

1 **The nuclear and mitochondrial genomes of the facultatively eusocial orchid**  
2 **bee *Euglossa dilemma*.**

3  
4 **Authors:**

5  
6 Philipp Brand<sup>\*†</sup>, Nicholas Saleh<sup>\*†</sup>, Hailin Pan<sup>‡</sup>, Cai Li<sup>‡</sup>, Karen M. Kapheim<sup>§</sup>, Santiago  
7 R. Ramírez<sup>\*</sup>

8  
9 **Affiliations:**

10  
11 \* Department for Evolution and Ecology, Center for Population Biology, University  
12 of California, Davis, California 95616

13  
14 † Graduate Group in Population Biology, University of California, Davis, California  
15 95616

16  
17 ‡ China National Genebank, BGI-Shenzhen, Shenzhen, 518083, China

18  
19 § Department of Biology, Utah State University, Logan, Utah 84322

20  
21 **Data availability**

22  
23 The *E. dilemma* genome assembly *Edil\_v1.0*, the annotation, and the original gene set  
24 *Edil\_OGS\_v1.0* are available for download via NCBI [XXX], Beebase [XXX], the i5k NAL  
25 workspace [xxx], and the Ramirez Lab website [URL]. The raw reads are available at  
26 the NCBI Sequence Read Archive [XXX]. The published raw transcriptome sequence  
27 reads are available at the NCBI Sequence Read Archive [SRA: SRX765918].

47 **Running Title:**

48

49 The genomes of the orchid bee *Euglossa dilemma*

50

51 **Key Words:**

52

53 whole-genome assembly; corbiculate bee; orchid bee; invasive species;  
54 mitochondrial genome

55

56 **Corresponding Author:**

57

58 Philipp Brand

59 Department of Evolution & Ecology

60 University of California, Davis

61 One Shields Ave

62 Davis, CA 95616

63

64 Office: (530) 752-7614

65 pbrand@ucdavis.edu

66

67

68

69

70

71

72

73

74

75

76

77

78

79

80

81

82

83

84

85

86

87

88

89

90

91

92

93 **Abstract**

94

95 Bees provide indispensable pollination services to both agricultural crops and wild  
96 plant populations, and several species of bees have become important models for  
97 the study of learning and memory, plant-insect interactions and social behavior.  
98 Orchid bees (Apidae: Euglossini) are especially important to the fields of pollination  
99 ecology, evolution, and species conservation. Here we report the nuclear and  
100 mitochondrial genome sequences of the orchid bee *Euglossa dilemma*. *Euglossa*  
101 *dilemma* was selected because it is widely distributed, highly abundant, and it was  
102 recently naturalized in the southeastern United States. We provide a high-quality  
103 assembly of the 3.3 giga-base genome, and an official gene set of 15,904 gene  
104 annotations. We find high conservation of gene synteny with the closely related  
105 honey bee. This genomic resource represents the first draft genome of the orchid  
106 bee genus *Euglossa*, and the first draft orchid bee mitochondrial genome, thus  
107 representing a valuable resource to the research community.

108

109 **Introduction**

110

111 Bees (Apoidea) are important models for the study of learning and memory (Menzel  
112 and Muller 1996), plant-insect interactions (Doetterl and Vereecken 2010) and the  
113 evolution of social behavior (Nowak *et al.* 2010; Woodard *et al.* 2011; Kapheim *et al.*  
114 2015). Among the >20,000 bee species worldwide, lineages have evolved varied  
115 degrees of specialization on floral resources such as pollen, resins, and oils (Wcislo  
116 and Cane 2003; Michener 2007; Litman *et al.* 2011). These relationships are wide-  
117 ranging and have substantial impact on the health and function of natural and  
118 agricultural systems (Klein *et al.* 2007). Furthermore, several transitions from an  
119 ancestral solitary to a derived eusocial behavior have occurred within bees  
120 (Danforth 2002; Cardinal and Danforth 2011; Branstetter *et al.* 2017). Thus, bees  
121 provide unique opportunities to investigate the genetic underpinnings of multiple  
122 major ecological and evolutionary transitions. The repeated evolution of different  
123 behavioral phenotypes in bees, including foraging and social behavior, provides a  
124 natural experiment that allows for the determination of general as well as species-  
125 specific molecular genomic changes underlying phenotypic transitions. In order to  
126 capitalize on this potential, whole-genome sequences of a divergent array of bee  
127 species with different life histories are needed (Kapheim *et al.* 2015).

128

129 Orchid bees (Apidae; Euglossini) are among the most important pollinators of  
130 thousands of diverse neotropical plant species (Ramírez *et al.* 2002). While female  
131 orchid bees collect nectar, pollen and resin for nest construction and provisioning,  
132 male bees collect perfume compounds from floral and non-floral sources (Vogel  
133 1966; Whitten *et al.* 1993; Eltz *et al.* 1999; Roubik and Hanson 2004). These volatile  
134 compounds are used to concoct a species-specific perfume blend that is  
135 subsequently used during courtship, presumably to attract conspecific females. This  
136 unique male scent-collecting behavior has recently been examined in a broad array  
137 of molecular ecological and evolutionary studies, focusing on phenotypic evolution,

138 chemical communication, plant-insect mutualisms, and speciation (Eltz *et al.* 2008;  
139 2011; Ramírez *et al.* 2011; Brand *et al.* 2015).

140

141 While most of the approximately 220 species of orchid bees appear to be solitary,  
142 several species have transitioned to living in coordinated social groups (Garófalo  
143 1985; Pech *et al.* 2008; Augusto and Garófalo 2009). Female *Euglossa dilemma*  
144 individuals, for example, can either form solitary nests and provision their own  
145 brood cells, or live in small groups where daughters remain in their natal nest and  
146 help their mother rear her offspring, instead of dispersing to found their own nest.  
147 The social *E. dilemma* nests (similar to the closely related *E. viridissima*; Pech *et al.*  
148 2008) exhibit true division of labor, with subordinate daughters foraging for  
149 resources and the reproductively dominant mothers laying eggs (Saleh & Ramirez  
150 pers. obs.). This facultative eusocial behavior represents an early stage in social  
151 evolution and makes *E. dilemma* well-suited for studying the genomic changes  
152 underlying the transition from solitary to eusocial behavior. While other facultative  
153 eusocial species evolved throughout the bee lineage, orchid bees have a unique  
154 taxonomic position (Cardinal and Danforth 2011). Orchid bees are part of the  
155 corbiculate bees, together with the honey bees, bumblebees and stingless bees,  
156 three obligately eusocial bee lineages (Figure 1a). As the sister group to all other  
157 corbiculate bee lineages (Romiguier *et al.* 2015; Peters *et al.* 2017; Branstetter *et al.*  
158 2017), orchid bees may provide key insights into the early stages of eusociality and  
159 the possible evolutionary trajectories that led to the obligate eusocial behavior  
160 exhibited by honey bees.

