1 **Evaluating multiplexed next-generation sequencing as a method in palynology**

2 **for mixed pollen samples**

3 Alexander Keller[1,2,$], Nadja Danner[1], Gudrun Grimmer[1,2], Markus Ankenbrand[3],

4 Katharina von der Ohe[4], Werner von der Ohe[4], Simone Rost[5], Stephan Härtel[1], Ingolf

5 Steffan-Dewenter[1]


6 [1] Department of Animal Ecology and Tropical Biology, University of Wuerzburg,

7 Biocenter, Am Hubland, Wuerzburg, Germany

8 [2] DNA Analytics Core Facility, University of Wuerzburg, Biocenter, Am Hubland,

9 Wuerzburg, Germany

10 [3] Department of Bioinformatics, University of Wuerzburg, Biocenter, Am Hubland,

11 Wuerzburg, Germany

12 [4] LAVES Institut für Bienenkunde, Herzogin-Eleonore-Allee 5, Celle, Germany

13 [5] Department of Human Genetics, University of Wuerzburg, Biocenter, Am Hubland,

14 Wuerzburg, Germany

15

16 Keywords: DNA barcoding, high-throughput, internal transcribed spacer 2, ITS2,

17 molecular ecology, meta-barcoding, pollination, plant-pollinator interactions,

18 phylotyping

19 [$] Corresponding author: Alexander Keller, Department of Animal Ecology and

20 Tropical Biology, Biocenter, Am Hubland, 97074 Würzburg, Germany, eMail:

21 a.keller@biozentrum.uni-wuerzburg.de, Fax: 0049 931 31 84352

22 Short title: Evaluating NGS-based palynology

## Abstract

The identification of pollen plays an important role in ecology, palaeo-climatology, honey quality control and other areas. Currently, expert knowledge and reference collections are essential to identify pollen origin through light microscopy. Pollen identification through molecular sequencing and DNA barcoding has been proposed as an alternative approach, but the assessment of mixed pollen samples originating from multiple plant species is still a tedious and error-prone task. Next-generation sequencing has been proposed to avoid this hindrance. In this study we assessed mixed pollen probes through next-generation sequencing of amplicons from the highly variable species-specific internal transcribed spacer 2 region of the nuclear ribosomal DNA. Further, we developed a bioinformatical workflow to analyse these high-throughput data with a newly created reference database. To evaluate the feasibility, we compared results from classical identification based on light microscopy from the same samples with our sequencing results. We assessed in total 16 mixed pollen samples, 14 originated from honey bee colonies and two from solitary bee nests. The sequencing technique resulted in higher taxon richness (deeper assignments and more identified taxa) compared to light microscopy. Abundance estimations from sequencing data were significantly correlated with counted abundances through light microscopy. Simulation analyses of taxon specificity and sensitivity indicate that 96% of taxa present in the database are correctly identifiable at the genus level and 70% at the species level. Next-generation sequencing thus presents a useful and efficient workflow to identify pollen at the genus and species level without requiring specialized palynological expert knowledge.

## Introduction

Palynology, the scientific study of pollen and identification of its origin, plays an important role in studying mechanisms of plant-pollinator interactions (Wilcock and Neiland, 2002), resource use of flower-visiting animals (Kleijn and Raemakers, 2008; Wcislo and Cane, 1996) and climate-related variation of plant communities through time (Marchant et al., 2001; Sugita, 1994; Tzedakis, 1993). Pollen grains often display a species-specific morphology with diverse structure and sculpture. However, it remains difficult to delineate between closely related species when using light microscopy (Mullins and Emberlin, 1997). As a result, many pollen types are simply grouped at genus or family level (Davies and Fall, 2001) and data analyses on pollen diversity are strongly limited (Bagella et al., 2013). DNA barcoding, i.e. to identify and classify organisms according to a nucleotide sequence was often and successfully applied to all major groups of organisms, also plants including pollen (Chen et al., 2010; Hebert et al., 2003; Zhou et al., 2007). Accordingly, molecular tools to analyze pollens have also substantially increased in their application and show great potential especially with difficult, also fossil taxa and those with low taxonomic knowledge (Bennett and Parducci, 2006; Wilson et al., 2010; Zhou et al., 2007).

It is further a promising new approach in ecology to directly determine the diversity of organisms in environmental samples (Sheffield et al., 2009; Valentini et al., 2009), i.e. samples that represent a mixture of species, e.g. faeces, soil or pollen collections, for which identification with classical methods is difficult or incomplete (Wilson et al., 2010). To analyze mixed sets of pollens originating from different plant organisms with DNA barcoding however is still a tedious and error-prone task, requiring manual separation of pollens to taxa, each to be amplified and sequenced individually. Studies evaluating applicability of high-throughput techniques to pollen materials are currently lacking (Taylor and Harris, 2012; Wilson et al., 2010) or are restricted to specific investigations using quantitative real time polymerase chain reaction (qrtPCR) where prior information about present organisms is

4

79    required (Agodi et al., 2006; Schnell et al., 2010). Palynology would therefore benefit

80    from species-level determination from mixed samples, larger counts, higher

81    processing speed, improved objectivity, and automation to be attractive for large

82    scale studies (Stillman and Flenley, 1996). Molecular methods based on high-

83    throughput DNA-sequencing could provide the requested features to extent and

84    improve classical pollen determination. Valentini et al. (2010) proposed next-

85    generation sequencing (NGS) as a suitable method for this task. We agree with this

86    idea and thus evaluated in this study the performance and reliability of the new

87    sequencing and bioinformatical strategies by directly comparing it with data

88    obtained by light microscopy.

89

90    Specifically we address the following challenges that emerge in DNA barcoding with

91    mixed pollen samples. (1) A laboratory routine has to be defined which can be

92    applied to all major plant clades, requiring universality of amplification priming

93    regions and adequate length to be suitable for next-generation sequencing while

94    holding enough sequence variation to differ between species. This routine includes

95    DNA extraction, amplification, sample multiplexing, library preparation, sequencing

96    with high-throughput devices and raw-data cleanup. Also, (2) a mapping algorithm

97    has to be developed which adequately maps obtained sequences in their full length

98    to references, preferably in a hierarchical progression with confidence values for

99    each level of taxonomy. Further, this algorithm has to be with good performance to

100   be able to process high-throughput data on a standard desktop computer and

101   produce results in reasonable time. (3) A comprehensive reference database is

102   required to derive the desired taxonomic annotations.

