

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
11 evolution. Here, we present a draft genome sequence for *B. anynana* to serve as a
12 genomics resource for current and future studies of this important model species.

13 **Findings:** Eight 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
15 and PacBio technology. 128 Gb raw Illumina data were filtered to 124 Gb and
16 assembled to a final size of 475 Mb (~260X assembly coverage). Contigs were
17 scaffolded using mate-pair, transcriptome and PacBio data into 10,800 sequences
18 with an N50 of 638 kb (longest scaffold 5 Mb). The genome is comprised of 26%
19 repetitive elements, and encodes a total of 22,642 predicted protein-coding genes.

1 20 Recovery of a BUSCO set of core metazoan genes was almost complete (98%).

2
3 21 Overall, these metrics compare well with other recently published lepidopteran
4
5 22 genomes.

6
7
8
9 23 **Conclusions:** We report a high-quality draft genome sequence for *Bicyclus anynana*.

10
11 24 The genome assembly and annotated gene models are available at LepBase
12
13 25 (<http://ensembl.lepbase.org/index.html>).

14
15
16
17 26 **Keywords:** *Bicyclus anynana*, Squinting bush brown, Nymphalidae, nymphalid,
18
19 27 satyrid, lepidopteran genome.

28 **Data description**

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

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 characteristics, are ongoing research questions to which the availability of a *B.*
43 *anynana* reference genome will contribute [10-12].

44 **Sampling and sequencing**

45 Genomic DNA was extracted from a *B. anynana* female that had been inbred via
46 seven generations of brother-sister matings. The captive laboratory stock population
47 from which these individuals originated was established in 1988 from 80 wild-caught
48 individuals, and has been maintained at large effective population sizes to minimise
49 the loss of genetic diversity [1]. Two short-insert libraries with insert sizes of 350 and
50 550 bp were constructed using Illumina TruSeq Nano reagents and sequenced (125
51 base, paired end) on an Illumina HiSeq2500 at Edinburgh Genomics (Edinburgh, UK).
52 DNA from a sister to this focal animal was used to construct four long-insert (mate-
53 pair) libraries with insert sizes of 3 and 5 kb (two of each) at the Centre for Genomic
54 Research, University of Liverpool (Liverpool, UK); libraries of both insert-sizes were
55 then sequenced on an Illumina HiSeq2500 and an Illumina MiSeq at Edinburgh
56 Genomics (Table 1). DNA from a female descendent of the same inbred line was used
57 to construct two long read libraries with insert sizes of 10 and 20 kb, sequenced on
58 the PacBio platform at the Genome Institute of Singapore at 20x coverage using 16
59 P6 SMRT cells. All raw data have been deposited in the Short Read Archive under the
60 accessions given in Table 1.

61 A total of 128.2 Gb raw Illumina data were filtered for low-quality bases and
62 adapter contamination using Skewer v0.2.2 [13], and both raw and trimmed reads
63 were inspected using FastQC v0.11.4 [14]. Only 4 Gb data (3.1%) were discarded,
64 indicating the high quality of the raw data. Kmer frequency distributions were

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
66
67
68
69
70
71
72
73

estimated using the “kmercountexact” program from the BMAP v36.02 package [15], and showed two major coverage peaks at ~105X and ~210X (Figure 1). The first peak (105X) represents the proportion of the genome that is heterozygous, and has an approximate span of 87.7 Mb (18.4% of the genome; calculated as one half of the area under the 105X curve, from 50X to 150X). The expected proportion of heterozygous sites given seven brother-sister (full-sib) matings is $0.75^7 = 13.3\%$, or 63.5 Mb. Thus, the greater than expected heterozygosity is likely to be due primarily to selection against highly inbred individuals during the course of the inbreeding regime [16].

74 **Contaminant filtering and assembly**

75 Short-insert libraries were screened for the presence of contaminant reads using
76 Taxon-Annotated GC-Coverage (TAGC) plots, or “blobplots” [17]. An initial draft
77 assembly was constructed using CLC assembler (CLCBio, Copenhagen) and compared
78 to the NCBI nucleotide database (nt) using Megablast v2.3.0+ [18], and against the
79 UniRef90 protein database using Diamond v0.7.10 [19]. Read coverage for each
80 contig was calculated by mapping both libraries to the CLC assembly using CLC
81 mapper (CLCBio, Copenhagen), and blobplots were generated using Blobtools
82 v0.9.19.4 [20] using the “bestsumorder” rule for taxonomic annotation of contigs
83 (Figure 2). Contigs that showed a substantially different coverage relative to that of
84 the main cluster of contigs and/or good hits to sequences annotated as non-
85 Arthropoda were classed as putative contaminants. A total of 237,394 pairs of reads
86 (~59 Mb) that were classed as either “mapped/mapped” or “mapped/unmapped” to
87 a putative contaminant were subsequently discarded from further analysis.

