

1 **Recovering individual haplotypes and a contiguous genome**
2 **assembly from pooled long-read sequencing of the**
3 **diamondback moth (Lepidoptera: Plutellidae)**

4
5
6 Samuel Whiteford*¹, Arjen E. van't Hof¹, Ritesh Krishna^{1,5}, Thea Marubbi²,
7 Stephanie Widdison³, Ilik J. Saccheri¹, Marcus Guest⁴, Neil I. Morrison², Alistair C.
8 Darby¹

9
10 ¹University of Liverpool, Crown Street, Liverpool, L69 7ZB, UK

11
12 ²Oxitec Ltd., 71 Innovation Drive, Milton Park, Abingdon, OX14 4RQ, UK

13
14 ³General Bioinformatics, Jealott's Hill International Research Centre, Bracknell,
15 Berkshire RG42 6EY, UK.

16
17 ⁴Syngenta, Jealott's Hill International Research Centre, Bracknell, Berkshire RG42
18 6EY, UK.

19
20 ⁵IBM Research UK, STFC Daresbury Laboratory, Warrington, WA4 4AD, UK

21
22 *Corresponding author: s.whiteford@liverpool.ac.uk

23
24
25
26 Running Title: *Plutella xylostella* genome validation

27
28 Keywords: Pool-seq, Haplotype, Assembly, *Plutella xylostella*

29
30
31
32
33
34
35
36
37
38 © The Author(s) (2022) . Published by Oxford University Press on behalf of the Genetics Society of America.
39 This is an Open Access article distributed under the terms of the Creative Commons Attribution License
40 (<http://creativecommons.org/licenses/by/4.0/>), which permits unrestricted reuse, distribution, and reproduction in any medium,
41 provided the original work is properly cited.

45 **Abstract**

46

47 The assembly of divergent haplotypes using noisy long-read data presents a
48 challenge to the reconstruction of haploid genome assemblies, due to overlapping
49 distributions of technical sequencing error, intra-locus genetic variation and inter-
50 locus similarity within these data. Here we present a comparative analysis of
51 assembly algorithms representing overlap-layout-consensus, repeat graph and de
52 Bruijn graph methods. We examine how post-processing strategies attempting to
53 reduce redundant heterozygosity interact with the choice of initial assembly
54 algorithm and ultimately produce a series of chromosome-level assemblies for an
55 agricultural pest, the diamondback moth, *Plutella xylostella* (L.). We compare
56 evaluation methods and show that BUSCO analyses may overestimate haplotig
57 removal processing in long-read draft genomes, in comparison to a k-mer method.
58 We discuss the trade-offs inherent in assembly algorithm and curation choices
59 and suggest that “best practice” is research question dependent. We demonstrate
60 a link between allelic divergence and allele-derived contig redundancy in final
61 genome assemblies and document the patterns of coding and non-coding
62 diversity between redundant sequences. We also document a link between an
63 excess of non-synonymous polymorphism and haplotigs that are unresolved by
64 assembly or post-assembly algorithms. Finally, we discuss how this phenomenon
65 may have relevance for the usage of noisy long-read genome assemblies in
66 comparative genomics.

67

68 **Introduction**

69

70 Technical and analytical advances in genomics have dramatically improved the
71 achievable standard of genome projects. The amount of high molecular weight
72 (HMW) DNA required to perform long-read sequencing has reduced significantly
73 and has been accompanied by a steady increase in sequencing read lengths and
74 read accuracy(Kingan *et al.*, 2019). Longer reads have aided genome assembly
75 efforts by providing information linking unique genomic sequences flanking
76 repetitive elements, which represented a challenge to algorithms reliant on short-

77 read data (Koren *et al.*, 2017). However, early long-read methods, such as Oxford
78 nanopore and SMRT sequencing contained higher error rates in raw sequence
79 reads (Derrington *et al.*, 2010; Chin *et al.*, 2013). Whilst this technical noise can be
80 tolerated by various means (Chin *et al.*, 2013, 2016; Koren *et al.*, 2017;
81 Kolmogorov *et al.*, 2019; Ruan and Li, 2020), it is ultimately confounded with the
82 real biological variation present in the underlying samples. The form this biological
83 variation takes and the way it is distributed across the genome of an organism can
84 influence the accuracy of a reconstructed haploid (or phased diploid) genome
85 assembly (Kajitani *et al.*, 2019). Various assembly pipelines and algorithms have
86 been explicitly designed to overcome challenges of heterozygosity (Chin *et al.*,
87 2016; Huang *et al.*, 2017; Roach *et al.*, 2018), repeat resolution (Kolmogorov *et al.*,
88 2019), speed (Ruan and Li, 2020) and integration of multiple data-types (Ye *et al.*,
89 2016; Zimin *et al.*, 2017). Furthermore, all assembly software has some level of
90 parameterisation available to optimise results, yielding a huge array of possible
91 outcomes.