103

104   Several genetic marker regions have been proposed for DNA barcoding in plants

105   that match the requirements, foremost presence and feasibility to be amplified in all

106   investigated taxa, as well as low intra-specific but high inter-specific variability to

107   succeed in being species-specific (Chen et al., 2010; Hebert et al., 2003;

108   Hollingsworth et al., 2011; Zhou et al., 2007). In this study, we use the internal

109   transcribed spacer 2 (ITS2) region, which has been shown to be suitable as a

110    barcode for plants (92.7% successful identifications in 6,600 samples, Chen et al.,

111    2010; Buchheim et al., 2011). Also the enclosed genetic regions (5.8S and 28S) are

112    highly conserved throughout the eukaryotes. Thus an universal primer for the

113    analysis of probes consisting of multiple organisms is applicable with a low risk to

114    exclude taxa from amplification (Chen et al., 2010; Keller et al., 2009; White et al.,

115    1990). A further reason for choosing this marker is that a comprehensive ITS2

116    database already exists (Koetschan et al., 2010) helping to prepare reference

117    sequences suitable for our needs.

118

119    We approached the targeted tasks by combining and adapting existing molecular

120    and bioinformatical tools to develop new functionalities for DNA barcoding of pollen

121    samples that consist of multiple taxa. We then evaluated the performance and

122    quality of the molecular and bioinformatical workflow by comparing our results

123    with data from classical light microscopy identification of pollen samples. Further,

124    we tested the applicability for samples with low pollen contents and performed

125    computer-based simulations to validate that the bioinformatical classification

126    pipeline is trustable.

127    ## Materials and Methods

128    **Pollen collection**

129    The honey-bee pollen samples were collected in twelve different landscapes in the

130    region around Bayreuth, Germany. The distance between landscapes was at least 3

131    km leading to diversified pollen inputs depending on the surrounding floral

132    resources. In the centre of each landscape we established a honey bee colony (*Apis*

133    *mellifera carnica* L.) with a pollen trap in front of the hive entrance. Returning

134    foragers had to pass a 5 mm grid taking off the pollen loads from their hind legs.

135    From 21.07.2009-12.08.2009 every one to three days accumulated pollen loads

136    were removed from the traps and stored as individual samples at -18 °C until the

137    end of the sampling period. Pollen samples were dried at 30 °C for one week.

138    Further, to assess variability in resource use of honeybees at one location, samples

139    from three colonies located at the same study site were separately analysed (in the

140    following designated as Samples 12a, 12b and 12c). From each of the fourteen

141    samples (one per colony) 20% of the collected pollen were randomly taken and

142    mixed for further analyses.

143

144    We performed NGS as well as microscopic assessment of the samples. The samples

145    were split into independent aliquots for these separate, blinded analyses. NGS was

146    performed by AK, GG and MA, whereas the samples were classified through classical

147    light microscopy by ND with expert guidance by KvO, without knowledge of the

148    other group's results.

149

150    Two further pollen samples were obtained from solitary bee nests (*Osmia bicornis*

151    L.) in October 2012 by swabbing the cell walls with cotton buds (Keller et al., 2013).

152    In contrast to the relatively pure pollen samples obtained from honey bees, this

153    experiment reflects samples strongly contaminated with nest building materials

154    (soil) and faeces, challenging to analyse with traditional methods. Solitary bee

155    samples were thus only processed with NGS.

156

157    **Classical pollen identification**

158    Pollen samples were first analyzed using light microscopy in the LAVES Institut für

159    Bienenkunde in Celle, Germany. For the microscopic pollen determination, 10mg

160    pollen loads of each sample were homogenized in 50ml demineralized water with a

161    magnetic stirrer for one hour. 15 µl of the solution and 30 µl demineralized water

162    were transferred to a slide, distributed equally over an area of the size of a cover

163    glass and embedded in glycerin gelatin after complete dehydration following the

164    method of Behm et al. (1996). From each sample 500 randomly selected pollen

165    grains were determined on genus level and where possible to species level. Very

166    rarely occurring pollen types were not determined (Behm et al., 1996).

167

168    **Molecular pollen identification**

169    Second pollen identification was done by DNA barcoding of the ITS2 region. The

170    main working steps described below were DNA extraction, amplification,

171    sequencing, bioinformatic cleanup and taxonomic classification.

172

173    <u>DNA extraction, amplification and sequencing:</u> For each sample, 2 g of pollens were

174    added to 4 ml of bidest $H_2O$ and homogenized with an electronic pistil within a

175    plastic tube. Of this emulsion, 200 µl (equaling approximately 50 mg of pollens)

176    were taken for the following extraction. We grinded the aliquot with the TissueLyser

177    LT (Qiagen, Hilden, Germany) and extracted DNA using the Machery-Nagel (Düren,

178    Germany) NucleoSpin Food Kit. We followed the special supplementary guidelines

179    for pollen samples provided by the manufacturer. For polymerase chain reaction

180    (PCR) amplification we used the primers S2F and ITS4R originally designed by Chen

181    et al. (2010) and White et al. (1990) to span a mean region of approximately 350bp.

182    This covers the complete ITS2 region. We adapted those primers to match 454

183    sequencing purposes and multiplexing by adding the 454 specific Adapters A and B,

184    the linker key, and a variable multiplex identifier (MID). Thus the forward "fusion"

185    primer was 5'-CGT ATC GCC TCC CTC GCG CCA TCA GAT GCG ATA CTT GGT GTG AAT

186    -3' and the reverse "fusion" primer 5'-5′CTA TGC GCC TTG CCA GCC CGC TCA GXX

187    XXX XXX XXT CCT CCG CTT ATT GAT ATG C-3', where the X-region designates a

188    variable MID. In total, 16 MIDs were taken from the official Roche technical bulletin

189    (454 Sequencing Technical Bulletin No. 005-2009, April 2009) to be able to process

190    all our samples with one sequencing chip.

191

192    PCR reaction mixes consisted of 0.25 µl of each forward and reverse primer (each

193    30µM molar), 3 µl of template DNA and 25µl of Phusion High-Fidelity DNA

194    polymerase PCR 2x MasterMix (Thermo Scientific, Waltham, MA, USA). Bidest $H_2O$

195    was added to a reaction volume of 50 µl. Samples were initially denaturated at 94 °C

196    for 4 min, then amplified by using 37 cycles of 95 °C for 40 s, 49 °C for 40 s and 72 °C

197    for 40 s. A final extension (72 °C) of 5 min was added at the end of the program to

198    ensure complete amplification. All samples were amplified in ten separate aliquots

199    to reduce random effects on the community during PCR amplification (Fierer et al.,