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

88 Filtered libraries were reassembled using the heterozygous-aware assembler
89 Platanus v1.2.4 [21], with default parameters. Contigs were further scaffolded with the
90 mate pair libraries using SSPACE v3.0 [22] and with 35,747 assembled *B. anynana*
91 transcripts using a combination of L_RNA_scaffolder [23] and SCUBAT v2 [24]. The
92 transcripts were assembled using Trinity v. 20140717 [25] from ca. 2×10^9 paired end
93 RNA-Seq reads sequenced from thorax and abdomen tissue of 72 outbred *B.*
94 *anynana* females of the standard captive laboratory stock population (Oostra et al., in
95 preparation). A final round of scaffolding was performed with PacBio long reads
96 (fastq files error-corrected using the RS_Preaassembler.2 protocol) using SSPACE-
97 LongRead v1.1 [26]. Finally, gaps between scaffolds were filled using GapFiller v1.10
98 [27] and PBJelly v15.8.24 [28].

99 Our final assembly (v1.2) comprised 10,800 scaffolds spanning a total of 475.4
100 Mb, with a scaffold N50 of 638 kb (Table 2). The genome-wide proportion of G+C
101 was 36.5%, while the number of undetermined bases (N's) was 5.8 Mb (~1.2% of the
102 total span). We determined assembly completeness by mapping both genomic and
103 transcriptomic reads from *B. anynana* (SRA whole genome sequencing accessions
104 ERR1102671-8, and transcriptome accessions ERR1022631, ERR1022635-7,
105 ERR1022640 and ERR1022644, downloaded October 2016) to the genome using BWA
106 mem v0.7.12 [29] and STAR v020201 [30] respectively. Over 99% of reads from the
107 two short-insert libraries mapped to the assembly, suggesting that the vast majority
108 of the genome represented by these data has been assembled. In addition, 94.9% of
109 RNA-Seq reads mapped to the assembly, suggesting that the majority of transcribed
110 genes are present. Gene-level completeness was assessed using CEGMA v2.5 [31] and

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

111 BUSCO v2.0 [32]. The proportion of CEGMA genes “completely” recovered ($n = 248$)
112 was 81%, increasing to 97% when partially recovered genes are included. The
113 recovery of BUSCO genes specific to the metazoa ($n = 978$) was higher, at 98% for
114 complete genes, increasing to 99% when partial genes are included. An almost
115 complete set (99.2%) of BUSCO genes specific to the Arthropoda ($n = 1,066$) was also
116 recovered. In addition, CEGMA indicated a duplication rate of 1.1 while BUSCO
117 estimated only ~2% genes were present in multiple copies. The high complete
118 CEGMA/BUSCO scores suggest a good assembly that has captured the majority of
119 core metazoan/Arthropod genes in full-length, and that the fragmentation of genes
120 across multiple scaffolds is low. In addition, the low duplication rates suggest that
121 most genes are present in single copy, and thus that the genome does not include
122 significant duplicated segments representing alternative haplotypes.

123 **Annotation**

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

134 assembly. Overall, approximately one quarter of the assembly (122.6 Mb) was masked
135 from gene prediction (Table 3).

136 **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%

137

138 Gene finding was performed following a two-pass approach [35]. Initial gene-
139 models were constructed with MAKER v2.31 [36], using HMMs derived from SNAP
140 [37] and GeneMark-ES v4.3 [38] in conjunction with a set of assembled *B. anynana*
141 transcripts as evidence (Oostra et al., in preparation). MAKER gene-models were then
142 passed to AUGUSTUS v3.0.3 [39] for refinement, resulting in an initial set of 26,722
143 predicted protein-coding genes. A set of basic filters was applied to remove likely
144 spurious gene models (Table 4), resulting in the deletion of 4,080 gene models.
145 Protein sequences from the filtered 22,642 genes were annotated using BLAST
146 searches versus UniRef90 and the NCBI non-redundant protein database (nr), and
147 domains/motifs were described using InterProScan5 [40]. Summary statistics for the
148 22,642 predicted gene models are given in Table 5.

