

Metagenomic DNA discoveries after sequencing everything (bacteria,
parasites, food, gut) inside a bee

Cong Liu

2021.8

A thesis submitted in partial fulfilment of the requirements for the degree of Master of Science/Research
at Imperial College London

Submitted for the MRes/MSc in Computational Methods in Ecology and Evolution

Declaration:

1 Introduction

The health of bees is a crucial issue that has received increasing concern (Amiri et al., 2020). As pollinators of numerous plants, bees play an important role in the stability of natural and agricultural systems (Hristov et al., 2020a, Földesi et al., 2021, Bänisch et al., 2021, Khalifa et al., 2021). However, diverse pathogens and environmental factors impose threaten on bee health. The trend of bee decline has been widely observed at a global scale and attributed to multiple factors including pesticide exposure, invasive species, habitat loss, disease prevalence and climate change (Brown and Paxton, 2009, Hristov et al., 2020b, Cheng and Ashton, 2021, Zattara and Aizen, 2021).

Bees are associated with a community of microorganisms influencing host health (Engel et al., 2016, Raymann and Moran, 2018). Bee-associated microbiome plays a role in food digestion and carbohydrate metabolism by mediating processes including breaking down plant-produced pectin and hemicellulose (Zheng et al., 2019); hydrolysis of sucrose (Engel et al., 2012, Lee et al., 2015), a predominant component of nectar (Nicolson and Thornburg, 2007) and metabolism of mannose (Engel et al., 2012, Lee et al., 2015) which presents in nectar (Adler, 2000) and is poisonous to honey bees (de la Fuente et al., 1986). Besides, it is also involved in immunity against pathogens, providing protection against pathogens including *Crithidia*, (Koch and Schmid-Hempel, 2011b, Cariveau et al., 2014), *Paenibacillus larvae*, (Ebeling et al., 2016, Forsgren et al., 2010) and *Nosema* (Cariveau et al., 2014, Maes et al., 2016).

Host-associated microbiome plays an important role in host homeostasis and metagenomics boosted by high throughput sequencing provides powerful tools for understanding its composition, function and dynamics. For example, amplicon sequencing is a powerful tool and has revealed incredible understanding in microbiome (Eckburg et al., 2005, Galloway-Peña and Hanson, 2020, New and Brito, 2020). However, it only focuses on a small part of the whole microbiome. Shotgun sequencing of microbiome, or metagenomics, unselectively captures DNA in a sample and is capable of providing comprehensive inventories of taxa and functional gene clusters (FGCs) with high resolution (Quince et al., 2017, New and Brito, 2020). However, utilization of shotgun metagenomics is hindered by challenges in bioinformatics. Typical goal of metagenomics is to provide taxonomic and functional profile of the community, and there is not a golden standard for bioinformatics of metagenome. Generaly, one of the first steps in metagenomic analysis is assembling short reads into long contigs, which can help improve accuracy of metagenomic annotation (Tran and Phan, 2020) and is necessary for discovery of novel taxa and genes (Culligan et al., 2014, Youngblut et al., 2020). However, metagenome assembly is complex, compromised by fragmental assembly, chimeras

(Mikheenko et al., 2016) and loss of taxon/function diversity due to unassembled reads (Vollmers et al., 2017, Ayling et al., 2020). Probably because of these shortcomings, assembly is skipped in some researches and short reads are directly proceeded for annotation (Tringe et al., 2005, Huson et al., 2007, Abubucker et al., 2012, Vermote et al., 2018, Bovo et al., 2018), although the accuracy may be compromised due to low information load of short reads (Wommack et al., 2008, Carr and Borenstein, 2014, Tran and Phan, 2020). A combination of both assembly-dependent and -free method might be helpful for overcoming the complexity of metagenomic assembly and inaccuracy of short read annotation (Becker et al., 2020).

Another challenge of metagenomics is the determination of sequencing depth. It is recommended to retrieve as many reads as possible (Quince et al., 2017), since insufficient sequencing causes imprecise estimation of taxon/FGC diversity (Cattonaro et al., 2018, Zaheer et al., 2018, Pereira-Marques et al., 2019, Gweon et al., 2019). However, deep sequencing is expensive, which hinders its utilization, especially in large-scale projects. Currently, there is hardly any published guideline for the sufficient sequencing depth of a given environment or study type in order to reach a trade-off between sequencing effort and reliable output.

The determination of sequencing depth can be conducted by rarefaction analysis, a method originated from traditional ecology (Sanders, 1968, Hurlbert, 1971, Heck Jr et al., 1975, Moreno and Halffter, 2000, Hortal and Lobo, 2005, Gómez-Anaya et al., 2014, Hughes et al., 2021). In field based surveys, a random sample of individuals is drawn from a community and assigned to species. Then a rarefaction curve illustrating expected biodiversity represented by given sampling effort (often measured by number of captured individuals) is generated by random subsampling the original sample without replacement, and quantified by model fitting (Hughes and Hellmann, 2005, Gotelli and Colwell, 2011). The slope of the rarefaction curve represents the expected rise of the curve if one more individual is captured. The sample is nearly complete if and only if the final slope is small, and the point at which the slope of rarefaction curve falls to a cut-off value represents the minimal sampling effort required for assessment of biodiversity (Heck Jr et al., 1975, Moreno and Halffter, 2000, Hortal et al., 2004, Hortal and Lobo, 2005, Chao and Jost, 2012, Gómez-Anaya et al., 2014, Roswell et al., 2021). This framework for determination of sampling effort can be used in shotgun metagenomics, since a metagenomic dataset can be viewed as a random sample of an assemblage of genomic sequences, and profiling is the process by which reads are assigned to taxa or FGCs. The concept corresponding to sampling effort is sequencing depth in metagenomics, which is measured by read number.

In macroecology, comprehensive analytical frameworks are available for computing expected biodiversity returned by given sample sizes using a reference sample (Heck Jr et al., 1975, Chao and Jost, 2012, Chao et al., 2014b). However, they are not suitable for metagenomics because of the assumptions that every

captured individual can be assuredly assigned to a species, and the present of rare species is as reliable as that of abundant species. In shotgun metagenomics, reads may be unannotated or annotated incorrectly due to low sequencing quality, host contamination, similar structures in genomes of different species, non-coding regions in genomes and limited sensitivity of profiling pipeline. As a result, taxa or FGCs of low relative abundance are more likely to be false positive than that of high abundance.

In this project, I aimed to (1) construct an integrated pipeline combining assembly-dependent and -independent methods for metagenomic profiling and (2) estimate minimal sequencing depth sufficient for covering species and FGC diversity of metagenome from three environmental types: honey bees, bumble bees and surface of flowers. In the integrated pipeline (Figure 1), assembly-dependent taxon profiling is conducted after quality filtering and removing host contamination. To address high negative rate of assembly-dependent search, a reference database of genomes from species present in the assembly is constructed and used for filtering non-host reads. Reads not recruited by the reference database are subjected to assembly-free taxonomic search. This integrated pipeline was used to profile metagenomic datasets from honey bees, bumble bees and environmental DNA (eDNA) washed from a flower (*Erigeron annuus*). It was shown that the construction of reference database and assembly-free search helped identified species not represented by the assembly. Then I simulated different sequencing depth by rarefaction and profiled subsampled datasets using the integrated pipeline, generating inventories of species and Kyoto Encyclopedia of Genes and Genomes (KEGG) orthologies (KOs) (Kanehisa and Goto, 2000). Each KO represents one FGC, *i.e.* an aggregate of genes with same function. Diversity of species and KOs was measured by Hill numbers of order 0, 1 and 2. Sequencing depth was measured by number of clean read pairs (150 bp read length). The relationship between Hill numbers and sequencing depth (number of 150 bp pair-end reads) was quantified by fitting and averaging asymptotic species accumulation models, and estimation of minimal sequencing depth was given by the point where the slope of rarefaction curve drops to cut-off values. Sample would be excluded from estimation of minimal sequencing depth if the original dataset was not sufficient for providing an almost complete inventory of species or KOs, which was characterized by small final slope of rarefaction curve for species or KO richness (Hill number of order 0). As a result, when measuring species diversity by Hill number of order 0, 1 and 2, the average minimal sequencing depth for honey bee samples are 40.33 million, 18.57 million and 17.45 million; for bumble bee samples, the averages are 42.49 million, 40.33 million and 24.77 million; and for flower surface, 0.88 million sequencing depth is not sufficient for providing an almost complete inventory of species, indicating higher sequencing depth is needed. As for functional diversity, no datasets involved in this project is big enough for providing an almost complete inventory of KOs, indicating the sequencing depth need to be higher than the maximum of each sample type, *i.e.* 45.06 million for honey bees, 53.61 million for bumble bees and 0.88 million for

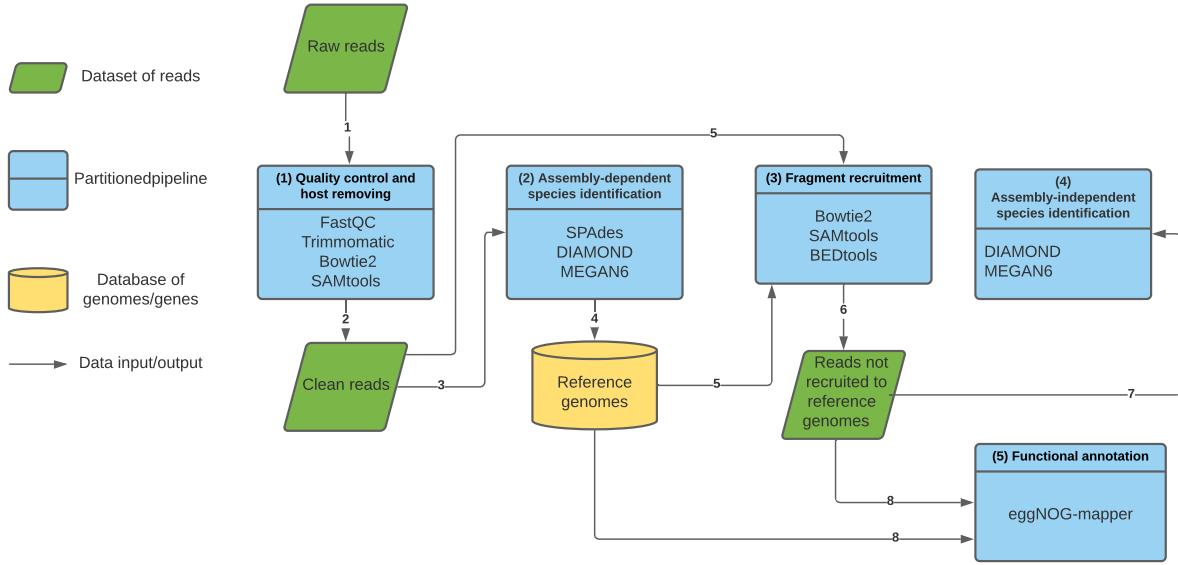


Figure 1: An overview of the integrated pipeline. The pipeline is separated into 5 modules, indicated by blue boxes. Green boxes indicate dataset of reads, while yellow boxes indicate database of genomes. Black arrows indicate the input and output of each step and the numbers on them indicate the order that each step is utilized.

2 Materials and Methods

2.1 Samples, DNA extraction and sequencing

2.2 Metagenomic profiling using integrated pipeline

2.2.1 Integrated pipeline

The integrated pipeline was designed for shotgun metagenomic profiling (assigning reads to taxa or FGCs). It is separated into five modules (Figure 1).

In the module of quality control and host removing, sequencing data quality is checked using FastQC v.0.11.5 (Andrews et al., 2010) and quality control is conducted by Trimmomatic v.0.39 (Bolger et al., 2014). Then clean reads are mapped to host genome using Bowtie2 v.2.4.2 (Langmead and Salzberg, 2012) and non-host reads are extracted by SAMtools v.1.11 (Li et al., 2009).

The non-host reads are subject to the module of assembly-dependent species identification. *De novo* assembly is conducted by SPAdes v.3.15.2 (Prjibelski et al., 2020). Assembled contigs are aligned to NCBI non-redundant (nr) database by DIAMOND v.2.0.7.145 (Buchfink et al., 2015), and assigned to

113 taxa by MEGAN6 (Huson et al., 2007).
 114 Then fragment recruitment is conducted. A reference database comprising reference genome dataset, *i.e.*
 115 genomic sequences in FASTA format and corresponding genome annotation in general feature format (gff),
 116 was constructed. For each species represented by assembly, its reference genome dataset, if available, is
 117 downloaded from NCBI using its *datasets* command-line tool and added to the reference database. Then
 118 non-host reads are mapped to the reference database by Bowtie2, and unmapped reads are extracted by
 119 SAMtools.
 120 Reads not recruited by the reference database are subjected to assembly-independent species identification.
 121 They are aligned to NCBI nr database through DIAMOND and assigned to taxa by MEGAN6.
 122 Finally, functional annotation is conducted by EggNOG-mapper v.2.1.2 (Huerta-Cepas et al., 2017). It
 123 takes coding sequences (CDSs) of genomes in the reference database and reads subject to assembly-
 124 independent species identification as input and assigns them to KOs.
 125 The integrated pipeline was used for analyzing metagenomic datasets involved in this study and details
 126 in parameter settings of each module are described in Supplementary 6.1.

127 **2.2.2 Taxon/function quantification and metabolic pathway reconstruction**

128 After profiling, identified species and KOs were quantified by calculating relative sequence abundance, *i.e.*
 129 proportion of reads assigned to a species/KO in all reads annotated. For species without available reference
 130 genomes, their abundances were calculated using reads assigned to them in assembly-independent search.
 131 As for taxa with available reference genomes, they may be identified in both assembly-dependent and
 132 -independent search due to strain-specific genomic structures that are not present in reference genomes.
 133 Their abundances were calculated by summing number of reads that (1) mapped to coding sequences
 134 (CDSs) of reference genomes and (2) assigned to them in assembly-independent search. Reads mapped to
 135 non-coding regions were not taken into consideration in order to avoid overestimation since the assembly-
 136 independent search was based on aligning reads to nr database, which is composed of proteins. As for
 137 KO quantification, CDSs with zero-coverage were excluded. Abundances of KOs were calculated by
 138 summing number of reads that (1) mapped to CDSs assigned to KOs and (2) assigned to KOs directly.
 139 Extraction of CDSs and calculation of their coverage were conducted by BEDtools v.2.30.0 (Quinlan and
 140 Hall, 2010).
 141 Metabolic pathways were inferred based on KOs. Reads assigned to plants and arthropods were not
 142 included since they were unlikely to represent living organisms. MinPath v.1.6 was used for pathway
 143 inference (Ye and Doak, 2009). It finds a minimal set of KEGG pathways that can explain all KOs
 144 provided as input.

145 2.3 Estimation of minimal sequencing depth required for metagenomic profiling

146 The impact of sequencing depth on taxon/function diversity of metagenome was simulated by rarefaction.
147 Diversity was measured by Hill numbers, which are calculated from inventory of relative abundances. For
148 taxon diversity, the inventory was generated by assigning reads to species; while for function diversity,
149 it was obtained by assigning reads to KOs. Then the rarefaction curve plotting Hill number against
150 sequencing depth was computed using a multimodel inference method, and estimation of minimal se-
151 quencing depth required for covering diversity were given by the point where the slope of rarefaction
152 curve drops to a cut-off value.

153 2.3.1 Simulating different sequencing depth by subsampling and measuring diversity by 154 Hill numbers

155 Species/KO inventories obtained from different sequencing depth was simulated by rarefaction, *i.e.* ran-
156 domly taking a proportion of the original dataset without replacement. Since the ratio between numbers
157 of raw and clean reads is dependent on sequencing process and is not influenced by sample type, sequenc-
158 ing depth here refers to number of clean read pairs to exclude variance caused by different proportion of
159 low quality reads in samples of same type. Besides, the expected ratio between host and non-host reads
160 in a metagenomic dataset is dependent on the DNA sample and not impacted significantly by sequencing
161 depth. Thus the proportion of non-host reads in each simulation is expected to be the same with that in
162 the original dataset.