92
93 Alongside these computational innovations, several experimental approaches and
94 supplemental data types can augment existing data. Trio-sequencing can partition
95 heterozygous variation in an F1 individual using information from parental
96 haplotypes (Koren *et al.*, 2018). Linked-reads utilise microfluidics to uniquely
97 barcode reads that derive from discrete large DNA fragments, thereby capturing
98 longer range information than standard short read preparations do not (Zheng *et al.*,
99 2016). Chromosome conformation capture (Hi-C) data crosslinks *in vivo*
100 chromatin molecules and recovers pairs of reads that derive from these crosslinks,
101 producing data that reflects the 3D organisation of the nucleus, and also long-
102 range cis-chromosome associations (Ghurye *et al.*, 2018). In addition to the
103 proliferation of supporting data-types, the efficacy and quality of core genomic
104 data has improved. Improvements to data quality predominantly come from
105 platform advancements such as high-fidelity (HiFi) long reads (Nurk *et al.*, 2020;
106 Cheng *et al.*, 2021) and updated nanopore proteins (Karst *et al.*, 2021). Meanwhile
107 concurrent developments in library preparation and whole genome amplification
108 have enabled the use of decreased input DNA amounts (Schneider *et al.*, 2021). A

109 recent study found that the best lepidopteran genome assembly available at the
110 time utilised a combination of HiFi data and Hi-C (Ellis et al., 2021).

111

112 Within genome assembly, accounting for genomic variation is largely a technical
113 consideration. However, this variation isn't uniformly or randomly distributed and
114 is shaped by a range of evolutionary and demographic processes. One particularly
115 challenging aspect of genome assembly is the resolution of highly divergent
116 regions (HDRs) (Kajitani *et al.*, 2019), which often cannot be determined as allelic
117 within the assembly process and requires supervised analysis (Roach et al., 2018).
118 Genome assembly projects often aim to pre-emptively avoid this problem by
119 severely inbreeding the source material to increase the proportion of genome
120 homozygosity (The Heliconius Genome Consortium, 2012; Nowell *et al.*, 2017).
121 However, previous studies indicate that high levels of heterozygosity are often
122 counter-intuitively maintained despite multiple generations of sib-sib inbreeding
123 (The Heliconius Genome Consortium, 2012; Nowell *et al.*, 2017). A candidate for
124 such an effect is the presence of overdominant or pseudo-overdominant loci.
125 These loci, by various mechanisms, produce severe fitness consequences in a
126 homozygous state. In the case of pseudo-overdominance, the presence of tightly
127 linked recessive lethal mutations on different alleles prevents either haplotype
128 from becoming homozygous (Charlesworth and Willis, 2009). Alternatively,
129 pseudo-overdominance may be produced by multiple linked mildly deleterious
130 alleles, of which the cumulative effect is functionally equivalent to a single
131 recessive lethal. Whatever the fundamental cause, these phenomena can also
132 accumulate linked neutral variation, particularly in recombination cold-spots
133 (Zhao and Charlesworth, 2016). These features appear to make pseudo-
134 overdominance blocks (PODs) a plausible candidate for the HDRs known to
135 interfere with genome assembly (Waller, 2021).

136

137 If HDRs can indicate regions experiencing particular forms of selection, failure to
138 properly resolve them could impact downstream analyses, particularly the
139 detection of balancing selection and overdominant loci. Since this bias is non-
140 random, it may also affect comparative genome analyses, for example in instances

141 of balancing selection predating speciation, or other forms of trans-species
142 polymorphism, such as the well-studied MHC locus (Azevedo *et al.*, 2015).
143 Similarly, there may be common features of genetic architecture that may be
144 more likely to produce effects like overdominance or pseudo-overdominance at
145 common ancestral regions. Nonetheless, the increasing quality of long-read
146 sequencing, read-lengths and supporting data, should help to mitigate the issue,
147 and enable evaluation of the scale of this problem across historic genome
148 datasets. One method that is used widely in evaluating the completeness of
149 genome assemblies is the use of highly conserved gene sequences that are
150 consistently present as single copy genes (Simão *et al.*, 2015). However, in the
151 context of HDRs and their putative sources, it is possible that these methods may
152 be biased against representing genome regions that are more likely to harbour
153 HDRs. One potential resolution to this is to find genome validation methods that
154 do not rely on such cross-species inferences (Rhie *et al.*, 2020).