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

Category	Description	Number of genes
(a)	Single-exon	7,112
(b)	Small exon (< 9 bp)	1,866
(c)	Small intron (\leq 40 bp)	45
(d)	Short (CDS < 120 bp)	127
(e)	No hit to <i>nr</i>	6,532
(f)	Duplicate (\geq 98% identity over \geq 98% query length)	822
Total¹		4,080

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).

150

151 **Comparison to other lepidopteran genomes**

152 To ascertain the relative quality of the *B. anynana* v1.2 assembly, we compared our
153 results to nine other published lepidopteran genomes available on LepBase
154 (<http://lepbase.org/>) [41]: *Bombyx mori* ASM15162v1 [42], *Danaus plexippus* v3 [43],
155 *Heliconius melpomene* Hmel2 [44,45], *Lerema accius* v1.1 [46], *Melitaea cinxia*
156 MelCinx1.0 [47], *Papilio glaucus* v1.1 [48], *Papilio polytes* Ppol 1.0 [49], *Papilio xuthus*
157 Pap_xu_1.0 [49] and *Plutella xylostella* DBM_FJ_v1.1 [50]. The *B. anynana* v1.2
158 assembly was of high quality compared to other published genomes, with the
159 majority of the genome represented in a relatively small number of scaffolds despite
160 being only marginally smaller than the longest lepidopteran genome, *B. mori* (Figure
161 3a). Interestingly, *B. anynana* v1.2 encodes the highest number of proteins of the 10
162 species compared (Figure 3b). Despite measures to eliminate potentially spurious
163 ORFs caused by annotation error or by duplication, *B. anynana* encodes ~3,250 more

164 genes that the diamondback moth *P. xylostella*, and ~10,400 more than the
165 swallowtail *P. polytes*. It is tempting to attribute the apparently high number of genes
166 to the developmental plasticity and alternative seasonal forms with divergent
167 morphologies and life histories in *B. anynana*. However, it remains to be determined
168 whether the number of genes predicted in *B. anynana* is a function of its larger
169 genome size or unusual life-history characteristics, or if further curation of the v1.2
170 gene models will reduce the number of inferred genes.

171 **Concluding remarks**

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

178 **Availability of supporting data**

179 All raw sequence data have been deposited in the Short Read Archive (SRA) and are
180 available for download using the accession numbers provided in Table 1. The *B.*
181 *anynana* v1.2 assembly, as well as final predicted gene-models and protein
182 annotations, are publicly available for viewing and download via LepBase [41], an
183 Ensembl [51] genome database for the Lepidoptera
184 (<http://ensembl.lepbase.org/index.html>). A previous *B. anynana* assembly (nBa.0.1) is
185 also available on LepBase.

186 **Abbreviations**

187 Gb: Gigabase; Mb: Megabase; kb: kilobase; bp: base pairs; CEGMA: Core Eukaryotic
188 Genes Mapping Approach; BUSCO: Benchmarking Universal Single-Copy Orthologs;
189 ORF: Open reading frame; CDS: Coding sequence

190 **Competing interests**

191 The authors declare that they have no competing interests.

192 **Author contributions**

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

199 **Acknowledgements**

200 We thank Edinburgh Genomics and Genome Institute of Singapore for genome
201 sequencing, initial QC and data delivery. Funding for the *Bicyclus anynana* genome
202 project was provided by the ERC Advanced Grant number 250325 (EMARES) to PMB
203 and by the South East Asian Biodiversity Genomics Center (NUS grant R-154-000-

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

204 648-646 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; ¹⁰Department
216 of Zoology, University of Cambridge, Cambridge, CB2 3EJ, United Kingdom

217 **Tables**

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

220 **Figure legends**

221 **Figure 1:** Kmer frequency distribution for *B. anynana* short-insert libraries ($k = 31$).
222 The bimodality of the distribution, with peaks at approximately 105X and again at
223 210X, is the result of heterozygosity in the sequence data.