161

162 Here we present the draft genome of the orchid bee species *Euglossa dilemma*. Using  
163 a combined paired-end and mate-pair library sequencing approach, we assembled  
164 18% of the predicted 3.2Gb genome, and annotated a high-quality gene set including  
165 15,904 genes. In addition, we reconstructed three quarters of the mitochondrial  
166 genome with the help of transcriptome data, representing the first orchid bee  
167 mitogenome. These genomic resources will facilitate the genetic study of  
168 outstanding ecological and evolutionary questions, such as the evolution of resource  
169 preferences and the evolution of eusociality. Moreover, it provides an important  
170 genomic resource for an endangered group of crucial neotropical pollinators, that  
171 are of specific concern for conservation biologists (Zimmermann *et al.* 2011; Suni  
172 and Brosi 2012; Suni 2016; Soro *et al.* 2016).

173

## 174 **Materials and Methods**

175

### 176 **Genome sequencing and assembly**

177

178 *Nuclear Genome.* Sequencing of the *E. dilemma* genome was based on six haploid  
179 male individuals collected at Fern Forest Nature Center in Broward County, FL in  
180 February 2011. This population was chosen due to its low nucleotide diversity  
181 resulting from a bottleneck during a single introduction to Southern Florida about  
182 15 years ago (Skov and Wiley 2005; Pemberton and Wheeler 2006; Zimmermann *et al.*  
183 *et al.* 2011). DNA was extracted from each bee independently and used for the

184 construction of four paired-end (two 170bp and 500bp libraries, respectively) and  
185 four mate-pair (two 2kb and 5kb libraries, respectively) sequencing libraries. Next,  
186 the paired-end libraries were sequenced in 90 cycles and the mate-pair libraries for  
187 49 cycles on an Illumina HiSeq2000. The resulting sequence data was run through  
188 fastuniq v1.1 (Xu *et al.* 2012) to remove PCR duplicates and quality trimmed using  
189 trim\_galore v0.3.7 (Babraham Bioinformatics). Subsequently, reads were used for *de*  
190 *novo* assembly with ALLPATHS-LG v51750 (Gnerre *et al.* 2011) and Soap-denovo2  
191 (Luo *et al.* 2012) with varying settings. Gaps within scaffolds were closed using  
192 GapCloser v1.12 (Luo *et al.* 2012) for each assembly. ALLPATHS-LG with default  
193 settings resulted in the highest-quality assembly, based on assessments of  
194 annotation completeness (see below). This assembly (*E. dilemma* genome assembly  
195 v1.0) was used for all subsequent analyses. All other assemblies were excluded from  
196 analysis, but are available upon request.

197

198 The pre-processed reads were used for k-mer based genome size estimates. We  
199 used ALLPATHS-LG to produce and analyze the k-mer frequency spectrum (k=25).  
200 Genome size was estimated on the basis of the consecutive length of all reads  
201 divided by the overall sequencing depth as  $(N \times (L - K + 1) - B)/D = G$ , where  $N$  is  
202 the total number of reads,  $L$  is the single-read length,  $K$  is the k-mer length,  $B$  is the  
203 total count of low-frequency (frequency  $\leq 3$ ) k-mers that are most likely due to  
204 sequencing errors,  $D$  is the k-mer depth estimated from the k-mer frequency  
205 spectrum, and  $G$  is the genome size. In addition, we used the ALLPATHS-LG k-mer  
206 frequency spectrum to predict the repetitive fraction of the genome.

207

208 The quality of the genome assembly was assessed using standard N statistics, and  
209 assembly completeness as measured by the CEGMA v2.5 (Parra *et al.* 2007) and  
210 BUSCO v1.1 (Simão *et al.* 2015) pipelines. CEGMA was run in default mode, whereas  
211 BUSCO was run with the arthropoda\_odb9 OrthoDB database (Zdobnov *et al.* 2017)  
212 in genome mode.

213

214 We estimated the mean per-base genome coverage on the basis of the pre-processed  
215 reads and the estimated genome size as  $\frac{\sum_{i=1}^4 (R_i * L_i)}{G} = C$ , where  $R$  is the number of  
216 reads and  $L$  the mean read length of sequence library  $i$ ,  $G$  is the estimated genome  
217 size and  $C$  the resulting per-base coverage.

218

219 *Mitogenome.* Initial attempts to reconstruct the mitochondrial genome from our  
220 whole-genome shotgun sequencing reads were only partially successful, due to high  
221 sequence variability of sequencing reads with similarity to mitochondrial loci (data  
222 not shown). In addition, we have observed that the amplification of mitochondrial  
223 DNA in standard polymerase chain reactions (PCR) leads to a high level of  
224 polymorphic sites in *E. dilemma* and other orchid bees (Brand & Ramírez pers. obs.).  
225 Together, this suggests the presence of nuclear copies of the mitochondrial genome  
226 (NUMTs) that interfere with the assembly process and PCR amplification. Therefore,  
227 we used available *E. dilemma* transcriptome assemblies in order to reconstruct the  
228 mitochondrial genome from cDNA (Brand *et al.* 2015). In order to find

229 mitochondrial genes in the transcriptome assembly of Brand et al. 2015, we used  
230 blastx with the honey bee mitochondrial genome as query (Crozier and Crozier  
231 1993) and an E-value cutoff of  $10E-12$  (Altschul *et al.* 1990; Camacho *et al.* 2009).  
232 The contigs and scaffolds that were detected with this approach were annotated  
233 following Dietz *et al.* 2015. Briefly, we performed tblastn and blastx searches with  
234 protein coding genes and rRNA genes of the honey bee mitochondrial genome,  
235 respectively. All hits were used for manual gene annotation using Geneious v8.0.5  
236 (Biomatters Ltd. 2012). Since the recovered mitochondrial mRNA scaffolds  
237 contained more than one gene, we searched and annotated intergenic tRNAs using  
238 ARWEN 1.2.3 (Laslett and Canbäck 2008) and tRNAscan-SE 1.21 (Lowe and Eddy  
239 1997).