155

156 Here we investigate the complex trade-offs that are made in the choice of
157 genome assembly algorithm using a long-read genomic dataset for the
158 diamondback moth - *Plutella xylostella*, which was the subject of a previous major
159 genome sequencing effort, culminating in the publication of an assembly in 2013
160 (GCA_000330985.1) (You *et al.*, 2013). The assembly strategy utilised the
161 sequencing of fosmid, in order to mitigate the short read lengths of Illumina
162 sequencing. The authors report extensive structural variation based on alignments
163 between their assembly and both the fosmids and a previously sequenced BAC
164 (GenBank accession GU058050). The genome of *P. xylostella* therefore represents
165 two distinct challenges to current long-read assembly methods, namely a large
166 proportion of structural variation and a small amount of extractable DNA per
167 individual. Our study includes the additional challenge of sequencing the
168 heterogametic sex, containing the W-chromosome which has been shown to be
169 highly repetitive and intractable to assembly (24).

170

171 **Methods**

172

173 *Insect material origin and DNA extraction*

174

175 Starting material was provided by Oxitec Ltd. (Abingdon, U.K.) from a lab colony
176 that has been continuously cultured on artificial diet and is derived from the Vero
177 Beach strain (Martins *et al.*, 2012). Several lines were inbred in parallel by mating
178 sib-sib pairs each generation. A 7 generation inbred family was selected for
179 genome sequencing. DNA was extracted by phenol-chloroform from a pool of 15
180 sisters of the final inbred generation and a single male and female (Saccheri and
181 Bruford, 1993).

182

183

184 *Library construction and sequencing*

185

186 The pooled DNA was sheared to 7Kbp or 10Kbp. A subset was size selected at
187 15Kbp on the BluePippin (Sage Science, Inc.). In total 66 SMRT cells were
188 sequenced with P5-C3 chemistry on the RSII platform (Pacific Biosciences, Inc.).
189 Reads were filtered according to subread length (>50bp), polymerase read quality
190 (>75bp) and polymerase read length(>50bp). Extracted DNA from the individual
191 male and female was sheared and used for individual libraries followed by 2x
192 100bp paired-end Illumina sequencing (Illumina, Inc).

193

194 *Genome assembly parameters*

195

196 We performed assembly using canu (version 2.1.1), flye (version 2.8.2-b1689) and
197 wtdbg2 (version 2.5). These assemblies were subsequently polished with the
198 same pacbio read-set for two iterations using quiver (version 2.3.3) As a
199 preliminary step, we applied author recommended parameters for producing
200 separated haplotypes in the presence of heterozygosity and subsequently used
201 the resulting assembly with the highest rate of duplicated BUSCO genes (run as
202 described below). For flye and wtdbg, we selected the default parameter result,
203 for canu we selected the assembly using the parameter set [genomeSize=340m

204 corOutCoverage=200 correctedErrorRate=0.040 "batOptions=-dg 3 -db 3 -dr 1 -ca
205 500 -cp 50"].

206

207 *Haplotype merging*

208

209 We trialled two post-assembly haplotype merging procedures, purge_dups and
210 Haplomerger 2. Genomes processed with Haplomerger2 were first masked using
211 windowmasker (version 20120730) . A species-specific scoring matrix was inferred
212 at 95% identity using the lastz_D_Wrapper.pl script included with Haplomerger2
213 (Huang et al., 2017). The masked genome and scoring matrix were then used to
214 run scripts B1-B5 of the Haplomerger2 pipeline (version 20161205)(Huang et al.,
215 2017).

216

217 *Scaffolding*

218

219 The preparation of HiC libraries was performed by Dovetail LLC. using pools of
220 starved larvae. HiC libraries were prepared as described by Kalhor *et al.* (Kalhor *et*
221 *al.*, 2012). Both library preparations used the restriction enzyme DpnII for
222 digestion after proximity ligation. Scaffolding and misassembly detection was
223 performed by running the 3D-DNA pipeline on each of the haplotype merged
224 assembly versions.