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

224 **Figure 2:** Taxon-annotated GC-coverage plots for **(a)** draft and **(b)** final *B. anynana*
225 genome assemblies. Each contig/scaffold in the assembly is represented as a circle,
226 coloured according to the best match to taxonomically annotated sequence
227 databases (see legends) and distributed according to the proportion GC (*x*-axis) and
228 read coverage (*y*-axis). The upper- and right-hand panels show the distribution of the
229 total span (kb) of contigs/scaffolds for a given coverage (upper panel) or GC (right
230 panel) bin. The heterozygosity in the sample is evident in the bimodal coverage
231 distribution seen in (a). The cluster of orange-coloured contigs at a lower coverage
232 and higher GC than the main cloud were likely derived from contaminant
233 *Enterococcus* present in the sample. The final assembly, (b), shows the effective
234 collapse of heterozygous regions, the removal of contaminant sequences and the
235 scaffolding of contigs into long contiguous sequences. Note that only taxon
236 annotations with a span > 1 Mb are shown in the legend for clarity.

237 **Figure 3:** Assembly and gene prediction comparison among 10 lepidopteran
238 genomes. **(a)** Cumulative assembly curves showing the relationship between the
239 number of scaffolds (*x*-axis) and the cumulative span of each assembly (*y*-axis),
240 coloured by species. Higher quality assemblies are represented by an almost-vertical
241 line (e.g., *H. melpomene* Hmel2 assembly in black), indicating a relatively small
242 number of scaffolds is required to reach the final genome span; conversely, a long tail
243 indicates the assembly includes a large number of smaller scaffolds. The curve for *B.*
244 *anynana* (brown and bold) suggests a good assembly for this species, with the
245 majority of the assembly comprised of relatively few scaffolds. **(b)** *B. anynana* v1.2
246 encodes the greatest number of genes of the 10 genomes, and is particularly

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

247 divergent from *B. mori*, which is of equivalent length. Species names/colours are as
248 follows: "bicyclus" (brown), *B. anynana*; "bombyx" (blue), *B. mori*; "danaus" (light
249 green), *D. plexippus*; "heliconius" (black), *H. melpomene*; "lerema" (dark green), *L.*
250 *accius*; "melitaea" (orange), *M. cinxia*; "glaucus" (red), *P. glaucus*; "polytes" (pink), *P.*
251 *polytes*; "xuthus" (violet), *P. xuthus*; "plutella" (grey), *P. xylostella*.

252

253 **References**

- 254 1. Brakefield PM, Beldade P, Zwaan BJ. The African butterfly *Bicyclus anynana*: a
255 model for evolutionary genetics and evolutionary developmental biology. Cold
256 Spring Harb Protoc. 2009; doi:10.1101/pdb.emo122-2.
- 257 2. Brakefield PM. Radiations of mycalesine butterflies and opening up their
258 exploration of morphospace. Am. Nat. 2010;176 Suppl 1:S77-87.
- 259 3. Prudic KL, Jeon C, Cao H, Monteiro A. Developmental plasticity in sexual roles of
260 butterfly species drives mutual sexual ornamentation. Science. 2011;331:73-5.
- 261 4. Westerman EL, Hodgins-Davis A, Dinwiddie A, Monteiro A. Biased learning affects
262 mate choice in a butterfly. Proc. Natl. Acad. Sci. U.S.A. 2012;109:10948-53.
- 263 5. Monteiro A. Origin, development, and evolution of butterfly eyespots. Annu. Rev.
264 Entomol. 2015;60:253-71.
- 265 6. Aduse-Poku K, Brakefield PM, Wahlberg N, Brattström O. Expanded molecular
266 phylogeny of the genus *Bicyclus* (Lepidoptera: Nymphalidae) shows the importance

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

267 of increased sampling for detecting semi-cryptic species and highlights potentials for
268 future studies. *Systematics and Biodiversity*. 2016;
269 doi:10.1080/14772000.2016.1226979.

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

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

276 9. Monteiro A, Tong X, Bear A, Liew SF, Bhardwaj S, Wasik BR, et al. Differential
277 Expression of Ecdysone Receptor Leads to Variation in Phenotypic Plasticity across
278 Serial Homologs. *PLoS Genet*. 2015;11:e1005529.