240

## 241 **Genome annotation**

242

243 *Gene annotation.* Genes were annotated based on sequence homology and *de novo*  
244 gene predictions. The homology approach was based on the recently updated high-  
245 quality official gene set of the honey bee (OGS v3.2; Elsiek *et al.* 2014). All honey bee  
246 original gene set (OGS) proteins were used in initial tblastn searches against all *E.*  
247 *dilemma* scaffolds with an E-value cutoff of  $10E-4$ . Proteins with a hit covering  $\geq$   
248 50% of the honeybee protein query were selected for accurate exon-intron  
249 boundary prediction for each scaffold using exonerate v2.42.1 (Slater and Birney  
250 2005) with the minimum fraction of the possible optimal similarity per query set to  
251 35%. In a second round, genes not annotated under the previous settings were  
252 rerun with minimum similarity set to 15%. In the case of multiple annotations, we  
253 discarded all but one annotation with the best hit to the honeybee OGS (based on  
254 completeness and similarity). This approach proved feasible due to the close  
255 relatedness of *E. dilemma* and the honey bee. For *de novo* gene prediction we used  
256 Augustus (Stanke *et al.* 2008) and SNAP (Korf 2004) trained on the honey bee, with  
257 the *E. dilemma* genome masked for repetitive regions (See below) as input. Only  
258 genes predicted by both programs were taken into account. Gene predictions with  
259  $\geq 85\%$  sequence similarity to each other were discarded, to prevent the inclusion of  
260 putative unmasked transposable element derived genes in the official gene set. *De*  
261 *novo* predictions were added to the *E. dilemma* OGS if not annotated by the  
262 homology-based approach.

263

264 *Repetitive element annotation.* Repetitive elements including tandem repeats,  
265 nuclear copies of the mitochondrial genome (NUMTs), and transposable elements  
266 (TEs) were annotated using multiple methods.

267

268 *Tandem repeats.* We searched for micro- and mini-satellites (1–6 bp and 7–1000 bp  
269 motif length, respectively) in all scaffolds using Phobos 3.3.12 (Mayer 2010). We  
270 performed two independent runs for each class of tandem repeats with Phobos  
271 parameter settings following Leese et al. 2012 (gap score and mismatch score set to  
272 -4 and a minimum repeat score of 12; Leese *et al.* 2012).

273

274 *NUMTs*. We annotated NUMTs using blastn runs with the partial mitochondrial  
275 genome (see above) as query and an E-value cutoff of 10E-4 as used in NUMT  
276 analyses of other insect genomes (Pamilo *et al.* 2007). This approach allowed us to  
277 find NUMTs with medium to high similarity to the actual transcriptome-based  
278 mitochondrial genome.

279

280 *TEs*. In order to annotate TEs, we used Repeatmasker v4.0.5 (Smit *et al.* 2016) with  
281 Crossmatch v. 0.990329 as search engine in sensitive mode. The Repbase  
282 invertebrate database v21.12 (Jurka 2000; Bao *et al.* 2015) was used as TE  
283 reference, due to the lack of a bee-specific database.

284

## 285 **Genome structure**

286

287 To analyze genome structure, we compared the genome wide gene synteny of *E.*  
288 *dilemma* and the honey bee. We used the genomic locations of homologous genes (as  
289 determined above) of the honey bee and *E. dilemma* scaffolds of at least 100kb  
290 length to build haplotype blocks with a minimum length of 1kb. Haplotype blocks  
291 included the entire gene span as well as intergenic regions whenever two or more  
292 adjacent genes were homologous in both species. We discarded gene annotations  
293 from downstream analysis that were recovered as homologous to multiple genomic  
294 locations in either species. Furthermore, we excluded *E. dilemma* genes that were  
295 recovered as homologous to honey bee scaffolds belonging to unknown linkage  
296 groups.

297

## 298 **Data availability**

299

300 The *E. dilemma* genome assembly *Edil\_v1.0*, the annotation, and the original gene set  
301 *Edil\_OGSv1.0* are available for download via NCBI [XXX], Beebase [XXX] (Elsik *et al.*  
302 2016), the i5k NAL workspace [xxx] (i5K Consortium 2013), and the Ramirez Lab  
303 website [URL]. The raw reads are available at the NCBI Sequence Read Archive  
304 [XXX]. The published raw transcriptome sequence reads are available at the NCBI  
305 Sequence Read Archive [SRA: SRX765918] (Brand *et al.* 2015).

306

## 307 **Results and Discussion**

308

### 309 **Whole-genome assembly**

310

311 The *E. dilemma* genome assembly resulted in 22,698 scaffolds with an N50 scaffold  
312 length of 144Kb and a total length of 588Mb (Table 1). This represents 18% of the k-  
313 mer based estimated genome size of 3.2Gb. Of all sequence reads, 68% aligned to  
314 the genome assembly, of which 56% aligned more than once. Further, the k-mer  
315 frequency spectrum based on all sequencing reads was strongly positively skewed  
316 indicating the presence of highly repetitive sequences in the read set (Figure 1b).  
317 Based on the k-mer frequency spectrum, 87.7% of the genome was estimated to be

318 repetitive. This suggests that the genome of *E. dilemma* consists largely of highly  
319 repetitive sequences, explaining the low consecutive assembly length and the high  
320 assembly fragmentation. The mean per-base coverage was estimated to be  
321 comparatively low in comparison to previous bee genome assemblies, with 19.7x  
322 based on the pre-processed reads and estimated genome size (Kocher *et al.* 2013;  
323 Kapheim *et al.* 2015). Total genomic GC content was 39.9%, and thus similar to  
324 previously sequenced bee genomes ranging between 32.7% and 41.5% (Table 1)  
325 (Kocher *et al.* 2013; Elsik *et al.* 2014; Kapheim *et al.* 2015).