225

226 *Validation procedures*

227

228 Here we quantify the the haplotype resolution processes using two independent
229 methods. Firstly, we utilised the gene-based BUSCO score (version 5.0.0) with the
230 "Lepidoptera odb10" database, consisting of 5286 gene groups and augustus
231 species model "heliconius_melpomene1". Secondly we utilised a combination of
232 the stacked k-mer coverage histograms (a.k.a. spectra-cn) plots generated with
233 KAT (version 2.4.2) and the read-based k-mer models produced by the
234 genomescope R script. In brief we used the R function "pmin" to intersect the
235 assembly copy number coverage distributions with the modelled distributions

236 from genomescope, specifically the error distribution, and the heterozygous and
237 homozygous components of the unique distribution. This provided a quantitative
238 k-mer based comparison of the haplotype resolution processes.

239

240 *Haplotype divergence assessment*

241

242 We quantified the divergence between duplicated genes identified by the BUSCO
243 analysis using an alignment-based and a supporting alignment-free method. The
244 amino-acid sequences of duplicated BUSCO gene copies were aligned with MAFFT
245 (version v6.864b), and subsequently translated into codon based alignments with
246 pal2nal.pl (version 14), followed by calculation of synonymous and non-
247 synonymous variants using the biopython function “cal_dn_ds”. For the
248 alignment-free comparison we used the full genomic sequence (including introns)
249 between the gene start and end coordinates identified by BUSCO and used the
250 python package “alfpy” with a word-size of two to calculate the Canberra distance
251 (see 28).

252

253 **Results**

254

255

256 *Insect materials, sequencing & assembly*

257

258

259 The material described in this study was inbred for 7 generations and is derived
260 from a long-term laboratory culture, itself derived from the “Vero Beach” strain
261 (U.S.A.). 15 sisters were pooled to meet minimum HMW DNA input requirements.
262 Initial assemblies showed that substantial genetic variation was retained and was
263 of sufficient complexity to produce multiple allelic sequence contigs from the
264 same locus, inflating the total size way beyond the expected size of 338.7 (+/-1.1)
265 Mbp, reported by Baxter *et al.* (2011). We subsequently trialled two approaches
266 to resolve the redundant heterozygosity Haplomerger2 and purge_dups (Huang et

267 al., 2017; Guan *et al.*, 2020). After filtering 2,655,788 PacBio subreads remained
268 (mean subread length=7,301bp, N50=10,398bp, total bases 19.4Gbp).

269

270

271 *Heterozygosity assessment*

272

273

274 All initial assembly strategies resulted in an over-inflated genome size, suggesting
275 differing amounts of redundant haplotig sequence (Fig. 1A). We determined
276 BUSCO results for the initial assemblies as evidence for the levels of allelic
277 redundancy (measured as duplicated BUSCO genes) and overall completeness
278 (Fig. 1B & 1C). We also utilised k-mer based methods. Firstly, a histogram of
279 corrected PacBio reads, provided an initial estimation of genome heterozygosity
280 as approximately 1.11%, which is moderately, but not exceedingly high for a North
281 American sample (see You *et al.*, 2020 for context) (Fig. 2A). Estimates from
282 related individuals (non-pooled) were 0.54% for a related inbred male and 1.00%
283 for a related inbred female (Sup. fig. 3). Stacked histograms coloured by assembly
284 coverage provided a qualitative assessment of genome assembly completeness
285 and redundant allelic variation (Fig. 2B). Secondly, we intersected the modelled
286 distributions of homozygous, heterozygous and sequencing error content from
287 genomescope with stacked histograms in order to make quantitative comparisons
288 of the initial assemblies (methodology illustrated in Fig. 1A, data shown in Fig. 2C)
289 (Vurture *et al.*, 2017).

290

291 The WTDBG2 assembly was the smallest in size (427Mbp) and contig number
292 (4023) (Fig. 1A). It had the lowest number of duplicated BUSCO genes (805) and
293 highest number of missing genes (110) (Fig. 1B). Consistent with the BUSCO
294 results, WTDGB2 had the lowest number of homozygous k-mers duplicated in the
295 assembly (Fig. 2B). But it also had the highest number of modelled error k-mers
296 present (Fig. 2C). In contrast, the Flye assembly was the largest in size (494Mbp)
297 and contig number (5985). It had the most duplicated BUSCO genes (1324) and
298 the least missing genes (88). Again the k-mer results show concordance with the

299 highest number of homozygous k-mers present duplicated in the assembly.
300 However, the number of modelled error k-mers present in the assembly was
301 comparable with the canu assembly. Canu produced intermediate values in total
302 size (448Mbp) and contig number (5341). Similarly, BUSCO results indicated an
303 intermediate number of duplicated genes (1105) and missing genes (106). k-mer
304 results followed the same pattern except for error k-mers in the final assembly.