163 Based on these considerations, I randomly subsampled non-host dataset of each sample, taking 10%-100%
164 of read pairs at interval of 10% by *reformat.sh* script of BBmap v.38.90 (Bushnell, 2014), and profiled
165 subsampled datasets by the integrated pipeline (Figure 1). The sequencing depth of each subsampled
166 dataset equals number of subsampled non-host read pairs divided by ratio between non-host and clean
167 read pairs. Thus, each subsampled dataset of non-host reads is corresponded to an imaginary dataset
168 of clean reads, whose proportion of non-host reads is the same with that of the original metagenomic
169 dataset.

170 After profiling subsampled datasets, species/KO diversity was measured by Hill numbers of order q ,
171 defined as Equation 1 (Hill, 1973).

$$D^{(q)} = \left(\sum_i (p_i)^q \right)^{\frac{1}{1-q}} \quad (1)$$

172 p_i represents the relative abundance of i th species/KO, and q determines sensitivity to relative abundances.
173 When $q = 0$, abundances are not taken into consideration and $D^{(0)}$ equals species/KO richness. When
174 $q = 1$, Hill number is defined as the limit of Equation 1 as q tends to 1 (2) and emphasis is given to

175 species/KOs with general abundances.

$$D^{(1)} = e^{-\sum_i p_i \log p_i} \quad (2)$$

176 When $q = 2$, high leverage is provided to abundant species/KOs and Hill number equals the inverse of
177 Simpson index (Equation 3).

$$D^{(2)} = \frac{1}{\sum_i (p_i)^2} \quad (3)$$

178 **2.3.2 Quantification of relationship between sequencing depth and Hill numbers by fitting** 179 **rarefaction curves**

180 Hill number of order q (Equation 1) measures diversity of an inventory as the number of equally abundant
181 categories in an imaginary inventory with the same diversity (Chao et al., 2014a, Roswell et al., 2021).
182 Order q determines leverage given to abundant categories. All Hill numbers behave in the following way:
183 if a proportion of categories in an inventory was removed randomly, all Hill numbers would decrease by
184 that proportion (Roswell et al., 2021). Thus, it can be hypothesized that as sequencing depth (number
185 of clean read pairs) increases, the detection of novel species/KOs leads to increase of Hill numbers, and
186 when sequencing depth is so big that all species/KOs present in the metagenomic DNA sample have been
187 detected, Hill numbers level off. Such a relationship can be fitted by asymptotic species accumulation
188 models.

189 Let Hill number of order q (Equation 1) be a function of sequencing depth x , which takes million read pairs
190 as the unit. This function was fitted using a multimodel inference method. First, a total of five candidate
191 models (Table S1) were fitted to rarefaction curves which plots Hill numbers against sequencing depth. R
192 package *minpack.lm* v.1.2.1, which employs Levenberg-Marquardt nonlinear least-square algorithm, was
193 used for model fitting. Then small sample unbiased Akaike information criterion (AICc) (Anderson, 2007)
194 of each candidate model was calculated (Equation 4):

$$AICc = -2L + 2k + \frac{2k(k+1)}{(n-k-1)} \quad (4)$$

195 where n is number of observed data points ($n = 10$ in this study), k is the number of fitted coefficients,
196 and L is maximized log-likelihood, given by Equation 5.

$$L = -0.5n \log\left(\frac{R_{ss}}{n}\right) \quad (5)$$

197 Rss represents residual sum of squares.

198 Then model averaging was conducted. First, differences of AICc scores between i th candidate models and
 199 the model with lowest AICc value were calculated using Equation 6.

$$\Delta_i = AICc_i - AICc_{min} \quad (6)$$

200 $AICc_i$ is the AICc score of i th plausible model and $AICc_{min}$ is the lowest AICc score among all candidate
 201 models. The Akaike weight of i th model is given by Equation 7 (Anderson, 2007).

$$w_i = \frac{e^{(-0.5\Delta_i)}}{\sum_i e^{(-0.5\Delta_i)}} \quad (7)$$

202 Denote i th candidate model by $D_i^{(q)} = D_i^{(q)}(x)$, the averaged model is given by Equation 8.

$$D^{(q)}(x) = \sum_i w_i D_i^{(q)}(x) \quad (8)$$

203 The slope of rarefaction curve was calculated by first derivative of averaged model (Equation 9). It reflects
 204 the increase rate of the curve.

$$\frac{dD^q}{dx} = \sum_i w_i \frac{dD_i^{(q)}}{dx} \quad (9)$$

205 The asymptote of rarefaction curve as sequencing depth tends to infinity is given by Equation 10. It pro-
 206 vides an estimation of the total diversity and is comparable among metagenomic DNA samples with
 207 different sequencing depth (Lamas et al., 1991, SoberónM and LlorenteB, 1993, Hortal et al., 2004,
 208 Jiménez-Valverde and Lobo, 2005, Hortal et al., 2006). However, the accuracy of asymptotic estima-
 209 tors is controversial in macroecology (SoberónM and LlorenteB, 1993, Colwell and Coddington, 1994,
 210 Chazdon et al., 1998, Jimenez-Valverde et al., 2006, Hortal et al., 2006), and their performance in shot-
 211 gun metagenomics is seldom evaluated.

$$\lim_{x \rightarrow +\infty} D^{(q)}(x) = \sum_i w_i \lim_{x \rightarrow +\infty} D_i^{(q)}(x) \quad (10)$$

212 **2.3.3 Estimating minimal sequencing depth using rarefaction curves**

213 Minimal sequencing depth is defined as the point at which diversity starts to level off as sequencing depth
 214 increases, and its precise estimation via rarefaction is based on the assumption that the original dataset
 215 is sufficient for detection of almost all species/KOs present. This assumption can be verified by looking at
 216 the rarefaction curve that plots species/KO richness (Hill number of order 0) against sequencing depth.
 217 The original dataset is sufficient for providing a reliable inventory if and only if the rarefaction curve of

richness is characterized by a small final slope (Heck Jr et al., 1975, Moreno and Halffter, 2000, Hortal et al., 2004, Hortal and Lobo, 2005, Gómez-Anaya et al., 2014). Then an estimation of minimal sequencing depth is provided by the point at which the slope of rarefaction curve decreases to a given cut-off value (Hortal and Lobo, 2005).

3 Results

3.1 Sequence reads and pipeline evaluation

Eight samples (four honey bees, three bumble bees and one flower eDNA) were sequenced. The quality reports of raw reads showed low-quality 3'-end (Figure S1a), uneven base content in 5'-end (Figure S1b) and the present of adaptors (Figure S1c). The quality control procedure covered these aspects and improved data quality (Figure S1d, S1e and S1f).

After quality control, read pairs aligned to host genome were removed. Table S2 reports numbers of raw, clean and non-host read pairs. Honey bee sample *Bee_Amellifera_13_1* was dropped for further analysis since its raw read pair number (1.10 million) is significantly lower than three other samples (about 59 million). After quality control, 62.08%-76.52% of raw read pairs were retained for these three honey bee samples. Then a different proportion of non-host read pairs (29.52%-86.19%) were retained. As for bumble bee samples, about 58 million raw read pairs were obtained for each sample and 82.1%-83.8% were retained after quality control. After host removing, 7.35%-10.08% of clean reads were retained. For flower eDNA sample, 1.44 million raw read pairs were obtained and 61.15% of them were retained.

In integrated pipeline (Figure 1), clean non-host reads are first assembled to into contigs and assigned to taxa. To address false negative results in assembly-dependent taxon search (Sharon et al., 2015, Vollmers et al., 2017), a reference database composed of genomes of species represented by contigs is constructed. Reads not aligned to the reference database are subjected to assembly-independent taxon search. The reference database and assembly-independent search helped improve species identification in all three sample types, especially in simulations of low sequencing depth (Figure 2).

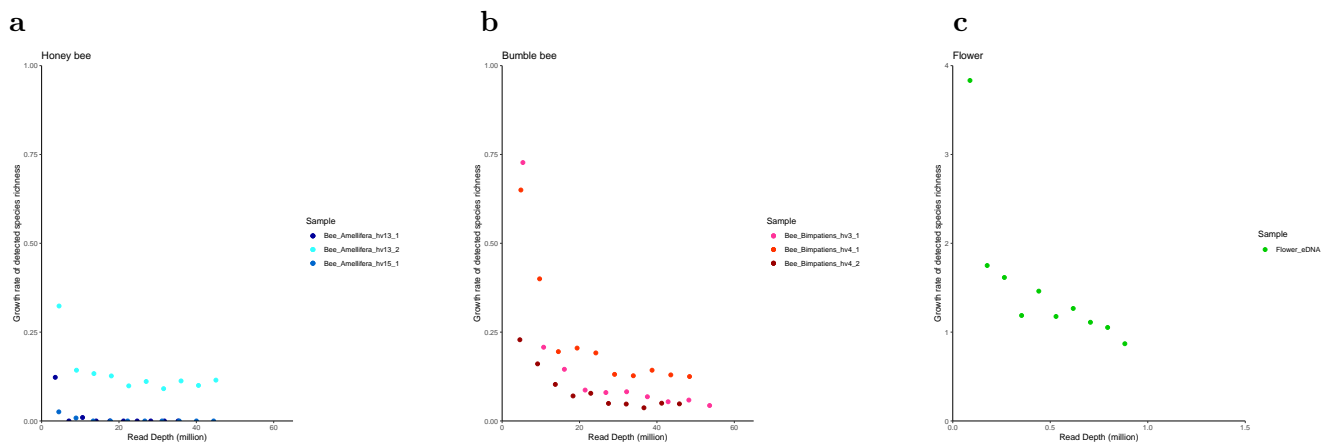


Figure 2: Integrated pipeline improves the detection of species richness. The horizontal axis represents sequencing depth, and the vertical axis represents the ratio between number of species only detected by assembly-independent search and number of species detected by assembly-dependent search. Sample type is shown in the top left of each subfigure.

3.2 Characterization of species composition

All seven samples were dominated by bacterial and/or arthropod species (Figure 3b). Honey bee samples *Bee_Amellifera_hv13_2* and *Bee_Amellifera_wild_1* were dominated by bacteria (approximate 90% relative abundance), while in *Bee_Amellifera_hv15_1*, arthropods, bacteria and viruses accounted for most annotated reads. As for bumble bee samples, they were all dominated by bacteria and arthropods. For the flower eDNA sample, arthropods were the most dominant and bacteria were the second.

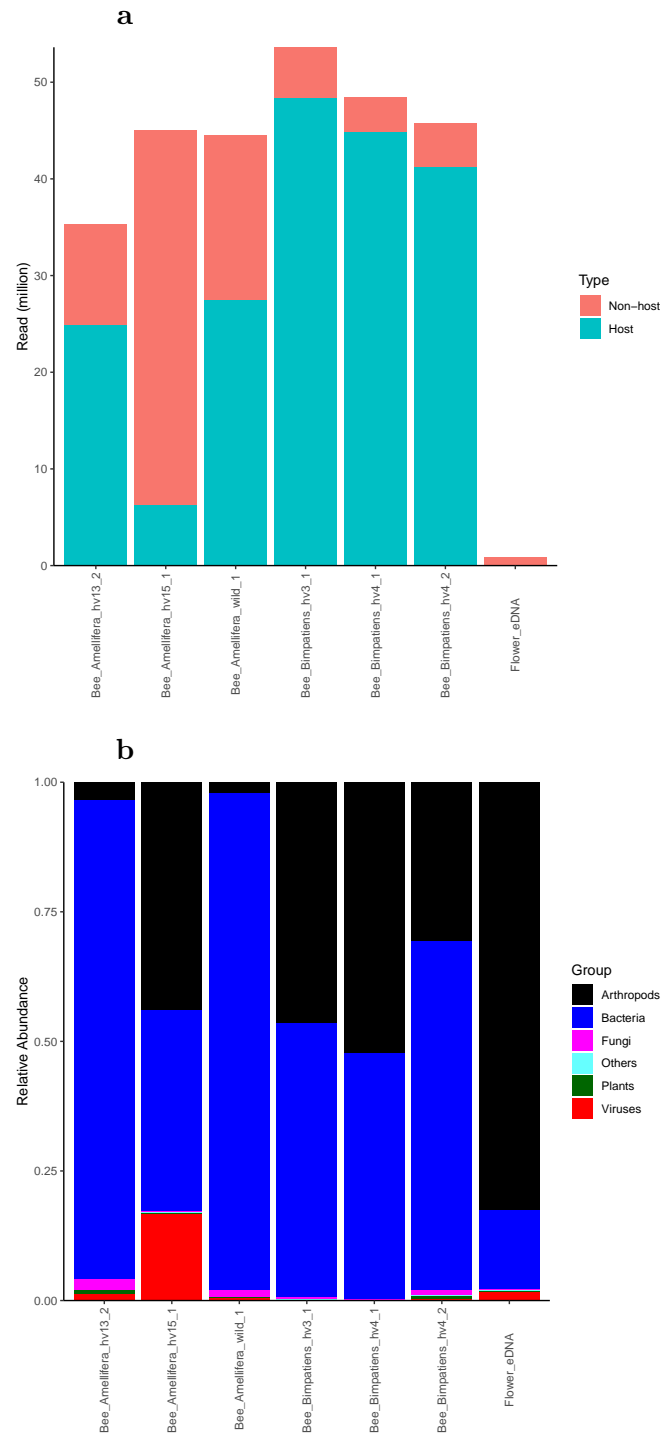


Figure 3: a. The number of host and non-host reads in each sample. b. The relative abundance of species under six taxonomic groups: superkingdom Viruses, superkingdom Bacteria, kingdom Viridiplantae (plants), kingdom Fungi, phylum Arthropoda and others (species that are not in the other five groups).

Figure 4 illustrates composition of bacteria, most of which are common bee-associated species within *Bifidobacterium*, *Frischella*, *Gilliamella*, *Snodgrassella* and *Lactobacillus* (Koch and Schmid-Hempel, 2011a, Moran, 2015, Kwong et al., 2017). *Apilactobacillus* and *Bombilactobacillus* are genera separated

251 from *Lactobacillus* recently (Zheng et al., 2020). *Fructobacillus* is often found in fructose-rich environ-
 252 ments like flowers (Endo and Dicks, 2014). *Bartonella apis* is related to animal pathogens (Kešnerová
 253 et al., 2016) and is widespread in honey bee workers (Raymann and Moran, 2018).
 254 Figure S2 shows composition of arthropods, and most of them are pollinators within *Apis* and *Bombus*.
 255 However, some of them might be considered as false positive. For example, *Apis cerana*, *Apis dorsata* and
 256 *Apis florea* are mainly found in Asia and unlikely to present in the area where samples were collected.
 257 These might derive by similarity between genomes of *Apis mellifera* and other *Apis* species.
 258 Figure S3 illustrates plant species, indicating foraging area of bees. Most samples were dominated by
 259 one or two plant species, except *Bee_Amellifera_hv13_2*. Some crops were identified, including *Brassica*
 260 *napus* (rape), *Brassica oleracea*, *Cicer arietinum* (chickpea), *Glycine max* (soybean), *Helianthus annuus*
 261 (sunflower), *Nicotiana glauca* (flowering tobacco) and *Raphanus sativus* (radish).
 262 Figure S4 shows compositions of fungal species. *Nosema ceranae*, a widespread bee pathogen, was the
 263 dominant fungal species in most samples including the flower eDNA. Three yeast species (*Clavispora*
 264 *lusitaniae*, *Saprochaete ingens* and *Wickerhamiella sorbophila*) were found in bees.
 265 Figure S5 illustrates viruses identified. Several phages were identified, including *Bifidobacterium phage*
 266 *BitterVaud1* that infects bee-commensal bacterium *Bifidobacterium asteroides* (Bonilla-Rosso et al., 2020),
 267 *Bacteriophage sp.* that infects *Pseudomonas aeruginosa* (Essou et al., 2015), an opportunistic pathogen
 268 that might contaminate bees (Bailey, 1968, Papadopoulou-Karabela et al., 1992, 1993), and species in
 269 Myoviridae and Siphoviridae. Most of the others are arthropod-associated viruses. *Apis mellifera fil-*
 270 *amentous virus* and *Bombus cryptarum densovirus* are bee-infecting viruses (Clark, 1978, Bailey et al.,
 271 1981, Schoonvaere et al., 2018). Several arthropod-infecting parvoviruses were found, including *Blattodean*
 272 *pefuambidensovirus 1*, *Hemipteran scindoambidensovirus 1*, *Hymenopteran scindoambidensovirus 1* and
 273 *Orthopteran scindoambidensovirus 1* (Pénzes et al., 2020).