326  
327 Despite the fragmentation of the genome assembly representing less than 20% of  
328 the estimated genome size, CEGMA analysis revealed complete assemblies of 231  
329 out of 248 core eukaryotic genes (93.2% completeness). Similarly, BUSCO analysis  
330 revealed that 1007 out of 1066 highly conserved arthropod genes were completely  
331 assembled (94.4% completeness). Our gene prediction approach generated a  
332 comprehensive official gene set including 15,904 protein-coding genes (Table 1). Of  
333 these gene models, 11,139 were derived from homology-based predictions,  
334 representing 73% of the 15,314 honey bee genes used for annotation. These  
335 annotations are well within or exceeding previous bee genome assemblies, and are  
336 similar to those reported for the other orchid bee genome available (Table 1)  
337 (Kocher *et al.* 2013; Elsik *et al.* 2014; Park *et al.* 2015; Sadd *et al.* 2015; Kapheim *et*  
338 *al.* 2015).

339  
340 The CEGMA and BUSCO analysis and the gene annotation results suggest that the  
341 gene-coding fraction of the *E. dilemma* genome was properly assembled, despite the  
342 large estimated genome size and comparatively low per-base sequencing coverage.  
343 Fragmentation of the assembly is thus likely to be primarily the result of highly  
344 repetitive genomic elements, and less the result of low coverage. Overall, the results  
345 suggest that our approach was sufficient to produce a high quality official gene set.  
346 Due to the chosen homology-based approach, the majority of annotated genes in the  
347 official gene set has known homology to honey bee genes (Table S1), which will  
348 greatly facilitate genome-wide expression studies including gene ontology analyses  
349 and comparative gene set analyses among insects.

### 350 351 **Mitochondrial Genome assembly**

352 The recently published transcriptome assembly used for the reconstruction of the  
353 mitochondrial genome contained four scaffolds between 1,222bp and 4,188bp long  
354 with a total consecutive length of 11,128bp (Figure 2). This corresponds to about  
355 75% of the estimated length of the mitochondrial genome, based on other  
356 corbiculate bee species (Crozier and Crozier 1993; Cha *et al.* 2007). The *E. dilemma*  
357 mitogenome fragments contained 5 out of 22 tRNAs, 11 out of 13 protein coding  
358 genes of which two were only partially recovered, and the 16S rRNA gene. Within  
359 scaffolds all genes showed the known hymenopteran gene order and orientation,  
360 while the orientation of the 5 tRNAs detected was identical to those in the honey bee  
361 (Crozier and Crozier 1993; Cha *et al.* 2007). Attempts to complete the mitochondrial  
362 genome using the nuclear genome assembly yielded no improvement of the  
363 assembly (data not shown).

364

365 The high success in mitochondrial gene reconstruction is likely due to the nature of  
366 the analyzed transcriptome data. Short intergenic regions as well as polycistronic  
367 mitochondrial mRNA likely lead to the assembly of multiple genes into single  
368 scaffolds. The A-T rich region is completely missing as well as the ND2 and 12S  
369 rRNA genes flanking the region in insect mitogenomes. This unrecovered region also  
370 contains a high number of tRNAs in the honeybee, which could explain the low  
371 number of recovered tRNAs in *E. dilemma*. While the partial mitochondrial genome  
372 assembly is only 75% complete, it represents the first mitogenome for the group of  
373 orchid bees and will thus be a valuable resource for future phylogenetic analyses  
374 within the lineage and between more distantly related bee taxa.

375

### 376 **Repetitive elements**

377

378 *Tandem Repeats.* We detected 76,001 microsatellite loci with a consecutive length of  
379 2,291,067 bp. Minisatellites with motif lengths from 7bp to 1000bp were less  
380 numerous in the genome (67,323 loci), totaling 13,343,515bp. Accordingly, tandem  
381 repeats represent 3.86% of the genome assembly, suggesting that they contribute  
382 only a small proportion to the overall genome size (Figure 1c).

383

384 *NUMTs.* We detected fragments with similarity to the draft mitochondrial genome  
385 on 129 scaffolds totaling a length of 150,670 bp. The fragments had a mean length of  
386 764.8 bp and a mean similarity of 91.5% to the mitogenome. This suggests that  
387 these fragments are not derived from the mitochondrial genome and represent  
388 actual NUMTs. A total of 39 scaffolds carried multiple fragments with high similarity  
389 to the mitogenome with a concatenated length of up to 6566bp, suggesting that  
390 respective NUMTs might have originated from larger fragments of the mitogenome.  
391 In total, only 0.04% of the whole-genome assembly had hits to the mitogenome  
392 (Figure 1c). This is likely an underestimate, due to the incompleteness of the  
393 reconstructed mitochondrial genome. Nevertheless, NUMTs likely represent only a  
394 small fraction of the whole nuclear *E. dilemma* genome. Previous analyses have  
395 shown a high density of NUMTs in the honey bee in comparison to other insect  
396 genomes totaling about 0.1% of the overall genome size (Pamilo *et al.* 2007).  
397 Accordingly, given the NUMT content detected in *E. dilemma*, it is possible that a  
398 comparatively high NUMT density is a common feature of corbiculate bee genomes.

399

400 *TEs.* In our RepeatMasker analysis we detected 47,553 interspersed repeats with a  
401 cumulative length of 7,747,824bp. This corresponds to 1.82% of the genome  
402 assembly (Figure 1c). This low number is surprising, since large genome sizes as in  
403 *E. dilemma* have been mainly associated with TE activity and content in genomes  
404 from unicellular eukaryotes to complex multicellular organisms including plants,  
405 invertebrates and vertebrates (Kidwell 2002). However, TEs are fast evolving and  
406 highly specific to their host lineages, which leads to large underestimates of genomic  
407 TE content in previously unstudied lineages (Chalopin *et al.* 2015; Platt *et al.* 2016).  
408 The only bee repeat content included in the Repbase database used for TE  
409 annotation is the honey bee, a species with a comparatively small genome (0.23Gb)

410 and low TE diversity and content (Weinstock *et al.* 2006; Kapheim *et al.* 2015). Thus,  
411 the low percentage of TEs detected in the *E. dilemma* genome is very likely an  
412 underestimate of the actual density. *De novo* TE reconstructions using the genomic  
413 resources presented here should be performed in the future to provide a better  
414 estimate of the actual TE density in the *E. dilemma* genome.