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

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

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

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

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

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

290 15. Bushnell B. BMap: BMap short read aligner. Available from:
291 sourceforge.net/projects/bbmap/

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

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

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

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

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

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

308 22. Boetzer M, Henkel CV, Jansen HJ, Butler D, Pirovano W. Scaffolding pre-

1
2
3
4 310 assembled contigs using SSPACE. *Bioinformatics*. 2011;27:578–9.
5
6 311 23. Xue W, Li J-T, Zhu Y-P, Hou G-Y, Kong X-F, Kuang Y-Y, et al. L_RNA_scaffolder:
7 scaffolding genomes with transcripts. *BMC Genomics*. 2013;14:604.
8
9 312 24. Koutsovoulos G. SCUBAT2. Available from: <https://github.com/GDKO/SCUBAT2>
10
11 313 25. Haas BJ, Papanicolaou A, Yassour M, Grabherr M, Blood PD, Bowden J, et al. De
12
13 314 novo transcript sequence reconstruction from RNA-seq using the Trinity platform for
14
15 315 reference generation and analysis. *Nat Protoc*. 2013;8:1494–512.
16
17 316 26. Boetzer M, Pirovano W. SSPACE-LongRead: scaffolding bacterial draft genomes
18
19 317 using long read sequence information. *BMC Bioinformatics*. 2014;15:211.
20
21 318 27. Boetzer M, Pirovano W. Toward almost closed genomes with GapFiller. *Genome*
22
23 319 *Biol*. 2012;13:R56.
24
25 320 28. English AC, Richards S, Han Y, Wang M, Vee V, Qu J, et al. Mind the Gap:
26
27 321 Upgrading Genomes with Pacific Biosciences RS Long-Read Sequencing Technology.
28
29 322 *PLoS ONE*. 2012;7.
30
31 323 29. Li H, Durbin R. Fast and accurate long-read alignment with Burrows-Wheeler
32
33 324 transform. *Bioinformatics*. 2010;26:589–95.
34
35 325 30. Dobin A, Davis CA, Schlesinger F, Drenkow J, Zaleski C, Jha S, et al. STAR: ultrafast
36
37 326 universal RNA-seq aligner. *Bioinformatics*. 2013;29:15–21.
38
39 327 31. Parra G, Bradnam K, Korf I. CEGMA: a pipeline to accurately annotate core genes
40
41 328 in eukaryotic genomes. *Bioinformatics*. 2007;23:1061–7.
42
43
44
45
46
47
48
49
50
51
52
53
54
55
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
- 329 32. Simão FA, Waterhouse RM, Ioannidis P, Kriventseva EV, Zdobnov EM. BUSCO:
330 assessing genome assembly and annotation completeness with single-copy
331 orthologs. *Bioinformatics*. 2015;31:3210–2.
- 332 33. Jurka J, Kapitonov VV, Pavlicek A, Klonowski P, Kohany O, Walichiewicz J. Repbase
333 Update, a database of eukaryotic repetitive elements. *Cytogenet Genome Res*.
334 2005;110:462–7.
- 335 34. Smit A, Hubley R, Green P. RepeatMasker. Available from
336 <http://www.repeatmasker.org>.
- 337 35. Koutsovoulos G. CGP-Pipeline. Available from: <https://gist.github.com/GDKO>
- 338 36. Holt C, Yandell M. MAKER2: an annotation pipeline and genome-database
339 management tool for second-generation genome projects. *BMC Bioinformatics*.
340 2011;12:491.
- 341 37. Korf I. Gene finding in novel genomes. *BMC Bioinformatics*. 2004;5:59.
- 342 38. Ter-Hovhannisyan V, Lomsadze A, Chernoff YO, Borodovsky M. Gene prediction in
343 novel fungal genomes using an ab initio algorithm with unsupervised training.
344 *Genome Res*. 2008;18:1979–90.
- 345 39. Stanke M, Diekhans M, Baertsch R, Haussler D. Using native and syntenically
346 mapped cDNA alignments to improve de novo gene finding. *Bioinformatics*.
347 2008;24:637–44.
- 348 40. Jones P, Binns D, Chang H-Y, Fraser M, Li W, McAnulla C, et al. InterProScan 5:
349 genome-scale protein function classification. *Bioinformatics*. 2014;30:1236–40.