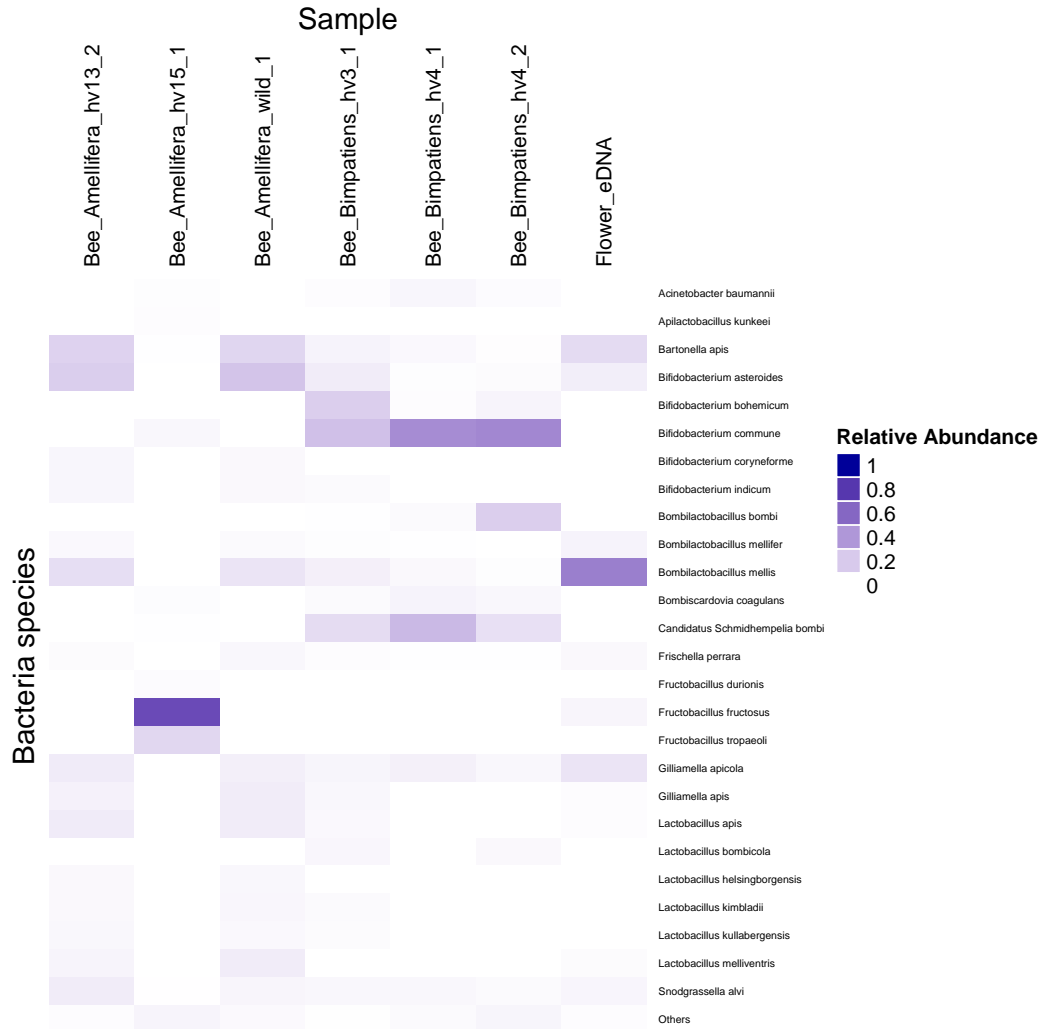


Figure 4: Heatmaps for bacterial species abundance distribution in all samples. The relative abundance takes reads assigned to bacterial species as background. Species with relative abundance smaller than 1% in all samples are collapsed as "others".

3.3 Characterization of functional profiling

In order to illustrate metabolic potentially of metagenomic samples, KEGG pathways were inferred using a parsimony approach. Reads assigned to plant and arthropod species were not involved in pathway inference. The coverage of a pathway was calculated by the ratio between number of annotated KOs and total number of KOs involved in that pathway.

Here concern is given to metabolism pathways of carbonhydrates and amino acids, which are crucial for bee health. Inferred pathways indicate potential capability of metabolism of sugars including fructose, sucrose, mannose and galactose (Figure S6), and all essential amino acids for honey bees (Figure 5)(Groot, 1953).

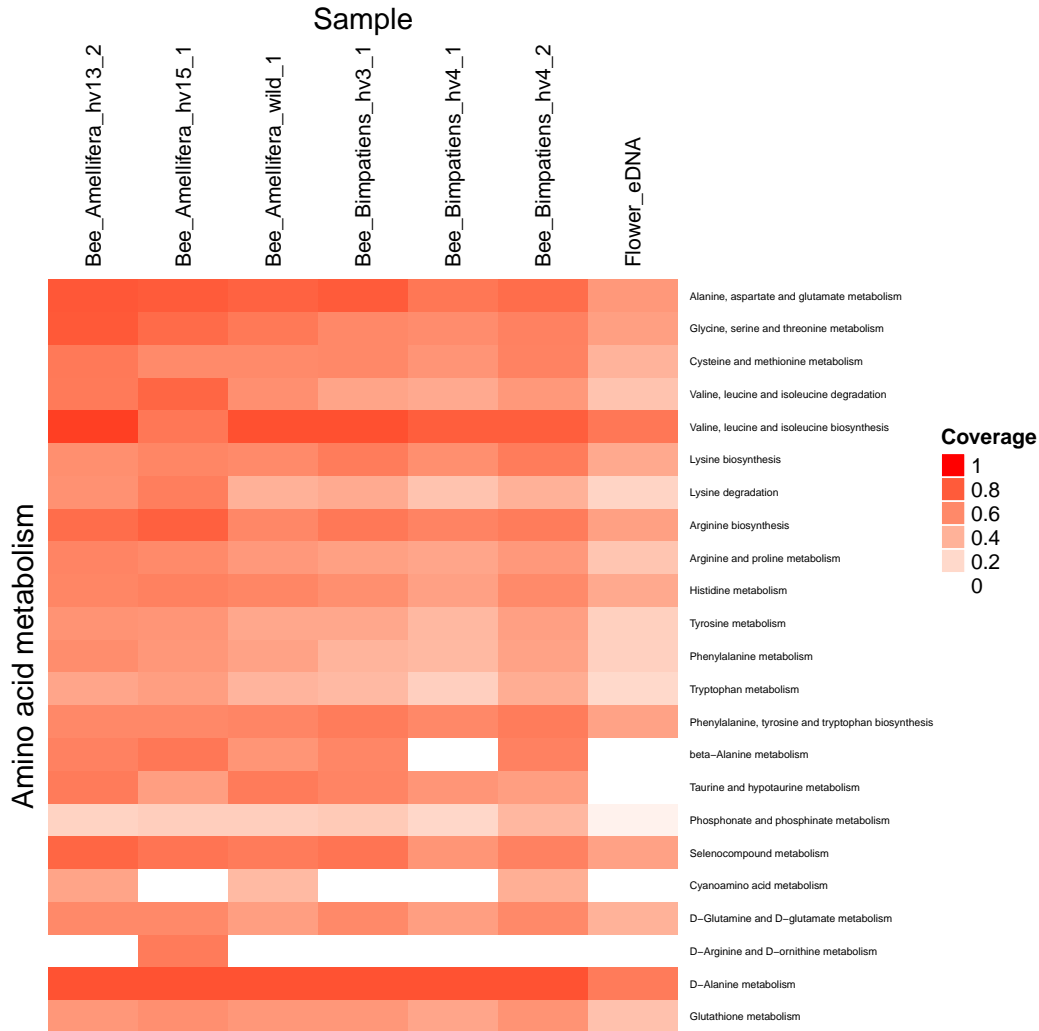


Figure 5: Heatmaps for pathways of amino acid metabolism.

3.4 Minimal sequencing depth required for detection of species and function diversity

In order to estimate minimal sequencing depth for combinations of sample type (honey bee, bumble bee and flower eDNA) and study type (species profiling and function profiling), shallow sequencing was simulated by rarefaction. The relationship between sequencing depth (clean read pairs) and species/KO diversity (Hill numbers of order 0, 1, 2) was quantified by fitting and averaging asymptotic species accumulation models. The asymptote of the model provides an estimation of total diversity and its slope reflects the increase rate of diversity. Minimal sequencing depth was estimated by the point at which the slope of the rarefaction curve drops to a cut-off value.

Rarefaction assumes that the original dataset provides an almost complete inventory, which can be verified by final slope of rarefaction curve for Hill number of order 0 (richness). Figure 6 shows rarefaction curves for species/KO richness and Table 1 summarizes their final slopes. For species diversity rarefaction,

all bumble bee samples are sufficient, with final slopes < 0.1 and completeness (ratio between final richness and asymptote) > 0.98 . As for honey bees, *Bee_Amellifera_hv15_1* and *Bee_Amellifera_wild_1* are sufficient, while *Bee_Amellifera_hv13_2* is insufficient, with final slope > 1 and completeness < 0.8 . For the flower eDNA sample, the final slope of species richness rarefaction curve is 10.8380 and its completeness is 1.32%, indicating more sequencing effort is needed for species profiling. As for function diversity rarefaction, the final slopes of all KO richness rarefaction curves are higher than 15, indicating no dataset can providing an almost complete inventory of KOs. Thus, estimation of minimal sequencing depth was conducted for the combinations of two sample types (honey bee and bumble bee) and one study type (species profiling), based on five datasets (two honey bees and three bumble bees).

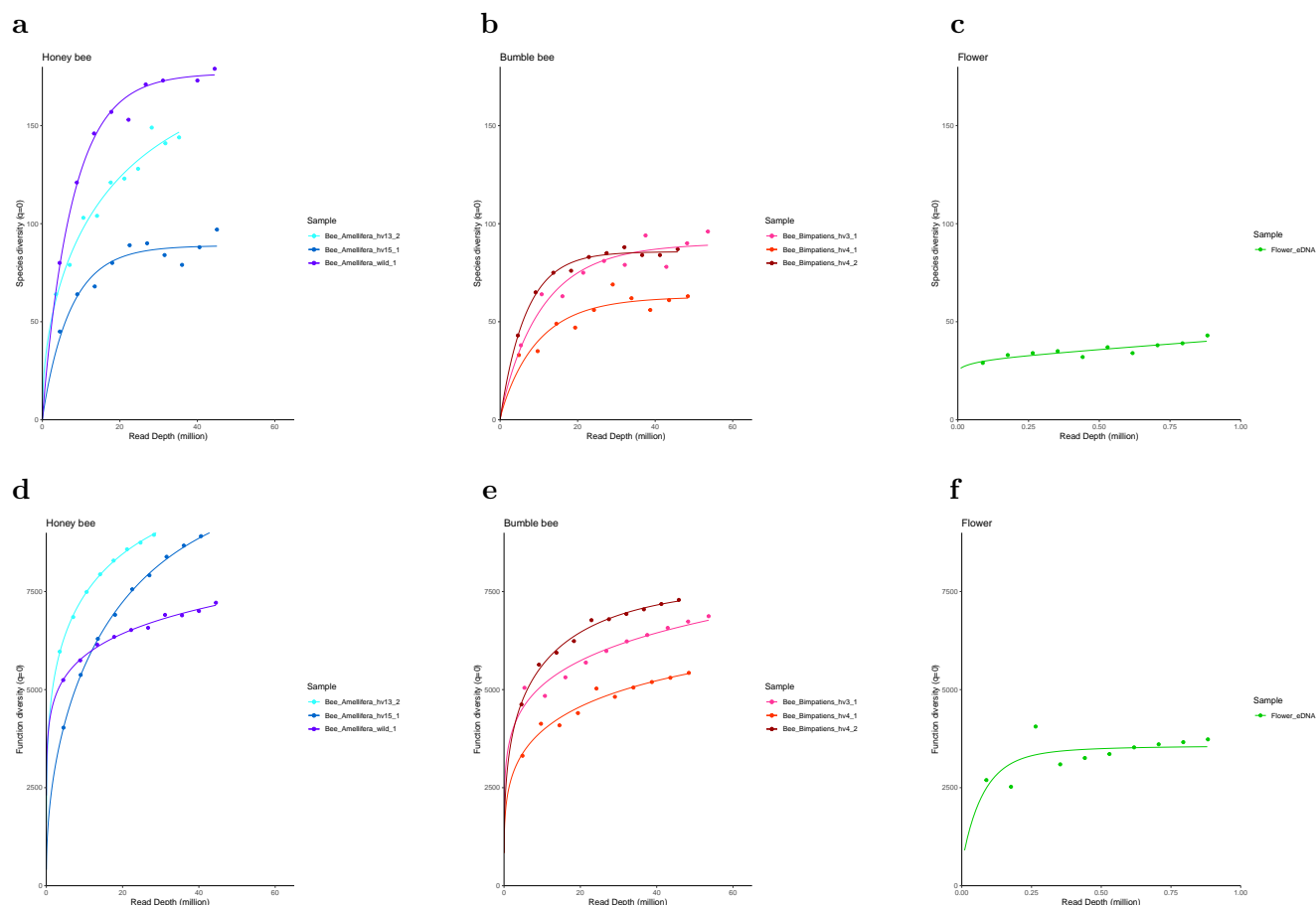


Figure 6: Rarefaction curves for species (a, b, c) or KO (d, e, f) richness (Hill number of order 0). The horizontal axis represents sequencing depth, and the vertical axis represents richness. Sample type is shown in the top left of each subfigure. Note that the scale of horizontal axis in subfigure c and f is much smaller than that in other subfigures.

Table 1: Summary of final point of rarefaction curve for species/KO richness (Hill number of order 0). Type: indicates whether this row reports rarefaction curve for species or KO richness. Depth: sequencing depth taking million read pairs as unit. OR: observed species/KO richness. ER: expected species/KO richness predicted by modeling rarefaction curve. FS: final slope of rarefaction curve. Asym: asymptote calculated by modeling rarefaction curve. Comp: completeness, represented by ratio between ER and Asym.

Sample	Type	Depth	OR	ER	FS	Asym	Comp
<i>Bee_Amellifera_hv13_2</i>	species	35.28	144	146.69	1.0958	184.94	0.7931
<i>Bee_Amellifera_hv15_1</i>	species	45.05	97	88.53	0.0344	88.84	0.9966
<i>Bee_Amellifera_wild_1</i>	species	44.46	179	175.89	0.0991	176.84	0.9946
<i>Bee_Bimpatiens_hv3_1</i>	species	53.61	96	89.00	0.0897	90.79	0.9803
<i>Bee_Bimpatiens_hv4_1</i>	species	48.42	63	62.06	0.0626	62.92	0.9862
<i>Bee_Bimpatiens_hv4_2</i>	species	45.80	87	85.59	0.0146	85.69	0.9988
<i>Flower_eDNA</i>	species	0.88	43	40.04	10.8380	3034.88	0.0132
<i>Bee_Amellifera_hv13_2</i>	KO	35.28	9289	9284.45	38.6844	10641.27	0.8725
<i>Bee_Amellifera_hv15_1</i>	KO	45.05	9046	9102.78	43.8717	10453.98	0.8707
<i>Bee_Amellifera_wild_1</i>	KO	44.46	7216	7151.72	21.4368	14779.89	0.4839
<i>Bee_Bimpatiens_hv3_1</i>	KO	53.61	6872	6769.80	21.1588	28111.79	0.2408
<i>Bee_Bimpatiens_hv4_1</i>	KO	48.42	5427	5410.71	19.6932	7665.58	0.7058
<i>Bee_Bimpatiens_hv4_2</i>	KO	45.8	7288	7259.56	15.6309	7873.64	0.9220
<i>Flower_eDNA</i>	KO	0.88	3732	3544.54	60.4367	3571.11	0.9926

Minimal sequencing depth is estimated by the point at which the slope of rarefaction curve drops to a cut-off value. Table 2 summarizes estimations of minimal sequencing depth for detection of species diversity, using 0.5, 0.1, 0.05 and 0.01 as cut-off values. When order q of Hill number equals 0, *i.e.* species abundances are not considered, cut-off value of 0.1 for slope of rarefaction curve is sufficient for providing completeness $> 97\%$ in all samples. The average minimal sequencing depth are 40.33 million for honey bees and 42.49 million for bumble bees. When species abundances are considered (order q equals 1 or 2), cut-off value of 0.01 can provide completeness $> 95\%$ in most pairs of sample and q value. For honey bees, the average minimal sequencing depth are 18.57 million ($q = 1$) and 17.45 million ($q = 2$). For bumble bee samples, the average minimal sequencing depth are 40.33 million ($q = 1$) and 24.77 million ($q = 2$).

