

1 **A high-coverage draft genome of the mycalesine**

2 **butterfly *Bicyclus anynana***

3 Reuben W. Nowell^{1,*}, Ben Elsworth¹, Vicencio Oostra², Bas J. Zwaan³, Christopher W.
4 Wheat⁴, Marjo Saastamoinen⁵, Ilik J. Saccheri⁶, Arjen E. van't Hof⁶, Bethany R. Wasik⁷,
5 Heidi Connahs⁸, Muhammad L. Aslam⁸, Sujai Kumar¹, Richard J. Challis¹, Antónia
6 Monteiro^{7,8,9}, Paul M. Brakefield¹⁰ and Mark Blaxter^{1,*}

7 *Correspondence: reubennowell@gmail.com; mark.blaxter@ed.ac.uk

8 **Abstract**

9 **Background:** The mycalesine butterfly *Bicyclus anynana*, the 'Squinting bush brown',
10 is a model organism in the study of lepidopteran ecology, development and evolution.
11 Here, we present a draft genome sequence for *B. anynana* to serve as a genomics
12 resource for current and future studies of this important model species.

13 **Findings:** Seven libraries with insert sizes ranging from 350 bp to 20 kb were
14 constructed using DNA from an inbred female and sequenced using both Illumina and
15 PacBio technology. 128 Gb raw Illumina data were filtered to 124 Gb and assembled
16 to a final size of 475 Mb (~260X assembly coverage). Contigs were scaffolded using
17 mate-pair, transcriptome and PacBio data into 10,800 sequences with an N50 of 638
18 kb (longest scaffold 5 Mb). The genome is comprised of 26% repetitive elements, and
19 encodes a total of 22,642 predicted protein-coding genes. Recovery of a BUSCO set
20 of core metazoan genes was almost complete (98%). Overall, these metrics compare
21 well with other recently published lepidopteran genomes.

22 **Conclusions:** We report a high-quality draft genome sequence for *Bicyclus anynana*.

23 The genome assembly and annotated gene models are available at LepBase
24 (<http://ensembl.lepbase.org/index.html>).

25 **Keywords:** *Bicyclus anynana*, Squinting bush brown, Nymphalidae, nymphalid,
26 satyrid, lepidopteran genome.

27 **Data description**

28 The squinting bush brown butterfly, *Bicyclus anynana*, is a member of the remarkably
29 speciose nymphalid subtribe Mycalesina, which is distributed across the Old World
30 tropics (Figure 1). *B. anynana* is an important model organism for the study of
31 lepidopteran ecology, development, speciation, behaviour, and evolution [1-6]. *B.*
32 *anynana* are found primarily in woodland habitats across East Africa (from southern
33 Sudan in the north to Swaziland in the south), and adults are typically observed flying
34 close to the ground where they feed on fallen fruit [1]. Strikingly, *B. anynana* exhibits
35 seasonal polyphenism, a form of phenotypic plasticity whereby individuals that
36 develop during the wet season differ in both behaviour, appearance and life history to
37 those that develop during the dry season [7-9]. Wet season butterflies are smaller,
38 have shorter lifespans, are more active, and show larger and more conspicuous
39 eyespots on their wings in comparison to dry season individuals. The genetic basis of
40 this plasticity, and its impacts on various other life-history and developmental
41 characteristics, are ongoing research questions to which the availability of a *B.*
42 *anynana* reference genome will contribute [10-12].

43 **Sampling and sequencing**

44 Genomic DNA was extracted from a *B. anynana* female that had been inbred via seven
45 generations of brother-sister matings. The captive laboratory stock population from
46 which these individuals originated was established in 1988 from 80 wild-caught

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47 individuals, and has been maintained at large effective population sizes to minimise
48 the loss of genetic diversity [1]. Two short-insert libraries with insert sizes of 350 and
49 550 bp were constructed using Illumina TruSeq Nano reagents and sequenced (125
50 base, paired end) on an Illumina HiSeq2500 at Edinburgh Genomics (Edinburgh, UK).
51 DNA from a sister to this focal animal was used to construct four long-insert (mate-
52 pair) libraries with insert sizes of 3 and 5 kb (two of each) at the Centre for Genomic
53 Research, University of Liverpool (Liverpool, UK); libraries of both insert-sizes were
54 then sequenced on an Illumina HiSeq2500 and an Illumina MiSeq at Edinburgh
55 Genomics (Table 1). DNA from a female descendent of the same inbred line was used
56 to construct two long read libraries with insert sizes of 10 and 20 kb, sequenced on the
57 PacBio platform at the Genome Institute of Singapore at 20x coverage using 16 P6
58 SMRT cells. All raw data have been deposited in the Short Read Archive under the
59 accessions given in Table 1.

60 A total of 128.2 Gb raw Illumina data were filtered for low-quality bases and
61 adapter contamination using Skewer v0.2.2 [13], and both raw and trimmed reads were
62 inspected using FastQC v0.11.4 [14]. Only 4 Gb data (3.1%) were discarded, indicating
63 the high quality of the raw data. Kmer frequency distributions were estimated using the
64 “kmercountexact” program from the BBMap v36.02 package [15], and showed two
65 major coverage peaks at ~105X and ~210X (Figure 2). The first peak (105X)
66 represents the proportion of the genome that is heterozygous, and has an approximate
67 span of 87.7 Mb (18.4% of the genome; calculated as one half of the area under the
68 105X curve, from 50X to 150X). The expected proportion of heterozygous sites given
69 seven brother-sister (full-sib) matings is $0.75^7 = 13.3\%$, or 63.5 Mb. Thus, the greater
70 than expected heterozygosity is likely to be due primarily to selection against highly
71 inbred individuals during the course of the inbreeding regime [16].