200 2008). PCR amplicons of these ten replicates were combined, gel-electrophoresed,
201 trimmed for amplicon length and cleaned with the HiYield PCR Clean-up Kit (Real
202 Biotech Corporation, Banqiao City, Taiwan) according to the manufacturers
203 description. Cleaned samples were quantified using a Qubit II Flurometer
204 (Invitrogen/Life Technologies, Carlsbad, CA, USA) and the dsDNA High-Sensitivity
205 Assay Kit (also Invitrogen/Life Technologies) as described in the vendors protocol.
206 We used the BioAnalyzer 2200 (Agilent, Santa Clara, CA, USA) with High Sensitivity
207 DNA Chips (also Agilent) for verification of fragment length distributions.
208 Pyrosequencing and library preparation was performed according to the guidelines
209 for the GS junior (Roche, Basel, Switzerland). Sequencing was performed in-house
210 with a GS junior device located in the Department of Human Genetics (University of
211 Würzburg, Germany) with original Roche GS junior titanium chemistry.
212
213 <u>Bioinformatic cleanup:</u> Data was demultiplexed into the different samples using the
214 MID adapter sequences and the QIIME software (Caporaso et al., 2010; Kuczynski et
215 al., 2011). During this step, only sequences spanning both priming regions were
216 further used, i.e. only completely sequenced amplicons. Primers, adapters and MIDs
217 were trimmed. Chimeric checking and quality filtering was also performed during
218 this step. We restricted data to high quality reads with a phred score ≥ 27, (Kunin et
219 al., 2010) and no reads with ambiguous characters were included in the following
220 downstream analyses.
221
222 <u>Hierarchic classification:</u> Taxonomic assignments were performed with the RDP
223 (Ribosomal Database Project) classifier (Wang et al., 2007) and an ITS2 specific,
224 novel reference set created and evaluated as described below. Further, we applied a
225 bootstrap cut-off at 85% as classification threshold with respect to the maximum f-
226 measure in the training database evaluation (also see below).
227
228 **Method comparison statistics**
229 Most of the analyses were performed on a generic level, as both methods yielded
230 some taxa only assignable to this level. With a generic analysis all identified taxa

231    were directly comparable. With this data we compared taxon richness and identified

232    species overlaps and differences obtained from the two methods. Rarefaction curves

233    for each plot were generated with R (R Development Core Team, 2010) in the NGS

234    data to evaluate species richness in relation to sequencing depth. Abundance was

235    assessed relatively in percent of total number of reads and in percent of 500 pollen

236    grains (Behm et al., 1996) for NGS and light microscopy, respectively. We used

237    overall and per-plot abundance of these relative accounts to compare between the

238    two methodologies by Pearson's product moment correlation using R (R

239    Development Core Team, 2010).

240

241    **Molecular reference database training**

242    Taxonomic classifications with DNA barcodes are currently mostly done via

243    phylogenetic analyses (Buchheim et al., 2011), pairwise alignments with specific

244    reference sequences (Chen et al., 2010) or BLAST searches (Basic Local Alignment

245    Search Tool) (Altschul et al., 1990) in GenBank (Benson et al., 2010) or other

246    nucleotide databases. The first both methods require that prior knowledge about

247    taxonomy is present to select suitable taxa included into the recalculated

248    phylogenetic tree or alignment. This is not feasible for mixed pollen collections,

249    where the included taxa are unknown prior to assessment or stem from very

250    different taxonomic groups. BLAST searches have to be performed very carefully, as

251    hits may include local alignments and identity calculations may thus be based only

252    on parts of the query and reference sequences. Further, the raw output of a BLAST

253    search is often obscured as a lot of hits are not taxonomically annotated or flagged

254    as "environmental samples". A novel approach to tackle these drawbacks has been

255    proposed with a Bayesian classification algorithm (Wang et al., 2007). It provides

256    hierarchical taxonomic assignments of DNA sequences and is well accepted in the

257    scientific community as especially high throughput analyses profit from the

258    efficiency and accuracy of the algorithm (Caporaso et al., 2010). Currently, the only

259    publicly available training sets are limited to bacterial 16S (Wang et al., 2007) and

260    fungal large ribosomal subunit (Liu et al., 2012).

261

262     In this study, a new ITS2 training set was designed for plants. We used the ITS2-

263     Database as an original database which is restricted to structure-validated

264     sequences (Koetschan et al., 2010). All ITS2 sequences matching the taxonomic

265     group "Viridiplantae" and with a sequence length between 200 bp and 400 bp were

266     downloaded, resulting in 73,853 sequences (accessed 3rd March 2013). The

267     taxonomy for each sequence was assigned using the GI (GenBank Identifier) and the

268     corresponding NCBI taxonomy (Federhen, 2012) by Perl scripting and reformated

269     to be usable with the python script "assign taxonomy.py" of the QIIME (Caporaso et

270     al., 2010) package. Additionally, RDP required formats of these preprocessed files

271     were generated. Training was performed with the RDP classifier v2.2 (Wang et al.,

272     2007) as implemented in QIIME. Before training of the final set, we evaluated the

273     performance by varying several parameters of the underlying data to maximize

274     effectiveness and allow quality estimations of the assignments as described in the

275     following.

276

277     Pre-clustering evaluation: Due to intraspecific variation (Song et al., 2012) and

278     sequencing errors in the underlying data (Kunin et al., 2010), pre-clustering of

279     reference sequences prior to training may prove useful to increase reliability of the

280     results (Lan et al., 2012). Thus, from the full data-set we generated eleven separate

281     training sets differing in the pre-clustering threshold of sequences before the actual

282     training. Clusters of sequences were generated at identity levels of 90%, 91% …

283     100%, and only the most abundant sequence of each cluster was picked. This also

284     generated an even distribution of taxonomic units in the sets. To assess the

285     assignment quality and depth, each sequence was reclassified to the training set.

286     Then starting from the root of the taxonomy of each sequence, every taxonomic

287     level of the assignment was compared to the correct taxonomy. If the bootstrap of

288     an assignment was less than 0.8, the level (and all sub-levels) was considered as

289     unassignable. If there was a mismatch between assigned taxonomy and expected

290     taxonomy, the number of remaining sub-levels (plus one), was called erroneous

291     levels. The number of assigned levels before the first mismatch or unassignable level

292     was called correct levels.

11

293

294 <u>Cut-off and assignment quality evaluation:</u> To estimate assignment qualities, the test

295 and training data had to be distinct sets. Further, we wanted to evaluate the

296 effectiveness to identify "new species" that do not have representatives in the

297 training data (Lan et al., 2012). The complete ITS2 reference data set was thus for

298 testing purposes artificially split into three sets representing "training data", "test

299 data A" with references, and "test data B" without references. This was achieved by

300 the following procedure: species with multiple sequences were separated into "test

301 data A" (one sequence) as well as "training data" (remaining sequences). Species

302 with only a single deposited sequence were assigned to category "test data B". For

303 this evaluation purpose, the algorithm was trained only with the set "training data"

304 (36,418 sequences). According to the measures for the RDP classifier evaluation

305 performed by Lan et al. (2012) for the original 16S dataset we estimated the number

306 of "true positive" (TP) and "false negative" (FN) assignments by classifying

307 sequences of "test data A" (10,635 sequences), where references were present in the

308 "training data", Only correct assignments were considered as TP, whereas wrong

309 assignments (to a different species) were added to the list of FNs. Similarly, we

310 classified sequences of "test data B" (26,800 sequences) to determine the number of

311 "true negative" (TN) and "false positive"(FP) hits. With that, we calculated

312 sensitivity $SN = \frac{TP}{TP+FN}$ to identify existing taxa and specificity $SP = \frac{TN}{TN+FP}$ to leave

313 sequences without references unclassified. Using these split data-sets, we were able

314 to estimate SN at species and genus level, whereas SP was only assessable at the

315 species level. We optimized our assignment bootstrap value for classification by

316 maximizing the f-measure as the harmonic mean of sensitivity and specificity at

317 species level $= \frac{2*SN*SP}{SN+SP}$ .