Table 2: Summary of minimal sequencing depth estimated from different cut-off values of slope. Minimal sequencing depth was estimated by the point at which the slope of rarefaction curve drops to a cut-off value (0.5, 0.1, 0.05 or 0.01), which is indicated in names of columns. For example, MinD_0.5 represents estimated minimal sequencing depth taking 0.5 as cut-off value, and ED_0.5 is expected Hill number from MinD_0.5. Comp_0.5 is the ratio between OptDiv_0.5 and the asymptote (represented by Asym). Unit of sequencing depth is million read pairs. q refers to order of Hill number determining sensitivity to species abundance distribution.

Sample	q	Asym	MinD_0.5	ED_0.5	Comp_0.5	MinD_0.1	ED_0.1	Comp_0.1	MinD_0.05	ED_0.05	Comp_0.05	MinD_0.01	ED_0.01	Comp_0.01
<i>Bee_Amellifera_hv15_1</i>	0	88.84	23.54	84.85	0.9551	36.27	87.99	0.9905	41.95	88.40	0.9951	55.75	88.74	0.9989
<i>Bee_Amellifera_wild_1</i>	0	176.84	30.53	172.46	0.9752	44.39	175.88	0.9946	50.57	176.33	0.9971	65.86	176.70	0.9992
<i>Bee_Amellifera_hv15_1</i>	1	12.28	3.18	10.67	0.8699	5.83	11.33	0.9224	7.21	11.43	0.9304	15.21	11.58	0.9429
<i>Bee_Amellifera_wild_1</i>	1	14.77	0.69	13.49	0.9137	3.30	14.03	0.9498	6.20	14.23	0.9635	21.93	14.58	0.9868
<i>Bee_Amellifera_hv15_1</i>	2	8.40	3.10	6.35	0.7561	8.17	7.49	0.8909	11.97	7.75	0.9229	27.99	8.11	0.9655
<i>Bee_Amellifera_wild_1</i>	2	9.88	0.18	9.49	0.9602	0.90	9.63	0.9747	1.73	9.69	0.9805	6.91	9.80	0.9916
<i>Bee_Bimpatiensi_hv3_1</i>	0	90.79	30.80	83.69	0.9218	52.03	88.85	0.9786	62.56	89.61	0.9870	92.35	90.32	0.9948
<i>Bee_Bimpatiensi_hv4_1</i>	0	62.92	24.83	57.22	0.9093	42.76	61.61	0.9791	51.24	62.22	0.9887	73.55	62.76	0.9973
<i>Bee_Bimpatiensi_hv4_2</i>	0	85.69	21.75	82.30	0.9604	32.67	85.01	0.9920	37.39	85.35	0.9960	48.40	85.62	0.9992
<i>Bee_Bimpatiensi_hv3_1</i>	1	21.59	3.63	10.50	0.4862	9.47	11.76	0.5449	15.37	12.18	0.5639	58.09	13.04	0.6040
<i>Bee_Bimpatiensi_hv4_1</i>	1	11.85	3.80	9.45	0.7973	10.15	10.87	0.9174	14.68	11.19	0.9445	31.22	11.57	0.9766
<i>Bee_Bimpatiensi_hv4_2</i>	1	12.66	4.17	10.24	0.8087	10.47	11.67	0.9213	14.85	11.98	0.9458	31.68	12.36	0.9760
<i>Bee_Bimpatiensi_hv3_1</i>	2	8.73	2.20	6.16	0.7051	6.28	7.06	0.8088	9.45	7.28	0.8344	23.89	7.60	0.8708
<i>Bee_Bimpatiensi_hv4_1</i>	2	8.33	3.12	6.28	0.7540	8.52	7.50	0.9002	12.30	7.76	0.9323	26.68	8.09	0.9718
<i>Bee_Bimpatiensi_hv4_2</i>	2	7.12	2.90	5.34	0.7493	7.47	6.38	0.8961	10.58	6.60	0.9270	23.73	6.90	0.9683

4 Discussion

Shotgun metagenomics provides powerful tools for microbiome investigations, yet its utilization is hindered by bioinformatical challenges and high cost of deep sequencing. Here, I constructed an integrated pipeline, which combines assembly-dependent and -independent profiling methods to help detect species richness and profiled real metagenomic datasets from honey bees, bumble bees and flower surface. Then I simulated the relationship between species/KO diversity measured by Hill numbers (order $q = 0, 1, 2$) and sequencing depth measured by clean read pair (2×150 bp) number. It is showed that to detect species diversity, the average minimal sequencing depth for honey bee samples are 40.33 million ($q = 0$), 18.57 million ($q = 1$) and 17.45 million ($q = 2$); for bumble bee samples, the averages are 42.49 million ($q = 0$), 40.33 million ($q = 1$) and 24.77 million ($q = 2$); and for flower surface, more than 0.88 million sequencing depth is needed. As for function diversity, no datasets in this project is sufficient for providing an almost complete inventory of KOs, indicating deeper sequencing is needed, *i.e.* sequencing depth need to be higher than 45.06 million for honey bees, 53.61 million for bumble bees and 0.88 million for flower surface. Additionally, bee symbionts and pathogens were found on the flower eDNA sample, and pathway inference indicates that bee-associated microbiome may have capability of metabolizing amino acids that are essential for host.

Shotgun metagenomics provides unique advantages over amplicon sequencing, which is vastly used in bee microbiome investigations (*e.g.* Geldert et al. (2021), Wang et al. (2021), Powell et al., Kapheim et al. (2021)). In amplicon sequencing, a species-specific barcode region is amplified and sequenced (Abdelfattah et al., 2018). Since it captures a small region of the whole genome, amplicon sequencing dataset is relatively small in size, which makes it low-cost both economically and computationally. Besides, comprehensive and streamlined softwares for amplicon sequencing analysis are available, including QIIME (Bolyen et al., 2019), Mothur (Schloss, 2020) and VSEARCH (Rognes et al., 2016). However, a single amplicon sequencing dataset does not represent all taxa in the community since different barcode region is used for identification of different taxonomic group, *e.g.* 16S ribosomal RNA (rRNA) for bacteria (Hayashi et al., 2002, Eckburg et al., 2005), internal transcribed spacer (ITS) for fungi (Nilsson et al., 2008), cytochrome c oxidase subunit I (COI) for Animalia (Hebert et al., 2003) and plastid genes for plants (Group et al., 2009). Besides, it does not provide information on FGCs. Metabolic capacity of microbiome need to be inferred based on reference genomes (Abhauer et al., 2015, Douglas et al., 2018). Shotgun metagenomics provides an alternative to overcome these drawbacks by capturing and sequencing DNA in a sample unselectively. It is capable of representing all taxonomic groups and providing information on FGC composition.

The advantages of shotgun metagenomics over amplicon sequencing are valuable for investigations in bee-associated microbiome for two reasons. First, bees visit numerous niches during environment exploration and foraging activities, and can get contact with diverse eDNA signatures, which provide insight into bee ecology by reflecting interactions between bees and other organisms including plants, arthropods, bacteria and viruses (Bovo et al., 2018, Ribani et al., 2020, Bovo et al., 2020, Matsuzawa et al., 2020). These materials represent a diverse assemblage which is difficult to be profiled comprehensively via amplicon sequencing. Besides, bee bacterial symbionts are diversified at strain level (Engel et al., 2012, Powell et al., 2016, Ellegaard et al., 2020). Bacterial strains are often highly variable in gene content (Cordero and Polz, 2014, Brockhurst et al., 2019), and thus different in metabolic capacity. However, resolving strain-level diversity by species-specific region is difficult (Rodriguez-R et al., 2018, Ciufo et al., 2018, Olm et al., 2020), and no information on gene content is provided by amplicon sequencing. Based on these considerations, shotgun metagenomics provides valuable tools for bee microbiome investigations. However, utilization of shotgun metagenomics is hindered by unique challenges. First, there is not a golden standard for bioinformatics of metagenomic annotation pipeline to provide abundances of taxa and FGCs. Second, shotgun metagenomics is often associated with high cost since it captures a big proportion of microbial genomes.

The integrated pipeline for taxonomic and functional profiling of metagenomic dataset provides several

361 advantages. First, through combination with assembly-free taxon search, it helps solve high false negative
 362 rate associated with assembly-dependent species identification, which is caused by reads left unassembled.
 363 This was shown by analyzing real metagenomic datasets from pollination system and simulating different
 364 sequencing depth by rarefaction (Figure 2). Second, the modularity of integrated pipeline provides flexi-
 365 bility for incorporation of alternative tools. For example, BWA aligner (Li and Durbin, 2009) serves as an
 366 alternative for Bowtie2, and SPAdes can be replaced by other metagenomic assemblers such as Megahit
 367 (Li et al., 2015) and IDBA-UD (Peng et al., 2012). Third, the output files generated by each step are
 368 recorded and can be inspected easily, which is important for troubleshooting.

369 Determination of sequencing depth is important for shotgun metagenomics since insufficient sequencing
 370 causes underestimation of taxonomic/functional diversity (Cattonaro et al., 2018, Zaheer et al., 2018,
 371 Gweon et al., 2019, Pereira-Marques et al., 2019), while deep sequencing is of high cost. Here, expected
 372 species/KO diversity provided by given sequencing depth was estimated by rarefaction and model fitting,
 373 using datasets from honey bees, bumble bees and flower surface. It was shown that increasing sequencing
 374 depth boosted identification of species/KO (Figure 6, Table 1), highlighting value of deep metagenomic
 375 sequencing. For function profiling, no dataset involved here is big enough to provide an almost complete
 376 inventory of KOs (Table 1) even though most samples were deeply sequenced (Table S1), indicating such
 377 task is demanding about sequencing depth. Therefore, when dealing with function potentiality of micro-
 378 biome associated with bees or flowers, retrieving as many reads as possible would be recommended. As
 379 for species profiling, although deep sequencing is still valuable, optimization of sequencing depth can be
 380 made when the budget is limited (Table 2). For honey bees, approximate 40 million sequencing depth
 381 can be sufficient for representing species richness. When assessing biodiversity by Hill numbers of order
 382 1 (Equation 2) or 2 (Equation 3), 17-19 million sequencing depth can provide robust estimation. As for
 383 bumble bees, about 40-43 million sequencing depth can provide reliable estimation for Hill numbers of
 384 order 0 (richness) or 1, and about 25 million sequencing depth can be sufficient for representing biodiver-
 385 sity measured by Hill number of order 2.

386 Additionally, the results of metagenomic profiling indicate overlap between microbiome of bees and flowers,
 387 and provide evidence on potentiality of bee microbiome in metabolism of carbohydrates and amino acids.
 388 In the flower eDNA sample, bee-associated microorganisms were identified, including both pathogens (*e.g.*
 389 *Nosema ceranae* and *Apis mellifera filamentous virus*) and typical bee symbionts (*e.g.* *Bifidobacterium*
 390 *asteroides*, *Bombilactobacillus mellis* and *Gilliamella apicola*) (Figure 4, S4 and S5). Shared microbiome
 391 between bees and flowers is revealed by growing evidence, although its role in pollination system remains
 392 an open question (Keller et al., 2020, Vannette, 2020). As for function profiling, it is indicated that
 393 bee-associated microbiome is capable of metabolizing carbohydrates such as glucose, fructose, sucrose

and mannose (Figure S6). Carbohydrates are main component of bee diet and microbiome-mediated carbohydrate-processing have been vastly investigated (Engel et al., 2012, Lee et al., 2015, 2018, Taylor et al., 2019). Besides, pathways for metabolism of all ten essential amino acids for honey bees (arginine, histidine, isoleucine, leucine, lysine, methionine, phenylalanine, threonine, tryptophan and valine) (Groot, 1953) were inferred (Figure 5). Essential amino acids are crucial for bee health (Simcock et al., 2014, Paoli et al., 2014, Stabler et al., 2015, Hendriksma et al., 2019) and influence feeding preference due to their potential deficiency in single pollen source (Cook et al., 2003, Hendriksma et al., 2014, Hendriksma and Shafir, 2016). Whether bee-associated microbiome influences host health by providing amino acids need to be further investigated.

5 Data and Code Availability

6 Supplementary

6.1 Parameter settings of integrated pipeline

6.1.1 Quality control and host removing

Quality control was conducted to reduce compromise from low quality reads. Data quality was checked using FastQC v.0.11.5 (Andrews et al., 2010) before filtering. FastQC reports of raw reads showed the following aspects need to be covered in quality control: (1) low base quality in 3'-end (Figure S1a); (2) uneven base content in 5'-end (Figure S1c) and (3) the present of Nextera adaptors (Figure S1e). Therefore, quality control was conducted using Trimmomatic v.0.39 (Bolger et al., 2014), which (1) trimmed adaptors; (2) cutted 15 bases from the 5'-ends of reads; (3) cutted bases off from 3'-ends of reads if Phred-33 quality is below 20; (4) dropped reads shorter than 50 bp; (5) dropped reads if average Phred-33 quality is below 20. Then unpaired reads were removed and quality of clean data was checked using FastQC (Figure S1d, S1e, S1f).

After quality control, host read pairs were removed. First, clean read pairs were mapped to host genome (GCA_003254395.2 for *Apis mellifera* and GCA_000188095 for *Bombus impatiens*, downloaded from NCBI) using Bowtie2 v.2.4.2 (Langmead and Salzberg, 2012) with flags *-end-to-end* and *-sensitive*. With flag *-end-to-end*, Bowtie2 requires the read aligned without any clipping from neither end, and *-sensitive* maintains a trade-off between speed and sensitivity. SAM files generated by Bowtie2 were converted to BAM format using SAMtools v.1.11 (Li et al., 2009). Then non-host read pairs were extracted from BAM files by SAMtools according to the present of SAM flag 12 (neither forward nor reverse read

424 in a pair of reads is mapped).

425 6.1.2 Assembly-dependent species identification

426 In order to identify taxa of metagenome, *de novo* assembly was conducted. Non-host read pairs were
427 assembled using SPAdes v.3.15.2 (Prijbelski et al., 2020). Values of k-mer ranged from 21 to 101 at interval
428 of 10. Flags *-only-assembler* and *-meta* were used. Through flag *-meta*, SPAdes runs metaSPAdes which
429 is developed for metagenomic assembly (Nurk et al., 2017). The *-only-assembler* flag skips read error
430 correction and runs assembly only. Its utilization is justified by the following facts. First, when *-only-*
431 *assembler* is not used, SPAdes conducts error correction before assembly. It is conducted by BayesHammie,
432 which is optimized for single cell sequencing instead of shotgun metagenomics (Nikolenko et al., 2013).
433 Besides, reads used for *de novo* assembly had been filtered to ensure their quality.
434 After assembly, taxon identification was conducted using DIAMOND v.2.0.7.145 (Buchfink et al., 2015)
435 and MEGAN6 (Huson et al., 2007). Assembled contigs with a length above 500 bp were aligned to
436 nr database using DIAMOND v.2.0.7.145 with *-long-reads* flag. This flag triggers frame-shift aware
437 alignment mode, which is optimized for long sequence alignment. Therefore, short contigs (length < 500
438 bp) were not retained. Besides, alignments with an E-value < $1e - 5$ or identity < 50% were removed,
439 and for each contig, only alignments above 10% of the best local bit score were retained. The output of
440 DIAMOND was analyzed by the *blast2rma* tool of MEGAN6 with *-lg* flag, which runs lowest-common-
441 ancestor (LCA)-based algorithm developed for long contigs and assigns each contig to a taxon (Huson
442 et al., 2018). The parameter *-supp* was 0, which means the present of a taxon would be identified as long
443 as at least one contig was assigned. This value was used because a contig is assembled from multiple short
444 reads and represents a strong signal for the present of a taxon.