72 **Contaminant filtering and assembly**

1
2
3 73 Short-insert libraries were screened for the presence of contaminant reads using
4
5 74 Taxon-Annotated GC-Coverage (TAGC) plots, or “blobplots” [17]. An initial draft
6
7 75 assembly was constructed using CLC assembler (CLCBio, Copenhagen) and
8
9
10 76 compared to the NCBI nucleotide database (nt) using Megablast v2.3.0+ [18], and
11
12 77 against the UniRef90 protein database using Diamond v0.7.10 [19]. Read coverage
13
14 78 for each contig was calculated by mapping both libraries to the CLC assembly using
15
16 79 CLC mapper (CLCBio, Copenhagen), and blobplots were generated using Blobtools
17
18 80 v0.9.19.4 [20] using the “bestsumorder” rule for taxonomic annotation of contigs
19
20
21 81 (Figure 3). Contigs that showed a substantially different coverage relative to that of the
22
23 82 main cluster of contigs and/or good hits to sequences annotated as non-Arthropoda
24
25 83 were classed as putative contaminants. A total of 237,394 pairs of reads (~59 Mb) that
26
27 84 were classed as either “mapped/mapped” or “mapped/unmapped” to a putative
28
29
30 85 contaminant were subsequently discarded from further analysis.

31
32
33 86 Filtered libraries were reassembled using the heterozygous-aware assembler
34
35 87 Platanus v1.2.4 [21], with default parameters. Contigs were further scaffolded with the
36
37 88 mate pair libraries using SSPACE v3.0 [22] and with 35,747 assembled *B. anynana*
38
39 89 transcripts using a combination of L_RNA_scaffolder [23] and SCUBAT v2 [24]. The
40
41
42 90 transcripts were assembled using Trinity v. 20140717 [25] from ca. 2×10^9 paired end
43
44 91 RNA-Seq reads sequenced from thorax and abdomen tissue of 72 outbred *B. anynana*
45
46 92 females of the standard captive laboratory stock population (Oostra et al., in
47
48
49 93 preparation). A final round of scaffolding was performed with PacBio long reads (fastq
50
51 94 files error-corrected using the RS_Preaassembler.2 protocol) using SSPACE-
52
53 95 LongRead v1.1 [26]. Finally, gaps between scaffolds were filled using GapFiller v1.10
54
55 96 [27] and PBJelly v15.8.24 [28].
56
57
58
59
60
61
62
63
64
65

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

97 Our final assembly (v1.2) comprised 10,800 scaffolds spanning a total of 475.4
98 Mb, with a scaffold N50 of 638 kb (Table 2). The genome-wide proportion of G+C was
99 36.5%, while the number of undetermined bases (N's) was 5.8 Mb (~1.2% of the total
100 span). We determined assembly completeness by mapping both genomic and
101 transcriptomic reads from *B. anynana* (SRA whole genome sequencing accessions
102 ERR1102671-8, and transcriptome accessions ERR1022636-7, ERR1022640-1, and
103 ERR1022644-5, downloaded October 2016) to the genome using BWA mem v0.7.12
104 [29] and STAR v020201 [30] respectively. Over 99% of reads from the two short-insert
105 libraries mapped to the assembly, suggesting that the vast majority of the genome
106 represented by these data has been assembled. In addition, 94.9% of RNA-Seq reads
107 mapped to the assembly, suggesting that the majority of transcribed genes are
108 present. Gene-level completeness was assessed using CEGMA v2.5 [31] and BUSCO
109 v2.0 [32]. The proportion of CEGMA genes “completely” recovered ($n = 248$) was 81%,
110 increasing to 97% when partially recovered genes are included. The recovery of
111 BUSCO genes specific to the metazoa ($n = 978$) was higher, at 98% for complete
112 genes, increasing to 99% when partial genes are included. An almost complete set
113 (99.2%) of BUSCO genes specific to the Arthropoda ($n = 1,066$) was also recovered.
114 In addition, CEGMA indicated a duplication rate of 1.1 while BUSCO estimated only
115 ~2% genes were present in multiple copies. The high complete CEGMA/BUSCO
116 scores suggest a good assembly that has captured the majority of core
117 metazoan/Arthropod genes in full-length, and that the fragmentation of genes across
118 multiple scaffolds is low. In addition, the low duplication rates suggest that most genes
119 are present in single copy, and thus that the genome does not include significant
120 duplicated segments representing alternative haplotypes.

121 **Annotation**

122 Prior to gene prediction, we masked the *B. anynana* assembly for repetitive elements
123 to minimise the number of spurious open-reading frames due to low-complexity repeat

124 regions or transposable elements. Repetitive motifs in the *B. anynana* assembly were
 125 modelled *ab initio* using RepeatModeler v1.0.5
 126 (<http://www.repeatmasker.org/RepeatModeler.html>). Repeats occurring within
 127 genuine coding regions were excluded by querying the proteins from a previous *B.*
 128 *anynana* assembly (v0.1) versus the RepeatModeler database using BLAST,
 129 removing any sequences showing a match at E -value $\leq 1e-10$ threshold. The filtered
 130 RepeatModeler database was combined with known repeats from the Lepidoptera
 131 using RepBase v20.05 [33] and input to RepeatMasker v4.0.5 [34] to mask the
 132 assembly. Overall, approximately one quarter of the assembly (122.6 Mb) was masked
 133 from gene prediction (Table 3).