415

## 416 **Genome structure**

417

418 Of the 22,698 *E. dilemma* scaffolds, 580 were at least 100kb in length and used for  
419 synteny analysis with the honey bee genome. A total of 356 of these scaffolds  
420 carried at least one gene annotation with known homology to the honey bee, and  
421 329 of these *E. dilemma* scaffolds were homologous to honey bee scaffolds with  
422 known linkage group (LG) association (Table S1). Of these scaffolds, 272 (83%)  
423 showed  $\geq 95\%$  syntenic homology to a single honey bee LG (Figure 1d). Overall, the  
424 detected syntenic linkage blocks cover 222MB of scaffold length with homology to  
425 the honeybee, representing 85% of the 329 filtered scaffolds. Mean syntenic block  
426 length is 206,033bp for *E. dilemma* (min: 1017bp; max: 4,458,123bp; median:  
427 47,154bp) and 143,221.4bp for the honey bee (min: 1029bp; max: 2,205,566bp;  
428 median: 41,726bp). This suggests that the genomic architecture is very similar  
429 between *E. dilemma* and the honey bee, representing a high level of conservation  
430 during the 80 million years since the two lineages diverged. Further, our results  
431 support a recent comparative analysis of the honey bee and the bumblebee  
432 genomes, which revealed high conservation of genomic synteny (Stolle *et al.* 2011).  
433 Together, these results support a general pattern of surprisingly slow genome  
434 evolution in gene coding regions in corbiculate bees, independent of the fraction of  
435 repetitive genome content.

436

## 437 **Conclusion**

438

439 The genome assembly of the orchid bee *E. dilemma* that we present here is of high  
440 quality, despite its large genome size (estimated to be 3.3Gb). The 15,904 gene  
441 annotations provide a comprehensive set of genes with known homology to the  
442 honey bee, facilitating future gene ontology and functional genomic analyses. While  
443 we were unable to annotate the mostly repetitive majority of the genome assembly  
444 with our approach, the provided sequence reads will be useful for future analyses of  
445 repetitive genetic elements in the genome. The nuclear and mitochondrial draft  
446 genomes represent a valuable genomic resource for the community of bee  
447 geneticists. This genomic resource will likely prove valuable in genetic and  
448 functional genomic analyses dealing with the ecology, evolution, and conservation of  
449 orchid bees. Furthermore, the genome of the facultatively eusocial *E. dilemma* will  
450 be helpful in the study of the evolution of eusociality, due to its taxonomic  
451 placement as the sister lineage to the three obligately eusocial corbiculate bee tribes  
452 including stingless bees, bumblebees, and honey bees.

453

## 454 **Acknowledgements**

455 We would like to thank Gene E. Robinson, Guojie Zhang, and the BGI for financial

456 support for genome sequencing. P.B. was supported by a fellowship of the German  
457 academic exchange service (Deutscher Akademischer Austauschdienst, DAAD) for  
458 parts of the project.

459

## 460 **References**

461

462 Altschul, S. F., W. Gish, W. Miller, E. W. Myers, and D. J. Lipman, 1990 Basic local  
463 alignment search tool. *J. Mol. Biol.* 215: 403–410.

464 Augusto, S. C., and C. A. Garófalo, 2009 Bionomics and sociological aspects of  
465 *Euglossa fimbriata* (Apidae, Euglossini). *Genet. Mol. Res.* 8: 525–538.

466 Bao, W., K. K. Kojima, and O. Kohany, 2015 Repbase Update, a database of repetitive  
467 elements in eukaryotic genomes. *Mobile DNA* 2015 6:1 6: 11.

468 Brand, P., S. R. Ramírez, F. Leese, J. J. Quezada-Euan, R. Tollrian *et al.*, 2015 Rapid  
469 evolution of chemosensory receptor genes in a pair of sibling species of orchid  
470 bees (Apidae: Euglossini). *BMC Evol. Biol.* 15: 176.

471 Branstetter, M. G., B. N. Danforth, J. P. Pitts, B. C. Faircloth, P. S. Ward *et al.*, 2017  
472 Phylogenomic Insights into the Evolution of Stinging Wasps and the Origins of  
473 Ants and Bees. *Current Biology* 27: 1019–1025.

474 Camacho, C., G. Coulouris, V. Avagyan, N. Ma, J. Papadopoulos *et al.*, 2009 BLAST plus  
475 : architecture and applications. *BMC Bioinformatics* 10:.

476 Cardinal, S., and B. N. Danforth, 2011 The Antiquity and Evolutionary History of  
477 Social Behavior in Bees (C. S. Moreau, Ed.). *PLoS ONE* 6:.

478 Cha, S. Y., H. J. Yoon, E. M. Lee, M. H. Yoon, J. S. Hwang *et al.*, 2007 The complete  
479 nucleotide sequence and gene organization of the mitochondrial genome of the  
480 bumblebee, *Bombus ignitus* (Hymenoptera: Apidae). *Gene* 392: 206–220.

481 Chalopin, D., M. Naville, F. Plard, D. Galiana, and J.-N. Volff, 2015 Comparative  
482 Analysis of Transposable Elements Highlights Mobilome Diversity and Evolution  
483 in Vertebrates. *Genome Biol Evol* 7: 567–580.

484 Crozier, R. H., and Y. C. Crozier, 1993 The mitochondrial genome of the honeybee  
485 *Apis mellifera*: complete sequence and genome organization. *Genetics* 133: 97–  
486 117.

487 Danforth, B. N., 2002 Evolution of sociality in a primitively eusocial lineage of bees.  
488 *Proc Natl Acad Sci USA* 99: 286–290.

489 Dietz, L., P. Brand, L. M. Eschner, and F. Leese, 2015 The mitochondrial genomes of  
490 the caddisflies *Sericostoma personatum* and *Thremma gallicum* (Insecta:  
491 Trichoptera). *Mitochondrial DNA* 1–2.