318

319 **Results**

320 **Pollen high-throughput sequencing and classification**

12

321   In total, our study produced 14,924 raw sequences for pollen samples passing

322   Roche's quality filtering of the 454 junior sequencing device. Of these, 9,310 ITS2

323   sequences matched our extended quality standards. The remainders were dismissed

324   as too short (<200 bp), with low quality score (<27), excessive homopolymers (>5

325   bp), chimeric or mismatches in primer regions (Caporaso et al., 2010; Kunin et al.,

326   2010). After removal of adapters and primers, mean sequence length was 348,3 bp

327   (± 28 bp standard deviation), spanning the complete ITS2 region. Individual

328   samples comprised 219-1,179 reads, with mean read length of 330,5 bp – 363,9 bp

329   (± 3,8 bp – 68,2 bp standard deviation). Beside plant sequences, we also found

330   several fungal sequences, belonging to *Issatchenkia occidentalis*, *Cochliobolus sativus*,

331   *Phoma* sp. and *Lewia infectoria*, which are regularly inhabiting or infecting plant

332   tissues.

333

334   **Honey bee pollen samples**

335   For the samples collected by honey-bees, 98.9% of all reads were assignable to

336   genus level with a bootstrap confidence higher or equal than 0.85. At the species

337   level we were able to classify 61.6% of our reads using the same bootstrap cut-off.

338   Reducing the filter's required sequence length to 150 bp did not produce any new

339   classifiable plant taxa. Taxon richness was not correlated with the number of reads

340   within a sample (Pearson's correlation, r = -0.099, df = 12, t = -0.3453, p-value >

341   0.05). Rarefaction showed that we reached a plateau regarding genera richness in all

342   samples (Fig. 1A). These observations suggest that the sequencing depth was

343   adequate to assess the underlying taxon richness.

344

345   We identified a total of 29 different genera of 16 families when we combined the

346   results from molecular sequencing and microscopy (Tab. 1). Further, 24 taxa were

347   also identifiable at the species level. With NGS we found 13 genera that were not

348   identified through microscopy, whereas four genera (*Heracleum, Carduus, Phacelia,*

349   *Convolvulus*) that were identified by light microscopy were missing in the NGS

350   results although having references in the database. One genus (*Vitis*) had no

351  trustable reference sequence in the database and was thus also not identifiable with
352  the NGS method.

353

354  From the phenology of the pollens and presence at plots, we assume that a
355  misidentification of very similar pollens happened with light microscopy which was
356  revealed by NGS: *Tanacetum* and *Scorzoneroides* were both manually misclassified
357  as *Taraxacum*. We observed higher intra-generic taxon richness for *Trifolium*,
358  *Hypochaeris*, *Chamerion* through NGS, yet lesser in *Centaurea* (Fig. 1B).
359  Improvement of the taxonomic assignment was found in four genera, where species
360  levels were obtainable only through NGS. However, *Helianthus* was only classified at
361  genus level, whereas microscopy was able to identify it as *Helianthus annuus*.

362

363  Based on NGS data, taxon richness within the samples ranged from 4 to 12 taxa that
364  were at least classifiable at genus level (Fig. 1B). Correspondingly, diversity ranged
365  from 4 to 12 taxa for the microscopy assessment. Pollen diversity collected by the
366  three colonies from site twelve was 12, 10 and 12 taxa, respectively. The
367  compositional profile was similar for the dominant pollen taxa in all three samples,
368  but still showed considerable variation (Fig. 1B).

369

370  Over all samples, we found a strong correlation of abundance estimations between
371  the two identification methods (Pearson's correlation, r = 0.86, t = 8.71, df = 26, p <
372  0.001***, Fig. 2). This relationship is also reflected on a per plot basis, yet with lower
373  correlation coefficient (Pearson's correlation, r = 0.66, t = 17.36, df = 390, p <
374  0.001***). These results indicate that the abundance estimates of taxa within plots
375  show relatively high similarity between the two methods.

376

377  **Pollens in solitary bee nests**
378  Pollen samples from both solitary bee nests were successfully processed with 100%
379  of reads identifiable at genus level despite high contamination of the samples with
380  nesting materials and faeces. Both samples harbored *Brassica* sp. and *Dioscorea* sp.

381    pollens, the latter most likely *Dioscorea* (*Tamus*) *communis* as the only

382    representative of the Dioscoreaceae present in the sampling region.

383

384    **Molecular reference database training**

385    Pre-clustering of data prior to training of the RDP classifier did not improve the

386    overall performance of classifications (Fig. 3). This was the case both for depth of

387    the assignment as well as the mean number of incorrectly assigned levels, which

388    respectively increase and decrease with higher pre-clustering thresholds. We thus

389    used a cut-off at 100% sequence identity, which equals unique sequences, for the

390    final training set. With that, of the 73,853 tested database sequences, 55,028 were

391    positively identifiable at the species and further 10,518 at the genus level.

392    Surprisingly, 6,104 sequences were assignable only to phylum level. They likely

393    represent contaminations in the reference database.

394

395    Regarding determination of the optimal cut-off threshold, specificity and sensitivity

396    of the novel/known classifications are shown with their dependency of the

397    bootstrap in Fig. 4. The best classification by means of f-measure is achieved with a

398    bootstrap cutoff of 0.85. Specificity and sensitivity at this threshold for species level

399    were both approximately 70%. At the genus level, sensitivity to correctly identify a

400    genus increased to 96%. We thus recommend this threshold when using the RDP

401    classifier with the generated training data.

402

403    Currently, all sequences in the reference data-set accumulate to 37,435 different

404    plant species and 6,162 genera according to NCBI taxonomy (Federhen, 2012). The

405    complete reference dataset is available for download and public usage at

406    http://www.dna-analytics.biozentrum.uni-wuerzburg.de.

407    **Discussion**

408    The demand for methods to identify pollen samples at a high-throughput level is

409    increasing for many applications in ecology and paleo-climatology (Bennett and

410    Parducci, 2006; Sheffield et al., 2009; Taylor and Harris, 2012; Valentini et al., 2009;

411  Wilson et al., 2010; Zhou et al., 2007). DNA barcoding is a frequently and

412  successfully applied method, yet pollens of mixed samples originating from more

413  than one source are currently not assessable through standard methods. Valentini et

414  al. (2010) proposed that next-generation sequencing may counter this deficiency,

415  i.e. to investigate such mixed samples by identifying all included plant organisms

416  together without manual separation. The goals of this study were thus to develop,

417  and moreover evaluate, a molecular laboratory procedure and bioinformatical

418  analysis for such a task. The complete workflow was applied to pollen samples from

419  two different studies (in total 16 samples). The resulting gene sequences allowed to

420  successfully identify taxon richness and abundance of the underlying samples. The

421  resulting taxonomic resolution is similar or better than results from classical light

422  microscopy. Details of the performance of each individual step of the workflow and

423  the resulting methodological and biological relevance are discussed in the following.