445 6.1.3 Fragment recruitment

446 To integrate individual genomic data of species identified by assembly-dependent search, a reference
447 database comprising reference genome dataset, *i.e.* genomic sequences in FASTA format and correspond-
448 ing gff file, was constructed. For each species represented by assembly, its reference genome dataset, if
449 available, was downloaded from NCBI using its *datasets* command-line tool and added to the reference
450 database.

451 Then fragment recruitment was conducted. The non-host read pairs were mapped to genomic sequences
452 in the reference database using Bowtie2. Read pairs that were not recruited were extracted using SAM-
453 tools. Settings for Bowtie2 and SAMtools were the same as that described in 6.1.1. Read pairs recruited
454 by the reference database were assigned to corresponding species, while the others were subjected to

assembly-independent search.

6.1.4 Assembly-independent species identification

In order to detect species not represented by assembly (Sharon et al., 2015, Vollmers et al., 2017), assembly-independent search was conducted, taking read pairs not recruited by the reference database as input. These reads were aligned to nr database through DIAMOND without using *-long-reads* flag, which triggers computing alignments for short metagenomic reads. Other settings were the same as described in 6.1.2. Then the output of DIAMOND was analysed by MEGAN6 (*blast2rma* tool), which assigns read pairs to taxa through LCA algorithm. Here the parameter *-supp* was 0.1, which means a taxon is reported after being represented by at least 0.1% of all assigned read pairs. It was used in order to avoid false positive results.

6.1.5 Functional annotation

Functional annotation was conducted by EggNOG-mapper v.2.1.2 (Huerta-Cepas et al., 2017). Sequences were searched against eggNOG database (Huerta-Cepas et al., 2019) for best seed orthologs using DIAMOND and fine-grained orthology assignments were retrieved from pre-computed eggNOG phylogenetic trees. Then functional descriptions of retrieved orthologs including Gene Ontology (GO) terms (Consortium, 2004), KOs, Enzyme Commission (EC) numbers Webb et al. (1992), Carbohydrate-Active Enzymes (CAZy) terms Cantarel et al. (2009) and Clusters of Orthologous Groups (COG) functional categories Tatusov et al. (2000) were transferred to query sequences.

6.2 Candidate models for fitting rarefaction curves

Table S1: Candidate species accumulation models. Dependent variable $D^{(q)}$ is Hill number of order q and independent variable x is sequencing depth. a, b, c, d are fitted coefficients.

Model	Parameter(k)	Derivative	Asymptote	Reference
$D^{(q)} = \frac{ax}{bx+1}$	2	$\frac{dD^{(q)}}{dx} = \frac{a}{(bx+1)^2}$	$\frac{a}{b}$	Clench (1979)
$D^{(q)} = a(1 - e^{-bx})$	2	$\frac{dD^{(q)}}{dx} = abe^{-bx}$	a	Miller and Wiegert (1989)
$D^{(q)} = a - bc^x$	3	$\frac{dD^{(q)}}{dx} = -bc^x \log(c)$	a	Ratkowsky (1983)
$D^{(q)} = a(1 - e^{-bx})^c$	3	$\frac{dD^{(q)}}{dx} = abce^{-bx}(1 - e^{-bx})^{c-1}$	a	Ratkowsky and Giles (1990)
$D^{(q)} = a(1 - (1 + (\frac{x}{c})^d)^{-b})$	4	$\frac{dD^{(q)}}{dx} = \frac{abd}{c}(\frac{x}{c})^{d-1}(1 + (\frac{x}{c})^d)^{-b-1}$	a	Mielke Jr and Johnson (1974)

474 6.3 Exemplification of effect of quality control

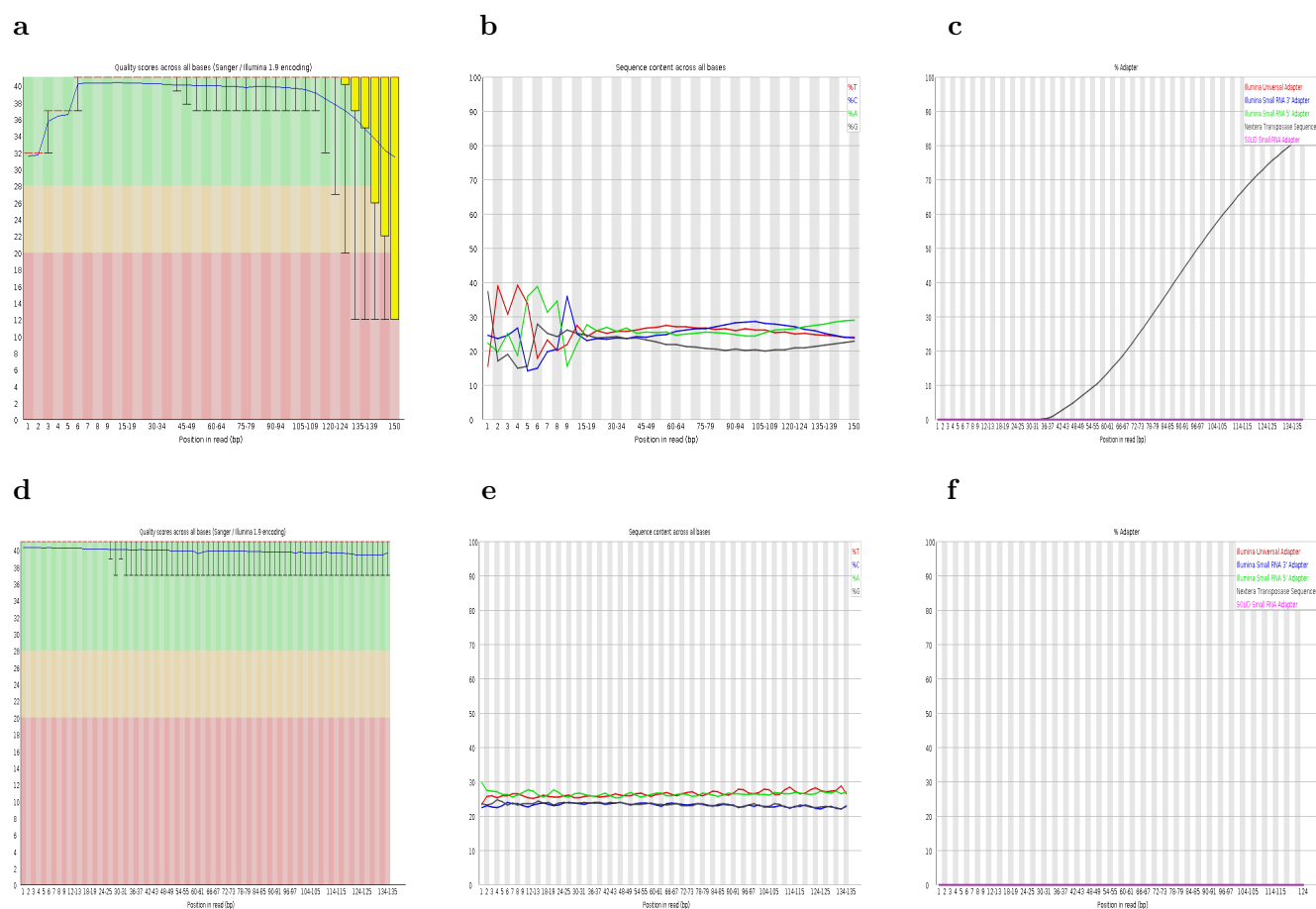


Figure S1: Data quality report of forward reads from bumble bee sample *Bee_Bimpatiens_hv4_1*. a, b and c shows low base quality in 3'-end, uneven base content in 5'-end and present of Nextera adaptors in raw data, respectively. d, e and f shows the same results from clean data.

6.4 Numbers of raw, clean and non-host read pairs

Table S2: Basic statistics of sequenced read pairs. Data in "Clean" considers as background the numbers of raw read pairs ("Raw"), while data in "Non-host" considers as background the numbers of clean read pairs ("Clean").

Sample	Host	Raw	Clean (%)	Non-host (%)
<i>Bee_Amellifera_hv13_1</i>	<i>Apis mellifera</i>	1104861	842095 (76.22%)	665507 (79.03%)
<i>Bee_Amellifera_hv13_2</i>	<i>Apis mellifera</i>	56836899	35282101 (62.08%)	10413704 (29.52%)
<i>Bee_Amellifera_hv15_1</i>	<i>Apis mellifera</i>	63159968	45059211 (71.34%)	38837759 (86.19%)
<i>Bee_Amellifera_wild_1</i>	<i>Apis mellifera</i>	58113227	44466776 (76.52%)	17014144 (38.26%)
<i>Bee_Bimpatiens_hv3_1</i>	<i>Bombus impatiens</i>	63973750	53612702 (83.80%)	5300592 (9.89%)
<i>Bee_Bimpatiens_hv4_1</i>	<i>Bombus impatiens</i>	58988182	48426748 (82.10%)	3557052 (7.35%)
<i>Bee_Bimpatiens_hv4_2</i>	<i>Bombus impatiens</i>	54955553	45805759 (83.35%)	4618023 (10.08%)
<i>Flower_eDNA</i>	None	1443107	882436 (61.15%)	882436 (100%)

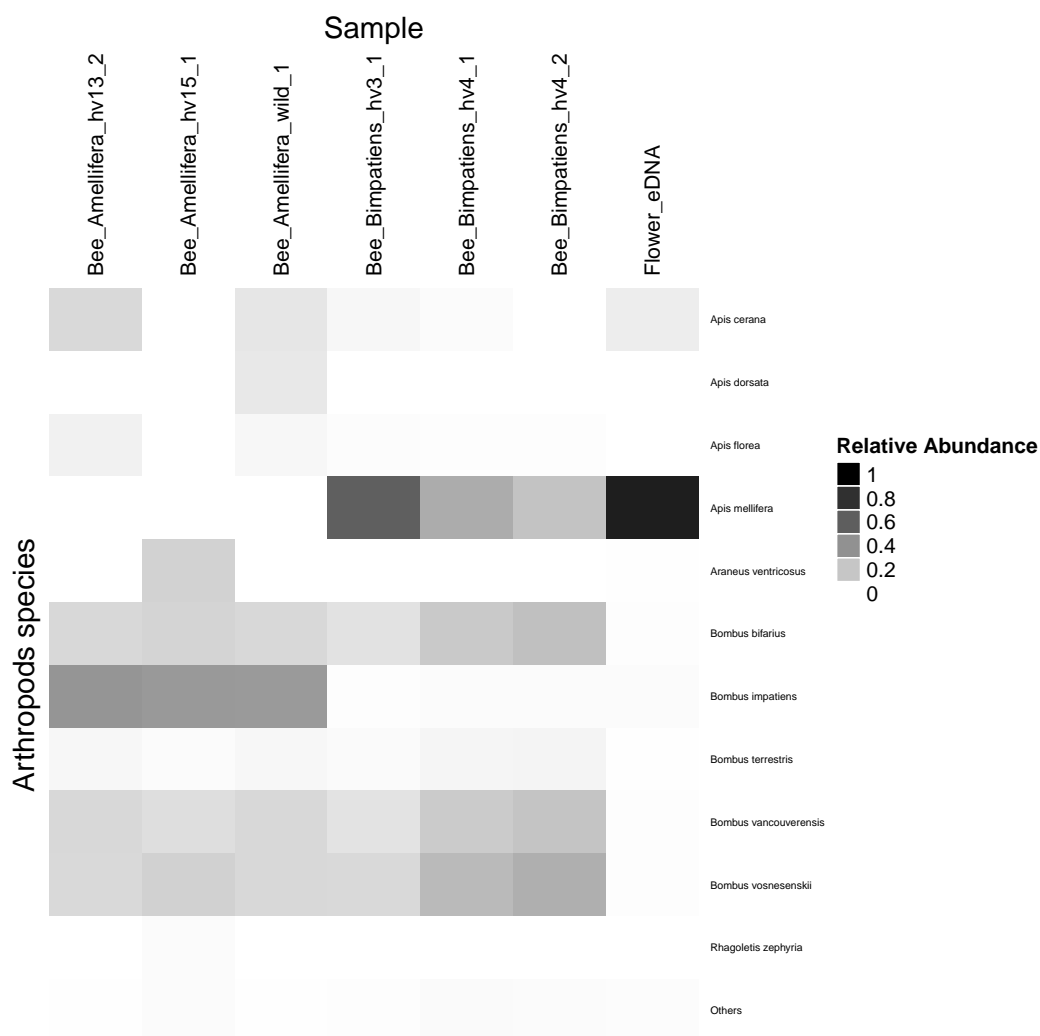


Figure S2: Heatmaps for arthropod species abundance distribution in all samples. The relative abundance takes reads assigned to arthropod species as background. Species with relative abundance smaller than 1% in all samples are collapsed as "others". It should be noted that for bee samples, host contamination was removed before taxon profiling. As a result, the relative abundances of honey bees are extremely low in three honey bee samples, and the same for bumble bees in three bumble bee samples.

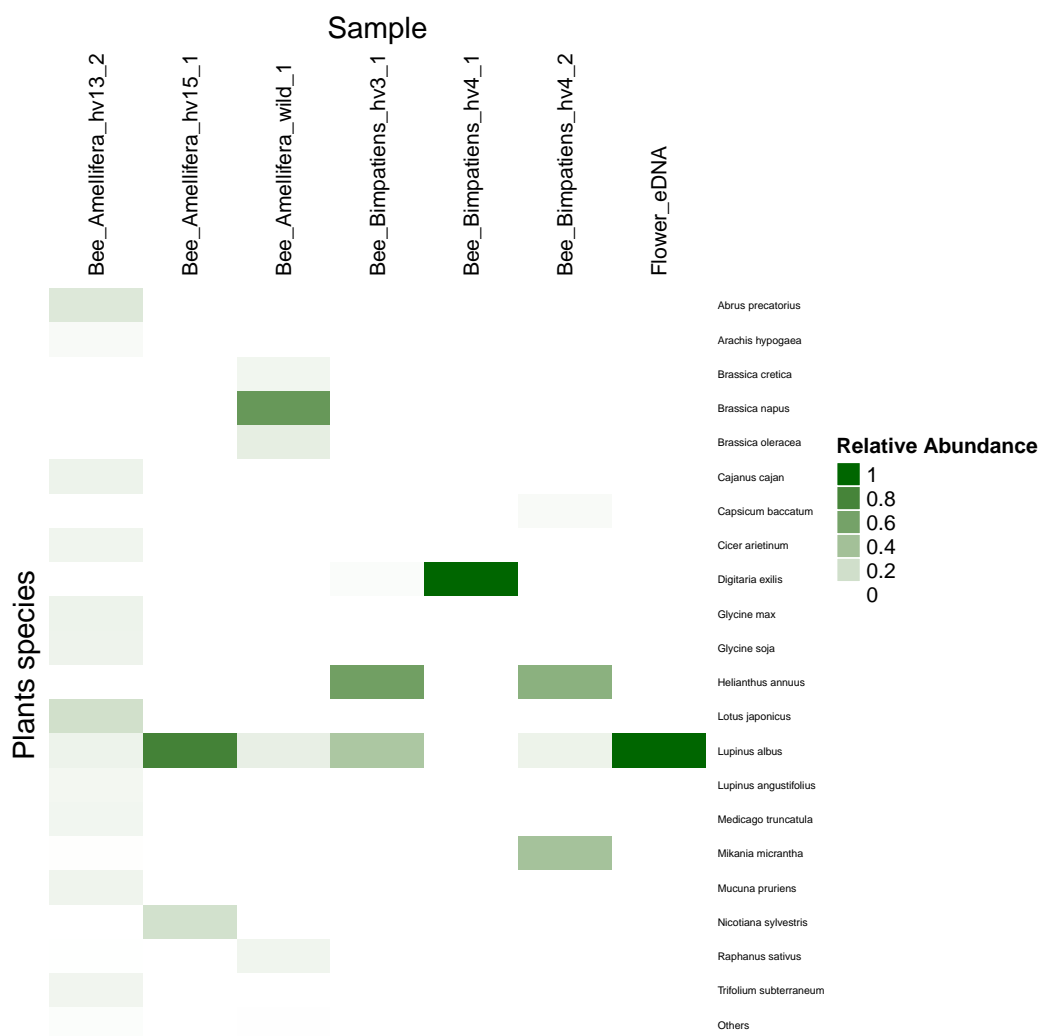


Figure S3: Heatmaps for plant species abundance distribution in all samples. The relative abundance takes reads assigned to plant species as background. Species with relative abundance smaller than 1% in all samples are collapsed as "others".