134 **Table 3:** Major types of repeat content for *B. anynana*.

Repeat type	Span (Mb)	Proportion of genome
SINE	10.8	2.3%
LINE	15.3	3.2%
LTR elements	1.1	0.2%
DNA elements	0.8	0.2%
Small RNA	10.8	2.3%
Unclassified	86.2	18.1%
Total	122.6	25.8%

135
 136 Gene finding was performed following a two-pass approach [35]. Initial gene-
 137 models were constructed with MAKER v2.31 [36], using HMMs derived from SNAP
 138 [37] and GeneMark-ES v4.3 [38] in conjunction with a recently published *B. anynana*
 139 transcriptome as evidence [39]. MAKER gene-models were then passed to
 140 AUGUSTUS v3.0.3 [40] for refinement, resulting in an initial set of 26,722 predicted
 141 protein-coding genes. A set of basic filters was applied to remove likely spurious gene
 142 models (Table 4), resulting in the deletion of 4,080 gene models. Protein sequences

143 from the filtered 22,642 genes were annotated using BLAST searches versus
 144 UniRef90 and the NCBI non-redundant protein database (nr), and domains/motifs
 145 were described using InterProScan5 [41]. Summary statistics for the 22,642 predicted
 146 gene models are given in Table 5.

147 **Table 4:** Number of genes in potential error categories.

Category	Description	Number of genes
(a)	Single-exon	7112
(b)	Small exon (< 9 bp)	1866
(c)	Small intron (\leq 40 bp)	45
(d)	Short (CDS < 120 bp)	127
(e)	No hit to <i>nr</i>	6532
(f)	Duplicate (\geq 98% identity over \geq 98% query length)	822
Total¹		4080

¹Defined as the non-redundant total of the intersection of each category (a) to (d) with category (e), plus the shorter of any duplicates identified in category (f).

149 Comparison to other lepidopteran genomes

150 To ascertain the relative quality of the *B. anynana* v1.2 assembly, we compared our
 151 results to nine other published lepidopteran genomes available on LepBase
 152 (<http://lepbase.org/>) [42]: *Bombyx mori* ASM15162v1 [43], *Danaus plexippus* v3 [44],
 153 *Heliconius melpomene* Hmel2 [45,46], *Lerema accius* v1.1 [47], *Melitaea cinxia*
 154 MelCinx1.0 [48], *Papilio glaucus* v1.1 [49], *Papilio polytes* Ppol 1.0 [50], *Papilio xuthus*
 155 Pap_xu_1.0 [50] and *Plutella xylostella* DBM_FJ_v1.1 [51]. The *B. anynana* v1.2
 156 assembly was of high quality compared to other published genomes, with the majority
 157 of the genome represented in a relatively small number of scaffolds despite being only
 158 marginally smaller than the largest lepidopteran genome, *B. mori* (Figure 4a).
 159 Interestingly, *B. anynana* v1.2 encodes the highest number of proteins of the 10

160 species compared (Figure 4b). Despite measures to eliminate potentially spurious
161 ORFs caused by annotation error or by duplication, *B. anynana* encodes ~3,250 more
162 genes than the diamondback moth *P. xylostella*, and ~10,400 more than the
163 swallowtail *P. polytes*. It is tempting to attribute the apparently high number of genes
164 to the developmental plasticity and alternative seasonal forms with divergent
165 morphologies and life histories in *B. anynana*. However, it remains to be determined
166 whether the number of genes predicted in *B. anynana* is a function of its larger genome
167 size or unusual life-history characteristics, or if further curation of the v1.2 gene models
168 will reduce the number of inferred genes.

169 **Concluding remarks**

170 We present a high-coverage, high quality draft assembly and annotation of the
171 mycalesine butterfly *B. anynana*. The assembly will be a core resource for ongoing
172 analyses of population genomics, discovery of *cis*-regulatory elements of wing
173 patterning and other genes, functional genetics and functional ecology of complex
174 gene families, and the evolution of novel and plastic lifecycle strategies in
175 lepidopterans and other arthropods.

176 **Availability of supporting data**

177 All raw sequence data have been deposited in the Short Read Archive (SRA) and are
178 available for download using the accession numbers provided in Table 1. The *B.*
179 *anynana* v1.2 assembly, as well as final predicted gene-models and protein
180 annotations, are publicly available for viewing and download via LepBase [42], an
181 Ensembl [52] genome database for the Lepidoptera
182 (<http://ensembl.lepbase.org/index.html>). Data supporting the manuscript, including
183 annotations as well as BUSCO and CEGMA results, are also available via the

184 *GigaScience* database, GigaDB [53]. A previous *B. anynana* assembly (nBa.0.1) is
185 also available on LepBase.

186 **Abbreviations**

187 BUSCO: Benchmarking Universal Single-Copy Orthologs; CEGMA: Core Eukaryotic
188 Genes Mapping Approach; CDS: Coding sequence; ORF: Open reading frame.