305

306 Both post-assembly allelic redundancy approaches reduced the overall sizes of the
307 assemblies and appears to follow the patterns observed in the initial assemblies,
308 such that WTDBG2 still retains the lowest number duplicated and highest number
309 of missing genes in contrast with Flye. For each of the three starting assemblies,
310 Haplomerger2 provided a greater reduction in total size and number of contigs
311 compared to purge_dups (Fig. 1A). When applied to the canu assembly we also
312 observe an increase in contiguity (Fig. 1A), due to a tiling effect produced when
313 corresponding redundant heterozygous regions are merged at the ends of contigs
314 (Supp. Fig. 1)

315

316 Post-assembly processing resolved most duplicated BUSCO genes to a single copy
317 regardless of the initial assembly algorithm, however in all cases, the numbers of
318 missing BUSCO genes also increased (Fig. 1B & 1C). We observe that purge_dups
319 resolved some duplications that are not resolved by Haplomerger2 and vice versa
320 (Fig. 1C). Similarly, genes that go from complete and single copy in the primary
321 assembly to fragmented or missing after post-processing are not necessarily the
322 same across the two methods (Fig. 1C). This suggests that removal of redundancy
323 is not simply, more or less “aggressive”, and that performance varies by algorithm
324 depending on specific sequence properties.

325

326 *Comparison of heterozygosity assessments*

327

328

329 Using the k-mer intersection approach described in Figure 2, we produce a k-mer
330 proxy of BUSCO genes for comparison. The proxy is calculated using the modelled
331 homozygous k-mers (analogous to single copy genes) and divides the occurrence

332 of duplicated assembly k-mers by the sum of the single copy and duplicated
333 assembly k-mers (Supp. Tab. 2). Levels of percent duplication in the initial
334 assemblies are remarkably concordant between the genic (BUSCO) and unbiased
335 (k-mer) methods (Supp. Tab. 2), with the exception of flye. This exception is likely
336 due to the relatively higher occurrence of three-copy redundancy observed with
337 the flye assembly algorithm (Fig. 2B), which are not captured in our k-mer proxy
338 measurement (BUSCO duplications do not distinguish 2 copy genes from >2 copy
339 genes). However, after assembly post-processing to remove redundant haplotigs,
340 BUSCO genes appear to overestimate the efficiency of the procedures in
341 comparison to the k-mer proxy. Across all methods and initial starting assemblies,
342 the k-mer proxy shows consistently higher residual duplication than suggested by
343 BUSCO genes (Supp. Tab. 2).

344

345

346 *HiC misassembly detection and scaffolding*

347

348

349 We observed the greatest overall number of detected misassembled region
350 candidates in ‘canu + purge_dups’ after two iterations of the HiC scaffolding
351 pipeline 3d-DNA. The least misassembled region candidates detected after two
352 iterations was in ‘wtdbg + HM2’, followed by ‘canu + HM2’ (Tab. 1). We find the
353 greatest disparity in total misassembled region candidates between post-
354 processing methods in the canu assemblies. Furthermore, we find that ‘canu +
355 purge_dups’ produced the lowest final N50 value, whereas all other assemblies
356 produced very similar results, though the metric is limited by karyotype at this
357 resolution (Tab. 1).

358

359

360 *Patterns of divergence between redundant alleles*

361

362

363 For BUSCO genes that were duplicated in the initial assembly and subsequently
364 reduced to single copy by the post-processing methods, we broadly describe the

365 variation between the copies using the ratio of non-synonymous and synonymous
366 nucleotide diversity and an alignment free method using the entire genomic
367 region (Zielezinski *et al.*, 2019). We observed that the distributions of genes
368 remaining after the application of `purge_dups` were more heavily weighted
369 toward a low k-mer based distance and a π_N/π_S ratio of 0 as compared to the
370 distributions of genes de-duplicated by Haplomerger2 (Fig. 3). We partitioned
371 duplicated BUSCOs depending on whether they are present on the same initial
372 assembly contig or not, as these genes may plausibly be real duplication events
373 rather than haplotypic redundancy. The distribution of π_N/π_S between putative
374 tandem duplicated BUSCOs (occurring on the same assembly contig) appears to
375 be somewhat inflated in comparison to putative haplotig BUSCO genes. This
376 pattern is also reflected in the overall genomic DNA k-mer distances (Fig. 3).
377

