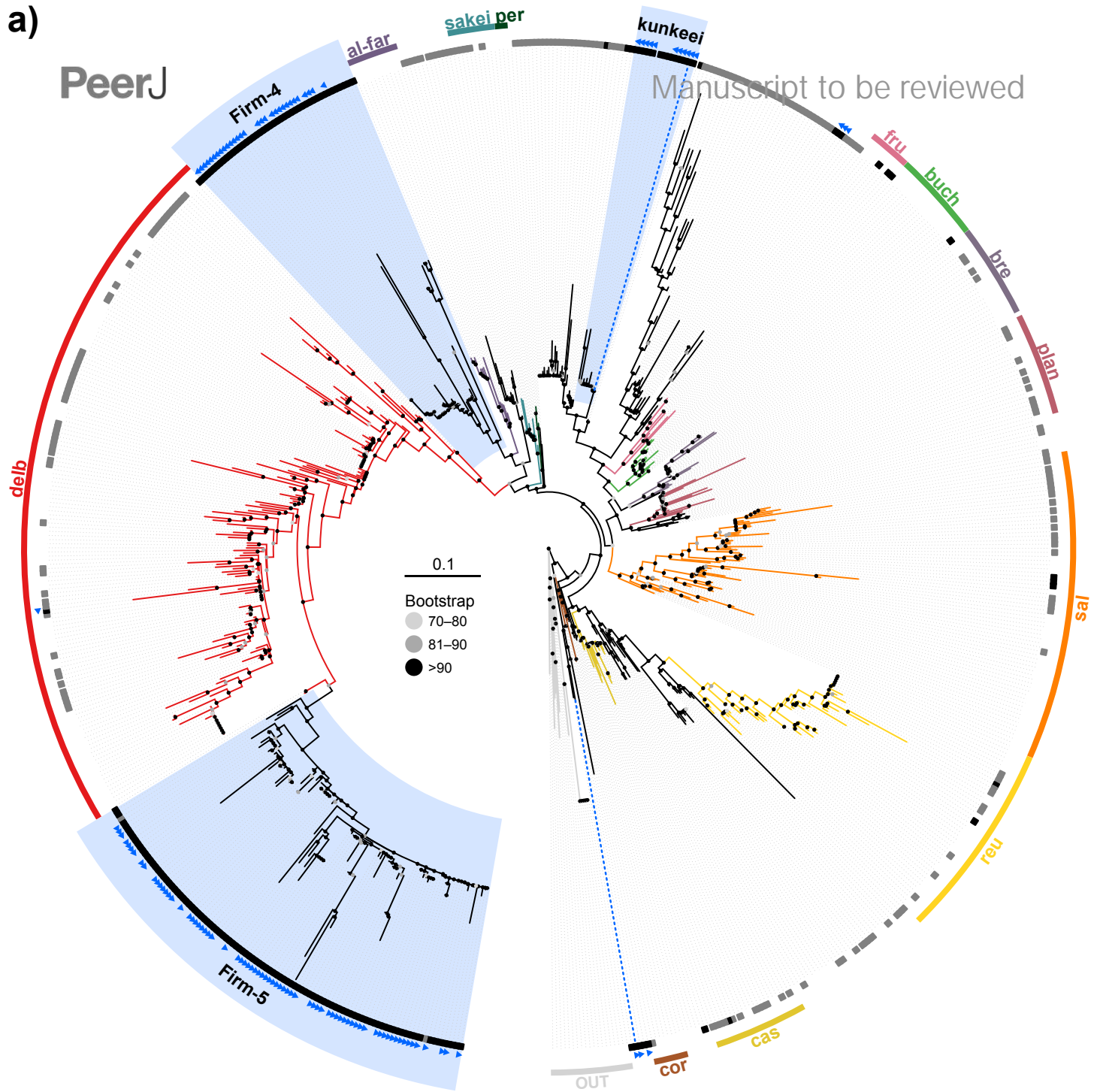
**Number of sequencing libraries**

# Figure 2(on next page)

'Contamination' from Lactobacilli in *Apis* short read libraries.

a) Maximum likelihood tree of 720 16S rRNA sequences from Lactobacilli. Branch colors and the color of the outer annotation circle correspond to *Lactobacillus* species groups according to Felis & Dellaglio (2007). Inner circle demarks taxa found Hymenoptera (grey squares) and in corbiculate apids (honey bees and relatives, black squares). *Lactobacillus* sequences recovered in this study from contaminated *Apis* libraries are labeled with blue triangles. The Lactobacilli typically associated with honey bees (Firm-4, Firm-5, *L. kunkeei*) are further highlighted with a blue background color. Two dotted blue lines denote the taxa of which whole draft genomes were recovered. See text for details. An interactive version of the tree containing all node labels is available under

http://www.evolgenius.info/evolview/#shared/wZcKHbwJuT . Abbreviations: al-far-alimentarius-farciminis, bre- brevis, buch- buchneri, cas- casei, cor- coryniformis, del-delbrueckii, fru- fructivorans, per- perolens, plan- plantarum, reu- reuteri, sak- sakei, sal-salivarius, OUT- outgroup. b) Phylogeny of *Lactobacillus kunkeei* strains based on maximum likelihood analyses of 947 concatenated single copy orthologs (290,774 amino acid positions). Tree is rooted with *Lactobacillus apinorum* Fhon13 (taxon not shown). Strain names correspond to the names used in Tamarit et al. (2015; see Table S3). Blue taxon label corresponds to the *L. kunkeei* strain recovered from 'contaminants' in library SRR1046114. Bootstrap values are given on nodes. See Table S3 for sources of genomes. c) Maximum likelihood tree of *Fructobacillus* (F.) and *Leuconostoc* (L.) species based on 435 concatenated single copy orthologs (145,069 amino acid positions). Tree is rooted with *Lactobacillus delbruecki*. Numbers on nodes correspond to bootstrap values. Again, blue taxon label denotes the *Fructobacillus* genome recovered from the 'contaminated' library SRR1046114. Note that the phylogenetic distance between *Fructobacillus fructosus* and the novel genome is similar to other between-species distances in this tree. See Table S3 for accession numbers of all genomes used for phylogenetic analysis. d) Assembly statistics for the two novel draft

genomes recovered from library SRR1046114. Abbreviations: CDS- coding sequences predicted with PROKKA, Comp. & Cont.- completeness and contamination as estimated with CheckM version 1.0.6 (Parks et al. 2015) based on the number of conserved marker loci. Phylogenetic affiliations of the two strains are depicted in Fig. 3b and 3c, respectively.

**Figure 3**(on next page)

Characteristics of *Spiroplasma melliferum* isolated from a 'contaminated' *Apis* sequencing library (SRR957082).

a) Venn diagram illustrating the number of orthologs shared between the novel strain and its closest sequenced relatives IBMB4A (Lo et al. 2013) and KC3 (Alexeev et al. 2012). b) Taxon-annotated GC-coverage plot of SRR951082 metaassembly created with Blobology. *Spiroplasma* and *Apis* contigs can be differentiated by coverage. c) Synteny across *Spiroplasma melliferum* genomes. Contigs from assemblies SRR957082 and IPMB4A were ordered against KC3, the most complete of the three *S. melliferum* genomes. d) Phylogenetic relationships within the genus *Spiroplasma*. Maximum likelihood tree is based on 206 concatenated loci (58,950 amino acid positions), numbers on branches correspond to bootstrap values. *Spiroplasma* groups are highlighted with colors. The taxon label of the novel genome is highlighted in bold. Accession numbers for all taxa are listed in Table S4.

**a)**



**b)**



**c)**



**d)**