189 **Competing interests**

190 The authors declare that they have no competing interests.

191 **Author contributions**

192 PMB and MB designed the study; AM and BRW collected samples and produced the
193 inbred line; AEVH, IJS and HC extracted DNA samples; RWN, BE and MB worked on
194 the genome assembly and annotation; VO, BJZ, CW and MS contributed
195 transcriptome data; AM, HC and MLA contributed PacBio data; SK and RJC uploaded
196 the assembly to LepBase. RWN, VO, AM, PMB and MB wrote the manuscript. All
197 authors have read and approved the final version of the manuscript.

198 **Acknowledgements**

199 We thank Edinburgh Genomics and Genome Institute of Singapore for genome
200 sequencing, initial QC and data delivery. We also thank two reviewers for helpful
201 comments on a previous version of this manuscript. Funding for the *Bicyclus anynana*
202 genome project was provided by the ERC Advanced Grant number 250325 (EMARES)
203 to PMB and by the South East Asian Biodiversity Genomics Center (NUS grant R-154-
204 000-648-646 and R-154-000-648-733) to AM. Funding for LepBase was provided by
205 BBSRC grant number BB/K020161.

206 **Author details**

207 ¹Institute of Evolutionary Biology, University of Edinburgh, Edinburgh EH9 3FL, United
208 Kingdom; ²Department of Genetics, Evolution and Environment, University College
209 London, United Kingdom; ³Laboratory of Genetics, Wageningen University, The
210 Netherlands; ⁴Department of Zoology, Stockholm University, Sweden;
211 ⁵Metapopulation Research Centre, Department of Biosciences, University of Helsinki,
212 Finland; ⁶Institute of Integrative Biology, University of Liverpool, Liverpool L69 7ZB,
213 United Kingdom; ⁷Department of Ecology and Evolutionary Biology, Yale University,
214 New Haven, CT 06511, USA; ⁸Department of Biological Sciences, National University
215 of Singapore, Singapore 117543; ⁹Yale-NUS College, Singapore 138609;
216 ¹⁰Department of Zoology, University of Cambridge, Cambridge, CB2 3EJ, United
217 Kingdom

218 **Tables**

219 Tables 1, 2 and 5 are in landscape orientation and can be found as additional files at
220 the end of this manuscript.

221 **Figure legends**

222 **Figure 1:** Wet-season morph of *Bicyclus anynana* (picture credit: William H. Piel and
223 Antónia Monteiro).

224 **Figure 2:** Kmer frequency distribution for *B. anynana* short-insert libraries ($k = 31$).
225 The bimodality of the distribution, with peaks at approximately 105X and again at 210X,
226 is the result of heterozygosity in the sequence data.

227 **Figure 3:** Taxon-annotated GC-coverage plots for **(a)** draft and **(b)** final *B. anynana*
228 genome assemblies. Each contig/scaffold in the assembly is represented as a circle,

229 coloured according to the best match to taxonomically annotated sequence databases
230 (see legends) and distributed according to the proportion GC (*x*-axis) and read
231 coverage (*y*-axis). The upper- and right-hand panels show the distribution of the total
232 span (kb) of contigs/scaffolds for a given coverage (upper panel) or GC (right panel)
233 bin. The heterozygosity in the sample is evident in the bimodal coverage distribution
234 seen in (a). The cluster of orange-coloured contigs at a lower coverage and higher GC
235 than the main cloud were likely derived from contaminant *Enterococcus* present in the
236 sample. The final assembly, (b), shows the effective collapse of heterozygous regions,
237 the removal of contaminant sequences and the scaffolding of contigs into long
238 contiguous sequences. Note that only taxon annotations with a span > 1 Mb are shown
239 in the legend for clarity.

240 **Figure 4:** Assembly and gene prediction comparison among 10 lepidopteran
241 genomes. **(a)** Cumulative assembly curves showing the relationship between the
242 number of scaffolds (*x*-axis) and the cumulative span of each assembly (*y*-axis),
243 coloured by species. Higher quality assemblies are represented by an almost-vertical
244 line (e.g., *H. melpomene* Hmel2 assembly in black), indicating a relatively small
245 number of scaffolds is required to reach the final genome span; conversely, a long tail
246 indicates the assembly includes a large number of smaller scaffolds. The curve for *B.*
247 *anyana* (brown and bold) suggests a good assembly for this species, with the majority
248 of the assembly comprised of relatively few scaffolds. **(b)** *B. anyana* v1.2 encodes
249 the greatest number of genes of the 10 genomes, and is particularly different from *B.*
250 *mori*, which is of equivalent length. Species names/colours are as follows: “bicyclus”
251 (brown), *B. anyana*; “bombyx” (blue), *B. mori*; “danaus” (light green), *D. plexippus*;
252 “heliconius” (black), *H. melpomene*; “lerema” (dark green), *L. accius*; “melitaea”
253 (orange), *M. cinxia*; “glaucus” (red), *P. glaucus*; “polytes” (pink), *P. polytes*; “xuthus”
254 (violet), *P. xuthus*; “plutella” (grey), *P. xylostella*.

255 References

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

256 1. Brakefield PM, Beldade P, Zwaan BJ. The African butterfly *Bicyclus anynana*: a
257 model for evolutionary genetics and evolutionary developmental biology. Cold Spring
258 Harb Protoc. 2009; doi:10.1101/pdb.emo122-2.

259 2. Brakefield PM. Radiations of mycalesine butterflies and opening up their
260 exploration of morphospace. Am. Nat. 2010;176 Suppl 1:S77-87.