378 Discussion

379

380 *Plutella xylostella* populations harbour large amounts of polymorphism (You *et al.*,
381 2013, 2020). We observed a relatively low heterozygosity in our pooled data
382 compared to other species, however this individual should not be considered
383 representative of the wild population due to severe inbreeding and prior
384 laboratory domestication. Despite this apparently reduced heterozygosity, a large
385 amount of redundant sequence remains after genome assembly, suggesting that
386 the heterozygosity is largely co-localised in highly divergent alleles. This pattern
387 may suggest regions of low recombination, enabling haplotypes to accumulate
388 linked neutral variation and persist through drift. Alternatively, in the case of
389 associative overdominance, neutral variation can accumulate alongside linked
390 overdominant or pseudo-overdominant loci (linked deleterious recessives with
391 opposing phase) (Ohta and Kimura, 1970). It is important that such regions are
392 represented appropriately in genome assemblies, as downstream analyses
393 involving mapping reads rely on both overall completeness and regions being
394 present in a haploid state (although see (Hickey *et al.*, 2020) for how this is
395 changing).

396

397 We tested two post-assembly redundancy reduction procedures (Haplomerger2 &
398 purge_dups) and found that Haplomerger2 generally appears to “resolve” more
399 redundant sequence, at the expense of erroneous removal of non-redundant
400 genome content and erroneous scaffolding of overlapping divergent regions. Both
401 programs utilise a self-alignment step to detect haplotigs, purge_dups then
402 implements a further QC step to these results by assessing the coverage of the
403 identified haplotigs. For self-alignment, Haplomerger2 utilises LASTZ and enables
404 users to calculate and use a sample specific scoring matrix, whilst purge_dups
405 utilises minimap2 with a fixed intra-species scoring parameter (asm5). The
406 parameterisation reflects a balance in differentiating intra-locus divergence, from
407 inter-locus paralogue similarity. To give specific examples; ancient balancing
408 selection vs relatively recent gene duplication or ancient balancing selection vs
409 genetic convergence. The idealised genome assembly or redundancy removal
410 pipeline can accurately differentiate these effects.

411

412 Genic analyses of assembly completeness such as BUSCO are widely used and
413 relatively straight forward to apply, however by definition they are limited to
414 genomic regions containing coding sequences (Simão *et al.*, 2015). The genes are
415 highly conserved at the amino-acid level, suggesting that non-synonymous
416 substitutions are largely deleterious. Because of this, the surrounding genomic
417 region (including non-coding variation) may be likely to harbour less variation
418 than a neutral region, due to the action of background selection (Gilbert *et al.*,
419 2020). In short, BUSCO genes are likely to inhabit (and help maintain) conserved
420 genomic regions. The practical implication is that BUSCO genes, when utilised to
421 assess the removal of redundant haplotigs, may systematically overestimate the
422 effectiveness of the procedure, as they are unlikely to represent HDRs. Indeed,
423 our results support the notion that before and after BUSCO duplication scores
424 overestimate the removal of redundant haplotig sequences when compared to an
425 analogous k-mer estimator. If BUSCO duplication results are liable to overestimate
426 the haploid nature of a given draft genome assembly, it may hamper comparative
427 genomic efforts to identify

428 balancing selection or overdominance (which may have either common or
429 independent origins).

430

431 Despite this potential limitation, BUSCO scores are still useful as a guide to
432 assembly completeness. BUSCO scores also provided insights into assembly post-
433 processing, showing that, despite resolving more duplications than `purge_dups`,
434 Haplomerger2 results do not completely overlap those of `purge_dups`. This
435 indicates that both underlying methods are sub-optimal and the results may be
436 complementary. We also note that BUSCO results, particularly missing genes, are
437 dependent on the optimisation of input parameters. For example, the “`--long`”
438 parameter can increase sensitivity at the cost of greater runtime. Similarly,
439 detection of BUSCO genes may differ between haplotypes, thereby
440 underestimating the number of duplicated genes.