- 492 Doetterl, S., and N. J. Vereecken, 2010 The chemical ecology and evolution of bee-  
493 flower interactions: a review and perspectives. *Canadian Journal of Zoology* 88:  
494 668–697.
- 495 Elsik, C. G., A. Tayal, C. M. Diesh, D. R. Unni, M. L. Emery *et al.*, 2016 Hymenoptera  
496 Genome Database: integrating genome annotations in HymenopteraMine.  
497 *Nucleic Acids Res.* 44: D793–D800.
- 498 Elsik, C. G., K. C. Worley, A. K. Bennett, M. Beye, F. Camara *et al.*, 2014 Finding the  
499 missing honey bee genes: lessons learned from a genome upgrade. *BMC*  
500 *Genomics* 15: 86.
- 501 Eltz, T., F. Fritsch, and J. R. Pech, 2011 Characterization of the orchid bee *Euglossa*  
502 *viridissima* (Apidae: Euglossini) and a novel cryptic sibling species, by  
503 morphological, chemical, and genetic characters. *Zool. J. Linn. Soc.*
- 504 Eltz, T., W. M. Whitten, D. W. Roubik, and K. E. Linsenmair, 1999 Fragrance  
505 Collection, Storage, and Accumulation by Individual Male Orchid Bees. *J. Chem.*  
506 *Ecol.* 25: 157–176.
- 507 Eltz, T., Y. Zimmermann, C. Pfeiffer, J. R. Pech, R. Twele *et al.*, 2008 An olfactory shift  
508 is associated with male perfume differentiation and species divergence in orchid  
509 bees. *Curr. Biol.* 18: 1844–1848.
- 510 Garófalo, C. A., 1985 Social Structure of *Euglossa cordata* Nests (Hymenoptera:  
511 Apidae: Euglossini). *Entomologia Generalis* 11: 77–83.
- 512 Gnerre, S., I. MacCallum, D. Przybylski, F. J. Ribeiro, J. N. Burton *et al.*, 2011 High-  
513 quality draft assemblies of mammalian genomes from massively parallel  
514 sequence data. *Proc. Natl. Acad. Sci. U.S.A.* 108: 1513–1518.
- 515 i5K Consortium, 2013 The i5K Initiative: Advancing Arthropod Genomics for  
516 Knowledge, Human Health, Agriculture, and the Environment. *Journal of*  
517 *Heredity* 104: 595–600.
- 518 Jurka, J., 2000 Repbase Update: a database and an electronic journal of repetitive  
519 elements. *Trends in Genetics* 16: 418–420.
- 520 Kapheim, K. M., H. Pan, C. Li, S. L. Salzberg, D. Puiu *et al.*, 2015 Genomic signatures of  
521 evolutionary transitions from solitary to group living. *Science* aaa4788.
- 522 Kidwell, M. G., 2002 Transposable elements and the evolution of genome size in  
523 eukaryotes. *Genetica* 115: 49–63.
- 524 Klein, A.-M., B. E. Vaissière, J. H. Cane, I. Steffan-Dewenter, S. A. Cunningham *et al.*,  
525 2007 Importance of pollinators in changing landscapes for world crops.  
526 *Proceedings of the Royal Society B: Biological Sciences* 274: 303–313.

- 527 Kocher, S. D., C. Li, W. Yang, H. Tan, S. V. Yi *et al.*, 2013 The draft genome of a socially  
528 polymorphic halictid bee, *Lasioglossum albipes*. *Genome Biology* 14: R142.
- 529 Korf, I., 2004 Gene finding in novel genomes. *BMC Bioinformatics* 5:.
- 530 Laslett, D., and B. Canbäck, 2008 ARWEN: a program to detect tRNA genes in  
531 metazoan mitochondrial nucleotide sequences. *Bioinformatics* 24: 172–175.
- 532 Leese, F., P. Brand, A. Rozenberg, C. Mayer, S. Agrawal *et al.*, 2012 Exploring  
533 Pandora's box: potential and pitfalls of low coverage genome surveys for  
534 evolutionary biology. (B. J. Mans, Ed.). *PLoS ONE* 7: e49202.
- 535 Litman, J. R., B. N. Danforth, C. D. Eardley, and C. J. Praz, 2011 Why do leafcutter bees  
536 cut leaves? New insights into the early evolution of bees. *Proc. Biol. Sci.* 278:  
537 3593–3600.
- 538 Lowe, T. M., and S. R. Eddy, 1997 tRNAscan-SE: A program for improved detection of  
539 transfer RNA genes in genomic sequence. *Nucleic Acids Res.* 25: 955–964.
- 540 Luo, R., B. Liu, Y. Xie, Z. Li, W. Huang *et al.*, 2012 SOAPdenovo2: an empirically  
541 improved memory-efficient short-read de novo assembler. *GigaScience* 2016 5:1  
542 1: 18.
- 543 Mayer, C., 2010 Phobos Version 3.3.12. A tandem repeat search program. 20 pp.
- 544 Menzel, R., and U. Muller, 1996 Learning and Memory in Honeybees: From Behavior  
545 to Neural Substrates. *Annu. Rev. Neurosci.* 19: 379–404.
- 546 Michener, C. D., 2007 *The Bees of the World*. 2nd. Ed. Johns Hopkins.
- 547 Nowak, M. A., C. E. Tarnita, and E. O. Wilson, 2010 The evolution of eusociality. 466:  
548 1057–1062.
- 549 Pamilo, P., L. Viljakainen, and A. Vihavainen, 2007 Exceptionally High Density of  
550 NUMTs in the Honeybee Genome. *Mol Biol Evol* 24: 1340–1346.
- 551 Park, D., J. W. Jung, B.-S. Choi, M. Jayakodi, J. Lee *et al.*, 2015 Uncovering the novel  
552 characteristics of Asian honey bee, *Apis cerana*, by whole genome sequencing.  
553 *BMC Genomics* 16:.
- 554 Parra, G., K. Bradnam, and I. Korf, 2007 CEGMA: a pipeline to accurately annotate  
555 core genes in eukaryotic genomes. *Bioinformatics* 23: 1061–1067.
- 556 Pech, M. E. C., W. de J May-Itzá, L. A. M. Medina, and J. J. G. Quezada-Euan, 2008  
557 Sociality in *Euglossa (Euglossa) viridissima* Friese (Hymenoptera, Apidae,  
558 *Euglossini*). 55: 428–433.
- 559 Pemberton, R. W., and G. S. Wheeler, 2006 Orchid bees don't need orchids: evidence