261 3. Prudic KL, Jeon C, Cao H, Monteiro A. Developmental plasticity in sexual roles of
262 butterfly species drives mutual sexual ornamentation. Science. 2011;331:73-5.

263 4. Westerman EL, Hodgins-Davis A, Dinwiddie A, Monteiro A. Biased learning affects
264 mate choice in a butterfly. Proc. Natl. Acad. Sci. 2012;109:10948-53.

265 5. Monteiro A. Origin, development, and evolution of butterfly eyespots. Annu. Rev.
266 Entomol. 2015;60:253-71.

267 6. Aduse-Poku K, Brakefield PM, Wahlberg N, Brattström O. Expanded molecular
268 phylogeny of the genus *Bicyclus* (Lepidoptera: Nymphalidae) shows the importance
269 of increased sampling for detecting semi-cryptic species and highlights potentials for
270 future studies. System. Biodivers. 2017;15:115-30.

271 7. Brakefield PM, Reitsma N. Phenotypic plasticity, seasonal climate and the
272 population biology of *Bicyclus* butterflies (Satyridae) in Malawi. Ecol. Entomol.
273 1991;16:291-303.

274 8. Brakefield PM, Gates J, Keys D, Kesbeke F, Wijngaarden PJ, Monteiro A, et al.
275 Development, plasticity and evolution of butterfly eyespot patterns. Nature.
276 1996;384:236-42.

277 9. Monteiro A, Tong X, Bear A, Liew SF, Bhardwaj S, Wasik BR, et al. Differential
278 expression of ecdysone receptor leads to variation in phenotypic plasticity across
279 serial homologs. PLoS Genet. 2015;11:e1005529.

280 10. Beldade P, Mateus ARA, Keller RA. Evolution and molecular mechanisms of
281 adaptive developmental plasticity. Mol. Ecol. 2011;20:1347-63.

282 11. Oostra V, Brakefield PM, Hiltemann Y, Zwaan BJ, Brattström O. On the fate of
283 seasonally plastic traits in a rainforest butterfly under relaxed selection. Ecol. Evol.
284 2014;4:2654-67.

285 12. Dion E, Monteiro A, Yew JY. Phenotypic plasticity in sex pheromone production
286 in *Bicyclus anynana* butterflies. Sci. Rep. 2016;6:39002.

287 13. Jiang H, Lei R, Ding S-W, Zhu S. Skewer: a fast and accurate adapter trimmer
288 for next-generation sequencing paired-end reads. BMC Bioinformatics. 2014;15:182.

289 14. Andrews S. FastQC: a quality control tool for high throughput sequence data.
290 Available from: <http://www.bioinformatics.babraham.ac.uk/projects/fastqc>

291 15. Bushnell B. BBMap short read aligner, and other bioinformatic tools. Available
292 from: sourceforge.net/projects/bbmap/

293 16. Saccheri IJ, Brakefield PM, Nichols RA. Severe inbreeding depression and rapid
294 fitness rebound in the butterfly *Bicyclus anynana* (Satyridae). Evolution.
295 1996;50:2000-13.

1 296 17. Kumar S, Jones M, Koutsovoulos G, Clarke M, Blaxter M. Blobology: exploring
2 297 raw genome data for contaminants, symbionts and parasites using taxon-annotated
3 298 GC-coverage plots. *Front Genet. Frontiers*; 2013;4:237.

4 299 18. Altschul SF, Madden TL, Schaffer AA, Zhang JH, Zhang Z, Miller W, et al.
5 300 Gapped BLAST and PSI-BLAST: a new generation of protein database search
6 301 programs. *Nucleic Acids Res.* 1997;25:3389–402.

7 302 19. Buchfink B, Xie C, Huson DH. Fast and sensitive protein alignment using
8 303 DIAMOND. *Nature Methods.* 2015;12:59–60.

9 304 20. Laetsch DR. Blobtools: application for the visualisation of draft genome
10 305 assemblies and general QC. Available from: <https://github.com/DRL/blobtools>

11 306 21. Kajitani R, Toshimoto K, Noguchi H, Toyoda A, Ogura Y, Okuno M, et al. Efficient
12 307 de novo assembly of highly heterozygous genomes from whole-genome shotgun
13 308 short reads. *Genome Res.* 2014;24:1384–95.

14 309 22. Boetzer M, Henkel CV, Jansen HJ, Butler D, Pirovano W. Scaffolding pre-
15 310 assembled contigs using SSPACE. *Bioinformatics.* 2011;27:578–9.

16 311 23. Xue W, Li J-T, Zhu Y-P, Hou G-Y, Kong X-F, Kuang Y-Y, et al.
17 312 L_RNA_scaffolder: scaffolding genomes with transcripts. *BMC Genomics.*
18 313 2013;14:604.

19 314 24. Koutsovoulos G. SCUBAT2. Available from: <https://github.com/GDKO/SCUBAT2>

20 315 25. Haas BJ, Papanicolaou A, Yassour M, Grabherr M, Blood PD, Bowden J, et al.
21 316 De novo transcript sequence reconstruction from RNA-seq using the Trinity platform
22 317 for reference generation and analysis. *Nat Protoc.* 2013;8:1494–512.

23 318 26. Boetzer M, Pirovano W. SSPACE-LongRead: scaffolding bacterial draft genomes
24 319 using long read sequence information. *BMC Bioinformatics.* 2014;15:211.