441

442 We demonstrate a supporting validation method, providing relative quantification
443 of assembly accuracy, using overlaps between modelled k-mer distributions and
444 the k-mer frequency histogram subdivided by numerical representation in the
445 final assembly version. Whilst this assessment applies to any non-repetitive
446 genome region it only offers a general comparative measure between assemblies
447 from the same read set and cannot determine appropriate representation at
448 specific sites, due to stochasticity in read coverage. However, the ability to
449 confidently determine truly heterozygous k-mers from homozygous will be
450 increased in low-error, high-coverage read datasets, such as those currently being
451 generated by projects like DToL (Darwin Tree of Life). This would offer an
452 independent and unbiased validation method, but only with sufficiently high
453 coverage to accurately partition the different k-mer peaks.

454

455 Our initial expectation was that a greater reduction in redundancy may
456 correspond with an increase in detected misassembled region candidates,
457 particularly in the case of Haplomerger2, which can join overlapping contigs (Sup.
458 Fig. 1). Instead, we find the opposite pattern, though this does not necessarily
459 imply more accurate assembly representation, since complex regions may be

460 absent from a final assembly altogether. For example, the lowest number of
461 misassemblies ('wtdbg + HM2'), occurred alongside both the lowest putative
462 allelic redundancy, but also the greatest values for missing BUSCO genes and
463 het/hom modelled k-mers with 0x assembly coverage).

464

465 After an appraisal of the results of haplotig resolution, we compared the overall
466 divergence of duplicated genes from the two methods. The results of purge_dups
467 retained a greater proportion of low divergence haplotigs and this also
468 corresponded to genes with a lower proportion of non-synonymous substitutions.
469 The remaining genes in both sets suggests that both methods did not resolve
470 more greatly diverged sequences and that this divergence corresponded to
471 elevated non-synonymous substitutions relative to synonymous substitutions.
472 Taken together with the high levels of coding sequence conservation intrinsic to
473 BUSCO genes, this pattern would appear consistent with pseudo-overdominant
474 regions generated by the linked arrangement of multiple deleterious non-
475 synonymous substitutions. However, additional investigations with
476 supplementary data will be required to establish this with confidence and
477 determine the processes responsible for these patterns.

478

479 **Conclusions**

480

481 Highly divergent alleles can pose a challenge to accurate haploid reconstruction
482 from noisy long read data. Post-processing can mitigate these problems
483 somewhat to produce mosaic resolved sequences for reference purposes,
484 however results in our case are largely imperfect and present a set of complex
485 trade-offs between assembly completeness, redundancy and mis-assembly.
486 Researchers producing or using genomes should be aware of these issues when
487 using genome assembly data derived from noisy long-reads, especially when
488 investigating genomic regions likely to harbour significant linked variation. Our
489 results lead us to the conclusion that unresolved HDRs may be widespread in draft
490 genomes assembled from noisy long-read data and that BUSCO analyses may

491 overestimate their resolution by post-processing methods. Plausible causes
492 include loci experiencing balancing-selection or overdominance effects that
493 originated prior to speciation events, or that exist within a genetic architecture
494 liable to parallel origins of these processes. This may impact comparative genomic
495 studies that aim to identify or describe these evolutionary processes, however
496 further investigation is required to examine this.

497

498 Finally as a recommendation to researchers utilising similar data, we suggest that
499 the optimal strategy is research question dependent. Our data shows that there
500 are complex trade-offs between gene set completeness, the presence and
501 abundance of redundant haplotigs, and overall genome contiguity. We list some
502 recommendations resulting from our dataset: 1) For comparative analyses of
503 large-scale shared-synteny or chromosome-level structures, we would
504 recommend wtdbg2 followed by Haplomerger2, however it should be stressed
505 that for this type of analysis researchers should supplement long-read data with
506 HiC, both to extend contigs into larger-scale scaffolds, but also to correct any
507 erroneous assembly, or post-processing mis-joins. When HiC data is available, the
508 contiguity differences between different assemblers becomes less important,
509 however wtdbg2 followed by Haplomerger2 should still reduce misleading inter-
510 specific alignment signals produced by residual haplotigs. 2) For comparative
511 analysis of orthologues, the decision is complicated, as there is a trade-off
512 between false-positive paralogues due to redundant haplotigs, vs false-negative
513 missing genome content that is eliminated by redundancy removal procedures. 3)
514 For analysis of a particular gene of interest, researchers can assemble their data
515 with flye or canu, and process the resulting assembly with purge_dups. Reference
516 to both the initial assembly and the post-processed data, should enable
517 researchers to recover their gene of interest and examine whether any allelic
518 variation is present. 4) For researchers wishing to produce a multi-purpose
519 reference assembly, with no specific research question, we would suggest
520 producing detailed and transparent methods, such that subsequent users can
521 understand the limitations and reanalyse for their specific purpose if necessary.