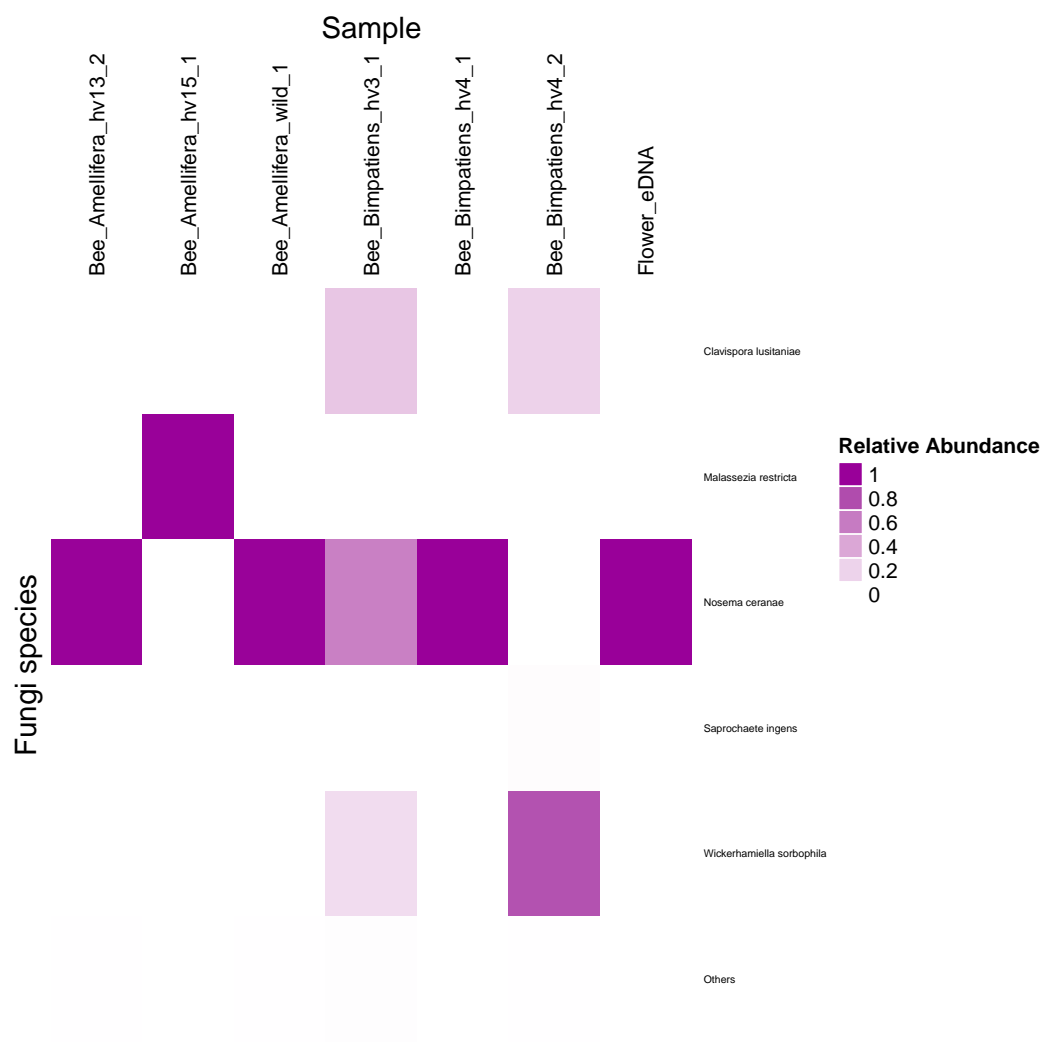


Figure S4: Heatmaps for fungal species abundance distribution in all samples. The relative abundance takes reads assigned to fungal species as background. Species with relative abundance smaller than 1% in all samples are collapsed as "others".

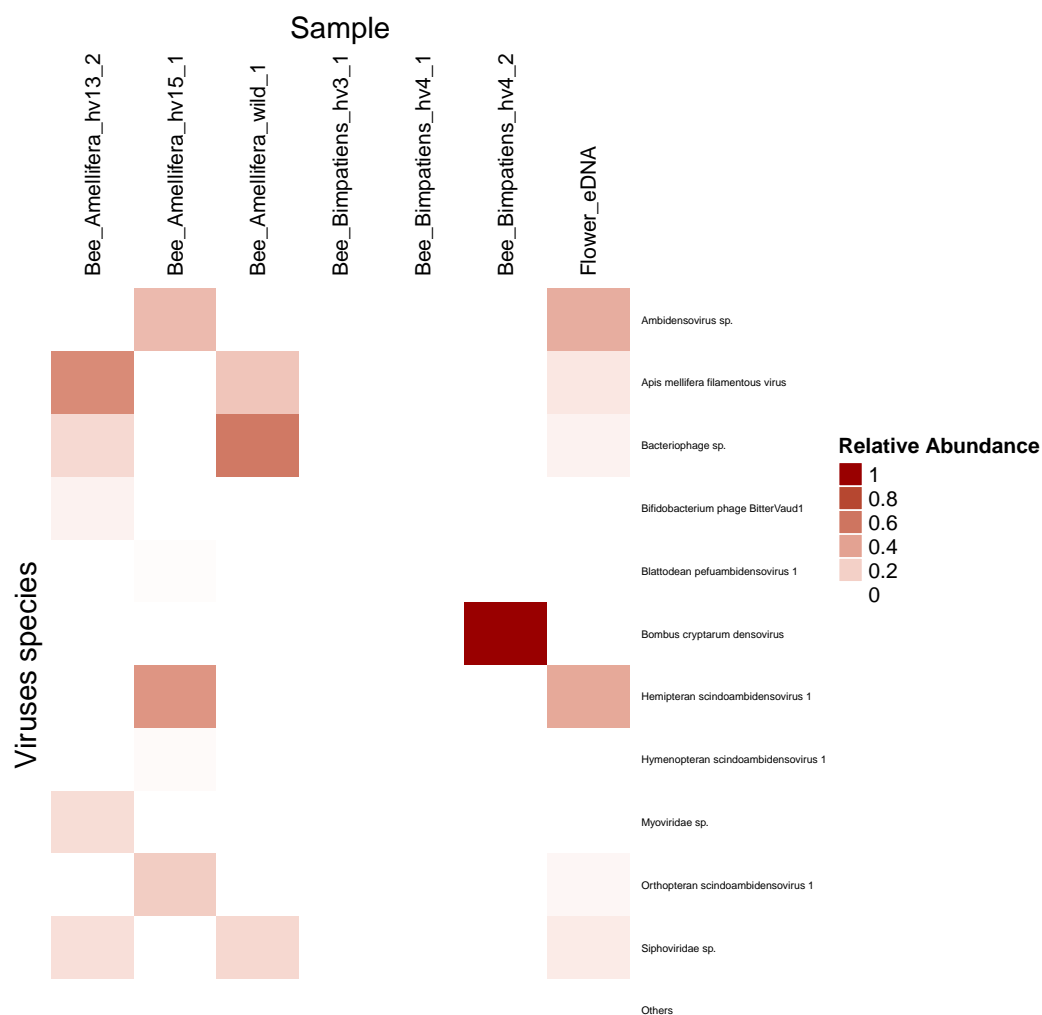


Figure S5: Heatmaps for virus species abundance distribution in all samples. The relative abundance takes reads assigned to virus species as background. Species with relative abundance smaller than 1% in all samples are collapsed as "others".

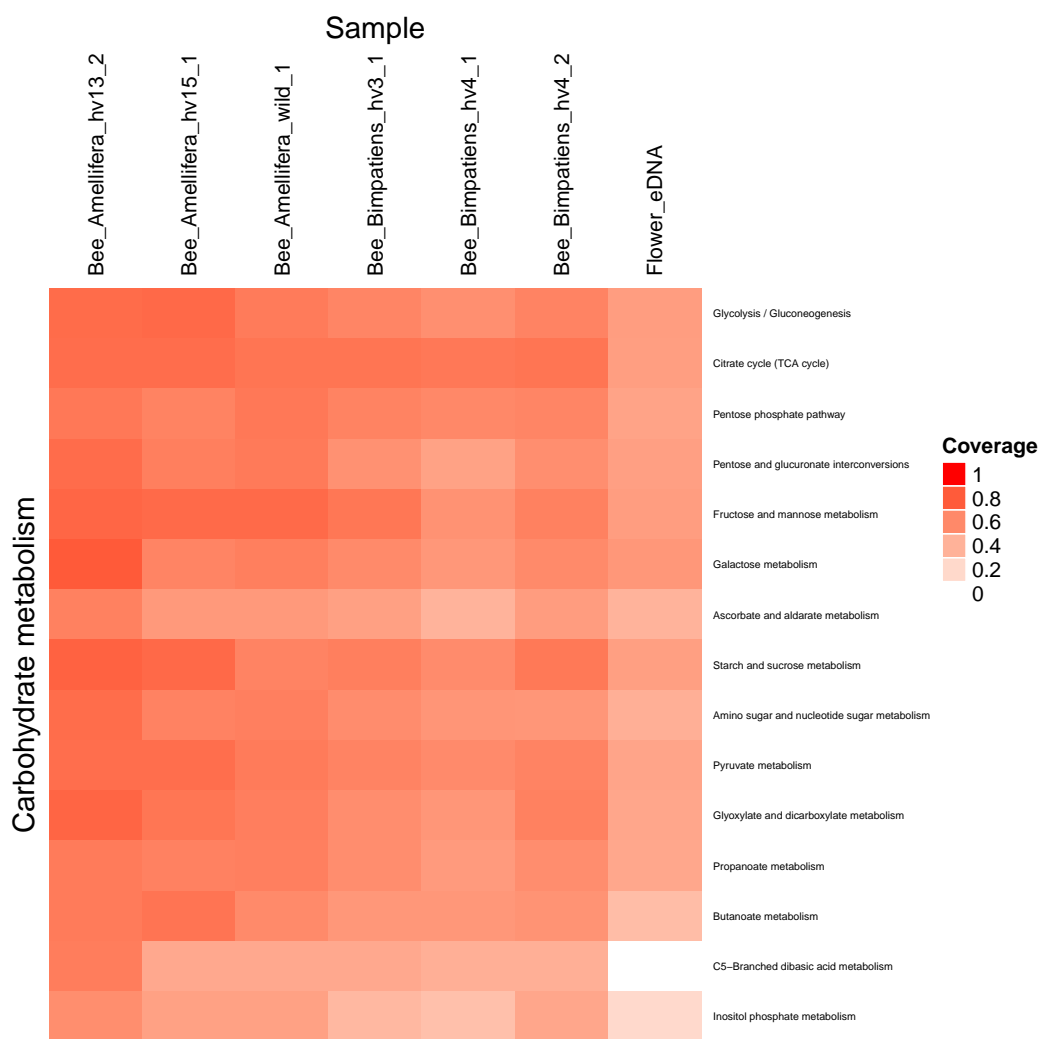


Figure S6: Heatmaps for pathways of carbohydrate metabolism.

References

- Ahmed Abdelfattah, Antonino Malacrino, Michael Wisniewski, Santa O Cacciola, and Leonardo Schena. Metabarcoding: A powerful tool to investigate microbial communities and shape future plant protection strategies. *Biological Control*, 120:1–10, 2018.
- Sahar Abubucker, Nicola Segata, Johannes Goll, Alyxandria M Schubert, Jacques Izard, Brandi L Cantarel, Beltran Rodriguez-Mueller, Jeremy Zucker, Mathangi Thiagarajan, Bernard Henrissat, et al. Metabolic reconstruction for metagenomic data and its application to the human microbiome. *PLoS computational biology*, 8(6):e1002358, 2012.
- Lynn S Adler. The ecological significance of toxic nectar. *Oikos*, 91(3):409–420, 2000.
- Esmaeil Amiri, Prashant Waiker, Olav Rueppell, and Prashanti Manda. Using manual and computer-based text-mining to uncover research trends for *apis mellifera*. *Veterinary sciences*, 7(2):61, 2020.
- David R Anderson. *Model based inference in the life sciences: a primer on evidence*. Springer Science & Business Media, 2007.
- Simon Andrews et al. Fastqc: a quality control tool for high throughput sequence data, 2010.
- Kathrin P Aßhauer, Bernd Wemheuer, Rolf Daniel, and Peter Meinicke. Tax4fun: predicting functional profiles from metagenomic 16s rRNA data. *Bioinformatics*, 31(17):2882–2884, 2015.
- Martin Ayling, Matthew D Clark, and Richard M Leggett. New approaches for metagenome assembly with short reads. *Briefings in bioinformatics*, 21(2):584–594, 2020.
- L Bailey, JM Carpenter, and RD Woods. Properties of a filamentous virus of the honey bee (*apis mellifera*). *Virology*, 114(1):1–7, 1981.
- Leslie Bailey. Honey bee pathology. *Annual review of entomology*, 13(1):191–212, 1968.
- Svenja Bänisch, Teja Tschardtke, Doreen Gabriel, and Catrin Westphal. Crop pollination services: complementary resource use by social vs solitary bees facing crops with contrasting flower supply. *Journal of Applied Ecology*, 58(3):476–485, 2021.
- Daniela Becker, Denny Popp, Hauke Harms, and Florian Centler. A modular metagenomics pipeline allowing for the inclusion of prior knowledge using the example of anaerobic digestion. *Microorganisms*, 8(5):669, 2020.

Anthony M Bolger, Marc Lohse, and Bjoern Usadel. Trimmomatic: a flexible trimmer for illumina
sequence data. *Bioinformatics*, 30(15):2114–2120, 2014.

Evan Bolyen, Jai Ram Rideout, Matthew R Dillon, Nicholas A Bokulich, Christian C Abnet, Gabriel A
Al-Ghalith, Harriet Alexander, Eric J Alm, Manimozhiyan Arumugam, Francesco Asnicar, et al. Re-
producible, interactive, scalable and extensible microbiome data science using qiime 2. *Nature biotech-*
nology, 37(8):852–857, 2019.

Germán Bonilla-Rosso, Théodora Steiner, Fabienne Wichmann, Evan Bexkens, and Philipp Engel. Honey
bees harbor a diverse gut virome engaging in nested strain-level interactions with the microbiota.
Proceedings of the National Academy of Sciences, 117(13):7355–7362, 2020.

Samuele Bovo, Anisa Ribani, Valerio Joe Utzeri, Giuseppina Schiavo, Francesca Bertolini, and Luca
Fontanesi. Shotgun metagenomics of honey dna: Evaluation of a methodological approach to describe
a multi-kingdom honey bee derived environmental dna signature. *PLoS One*, 13(10):e0205575, 2018.