25 320 27. Boetzer M, Pirovano W. Toward almost closed genomes with GapFiller. *Genome*
26 321 *Biol.* 2012;13:R56.

27 322 28. English AC, Richards S, Han Y, Wang M, Vee V, Qu J, et al. Mind the Gap:
28 323 Upgrading Genomes with Pacific Biosciences RS Long-Read Sequencing
29 324 Technology. *PLoS ONE.* 2012;7.

30 325 29. Li H, Durbin R. Fast and accurate long-read alignment with Burrows-Wheeler
31 326 transform. *Bioinformatics.* 2010;26:589–95.

32 327 30. Dobin A, Davis CA, Schlesinger F, Drenkow J, Zaleski C, Jha S, et al. STAR:
33 328 ultrafast universal RNA-seq aligner. *Bioinformatics.* 2013;29:15–21.

34 329 31. Parra G, Bradnam K, Korf I. CEGMA: a pipeline to accurately annotate core
35 330 genes in eukaryotic genomes. *Bioinformatics.* 2007;23:1061–7.

36 331 32. Simão FA, Waterhouse RM, Ioannidis P, Kriventseva EV, Zdobnov EM. BUSCO:
37 332 assessing genome assembly and annotation completeness with single-copy
38 333 orthologs. *Bioinformatics.* 2015;31:3210–2.

39 334 33. Jurka J, Kapitonov VV, Pavlicek A, Klonowski P, Kohany O, Walichiewicz J.
40 335 Repbase Update, a database of eukaryotic repetitive elements. *Cytogenet Genome*

- 336 Res. 2005;110:462–7.
- 337 34. Smit A, Hubley R, Green P. RepeatMasker. Available from
338 <http://www.repeatmasker.org>.
- 339 35. Koutsovoulos G. CGP-Pipeline. Available from: <https://gist.github.com/GDKO/>.
- 340 36. Holt C, Yandell M. MAKER2: an annotation pipeline and genome-database
341 management tool for second-generation genome projects. BMC Bioinformatics.
342 2011;12:491.
- 343 37. Korf I. Gene finding in novel genomes. BMC Bioinformatics. 2004;5:59.
- 344 38. Ter-Hovhannisyan V, Lomsadze A, Chernoff YO, Borodovsky M. Gene prediction
345 in novel fungal genomes using an ab initio algorithm with unsupervised training.
346 Genome Res. 2008;18:1979–90.
- 347 39. Oostra V, Saastamoinen M, Zwaan BJ, Wheat CW. Extensive phenotypic
348 plasticity in a seasonal butterfly limits potential for evolutionary responses to
349 environmental change. bioRxiv. 2017; doi: <https://doi.org/10.1101/126177>.
- 350 40. Stanke M, Diekhans M, Baertsch R, Haussler D. Using native and syntenically
351 mapped cDNA alignments to improve de novo gene finding. Bioinformatics.
352 2008;24:637–44.
- 353 41. Jones P, Binns D, Chang H-Y, Fraser M, Li W, McAnulla C, et al. InterProScan 5:
354 genome-scale protein function classification. Bioinformatics. 2014;30:1236–40.
- 355 42. Challis RJ, Kumar S, Dasmahapatra KKK, Jiggins CD, Blaxter M. Lepbase: the
356 Lepidopteran genome database. bioRxiv. 2016; doi:
357 <http://dx.doi.org/10.1101/056994>.
- 358 43. Duan J, Li R, Cheng D, Fan W, Zha X, Cheng T, et al. SilkDB v2.0: a platform for
359 silkworm (*Bombyx mori*) genome biology. Nucleic Acids Res. 2010;38:D453–6.
- 360 44. Zhan S, Merlin C, Boore JL, Reppert SM. The monarch butterfly genome yields
361 insights into long-distance migration. Cell. 2011;147:1171–85.
- 362 45. Heliconius Genome Consortium. Butterfly genome reveals promiscuous
363 exchange of mimicry adaptations among species. Nature. 2012;487:94–8.
- 364 46. Davey JW, Chouteau M, Barker SL, Maroja L, Baxter SW, Simpson F, et al.
365 Major improvements to the *Heliconius melpomene* genome assembly used to
366 confirm 10 chromosome fusion events in 6 million years of butterfly evolution. G3.
367 2016;6:695–708.
- 368 47. Cong Q, Borek D, Otwinowski Z, Grishin NV. Skipper genome sheds light on
369 unique phenotypic traits and phylogeny. BMC Genomics. 2015;16:639.
- 370 48. Ahola V, Lehtonen R, Somervuo P, Salmela L, Koskinen P, Rastas P, et al. The
371 Glanville fritillary genome retains an ancient karyotype and reveals selective
372 chromosomal fusions in Lepidoptera. Nature Communications. 2014;5:1–9.
- 373 49. Cong Q, Borek D, Otwinowski Z, Grishin NV. Tiger swallowtail genome reveals
374 mechanisms for speciation and caterpillar chemical defense. Cell Rep. 2015;10:910–
375 9.