- 560 from the naturalization of an orchid bee in Florida. *Ecology* 87: 1995–2001.
- 561 Peters, R. S., L. Krogmann, C. Mayer, A. Donath, S. Gunkel *et al.*, 2017 Evolutionary  
562 History of the Hymenoptera. *Current Biology* 0:
- 563 Platt, R. N., L. Blanco-Berdugo, and D. A. Ray, 2016 Accurate Transposable Element  
564 Annotation Is Vital When Analyzing New Genome Assemblies. *Genome Biol Evol*  
565 8: 403–410.
- 566 Ramírez, S., R. L. Dressler, and M. Ospina, 2002 Abejas euglosinas (Hymenoptera:  
567 Apidae) de la Región Neotropical: Listado de especies con notas sobre su  
568 biología. 3: 7.
- 569 Ramírez, S. R., T. Eltz, M. K. Fujiwara, G. Gerlach, B. Goldman-Huertas *et al.*, 2011  
570 Asynchronous Diversification in a Specialized Plant-Pollinator Mutualism.  
571 *Science* 333: 1742–1746.
- 572 Romiguier, J., S. A. Cameron, S. H. Woodard, B. J. Fischman, L. Keller *et al.*, 2015  
573 Phylogenomics controlling for base compositional bias reveals a single origin of  
574 eusociality in corbiculate bees. *Mol Biol Evol* msv258.
- 575 Roubik, D. W., and P. E. Hanson, 2004 *Orchid bees of tropical America: biology and*  
576 *field guide*. Instituto Nacional de Biodiversidad (INBio), Santo Domingo De  
577 Heredia.
- 578 Sadd, B. M., S. M. Barribeau, G. Bloch, D. C. de Graaf, P. Dearden *et al.*, 2015 The  
579 genomes of two key bumblebee species with primitive eusocial organization.  
580 *Genome Biology* 16: 76.
- 581 Simão, F. A., R. M. Waterhouse, P. Ioannidis, E. V. Kriventseva, and E. M. Zdobnov,  
582 2015 BUSCO: assessing genome assembly and annotation completeness with  
583 single-copy orthologs. *Bioinformatics* 31: 3210–3212.
- 584 Skov, C., and J. Wiley, 2005 Establishment of the neotropical orchid bee *Euglossa*  
585 *viridissima* (Hymenoptera: Apidae) in Florida. *Florida Entomologist* 88: 225–  
586 227.
- 587 Slater, G. S. C., and E. Birney, 2005 Automated generation of heuristics for biological  
588 sequence comparison. *BMC Bioinformatics* 6: 31.
- 589 Soro, A., J. J. G. Quezada-Euan, P. Theodorou, R. F. A. Moritz, and R. J. Paxton, 2016  
590 The population genetics of two orchid bees suggests high dispersal, low diploid  
591 male production and only an effect of island isolation in lowering genetic  
592 diversity. *Conserv Genet* 1–13.
- 593 Stanke, M., M. Diekhans, R. Baertsch, and D. Haussler, 2008 Using native and  
594 syntenically mapped cDNA alignments to improve de novo gene finding.

- 595        Bioinformatics 24: 637–644.
- 596        Stolle, E., L. Wilfert, R. Schmid-Hempel, P. Schmid-Hempel, M. Kube *et al.*, 2011 A  
597        second generation genetic map of the bumblebee *Bombus terrestris* (Linnaeus,  
598        1758) reveals slow genome and chromosome evolution in the Apidae. BMC  
599        Genomics 12: 48.
- 600        Suni, S. S., 2016 Dispersal of the orchid bee *Euglossa imperialis* over degraded  
601        habitat and intact forest. *Conserv Genet* 1–10.
- 602        Suni, S. S., and B. J. Brosi, 2012 Population genetics of orchid bees in a fragmented  
603        tropical landscape. *Conserv Genet* 13: 323–332.
- 604        Vogel, S., 1966 Parfümsammelnde Bienen als Bestäuber von Orchidaceen und  
605        Gloxinia. *Österr bot Z* 113: 302–361.
- 606        Wcislo, W. T., and J. H. Cane, 2003 Floral Resource Utilization by Solitary Bees  
607        (Hymenoptera: Apoidea) and Exploitation of Their Stored Foods by Natural  
608        Enemies. *41*: 257–286.
- 609        Weinstock, G. M., G. E. Robinson, R. A. Gibbs, K. C. Worley, J. D. Evans *et al.*, 2006  
610        Insights into social insects from the genome of the honeybee *Apis mellifera*.  
611        *Nature* 443: 931–949.
- 612        Whitten, W. M., A. M. Young, and D. L. Stern, 1993 Nonfloral sources of chemicals  
613        that attract male euglossine bees (Apidae: Euglossini). *J. Chem. Ecol.* 19: 3017–  
614        3027.
- 615        Woodard, S. H., B. J. Fischman, A. Venkat, M. E. Hudson, K. Varala *et al.*, 2011 Genes  
616        involved in convergent evolution of eusociality in bees. *Proc Natl Acad Sci USA*  
617        108: 7472–7477.
- 618        Xu, H., X. Luo, J. Qian, X. Pang, J. Song *et al.*, 2012 FastUniq: a fast de novo duplicates  
619        removal tool for paired short reads. (D. Doucet, Ed.). *PLoS ONE* 7: e52249.
- 620        Zdobnov, E. M., F. Tegenfeldt, D. Kuznetsov, R. M. Waterhouse, F. A. Simão *et al.*,  
621        2017 OrthoDB v9.1: cataloging evolutionary and functional annotations for  
622        animal, fungal, plant, archaeal, bacterial and viral orthologs. - PubMed - NCBI.  
623        *Nucleic Acids Res.* 45: D744–D749.
- 624        Zimmermann, Y., D. L. P. Schorkopf, R. F. A. Moritz, R. W. Pemberton, J. J. G. Quezada-  
625        Euan *et al.*, 2011 Population genetic structure of orchid bees (Euglossini) in  
626        anthropogenically altered landscapes. *Conserv Genet* 12: 1183–1194.
- 627
- 628

629 **Figures**

630

631 **Figure 1. Genomic features. (A)** Phylogeny of the four corbiculate bee tribes with  
 632 orchid bees as sistergroup to honey bees, stingless bees, and bumblebees (Romiguer  
 633 et al. 2015). **(B)** K-mer distribution spectrum (k=25) of genomic sequence reads.