424

425  **High-throughput pollen sequencing**

426  In general, our laboratory workflow was suitable in processing mixed pollen probes

427  through next-generation sequencing. However, quality filtering according to our

428  rigorous restrictions reduced the obtained sequences from approximately 15,000

429  sequences to 10,000. Most of them were removed due to failure to include both

430  primer regions and/or multiplex identifier due to low quality scores towards the

431  end of sequences or short read lengths (Caporaso et al., 2010). The first indicates

432  that a large proportion of reads was not fully sequenced with sufficient quality,

433  whereas the latter shows that the primers also amplified shorter fragments than the

434  intended plant ITS2 region. Not fully sequenced reads are a technical issue that is

435  regularly improved by increase of read length and quality through new generations

436  of sequencing devices and chemistry (Metzker, 2009). Improvements are also

437  expectable by applying paired-end strategies, as quality near the ends will increase,

438  or to use technologies with general lower sequencing error rates. Shorter, fully

439  sequenced sequences are project specific problems, but also expectable: as a

440  drawback of universal primers, they will as well amplify fungal ITS2 (White et al.,

441  1990) ranging from ~100 to 250 bp and even other eukaryotic protists with far

442    shorter ITS2 regions (Keller et al., 2009). Further, the existence of non-functional

443    pseudo-genes is known (Harpke and Peterson, 2008). Thus studies investigating

444    plant ITS2 sequences should account for a sufficient overhead of estimated

445    sequences per sample during project design due to sequencing technology and

446    potential contamination through unwanted organisms (Parameswaran et al., 2007).

447    The remaining high quality reads showed a high proportion of classifiable

448    sequences (~99%), whereas reduction of the minimum sequence length had no

449    impact on plant species diversity. Both observations suggest that the filters are

450    adequate to concentrate on the data of interest, i.e. plant sequences.

451

452    **Classification pipeline**

453    To be able to use the RDP classifier (Wang et al., 2007) for taxonomic assignments

454    with plants and with the ITS2 marker, we re-trained the algorithm with structurally

455    verified sequences obtained from the ITS2 database (Koetschan et al., 2010). The

456    underlying dataset incorporates more than 70,000 different plant sequences and

457    represents a cross-section throughout the Viridiplantae. Sequences originate from

458    all biogeographic regions of the world since the primary database is GenBank

459    (Benson et al., 2010). Currently, all sequences in the reference data-set accumulate

460    to 37,435 different plant species and 6,162 genera according to NCBI taxonomy

461    (Federhen, 2012). Exemplarily for the data analysed in this study, the dataset covers

462    79% of all vascular plant genera and 54% of species known to exist within the

463    Federal state Bavaria, Germany, where our samples were obtained (comprehensive

464    plant database http://www.bayernflora.de, accessed 6th November 2013,

465    Staatliche Naturwissenschaftliche Sammlungen Bayern, 2013). As 99% of reads

466    were classifiable to genus level and only one genus (*Vitis*) of the assessed 29 genera

467    in total was missing in the reference database, most of abundant and bee relevant

468    plant genera seem to be included. Further, the classifier's dataset is updateable to

469    match the constantly increasing numbers of sequences deposited in GenBank and

470    the ITS2 database in the future (Wang et al., 2007).

471

| 472 | In the computational evaluation of database and classifier for an ITS2 dataset, we |
| 473 | obtained values comparable to those of existing datasets published for bacteria |
| 474 | (Wang et al., 2007) and fungi (Liu et al., 2012). Taxonomic classifications performed |
| 475 | best regarding sensitivity, i.e. to identify taxa existing in the database, and |
| 476 | specificity, i.e. to restrain from classifying organisms without references, at a |
| 477 | bootstrap threshold level of approximately 0.85 (Lan et al., 2012). Species and genus |
| 478 | level sensitivity to correctly identify sequences with this bootstrap were 70% and |
| 479 | 96%, respectively. This is similar to the classifier's preferred level used to classify |
| 480 | microbial organisms (0.80, Lan et al., 2012; Wang et al., 2007). From a technical |
| 481 | perspective it is thus valid to apply the classification algorithm also for ITS2 |
| 482 | sequences of plants. |
| 483 | |
| 484 | **Comparison of assessment methods** |
| 485 | Using next-generation sequencing, we were clearly able to improve palynology |
| 486 | diversity assessments in comparison with traditional optical microscopy. This |
| 487 | appears in novel taxa that were identified, as well as improvement of classification |
| 488 | of taxa and better possibilities to distinguish species within a genus. Further, some |
| 489 | misidentifications of pollen through microscopy were revealed that were caused by |
| 490 | very similar morphological appearance of closely related species. Also, molecular |
| 491 | assessments were successful for solitary bee nest samples, where swabs included |
| 492 | pollens as well as contaminating materials. Sequencing assessments were |
| 493 | repeatable, identifying similar diversity in samples obtained from different bee |
| 494 | colonies placed within the same landscape. |
| 495 | |
| 496 | However, using the high-throughput approach we also encountered limitations, |
| 497 | which are partly related to the data used for training of the classifier. Regarding the |
| 498 | Vitaceae, the ITS2 database is currently lacking trustable reference sequences. We |
| 499 | validated the only existing sequence, which was considerably short (~200 bp) and |
| 500 | derived from a whole genome shotgun sequencing study (assembled sequence from |
| 501 | short length reads, GenBank ID: AM462492.2, Velasco et al., 2007). Due to intra- |
| 502 | genomic variation of the ITS2 (Song et al., 2012), we assume the assembly yielded a |

503 consensus, stacked ITS2 sequence, unusable for barcoding purposes or that a non-

504 ITS2 region was falsely identified as such by the ITS2 database annotation algorithm

505 (Keller et al., 2009). We therefore dismissed the sequence as missing within the

506 reference database. In general, taxa missing or with inadequate sequences in the

507 underlying database are not identifiable. As shown exemplarily for the geographic

508 region Bavaria, 22% of known plant genera are missing and thus the current

509 coverage is far from complete (Staatliche Naturwissenschaftliche Sammlungen

510 Bayern, 2013). Also, valid sequences with wrong taxonomic annotations may lead to

511 mis-training of the classification model regarding the respective taxa (Bridge et al.,

512 2003). This is exemplified by a proportion of sequences re-classified in the

513 evaluation to a different phylum, suggesting wrong taxonomic annotation of

514 GenBank database sequences. To address limitations of the underlying database

515 (missing or misclassified sequences) in a given research question, we suggest that

516 applied studies should consider also reviewing one cross-section pool of all samples

517 in parallel through optical means to verify the overall richness of taxa relevant for

518 the study. This will also maintain comparability between studies applying

519 traditional and molecular approaches. Despite these database-specific drawbacks,

520 the classifier produced taxonomic assignments that are congruent with light

521 microscopy, and thus corroborating the positive technical evaluation of the pipeline

522 above with a direct comparison of biological data.