1 376 50. Nishikawa H, Iijima T, Kajitani R, Yamaguchi J, Ando T, Suzuki Y, et al. A genetic
2 377 mechanism for female-limited Batesian mimicry in *Papilio* butterfly. Nat Genet.
3 378 2015;47:405–9.
4 379 51. You M, Yue Z, He W, Yang X, Yang G, Xie M, et al. A heterozygous moth
5 380 genome provides insights into herbivory and detoxification. Nat. Genet.
6 381 2013;45:220–5.
7
8 382 52. Yates A, Akanni W, Amode MR, Barrell D, Billis K, Carvalho-Silva D, et al.
9 383 Ensembl 2016. Nucleic Acids Res. 2016;44:D710–6.
10
11 384 53. Nowell RW, Elsworth B, Oostra V, Zwaan BJ, Wheat CW, Saastamoinen M, et al.
12 385 Supporting data for "A high-coverage draft genome of the mycalesine
13 386 butterfly *Bicyclus anynana*". *GigaScience* Database. 2017.
14 387 <http://dx.doi.org/10.5524/100280>
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

Tables

Table 1: Data counts and library information.

Library type	Platform	Read length	Insert size (expected)	Number of reads (raw)	Number of reads (trimmed)	Number of bases (trimmed)	SRA run accessions
Short insert	Illumina HiSeq2500	125 bp paired-end	350 bp	271808057 pairs	267241712 (98.3%)	66334099834 (97.6%)	ERR1102671-2, ERR1102675-6
Short insert	Illumina HiSeq2500	125 bp paired-end	550 bp	241050065 pairs	234269871 (97.2%)	57913474128 (96.1%)	ERR1102673-4, ERR1102677-8
Mate pair	Illumina HiSeq2500	100 bp paired-end	3 kb	77105680 pairs	31848200 (41.3%)	5758856502 (37.3%)	ERR1750945
Mate pair	Illumina MiSeq	100 bp paired-end	3 kb	5641764 pairs	2170610 (38.5%)	397993018 (35.3%)	ERR754051
Mate pair	Illumina HiSeq2500	100 bp paired-end	5 kb	77614870 pairs	45676725 (58.9%)	8203769131 (52.8%)	ERR1750946
Mate pair	Illumina	100 bp	5 kb	7939601 pairs	4734000 (59.6%)	861352793 (54.2%)	ERR754052

	MiSeq	paired-end					
Long read	PacBio P6	0.80-50 kb	10 kb	1388796	1199064 (86.3%)	4086394966	ERR1797559-74

Table 2: Summary of *B. anynana* genome assembly and comparison to selected lepidopteran genomes.

	<i>B. anynana</i>	<i>B. mori</i>	<i>D. plexippus</i>	<i>H. melpomene</i>	<i>M. cinxia</i>
Assembly version	v1.2	ASM15162v1	v3	Hmel2	MelCinx1.0
Span	475.4 Mb	481.8 Mb	248.6 Mb	275.2 Mb	389.9 Mb
Contigs					
Number	23699	88673	10682	3100	48180
N50 ¹	78.7 kb	15.5 kb	111.0 kb	328.9 kb	14.1 kb
NumN50 ²	1543	8075	548	214	7366
Scaffolds					
Number	10800	43379	5397	795	8261
N50	638.3 kb	4008.4 kb	715.6 kb	2102.7 kb	119.3 kb
NumN50	194	38	101	34	970
N90	99.3 kb	61.1 kb	160.5 kb	273.1 kb	29.6 kb
NumN90	909	258	366	176	3396
Shortest / longest	201 b / 5 Mb	53 b / 16.2 Mb	300 b / 6.2 Mb	394 b / 9.4 Mb	1.5 kb / 668 kb
G+C content	36.5%	37.7%	31.6%	32.8%	32.6%

NNNs					
Span	5.8 Mb (1.2%)	50.1 Mb (10.4%)	6.7 Mb (2.7%)	986 kb (0.4%)	28.9 Mb (7.4%)
N50	1.4 kb	5.0 kb	2.5 kb	2.4 kb	1.4 kb
CEGMA ³ (<i>n</i> = 248)	C :81.1%; D :1.1; F :97.2%	C :76.6%; F :96.8%	C :90.3%; F :96%	C :88.7%; F :96.8%	NA
BUSCO ³ (<i>n</i> = 1066)	C :98.3%; D :1%; F :99.2%	C :97.5%; D :0.5%; F :98.4%	C :97.4%; D :8.6%; F :98.5%	C :98.8%; D :0.7%; F :99.3%	C :85.7%; D :0.2%; F :91.8%

¹N50: the length of the contig/scaffold at which 50% of the genome span is accounted, given a list of sequences sorted by length. ²numN50: the number of sequences required to reach the N50 sequence. ³CEGMA / BUSCO notation: **C**, proportion (%) genes completely recovered; **D**, duplication rate; **F**, proportion (%) genes at least partially recovered (including complete genes); *n*, number of queries. Note that duplication rate (D) for CEGMA is given as the average number of (complete) genes recovered, whereas for BUSCO it is the proportion of complete genes recovered multiple times. BUSCO values are based on comparisons to the Arthropoda gene set.

Table 5: Summary of *B. anynana* gene prediction.