634 The positively skewed spectrum reveals a high abundance of a few k-mers, leading  
 635 to an estimate of 87.7% repetitiveness of the *E. dilemma* genome. Red shows the k-  
 636 mer spectrum before, and blue after error correction.

637 **(C)** Genomic element density including genic and non-genic features as a fraction of the overall genome assembly  
 638 length. Over 85% of the assembly could not be annotated with the selected methods.

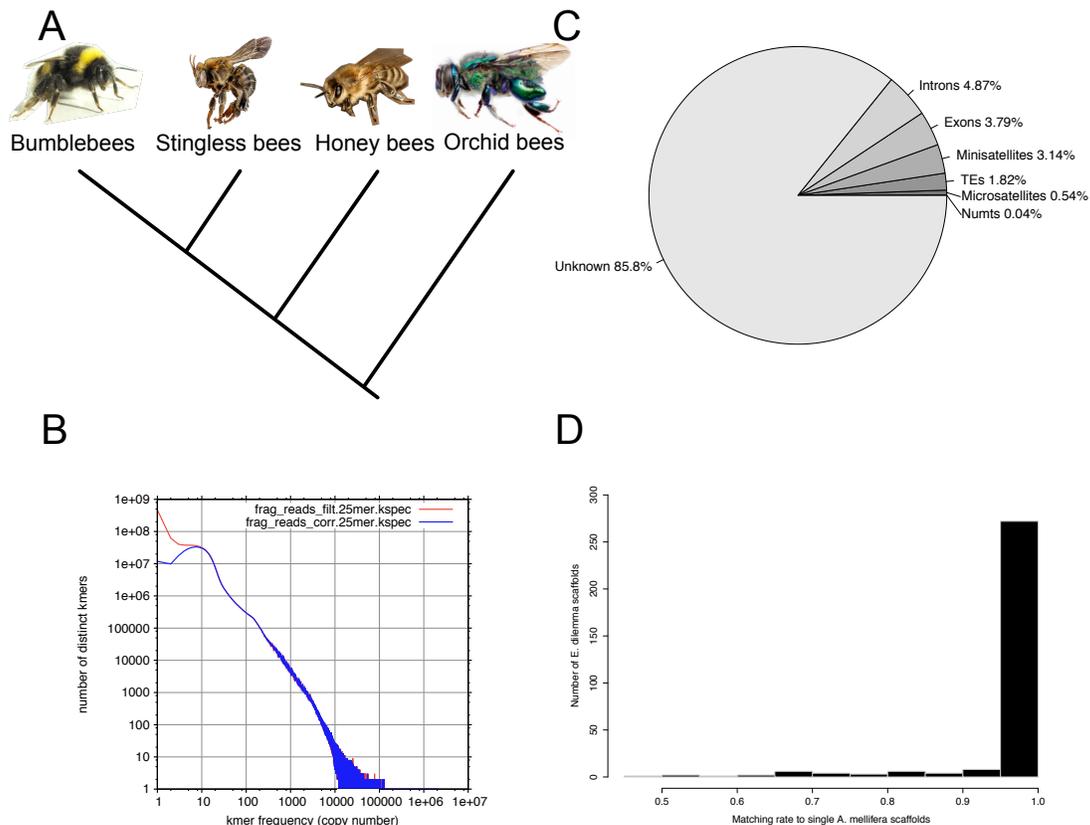
639 **(D)** Synteny between the *E. dilemma* and the honey bee (*Apis mellifera*) genome. In  
 640 an analysis including *E. dilemma* scaffolds of  $\geq 100\text{kb}$  length, 83% showed  $\geq 95\%$   
 641 synteny to a single honeybee scaffold. Photographs in **(A)** are reproduced from  
 642 Wikimedia under the CC BY-SA 3.0 license.

643

644 **Figure 2. Mitochondrial genome reconstruction.** The structure of the honey bee  
 645 mitochondrial genome and information of the homologous reconstructed parts of  
 646 the *E. dilemma* mitochondrial genome. Non-reconstructed parts of incompletely  
 647 reconstructed genes are hatched.

648

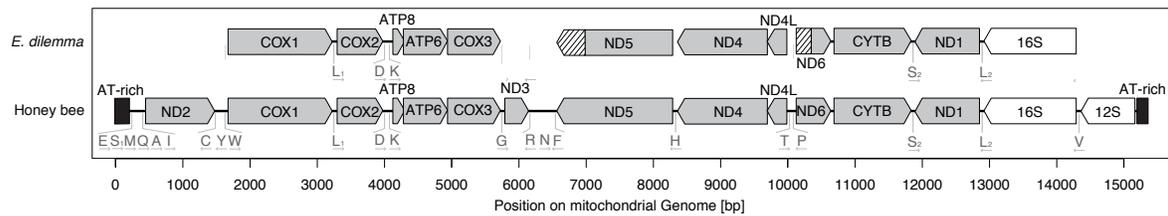
649



650  
 651

**Figure 1**

652



653  
654

Figure 2

655 Tables

656

657 Table 1. *E. dilemma* genome assembly statistics in comparison to previously  
658 published bee genomes.

659

660

Table 1

Species	N50	N25	Longest scaffold	Scaffolds	Assembly length	%GC	Predicted genes	Ref.
<i>Euglossa dilemma</i>	143,590	1,417,006	10,108,120	22,698	588,199,720	39.94	15,904	1
<i>Eufriesea mexicana</i>	2,427	443,231	4,677,300	3,522,543	1,031,837,970	41.38	12,022	2
<i>Apis mellifera</i>	997,192	1,922,192	4,736,299	5,644	234,070,657	32.70	15,314	3
<i>Melipona quadrifasciata</i>	68,085	1,896,322	12,087,087	38,604	507,114,161	38.88	14,257	2
<i>Bombus impatiens</i>	1,399,493	2,389,513	5,466,090	5,559	249,185,056	37.75	15,896	4
<i>Lasioglossum albipes</i>	616,426	1,130,413	3,533,895	41,433	341,616,641	41.50	13,448	5

661

662

663

N50 and N25 indicate the length of the shortest scaffold of those including 50% and 25% of the base pairs in a genome assembly. References (Ref.): 1: This study, 2: Kapheim et al. 2015, 3: Elsik et al. 2014, 4: Sadd et al. 2015, 5: Kocher et al. 2012.