- 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
- 350 41. Challis RJ, Kumar S, Dasmahapatra KKK, Jiggins CD, Blaxter M. Lepbase: the
351 Lepidopteran genome database. bioRxiv. 2016; doi: <http://dx.doi.org/10.1101/056994>
- 352 42. Duan J, Li R, Cheng D, Fan W, Zha X, Cheng T, et al. SilkDB v2.0: a platform for
353 silkworm (*Bombyx mori*) genome biology. Nucleic Acids Res. 2010;38:D453–6.
- 354 43. Zhan S, Merlin C, Boore JL, Reppert SM. The monarch butterfly genome yields
355 insights into long-distance migration. Cell. 2011;147:1171–85.
- 356 44. Heliconius Genome Consortium. Butterfly genome reveals promiscuous exchange
357 of mimicry adaptations among species. Nature. 2012;487:94–8.
- 358 45. Davey JW, Chouteau M, Barker SL, Maroja L, Baxter SW, Simpson F, et al. Major
359 improvements to the *Heliconius melpomene* genome assembly used to confirm 10
360 chromosome fusion events in 6 million years of butterfly evolution. G3. 2016;6:695–
361 708.
- 362 46. Cong Q, Borek D, Otwinowski Z, Grishin NV. Skipper genome sheds light on
363 unique phenotypic traits and phylogeny. BMC Genomics. 2015;16:639.
- 364 47. Ahola V, Lehtonen R, Somervuo P, Salmela L, Koskinen P, Rastas P, et al. The
365 Glanville fritillary genome retains an ancient karyotype and reveals selective
366 chromosomal fusions in Lepidoptera. Nature Communications. 2014;5:1–9.
- 367 48. Cong Q, Borek D, Otwinowski Z, Grishin NV. Tiger swallowtail genome reveals
368 mechanisms for speciation and caterpillar chemical defense. Cell Rep. 2015;10:910–9.
- 369 49. Nishikawa H, Iijima T, Kajitani R, Yamaguchi J, Ando T, Suzuki Y, et al. A genetic
370 mechanism for female-limited Batesian mimicry in *Papilio* butterfly. Nat Genet.

1
2
3
4 372 50. You M, Yue Z, He W, Yang X, Yang G, Xie M, et al. A heterozygous moth genome
5
6 373 provides insights into herbivory and detoxification. *Nat Genet.* 2013;45:220–5.
7
8
9
10 374 51. Yates A, Akanni W, Amode MR, Barrell D, Billis K, Carvalho-Silva D, et al. Ensembl
11
12 375 2016. *Nucleic Acids Res.* 2016;44:D710–6.
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

Tables

Table 1: Data counts and library information.

Library type	Platform	Read length	Insert (mean)	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	271,808,057 pairs	267,241,712 (98.3%)	66,334,099,834 (97.6%)	ERR1102671-2, ERR1102675-6
Short insert	Illumina HiSeq2500	125 bp paired-end	550 bp	241,050,065 pairs	234,269,871 (97.2%)	57,913,474,128 (96.1%)	ERR1102673-4, ERR1102677-8
Mate pair	Illumina HiSeq2500	100 bp paired-end	3 kb	77,105,680 pairs	31,848,200 (41.3%)	5,758,856,502 (37.3%)	ERR1750945
Mate pair	Illumina MiSeq	100 bp paired-end	3 kb	5,641,764 pairs	2,170,610 (38.5%)	397,993,018 (35.3%)	ERR754051
Mate pair	Illumina HiSeq2500	100 bp paired-end	5 kb	77,614,870 pairs	45,676,725 (58.9%)	8,203,769,131 (52.8%)	ERR1750946
Mate pair	Illumina	100 bp	5 kb	7,939,601 pairs	4,734,000 (59.6%)	861,352,793 (54.2%)	ERR754052

MiSeq

Long read	PacBio P6	0.80-50 kb	10kb	1,388,796	1199064 (86.3%)	4,086,394,966	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	23,699	10,682	10,682	3,100	48,180
N50 ¹	78.7 kb	111 kb	111 kb	328.9 kb	14.1 kb
NumN50 ²	1,543	8,075	548	214	7,366
Scaffolds					
Number	10,800	43,379	5,397	795	8,261
N50	638.3 kb	4,008.4 kb	715.6 kb	2,102.7 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	3,396
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> = 1,066)	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 genes	22642	15488	15130	15261	16751
Number of CDS	22642	19618	15130	13178	16668
Mean length	1.4 kb	1.6 kb	1.4 kb	1.3 kb	958 b
Median length	1.2 kb	1.2 kb	981 b	927 b	693 b
Min/max	84 b / 28.3 kb	23 b / 60.3 kb	9 b / 58.9 kb	45 b / 46.4 kb	6 b / 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 b	960/416 b	NA
Exons					
Length (mean/median)	208/126 b	283/161 b	206/149 b	284/157 b	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).