523

524 Abundance estimations of both methods showed a strong correlation, suggesting

525 that abundance estimates based on high-throughput sequencing regarding high or

526 low sequence frequency of taxa within the sample are valid. In our study, we took

527 care to reduce amplification biases through PCR with ten aliquots of each sample

528 simultaneously (typical in microbiota studies: three, Fierer et al., 2008)  and a low

529 number of amplification cycles (Suzuki and Giovannoni, 1996). Still, abundances

530 retained from PCR amplified DNA samples have to be regarded critically, as

531 amplification biases through priming preference of specific taxonomic groups,

532 random effects and the exponential nature of the amplification process are not

533 excludable (Spooner, 2009; Suzuki and Giovannoni, 1996). Abundances are thus

534 likely better interpreted categorical (e.g. high abundance, low abundance) than with

535 linear association. With the advent of increased sequencing throughput and third-

536 generation single molecule sequencers without need for amplification (Metzker,

537 2009; Roberts et al., 2013), improved abundance estimations by sequencing are

538 likely in the near future.

539

540 Expenses per sample were almost equal for both applied methods when considering

541 time consumption and consumables. As the trend of sequencing technologies goes

542 rapidly toward higher throughput and resulting multiplexing possibilities (Kozich et

543 al., 2013; Metzker, 2009), we expect price efficiency per sample with next-

544 generation sequencing to outpace optical assessments in the near future.

545

546 **Fields of application**

547 Various applications arise for the proposed method. These include studies of pollen

548 material from various origins, including plants themselves, pollinators, soil samples

549 and wind collections. The results of such assessments are of great importance in

550 identifying the diversity and specialization of plant-pollinator interaction networks

551 (Bosch et al., 2009) and also in supporting agricultural and ecological management

552 decisions (e.g. Girard et al., 2012; Odoux et al., 2012). Further, paleo-ecological and

553 climate-change associated studies investigating fossil pollens may also largely profit

554 (Bennett and Parducci, 2006).

555

556 Special attention is currently required in quality control of honey-bee products,

557 including the geographical origin, correct labeling of different varieties based on the

558 used floral resources and detection of contaminations from genetically modified

559 (GM) crops (Hemmer, 1997; Picard-Nizou et al., 1995). As pollen is naturally

560 incorporated into honey and protocols to isolate them are common usage

561 (Sowunmi, 1976), high-throughput sequencing and classification may contribute

562 largely to this endeavor by facilitating the analytical process and inclusion of

563 references from plant taxa throughout the world (Ruoff et al., 2007; Sowunmi,

564 1976).

565

566  Furthermore, the methodology may be equivalently applied to other questions not

567  only related to pollens. Other target samples are naturally occurring communities of

568  plants, (e.g. green algae), or artificially mixed probes of plant tissue fragments

569  (Schlumbaum et al., 2008). As the primers used in this study also efficiently amplify

570  fungal ITS2 sequences, ancillary information is automatically gained about this

571  group including pathogens as *Ascosphaera* spp. that may be present in collected

572  pollen samples and vectorised through harvesting flights of worker bees (Gilliam,

573  1990; White et al., 1990).

## Conclusions

575  Expert knowledge is essential to identify pollens adequately through traditional

576  light microscopy and taxonomic expertise is also often restricted to specific plant

577  groups or geographical regions. Further, mixed samples of pollens from several

578  plant origins present a problem in current palynology. With this study we evaluated

579  next-generation sequencing to approach pollen assessments through molecular

580  techniques including their bioinformatical analysis. The analytical pipeline is

581  designed for high-throughput data, but also adaptable to single sequences. It is a

582  useful technique broadening the assessment capabilities from expert labs to all

583  workgroups with access to standard molecular laboratory equipment. Further, our

584  results show that this assessment method improves the standard technique with

585  regard to taxonomical deepness, overall diversity and rectifying misidentifications.

## Data Accessibility

Sequences have been deposited at the ENA:SRA (https://www.ebi.ac.uk/ena) and are accessible under study accession number PRJEB5016. The used training set alongside installation and application notes is available for download at http://www.dna-analytics.biozentrum.uni-wuerzburg.de.

## References

Agodi, A., Barchitta, M., Grillo, A., Sciacca, S. (2006) Detection of genetically modified DNA sequences in milk from the Italian market. International Journal of Hygiene and Environmental Health, **209**(1), 81-88.

Altschul, S.F., Gish, W., Miller, W., Myers, E.W., Lipman, D.J. (1990) Basic local alignment search tool. Journal of Molecular Biology, **215**(3), 403-410.

Bagella, S., Satta, A., Floris, I., Caria, M.C., Rossetti, I., Podani, J. (2013) Effects of plant community composition and flowering phenology on honeybee foraging in Mediterranean sylvo-pastoral systems. Applied Vegetation Science.

Behm, F., von der Ohe, K., Henrich, W. (1996) Zuverlässigkeit der Pollenanalyse von Honig: Bestimmung der Pollenhäufigkeit. Deutsche Lebensmittel-Rundschau, **92**(6), 183-188.

Bennett, K.D. and Parducci, L. (2006) DNA from pollen: principles and potential. The Holocene, **16**(8), 1031-1034.

Benson, D.A., Karsch-Mizrachi, I., Lipman, D.J., Ostell, J., Sayers, E.W. (2010) GenBank. Nucleic Acids Research, **38**(suppl 1), D46-D51.

Bosch, J., Martín González, A.M., Rodrigo, A., Navarro, D. (2009) Plant–pollinator networks: adding the pollinator's perspective. Ecology Letters, **12**(5), 409-419.

Bridge, P.D., Roberts, P.J., Spooner, B.M., Panchal, G. (2003) On the unreliability of published DNA sequences. New Phytologist, **160**(1), 43-48.

Buchheim, M., Keller, A., Koetschan, C., Forster, F., Merget, B., Wolf, M. (2011) Internal Transcribed Spacer 2 (nu ITS2 rRNA) sequence-structure phylogenetics: towards an automated reconstruction of the green algal tree of life. PloS one, **6**(2).

Caporaso, J.G., Kuczynski, J., Stombaugh, J., Bittinger, K., Bushman, F.D., Costello, E.K., Fierer, N., Pena, A.G., Goodrich, J.K., Gordon, J.I. (2010) QIIME allows analysis of high-throughput community sequencing data. Nature Methods, **7**(5), 335-336.