Samuele Bovo, Valerio Joe Utzeri, Anisa Ribani, Riccardo Cabbri, and Luca Fontanesi. Shotgun se-
quencing of honey dna can describe honey bee derived environmental signatures and the honey bee
hologenome complexity. *Scientific reports*, 10(1):1–17, 2020.

Michael A Brockhurst, Ellie Harrison, James PJ Hall, Thomas Richards, Alan McNally, and Craig
MacLean. The ecology and evolution of pangenomes. *Current Biology*, 29(20):R1094–R1103, 2019.

Mark JF Brown and Robert J Paxton. The conservation of bees: a global perspective. *Apidologie*, 40(3):
410–416, 2009.

Benjamin Buchfink, Chao Xie, and Daniel H Huson. Fast and sensitive protein alignment using diamond.
Nature methods, 12(1):59–60, 2015.

Brian Bushnell. Bbmap: a fast, accurate, splice-aware aligner. Technical report, Lawrence Berkeley
National Lab.(LBNL), Berkeley, CA (United States), 2014.

Brandi L Cantarel, Pedro M Coutinho, Corinne Rancurel, Thomas Bernard, Vincent Lombard, and
Bernard Henrissat. The carbohydrate-active enzymes database (cazy): an expert resource for glycoge-
nomics. *Nucleic acids research*, 37(suppl_1):D233–D238, 2009.

Daniel P Cariveau, J Elijah Powell, Hauke Koch, Rachael Winfree, and Nancy A Moran. Variation in
gut microbial communities and its association with pathogen infection in wild bumble bees (*bombus*).
The ISME journal, 8(12):2369–2379, 2014.

534 Rogan Carr and Elhanan Borenstein. Comparative analysis of functional metagenomic annotation and
535 the mappability of short reads. *PloS one*, 9(8):e105776, 2014.

536 Federica Cattonaro, Alessandro Spadotto, Slobodanka Radovic, and Fabio Marroni. Do you cov me?
537 effect of coverage reduction on metagenome shotgun sequencing studies. *F1000Research*, 7, 2018.

538 Anne Chao and Lou Jost. Coverage-based rarefaction and extrapolation: standardizing samples by
539 completeness rather than size. *Ecology*, 93(12):2533–2547, 2012.

540 Anne Chao, Chun-Huo Chiu, and Lou Jost. Unifying species diversity, phylogenetic diversity, functional
541 diversity, and related similarity and differentiation measures through hill numbers. *Annual review of*
542 *ecology, evolution, and systematics*, 45:297–324, 2014a.

543 Anne Chao, Nicholas J Gotelli, TC Hsieh, Elizabeth L Sander, KH Ma, Robert K Colwell, and Aaron M
544 Ellison. Rarefaction and extrapolation with hill numbers: a framework for sampling and estimation in
545 species diversity studies. *Ecological monographs*, 84(1):45–67, 2014b.

546 Robin L Chazdon, Robert K Colwell, Julie S Denslow, and Manuel R Guariguata. Statistical methods for
547 estimating species richness of woody regeneration in primary and secondary rain forests of northeastern
548 costa rica. 1998.

549 Wenda Cheng and Louise Ashton. Ecology: What affects the distribution of global bee diversity. *Current*
550 *Biology*, 31(3):R127–R128, 2021.

551 Stacy Ciufo, Sivakumar Kannan, Shobha Sharma, Azat Badretdin, Karen Clark, Seán Turner, Slava
552 Brover, Conrad L Schoch, Avi Kimchi, and Michael DiCuccio. Using average nucleotide identity to
553 improve taxonomic assignments in prokaryotic genomes at the ncbi. *International journal of systematic*
554 *and evolutionary microbiology*, 68(7):2386, 2018.

555 Truman B Clark. A filamentous virus of the honey bee. *Journal of Invertebrate pathology*, 32(3):332–340,
556 1978.

557 Harry K Clench. How to make regional lists of butterflies: some thoughts. *Journal of the Lepidopterists’*
558 *Society*, 1979.

559 Robert K Colwell and Jonathan A Coddington. Estimating terrestrial biodiversity through extrapolation.
560 *Philosophical Transactions of the Royal Society of London. Series B: Biological Sciences*, 345(1311):
561 101–118, 1994.

Gene Ontology Consortium. The gene ontology (go) database and informatics resource. *Nucleic acids research*, 32(suppl_1):D258–D261, 2004.

Samantha M Cook, Caroline S Awmack, Darren A Murray, and Ingrid H Williams. Are honey bees’ foraging preferences affected by pollen amino acid composition? *Ecological Entomology*, 28(5):622–627, 2003.

Otto X Cordero and Martin F Polz. Explaining microbial genomic diversity in light of evolutionary ecology. *Nature Reviews Microbiology*, 12(4):263–273, 2014.

Eamonn P Culligan, Roy D Sleator, Julian R Marchesi, and Colin Hill. Metagenomics and novel gene discovery: promise and potential for novel therapeutics. *Virulence*, 5(3):399–412, 2014.

Mónica de la Fuente, Pablo F Penas, and Alberto Sols. Mechanism of mannose toxicity. *Biochemical and biophysical research communications*, 140(1):51–55, 1986.

Gavin M Douglas, Robert G Beiko, and Morgan GI Langille. Predicting the functional potential of the microbiome from marker genes using picrust. In *Microbiome Analysis*, pages 169–177. Springer, 2018.

Julia Ebeling, Henriette Knispel, Gillian Hertlein, Anne Fünfhaus, and Elke Genersch. Biology of paenibacillus larvae, a deadly pathogen of honey bee larvae. *Applied microbiology and biotechnology*, 100(17):7387–7395, 2016.

Paul B Eckburg, Elisabeth M Bik, Charles N Bernstein, Elizabeth Purdom, Les Dethlefsen, Michael Sargent, Steven R Gill, Karen E Nelson, and David A Relman. Diversity of the human intestinal microbial flora. *science*, 308(5728):1635–1638, 2005.

Kirsten M Ellegaard, Shota Suenami, Ryo Miyazaki, and Philipp Engel. Vast differences in strain-level diversity in the gut microbiota of two closely related honey bee species. *Current Biology*, 30(13):2520–2531, 2020.

Akihito Endo and Leon MT Dicks. The genus fructobacillus. *Lactic acid bacteria: biodiversity and taxonomy*, pages 381–389, 2014.

Philipp Engel, Vincent G Martinson, and Nancy A Moran. Functional diversity within the simple gut microbiota of the honey bee. *Proceedings of the National Academy of Sciences*, 109(27):11002–11007, 2012.

Philipp Engel, Waldan K Kwong, Quinn McFrederick, Kirk E Anderson, Seth Michael Barribeau, James Angus Chandler, R Scott Cornman, Jacques Dainat, Joachim R De Miranda, Vincent Doublet,

et al. The bee microbiome: impact on bee health and model for evolution and ecology of host-microbe interactions. *MBio*, 7(2):e02164–15, 2016.

Christiane Essoh, Libera Latino, Cédric Midoux, Yann Blouin, Guillaume Loukou, Simon-Pierre A Nguetta, Serge Lathro, Arsher Cablanmian, Athanase K Kouassi, Gilles Vergnaud, et al. Investigation of a large collection of pseudomonas aeruginosa bacteriophages collected from a single environmental source in abidjan, côte d’ivoire. *PLoS one*, 10(6):e0130548, 2015.

Rita Földesi, Brad G Howlett, Ingo Grass, and Péter Batáry. Larger pollinators deposit more pollen on stigmas across multiple plant species a meta-analysis. *Journal of Applied Ecology*, 58(4):699–707, 2021.

Eva Forsgren, Tobias C Olofsson, Alejandra Váasquez, and Ingemar Fries. Novel lactic acid bacteria inhibiting paenibacillus larvae in honey bee larvae. *Apidologie*, 41(1):99–108, 2010.

Jessica Galloway-Peña and Blake Hanson. Tools for analysis of the microbiome. *Digestive diseases and sciences*, 65(3):674–685, 2020.

Christina Geldert, Zaid Abdo, Jane E Stewart, and Arathi HS. Dietary supplementation with phytochemicals improves diversity and abundance of honey bee gut microbiota. *Journal of Applied Microbiology*, 130(5):1705–1720, 2021.

José Antonio Gómez-Anaya, Rodolfo Novelo-Gutiérrez, Alonso Ramírez, and Roberto Arce-Pérez. Using empirical field data of aquatic insects to infer a cut-off slope value in asymptotic models to assess inventories completeness. *Revista mexicana de biodiversidad*, 85(1):218–227, 2014.

Nicholas J Gotelli and Robert K Colwell. Estimating species richness. *Biological diversity: frontiers in measurement and assessment*, 12(39-54):35, 2011.

Antonius Petrus de Groot. Protein and amino acid requirements of the honeybee (*apis mellifica* l.). 1953.

CBOL Plant Working Group, Peter M Hollingsworth, Laura L Forrest, John L Spouge, Mehrdad Hajibabaei, Sujeevan Ratnasingham, Michelle van der Bank, Mark W Chase, Robyn S Cowan, David L Erickson, et al. A dna barcode for land plants. *Proceedings of the National Academy of Sciences*, 106(31):12794–12797, 2009.

H Soon Gweon, Liam P Shaw, Jeremy Swann, Nicola De Maio, Manal AbuOun, Rene Niehus, Alasdair TM Hubbard, Mike J Bowes, Mark J Bailey, Tim EA Peto, et al. The impact of sequencing depth on the inferred taxonomic composition and amr gene content of metagenomic samples. *Environmental Microbiome*, 14(1):1–15, 2019.

620 Hidenori Hayashi, Mitsuo Sakamoto, and Yoshimi Benno. Phylogenetic analysis of the human gut mi-
621 crobiota using 16s rdna clone libraries and strictly anaerobic culture-based methods. *Microbiology and*
622 *immunology*, 46(8):535–548, 2002.

623 Paul DN Hebert, Sujeevan Ratnasingham, and Jeremy R De Waard. Barcoding animal life: cytochrome c
624 oxidase subunit 1 divergences among closely related species. *Proceedings of the Royal Society of London.*
625 *Series B: Biological Sciences*, 270(suppl.1):S96–S99, 2003.

626 Kenneth L Heck Jr, Gerald van Belle, and Daniel Simberloff. Explicit calculation of the rarefaction
627 diversity measurement and the determination of sufficient sample size. *Ecology*, 56(6):1459–1461, 1975.

628 Harmen P Hendriksma and Sharoni Shafir. Honey bee foragers balance colony nutritional deficiencies.
629 *Behavioral Ecology and Sociobiology*, 70(4):509–517, 2016.

630 Harmen P Hendriksma, Karmi L Oxman, and Sharoni Shafir. Amino acid and carbohydrate tradeoffs
631 by honey bee nectar foragers and their implications for plant–pollinator interactions. *Journal of insect*
632 *physiology*, 69:56–64, 2014.

633 Harmen P Hendriksma, Collin D Pachow, and James C Nieh. Effects of essential amino acid supplemen-
634 tation to promote honey bee gland and muscle development in cages and colonies. *Journal of insect*
635 *physiology*, 117:103906, 2019.

636 Mark O Hill. Diversity and evenness: a unifying notation and its consequences. *Ecology*, 54(2):427–432,
637 1973.

638 Joaquín Hortal and Jorge M Lobo. An ed-based protocol for optimal sampling of biodiversity. *Biodiversity*
639 *& Conservation*, 14(12):2913–2947, 2005.

640 Joaquín Hortal, Patrícia Garcia-Pereira, and Enrique García-Barros. Butterfly species richness in main-
641 land portugal: predictive models of geographic distribution patterns. *Ecography*, 27(1):68–82, 2004.

642 Joaquín Hortal, Paulo AV Borges, and Clara Gaspar. Evaluating the performance of species richness
643 estimators: sensitivity to sample grain size. *Journal of animal ecology*, 75(1):274–287, 2006.

644 Peter Hristov, Boyko Neov, Rositsa Shumkova, and Nadezhda Palova. Significance of apoidea as main
645 pollinators. ecological and economic impact and implications for human nutrition. *Diversity*, 12(7):280,
646 2020a.

647 Peter Hristov, Rositsa Shumkova, Nadezhda Palova, and Boyko Neov. Factors associated with honey bee
648 colony losses: a mini-review. *Veterinary Sciences*, 7(4):166, 2020b.

Jaime Huerta-Cepas, Kristoffer Forslund, Luis Pedro Coelho, Damian Szklarczyk, Lars Juhl Jensen, Christian Von Mering, and Peer Bork. Fast genome-wide functional annotation through orthology assignment by eggno-mapper. *Molecular biology and evolution*, 34(8):2115–2122, 2017.

Jaime Huerta-Cepas, Damian Szklarczyk, Davide Heller, Ana Hernández-Plaza, Sofia K Forslund, Helen Cook, Daniel R Mende, Ivica Letunic, Thomas Rattei, Lars J Jensen, et al. eggno 5.0: a hierarchical, functionally and phylogenetically annotated orthology resource based on 5090 organisms and 2502 viruses. *Nucleic acids research*, 47(D1):D309–D314, 2019.

Jennifer B Hughes and Jessica J Hellmann. The application of rarefaction techniques to molecular inventories of microbial diversity. *Methods in enzymology*, 397:292–308, 2005.

Robert M Hughes, Alan T Herlihy, and David V Peck. Sampling efforts for estimating fish species richness in western usa river sites. *Limnological*, 87:125859, 2021.

Stuart H Hurlbert. The nonconcept of species diversity: a critique and alternative parameters. *Ecology*, 52(4):577–586, 1971.

Daniel H Huson, Alexander F Auch, Ji Qi, and Stephan C Schuster. Megan analysis of metagenomic data. *Genome research*, 17(3):377–386, 2007.

Daniel H Huson, Benjamin Albrecht, Caner Bağcı, Irina Bessarab, Anna Gorska, Dino Jolic, and Rohan BH Williams. Megan-lr: new algorithms allow accurate binning and easy interactive exploration of metagenomic long reads and contigs. *Biology direct*, 13(1):1–17, 2018.

Alberto Jiménez-Valverde and Jorge M Lobo. Determining a combined sampling procedure for a reliable estimation of araneidae and thomisidae assemblages (arachnida, araneae). *The Journal of Arachnology*, 33(1):33–42, 2005.

Alberto Jimenez-Valverde, Silvia Jimenez Mendoza, Jose Martin Cano, and Miguel L Munguira. Comparing relative model fit of several species-accumulation functions to local papilionoidea and hesperioidea butterfly inventories of mediterranean habitats. In *Arthropod diversity and conservation*, pages 163–176. Springer, 2006.

Minoru Kanehisa and Susumu Goto. Kegg: kyoto encyclopedia of genes and genomes. *Nucleic acids research*, 28(1):27–30, 2000.

Karen M Kapheim, Makenna M Johnson, and Maggi Jolley. Composition and acquisition of the microbiome in solitary, ground-nesting alkali bees. *Scientific reports*, 11(1):1–11, 2021.

678 Alexander Keller, Quinn S McFrederick, Prarthana Dharampal, Shawn Steffan, Bryan N Danforth, and
679 Sara D Leonhardt. (more than) hitchhikers through the network: The shared microbiome of bees and
680 flowers. *Current Opinion in Insect Science*, 2020.

681 Lucie Kešnerová, Roxane Moritz, and Philipp Engel. *Bartonella apis* sp. nov., a honey bee gut symbiont
682 of the class alphaproteobacteria. *International journal of systematic and evolutionary microbiology*, 66
683 (1):414–421, 2016.

684 Shaden AM Khalifa, Esraa H Elshafiey, Aya A Shetaia, Aida A Abd El-Wahed, Ahmed F Algethami,
685 Syed G Musharraf, Mohamed F AlAjmi, Chao Zhao, Saad HD Masry, Mohamed M Abdel-Daim, et al.
686 Overview of bee pollination and its economic value for crop production. *Insects*, 12(8):688, 2021.