	<i>B. anynana</i>	<i>B. mori</i>	<i>D. plexippus</i>	<i>H. melpomene</i>	<i>M. cinxia</i>
Assembly version	v1.2	ASM15162v1	v3	Hmel2	MelCinx1.0
Number of CDS	22642	19618	15130	13178	16668
Mean length	1.4 kb	1.6 kb	1.4 kb	1.3 kb	958 bp
Median length	1.2 kb	1.2 kb	981 bp	927 bp	693 bp
Min/max	84 bp / 28.3 kb	23 bp / 60.3 kb	9 bp / 58.9 kb	45 bp / 46.4 kb	6 bp / 45.4 kb
Introns					
Mean number per gene	4.4	9.9	5.7	5	NA ¹
Length (mean/median)	1.3/0.6 kb	2.4/0.8 kb	795/280 bp	960/416 bp	NA
Exons					
Length (mean/median)	208/126 bp	283/161 bp	206/149 bp	284/157 bp	NA
Number of single-exon genes	3571	1744	1461	3113	NA
Transcript GC	49.2%	48.3%	46.5%	43%	41.7%
Gene frequency² (genes per Mb)	47.7	32.1	60.9	55.5	NA

¹GFF for *M. cinxia* not available; ²Defined as the number of genes divided by the total genome span (Mb).



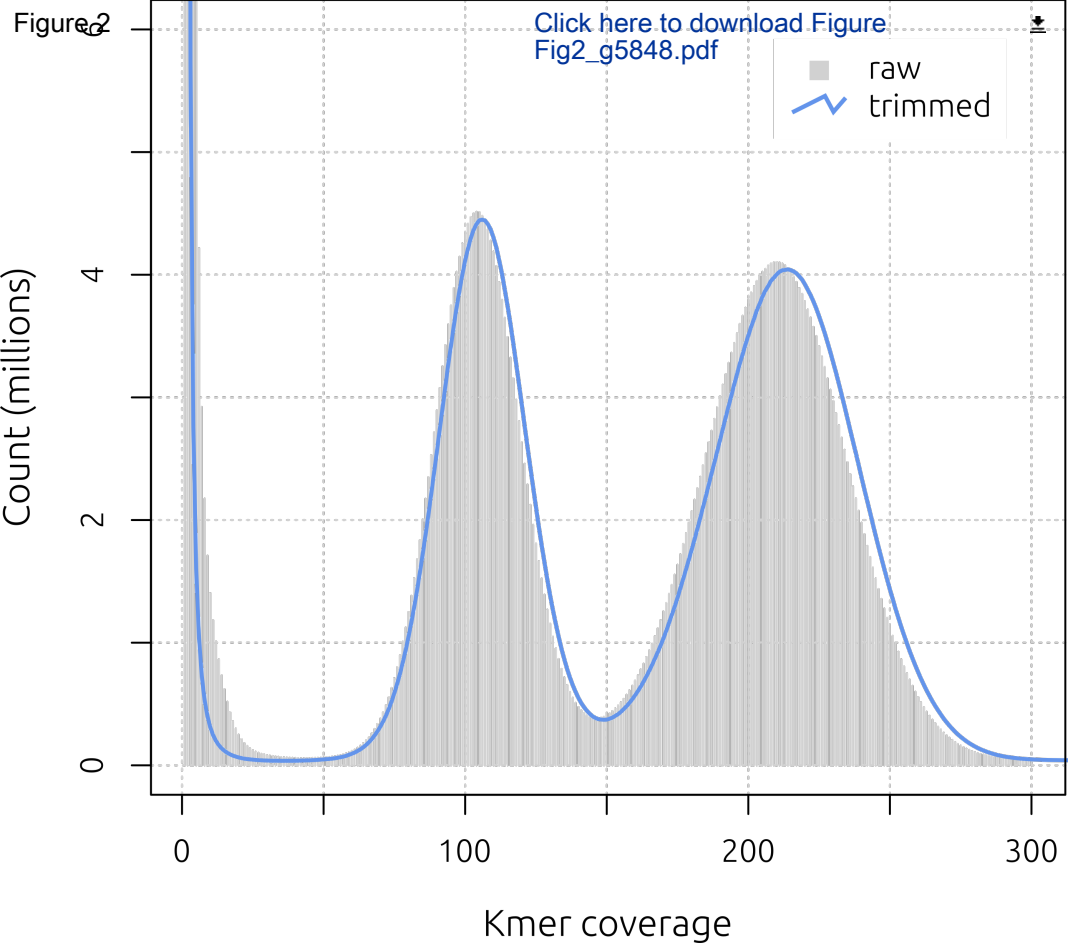
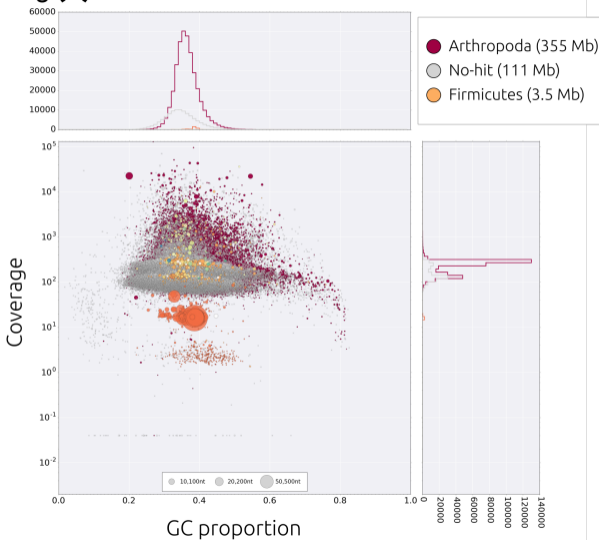


Fig 3

Click **(b)** [here](#) to download Figure Fig3_g4698.pdf