Chen, S., Yao, H., Han, J., Liu, C., Song, J., Shi, L., Zhu, Y., Ma, X., Gao, T., Pang, X. (2010) Validation of the ITS2 region as a novel DNA barcode for identifying medicinal plant species. PloS one, **5**(1), e8613.

Davies, C.P. and Fall, P.L. (2001) Modern pollen precipitation from an elevational transect in central Jordan and its relationship to vegetation. Journal of Biogeography, **28**(10), 1195-1210.

633    Federhen, S. (2012) The NCBI taxonomy database. Nucleic Acids Research, **40**(D1),
634        D136-D143.
635    Fierer, N., Hamady, M., Lauber, C.L., Knight, R. (2008) The influence of sex,
636        handedness, and washing on the diversity of hand surface bacteria.
637        Proceedings of the National Academy of Sciences, **105**(46), 17994-17999.
638    Gilliam, M. (1990) Chalkbrood disease of honey bees, Apis mellifera, caused by the
639        fungus, *Ascosphaera apis*: A review of past and current research. Vth
640        International Colloquium on Invertebrate Pathology and Microbial Control,
641        Adelaide, Australia, 398-402.
642    Girard, M., Chagnon, M., Fournier, V. (2012) Pollen diversity collected by honey bees
643        in the vicinity of Vaccinium spp. crops and its importance for colony
644        development 1 1 This article is part of a Special Issue entitled "Pollination
645        biology research in Canada: Perspectives on a mutualism at different scales".
646        Botany, **90**(7), 545-555.
647    Harpke, D. and Peterson, A. (2008) 5.8 S motifs for the identification of pseudogenic
648        ITS regions. Botany, **86**(3), 300-305.
649    Hebert, P.D., Cywinska, A., Ball, S.L. (2003) Biological identifications through DNA
650        barcodes. Proceedings of the Royal Society of London. Series B: Biological
651        Sciences, **270**(1512), 313-321.
652    Hemmer, W. (1997) *Foods derived from genetically modified organisms and detection*
653        *methods* Agency for Biosafety Research and Assessment of Technology
654        Impacts of the Swiss Priority Programme Biotechnology of the Swiss National
655        Science Foundation.
656    Hollingsworth, P.M., Graham, S.W., Little, D.P. (2011) Choosing and using a plant
657        DNA barcode. PloS one, **6**(5), e19254.
658    Keller, A., Grimmer, G., Steffan-Dewenter, I. (2013) Diverse microbiota identified in
659        whole intact nest chambers of the red mason bee *Osmia bicornis* (Linnaeus
660        1758). PloS one, e78296.
661    Keller, A., Schleicher, T., Schultz, J., Müller, T., Dandekar, T., Wolf, M. (2009) 5.8 S-28S
662        rRNA interaction and HMM-based ITS2 annotation. Gene, **430**(1-2), 50-57.
663    Kleijn, D. and Raemakers, I. (2008) A retrospective analysis of pollen host plant use
664        by stable and declining bumble bee species. Ecology, **89**(7), 1811-1823.
665    Koetschan, C., Forster, F., Keller, A., Schleicher, T., Ruderisch, B., Schwarz, R., Muller,
666        T., Wolf, M., Schultz, J. (2010) The ITS2 Database III -sequences and
667        structures for phylogeny. Nucleic Acids Research, **38**(Database issue), D275-
668        279.
669    Kozich, J.J., Westcott, S.L., Baxter, N.T., Highlander, S.K., Schloss, P.D. (2013)
670        Development of a dual-index sequencing strategy and curation pipeline for
671        analyzing amplicon sequence data on the MiSeq Illumina sequencing
672        platform. Applied and Environmental Microbiology.
673    Kuczynski, J., Stombaugh, J., Walters, W.A., González, A., Caporaso, J.G., Knight, R.
674        (2011) Using QIIME to analyze 16S rRNA gene sequences from microbial
675        communities. Current Protocols in Bioinformatics, **10**, 1-10.17.
676    Kunin, V., Engelbrektson, A., Ochman, H., Hugenholtz, P. (2010) Wrinkles in the rare
677        biosphere: pyrosequencing errors can lead to artificial inflation of diversity
678        estimates. Environmental Microbiology, **12**(1), 118-123.

679    Lan, Y., Wang, Q., Cole, J.R., Rosen, G.L. (2012) Using the RDP classifier to predict
680        taxonomic novelty and reduce the search space for finding novel organisms.
681        PloS one, **7**(3), e32491.
682    Liu, K.-L., Porras-Alfaro, A., Kuske, C.R., Eichorst, S.A., Xie, G. (2012) Accurate, rapid
683        taxonomic classification of fungal large-subunit rRNA genes. Applied and
684        Environmental Microbiology, **78**(5), 1523-1533.
685    Marchant, R., Behling, H., Berrio, J.C., Cleef, A., Duivenvoorden, J., Hooghiemstra, H.,
686        Kuhry, P., Melief, B., Geel, B.V., Hammen, T.V.d. (2001) Mid-to Late-Holocene
687        pollen-based biome reconstructions for Colombia. Quaternary Science
688        Reviews, **20**(12), 1289-1308.
689    Metzker, M.L. (2009) Sequencing technologies—the next generation. Nature
690        Reviews Genetics, **11**(1), 31-46.
691    Mullins, J. and Emberlin, J. (1997) Sampling pollens. Journal of Aerosol Science,
692        **28**(3), 365-370.
693    Odoux, J.-F., Feuillet, D., Aupinel, P., Loublier, Y., Tasei, J.-N., Mateescu, C. (2012)
694        Territorial biodiversity and consequences on physico-chemical
695        characteristics of pollen collected by honey bee colonies. Apidologie, **43**(5),
696        561-575.
697    Parameswaran, P., Jalili, R., Tao, L., Shokralla, S., Gharizadeh, B., Ronaghi, M., Fire,
698        A.Z. (2007) A pyrosequencing-tailored nucleotide barcode design unveils
699        opportunities for large-scale sample multiplexing. Nucleic Acids Research,
700        **35**(19), e130.
701    Picard-Nizou, A., Pham-Delegue, M., Kerguelen, V., Douault, P., Marilleau, R., Olsen, L.,
702        Grison, R., Toppan, A., Masson, C. (1995) Foraging behaviour of honey bees
703        (*Apis mellifera* L.) on transgenic oilseed rape (B*rassica napus* L. var. *oleifera*).
704        Transgenic Research, **4**(4), 270-276.
705    R Development Core Team. (2010) R: A language and environment for statistical
706        computing. R Foundation for Statistical Computing Vienna Austria(01/19).
707    Roberts, R.J., Carneiro, M.O., Schatz, M.C. (2013) The advantages of SMRT
708        sequencing. Genome Biology, **14**(6), 405.
709    Ruoff, K., Luginbühl, W., Kilchenmann, V., Bosset, J.O., von der Ohe, K., von der Ohe,
710        W., Amadò, R. (2007) Authentication of the botanical origin of honey using
711        profiles of classical measurands and discriminant analysis. Apidologie, **38**(5),
712        438-452.
713    Schlumbaum, A., Tensen, M., Jaenicke-Després, V. (2008) Ancient plant DNA in
714        archaeobotany. Vegetation History and Archaeobotany, **17**(2), 233-244.
715    Schnell, I.B., Fraser, M., Willerslev, E., Gilbert, M.T.P. (2010) Characterisation of
716        insect and plant origins using DNA extracted from small volumes of bee
717        honey. Arthropod-Plant Interactions, **4**(2), 107-116.
718    Sheffield, C.S., Hebert, P.D., Kevan, P.G., Packer, L. (2009) DNA barcoding a regional
719        bee (Hymenoptera: Apoidea) fauna and its potential for ecological studies.
720        Molecular Ecology Resources, **9 Suppl s1**, 196-207.
721    Song, J., Shi, L., Li, D., Sun, Y., Niu, Y., Chen, Z., Luo, H., Pang, X., Sun, Z., Liu, C. (2012)
722        Extensive pyrosequencing reveals frequent intra-genomic variations of
723        internal transcribed spacer regions of nuclear ribosomal DNA. PloS one, **7**(8),
724        e43971.