687 Hauke Koch and Paul Schmid-Hempel. Bacterial communities in central european bumblebees: low
688 diversity and high specificity. *Microbial Ecology*, 62(1):121–133, 2011a.

689 Hauke Koch and Paul Schmid-Hempel. Socially transmitted gut microbiota protect bumble bees against
690 an intestinal parasite. *Proceedings of the National Academy of Sciences*, 108(48):19288–19292, 2011b.

691 Waldan K Kwong, Luis A Medina, Hauke Koch, Kong-Wah Sing, Eunice Jia Yu Soh, John S Ascher,
692 Rodolfo Jaffé, and Nancy A Moran. Dynamic microbiome evolution in social bees. *Science Advances*,
693 3(3):e1600513, 2017.

694 Gerardo Lamas, Robert K Robbins, and Donald J Harvey. A preliminary survey of the butterfly fauna
695 of pakitza, parque nacional del manu, peru, with an estimate of its species richness. *Publicaciones del*
696 *Museo de Historia Natural Universidad Nacional Mayor de San Marcos*, 1991.

697 Ben Langmead and Steven L Salzberg. Fast gapped-read alignment with bowtie 2. *Nature methods*, 9(4):
698 357, 2012.

699 Fredrick J Lee, Douglas B Rusch, Frank J Stewart, Heather R Mattila, and Irene LG Newton. Saccharide
700 breakdown and fermentation by the honey bee gut microbiome. *Environmental microbiology*, 17(3):
701 796–815, 2015.

702 Fredrick J Lee, Kayla I Miller, James B McKinlay, and Irene LG Newton. Differential carbohydrate
703 utilization and organic acid production by honey bee symbionts. *FEMS microbiology ecology*, 94(8):
704 fy113, 2018.

705 Dinghua Li, Chi-Man Liu, Ruibang Luo, Kunihiro Sadakane, and Tak-Wah Lam. Megahit: an ultra-
706 fast single-node solution for large and complex metagenomics assembly via succinct de bruijn graph.
707 *Bioinformatics*, 31(10):1674–1676, 2015.

708 Heng Li and Richard Durbin. Fast and accurate short read alignment with burrows-wheeler transform.
709 *bioinformatics*, 25(14):1754–1760, 2009.

710 Heng Li, Bob Handsaker, Alec Wysoker, Tim Fennell, Jue Ruan, Nils Homer, Gabor Marth, Goncalo
711 Abecasis, and Richard Durbin. The sequence alignment/map format and samtools. *Bioinformatics*, 25
712 (16):2078–2079, 2009.

713 Patrick W Maes, Pedro AP Rodrigues, Randy Oliver, Brendon M Mott, and Kirk E Anderson. Diet-
714 related gut bacterial dysbiosis correlates with impaired development, increased mortality and nosema
715 disease in the honeybee (*apis mellifera*). *Molecular Ecology*, 25(21):5439–5450, 2016.

716 Tomonori Matsuzawa, Ryo Kohsaka, and Yuta Uchiyama. Application of environmental dna: Honey
717 bee behavior and ecosystems for sustainable beekeeping. In *Modern beekeeping-bases for sustainable*
718 *production*. IntechOpen, 2020.

719 Paul W Mielke Jr and Earl S Johnson. Some generalized beta distributions of the second kind having
720 desirable application features in hydrology and meteorology. *Water Resources Research*, 10(2):223–226,
721 1974.

722 Alla Mikheenko, Vladislav Saveliev, and Alexey Gurevich. Metaquast: evaluation of metagenome assem-
723 blies. *Bioinformatics*, 32(7):1088–1090, 2016.

724 Ronald I Miller and Richard G Wiegert. Documenting completeness, species-area relations, and the species-
725 abundance distribution of a regional flora. *Ecology*, 70(1):16–22, 1989.

726 Nancy A Moran. Genomics of the honey bee microbiome. *Current opinion in insect science*, 10:22–28,
727 2015.

728 Claudia E Moreno and Gonzalo Halffter. Assessing the completeness of bat biodiversity inventories using
729 species accumulation curves. *Journal of Applied ecology*, 37(1):149–158, 2000.

730 Felicia N New and Ilana L Brito. What is metagenomics teaching us, and what is missed? *Annual Review*
731 *of Microbiology*, 74:117–135, 2020.

732 Susan W Nicolson and Robert W Thornburg. Nectar chemistry. In *Nectaries and nectar*, pages 215–264.
733 Springer, 2007.

734 Sergey I Nikolenko, Anton I Korobeynikov, and Max A Alekseyev. Bayeshammer: Bayesian clustering for
735 error correction in single-cell sequencing. In *BMC genomics*, volume 14, pages 1–11. Springer, 2013.

736 R Henrik Nilsson, Erik Kristiansson, Martin Ryberg, Nils Hallenberg, and Karl-Henrik Larsson. Intraspe-
737 cific its variability in the kingdom fungi as expressed in the international sequence databases and its
738 implications for molecular species identification. *Evolutionary bioinformatics*, 4:EBO–S653, 2008.

739 Sergey Nurk, Dmitry Meleshko, Anton Korobeynikov, and Pavel A Pevzner. metaspades: a new versatile
740 metagenomic assembler. *Genome research*, 27(5):824–834, 2017.

741 Matthew R Olm, Alexander Crits-Christoph, Spencer Diamond, Adi Lavy, Paula B Matheus Carnevali,
742 and Jillian F Banfield. Consistent metagenome-derived metrics verify and delineate bacterial species
743 boundaries. *Msystems*, 5(1):e00731–19, 2020.

744 Pier P Paoli, Luisa A Wakeling, Geraldine A Wright, and Dianne Ford. The dietary proportion of essential
745 amino acids and sir2 influence lifespan in the honeybee. *Age*, 36(3):1239–1247, 2014.

746 K Papadopoulou-Karabela, N Iliadis, V Liakos, and E Bourdzy-Hatzopoulou. Experimental infection of
747 honeybees by pseudomonas aeruginosa. *Apidologie*, 23(5):393–397, 1992.

748 K Papadopoulou-Karabela, N Iliadis, and V Liakos. Haemocyte changes in honeybee (*apis mellifera* l)
749 artificially infected by pseudomonas aeruginosa. *Apidologie*, 24(1):81–86, 1993.

750 Yu Peng, Henry CM Leung, Siu-Ming Yiu, and Francis YL Chin. Idba-ud: a de novo assembler for single-
751 cell and metagenomic sequencing data with highly uneven depth. *Bioinformatics*, 28(11):1420–1428,
752 2012.

753 Judit J Péntzes, Maria Söderlund-Venermo, Marta Canuti, Anna Maria Eis-Hübinger, Joseph Hughes,
754 Susan F Cotmore, and Balázs Harrach. Reorganizing the family parvoviridae: a revised taxonomy
755 independent of the canonical approach based on host association. *Archives of Virology*, 165(9):2133–
756 2146, 2020.

757 Joana Pereira-Marques, Anne Hout, Rui M Ferreira, Michiel Weber, Ines Pinto-Ribeiro, Leen-Jan van
758 Doorn, Cornelis Willem Knetsch, and Ceu Figueiredo. Impact of host dna and sequencing depth
759 on the taxonomic resolution of whole metagenome sequencing for microbiome analysis. *Frontiers in*
760 *microbiology*, 10:1277, 2019.

761 Elijah Powell, Nalin Ratnayeke, and Nancy A Moran. Strain diversity and host specificity in a specialized
762 gut symbiont of honeybees and bumblebees. *Molecular ecology*, 25(18):4461–4471, 2016.

763 J Elijah Powell, Zac Carver, Sean P Leonard, and Nancy A Moran. Field-realistic tylosin exposure
764 impacts honey bee microbiota and pathogen susceptibility, which is ameliorated by native gut probiotics.
765 *Microbiology Spectrum*, 9(1):e00103–21.

766 Andrey Prjibelski, Dmitry Antipov, Dmitry Meleshko, Alla Lapidus, and Anton Korobeynikov. Using
767 spades de novo assembler. *Current Protocols in Bioinformatics*, 70(1):e102, 2020.

768 Christopher Quince, Alan W Walker, Jared T Simpson, Nicholas J Loman, and Nicola Segata. Shotgun
769 metagenomics, from sampling to analysis. *Nature biotechnology*, 35(9):833–844, 2017.

770 Aaron R Quinlan and Ira M Hall. Bedtools: a flexible suite of utilities for comparing genomic features.
771 *Bioinformatics*, 26(6):841–842, 2010.

772 David Ratkowsky. Nonlinear regression modelling. 1983.

773 David A Ratkowsky and David EA Giles. *Handbook of nonlinear regression models*. Number 04; QA278.
774 2, R3. M. Dekker New York, 1990.

775 Kasie Raymann and Nancy A Moran. The role of the gut microbiome in health and disease of adult honey
776 bee workers. *Current opinion in insect science*, 26:97–104, 2018.

777 Anisa Ribani, Valerio Joe Utzeri, Valeria Taurisano, and Luca Fontanesi. Honey as a source of environ-
778 mental dna for the detection and monitoring of honey bee pathogens and parasites. *Veterinary sciences*,
779 7(3):113, 2020.

780 Luis M Rodriguez-R, Juan C Castro, Nikos C Kyrpides, James R Cole, James M Tiedje, and Kon-
781 stantinos T Konstantinidis. How much do rrna gene surveys underestimate extant bacterial diversity?
782 *Applied and Environmental Microbiology*, 84(6):e00014–18, 2018.

783 Torbjørn Rognes, Tomáš Flouri, Ben Nichols, Christopher Quince, and Frédéric Mahé. Vsearch: a versatile
784 open source tool for metagenomics. *PeerJ*, 4:e2584, 2016.

785 Michael Roswell, Jonathan Dushoff, and Rachael Winfree. A conceptual guide to measuring species
786 diversity. *Oikos*, 130(3):321–338, 2021.

787 Howard L Sanders. Marine benthic diversity: a comparative study. *The American Naturalist*, 102(925):
788 243–282, 1968.

789 Patrick D Schloss. Reintroducing mothur: 10 years later. *Applied and environmental microbiology*, 86(2):
790 e02343–19, 2020.

791 Karel Schoonvaere, Guy Smagghe, Frédéric Francis, and Dirk C de Graaf. Study of the metatranscrip-
792 tome of eight social and solitary wild bee species reveals novel viruses and bee parasites. *Frontiers in*
793 *microbiology*, 9:177, 2018.

794 Itai Sharon, Michael Kertesz, Laura A Hug, Dmitry Pushkarev, Timothy A Blauwkamp, Cindy J Castelle,
795 Mojgan Amirebrahimi, Brian C Thomas, David Burstein, Susannah G Tringe, et al. Accurate, multi-kb
796 reads resolve complex populations and detect rare microorganisms. *Genome research*, 25(4):534–543,
797 2015.

798 Nicola K Simcock, Helen E Gray, and Geraldine A Wright. Single amino acids in sucrose rewards modulate
799 feeding and associative learning in the honeybee. *Journal of insect physiology*, 69:41–48, 2014.

800 Jorge SoberónM and Jorge LlorenteB. The use of species accumulation functions for the prediction of
801 species richness. *Conservation biology*, 7(3):480–488, 1993.

802 Daniel Stabler, Pier P Paoli, Susan W Nicolson, and Geraldine A Wright. Nutrient balancing of the adult
803 worker bumblebee (*bombus terrestris*) depends on the dietary source of essential amino acids. *The*
804 *Journal of experimental biology*, 218(5):793–802, 2015.

805 Roman L Tatusov, Michael Y Galperin, Darren A Natale, and Eugene V Koonin. The cog database: a
806 tool for genome-scale analysis of protein functions and evolution. *Nucleic acids research*, 28(1):33–36,
807 2000.

808 Michelle A Taylor, Alastair W Robertson, Patrick J Biggs, Kate K Richards, Daniel F Jones, and Shan-
809 thi G Parkar. The effect of carbohydrate sources: Sucrose, invert sugar and components of mānuka
810 honey, on core bacteria in the digestive tract of adult honey bees (*apis mellifera*). *PloS one*, 14(12):
811 e0225845, 2019.

812 Quang Tran and Vinhthuy Phan. Assembling reads improves taxonomic classification of species. *Genes*,
813 11(8):946, 2020.

814 Susannah Green Tringe, Christian Von Mering, Arthur Kobayashi, Asaf A Salamov, Kevin Chen, Hwai W
815 Chang, Mircea Podar, Jay M Short, Eric J Mathur, John C Detter, et al. Comparative metagenomics
816 of microbial communities. *Science*, 308(5721):554–557, 2005.

817 Rachel L Vannette. The floral microbiome: plant, pollinator, and microbial perspectives. *Annual Review*
818 *of Ecology, Evolution, and Systematics*, 51:363–386, 2020.

819 Louise Vermote, Marko Verce, Luc De Vuyst, and Stefan Weckx. Amplicon and shotgun metagenomic
820 sequencing indicates that microbial ecosystems present in cheese brines reflect environmental inoculation
821 during the cheese production process. *International Dairy Journal*, 87:44–53, 2018.

822 John Vollmers, Sandra Wiegand, and Anne-Kristin Kaster. Comparing and evaluating metagenome
823 assembly tools from a microbiologists perspective-not only size matters! *PloS one*, 12(1):e0169662,
824 2017.

825 Kai Wang, Jiahuan Li, Liuwei Zhao, Xiyan Mu, Chen Wang, Miao Wang, Xiaofeng Xue, Suzhen Qi, and
826 Liming Wu. Gut microbiota protects honey bees (*apis mellifera* l.) against polystyrene microplastics
827 exposure risks. *Journal of Hazardous Materials*, 402:123828, 2021.

828 Edwin C Webb et al. *Enzyme nomenclature 1992. Recommendations of the Nomenclature Committee of*
829 *the International Union of Biochemistry and Molecular Biology on the Nomenclature and Classification*
830 *of Enzymes*. Number Ed. 6. Academic Press, 1992.

831 K Eric Wommack, Jaysheel Bhavsar, and Jacques Ravel. Metagenomics: read length matters. *Applied*
832 *and environmental microbiology*, 74(5):1453–1463, 2008.

833 Yuzhen Ye and Thomas G Doak. A parsimony approach to biological pathway reconstruction/inference
834 for genomes and metagenomes. *PLoS computational biology*, 5(8):e1000465, 2009.

835 Nicholas D Youngblut, Jacobo De la Cuesta-Zuluaga, Georg H Reischer, Silke Dauser, Nathalie Schuster,
836 Chris Walzer, Gabrielle Stalder, Andreas H Farnleitner, and Ruth E Ley. Large-scale metagenome as-
837 sembly reveals novel animal-associated microbial genomes, biosynthetic gene clusters, and other genetic
838 diversity. *Msystems*, 5(6):e01045–20, 2020.

839 Rahat Zaheer, Noelle Noyes, Rodrigo Ortega Polo, Shaun R Cook, Eric Marinier, Gary Van Domselaar,
840 Keith E Belk, Paul S Morley, and Tim A McAllister. Impact of sequencing depth on the characterization
841 of the microbiome and resistome. *Scientific reports*, 8(1):1–11, 2018.

842 Eduardo E Zattara and Marcelo A Aizen. Worldwide occurrence records suggest a global decline in bee
843 species richness. *One Earth*, 4(1):114–123, 2021.

844 Hao Zheng, Julie Perreau, J Elijah Powell, Benfeng Han, Zijing Zhang, Waldan K Kwong, Susannah G
845 Tringe, and Nancy A Moran. Division of labor in honey bee gut microbiota for plant polysaccharide
846 digestion. *Proceedings of the National Academy of Sciences*, 116(51):25909–25916, 2019.

847 Jinshui Zheng, Stijn Wittouck, Elisa Salvetti, Charles MAP Franz, Hugh Harris, Paola Mattarelli, Paul W
848 O’Toole, Bruno Pot, Peter Vandamme, Jens Walter, et al. A taxonomic note on the genus *Lactobacillus*:
849 Description of 23 novel genera, emended description of the genus *Lactobacillus* Beijerinck 1901, and union
850 of *Lactobacillaceae* and *Leuconostocaceae*. 2020.