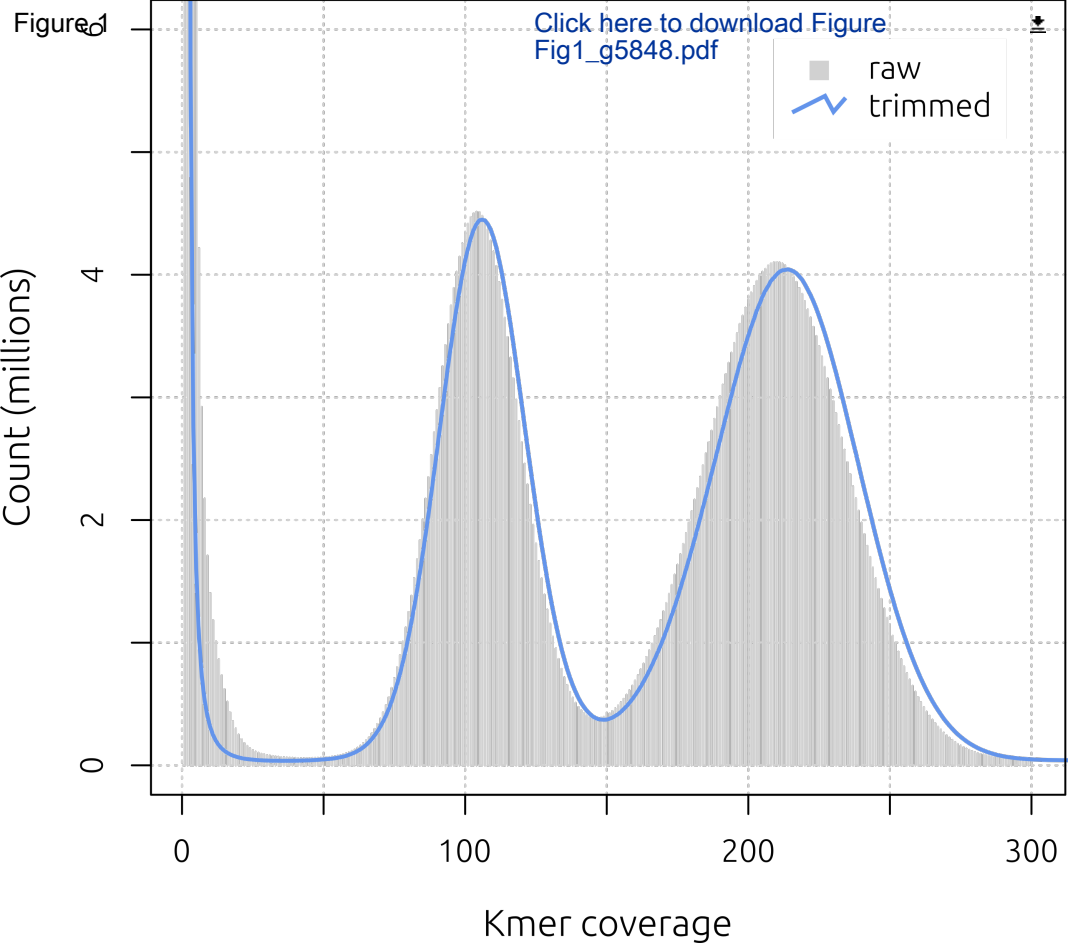
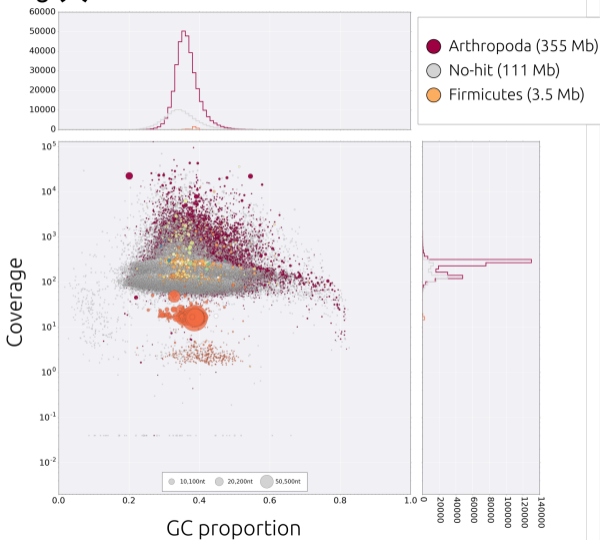


Fig (a) 2

Click (b) here to download Figure Fig2_g4698.pdf 