725     Sowunmi, M. (1976) The potential value of honey in palaeopalynology and
726          archaeology. Review of Palaeobotany and Palynology, **21**(2), 171-185.
727     Spooner, D.M. (2009) DNA barcoding will frequently fail in complicated groups: an
728          example in wild potatoes. American Journal of Botany, **96**(6), 1177-1189.
729     Staatliche Naturwissenschaftliche Sammlungen Bayern. (2013) Botanischer
730          Informationsknoten Bayern.
731     Stillman, E. and Flenley, J.R. (1996) The needs and prospects for automation in
732          palynology. Quaternary Science Reviews, **15**(1), 1-5.
733     Sugita, S. (1994) Pollen representation of vegetation in Quaternary sediments:
734          theory and method in patchy vegetation. Journal of Ecology, 881-897.
735     Suzuki, M.T. and Giovannoni, S.J. (1996) Bias caused by template annealing in the
736          amplification of mixtures of 16S rRNA genes by PCR. Applied and
737          Environmental Microbiology, **62**(2), 625-630.
738     Taylor, H.R. and Harris, W.E. (2012) An emergent science on the brink of
739          irrelevance: a review of the past 8 years of DNA barcoding. Molecular Ecology
740          Resources, **12**(3), 377-388.
741     Tzedakis, P. (1993) Long-term tree populations in northwest Greece through
742          multiple Quaternary climatic cycles. Nature, **364**(6436), 437-440.
743     Valentini, A., Miquel, C., Taberlet, P. (2010) DNA barcoding for honey biodiversity.
744          Diversity, **2**(4), 610-617.
745     Valentini, A., Pompanon, F., Taberlet, P. (2009) DNA barcoding for ecologists. Trends
746          in Ecology & Evolution, **24**(2), 110-117.
747     Velasco, R., Zharkikh, A., Troggio, M., Cartwright, D.A., Cestaro, A., Pruss, D., Pindo,
748          M., FitzGerald, L.M., Vezzulli, S., Reid, J. (2007) A high quality draft consensus
749          sequence of the genome of a heterozygous grapevine variety. PloS one, **2**(12),
750          e1326.
751     Wang, Q., Garrity, G.M., Tiedje, J.M., Cole, J.R. (2007) Naive Bayesian classifier for
752          rapid assignment of rRNA sequences into the new bacterial taxonomy.
753          Applied and Environmental Microbiology, **73**(16), 5261-5267.
754     Wcislo, W.T. and Cane, J.H. (1996) Floral resource utilization by solitary bees
755          (Hymenoptera: Apoidea) and exploitation of their stored foods by natural
756          enemies. Annual Review in  Entomology, **41**, 257-286.
757     White, T., Bruns, T., Lee, S., Taylor, J. (1990) Amplification and direct sequencing of
758          fungal ribosomal RNA genes for phylogenetics. PCR-protocols: a Guide to
759          Methods and Applications, 315 - 322.
760     Wilcock, C. and Neiland, R. (2002) Pollination failure in plants: why it happens and
761          when it matters. Trends in Plant Science, **7**(6), 270-277.
762     Wilson, E.E., Sidhu, C.S., LeVan, K.E., Holway, D.A. (2010) Pollen foraging behaviour
763          of solitary Hawaiian bees revealed through molecular pollen analysis.
764          Molecular Ecology, **19**(21), 4823-4829.
765     Zhou, L.J., Pei, K.Q., Zhou, B., Ma, K.P. (2007) A molecular approach to species
766          identification of Chenopodiaceae pollen grains in surface soil. American
767          Journal of Botany, **94**(3), 477-481.
768
769

## Tables

Tab. 1: Plant families with their number of genera and number of species assessed by Next Generation Sequencing (NGS) and optical microscopy.

## Figures

Figure 1: A) Rarefaction of genera richness obtained for each honey bee sample with respect to sequencing depth. B) Plot-based comparison of pollen identification through optical microscopy and NGS. Taxonomic assignments are illustrated at the genus level. Positive identification of a taxonomic unit within a sample is indicated in the community matrix as dark gray for microscopy and light gray for NGS. Relative abundance estimations are indicated by size at two levels, i.e. >=5% (fully-filled box) and <5% (half-filled box) of total abundance within a sample. Genera misidentified in optical microscopy were combined for direct comparison and are indicated by quote marks in abbreviated form (Tar = *Taraxacum*, Sco = *Scorzoneroides*, Tan = *Tanacetum*). Availability in the reference database is indicated in the column DB. *For sample 12, three samples were taken from the same study site but different colonies. All three samples were analyzed using NGS to evaluate repeatability, yet optical microscopy was only performed for 12a.

Figure 2: Overall log-scaled relative abundance comparison of genera between the two classification strategies. Rectangles at the axes represent genera only found with one of the two sampling techniques. Pearson's correlation r = 0.86, t = 8.71, df = 26, p < 0.001***.

Figure 3: Pre-clustering evaluation: Starting from the root of the taxonomy of each sequence, every taxonomic level of the assignment was compared to its correct lineage.  The overall mean of correct assignments according to the different pre-clustering levels is presented as green dots in the figure (left scale). Similarly, each sequence was tested for erroneous levels of classification with means displayed as red squares and the scale on the right side.

800    Figure 4: Dependence of sensitivity and specificity by the bootstrap threshold.

801    Sensitivity to identify at species level is illustrated with a red and single-dashed line,

802    whereas generic identification as a red two-dashed line. Specificity is displayed as a

803    green dotted line. The harmonic mean of both species level measures is displayed by

804    a solid black curve, maximized at approximately 0.85 as the suggested optimal

805    classification threshold.