522

524 Declarations

525

526 Ethics approval and consent to participate:

527 N/A

528

529 Consent for publication:

530 All authors agreed to the publication of this manuscript.

531

532 Availability of data and materials:

533 The read datasets generated during the current study are available in the ENA
534 database under accession PRJEB34571 (see Supp. Tab. 1). All assembly versions535 are available at **10.5281/zenodo.5647466**. Code provided at536 https://github.com/swomics/Plutella_genomes

537

538

539 Competing interests:

540 Samuel Whiteford's studentship was co-funded by Oxitec

541 Thea Marubbi and Neil I. Morrison are employees of Oxitec

542 Ritesh Krishna is an employee of IBM U.K.

543 Stephanie Widdison is an employee of General Bioinformatics

544 Marcus Guest is an employee of Syngenta

545 Alistair Darby has held multiple grants with Oxitec

546

547 Funding:

548 NIM and ACD received funding from the BBSRC and Innovate UK grants

549 BB/M001512/1 & BB/M503472/1.

550

551

552 Authors' contributions:

553 SW, NIM and ACD designed and planned the project. SW performed analyses with

554 ACD. AH, SH, IS and MG prepared material for sequencing. The manuscript was

555 written by SW, IJS and ACD with input from all authors.

556

557 Acknowledgements: PacBio and Illumina sequencing were performed at the

558 Centre for Genomics Research, University of Liverpool by Margaret Hughes and

559 Lucille Rainbow. We thank Nerys Humphrey-Jones & Adam S. Walker for their

560 work maintaining insect strains. We also thank Matthew R. Gemmell for providing

561 a template script that was modified for the quiver polishing procedure.

562

563 Authors' information:

564

565 Samuel Whiteford*¹, Arjen E. van't Hof^{1,5}, Ritesh Krishna^{1,6}, Thea Marubbi²,566 Stephanie Widdison³, Ilik J. Saccheri¹, Marcus Guest⁴, Neil I. Morrison², Alistair C.567 Darby¹

- 568
569 1. Institute of Integrative Biology, University of Liverpool, Crown Street,
570 Liverpool, L69 7ZB, UK
571 2. Oxitec Ltd., 71 Innovation Drive, Milton Park, Abingdon, OX14 4RQ, UK
572 3. General Bioinformatics, Jealott's Hill International Research Centre,
573 Bracknell, Berkshire RG42 6EY, UK.
574 4. Syngenta, Jealott's Hill International Research Centre, Bracknell, Berkshire
575 RG42 6EY, UK.
576 5. IBM Research UK, STFC Daresbury Laboratory, Warrington, WA4 4AD, UK
577

578 *corresponding author

579
580 s.whiteford@liverpool.ac.uk
581 arjenvth@hotmail.com
582 ritesh.krishna@uk.ibm.com
583 thea.marubbi@oxitec.com
584 stephanie.widdison@syngenta.com
585 saccheri@liverpool.ac.uk
586 marcus.guest@syngenta.com
587 neil.morrison@oxitec.com
588 acdarby@liverpool.ac.uk
589

590 **References**

- 591
592
593
594 Armstrong, J. *et al.* (2020) 'Progressive Cactus is a multiple-genome aligner for the
595 thousand-genome era', *Nature*. Springer US, 587(7833), pp. 246–251. doi:
596 10.1038/s41586-020-2871-y.
597 Azevedo, L. *et al.* (2015) 'Trans-species polymorphism in humans and the great
598 apes is generally maintained by balancing selection that modulates the host
599 immune response', *Human Genomics*. Human Genomics, 9(1), pp. 4–9. doi:
600 10.1186/s40246-015-0043-1.
601 Baxter, S. W. *et al.* (2011) 'Linkage mapping and comparative genomics using next-
602 generation RAD sequencing of a non-model organism', *PLoS one*, 6(4), p. e19315.
603 doi: 10.1371/journal.pone.0019315.
604 Charlesworth, D. and Willis, J. H. (2009) 'The genetics of inbreeding depression',
605 *Nature Reviews Genetics*, 10(11), pp. 783–796. doi: 10.1038/nrg2664.
606 Cheng, H. *et al.* (2021) 'Haplotype-resolved de novo assembly using phased
607 assembly graphs with hifiasm', *Nature Methods*. Springer US, 18(2), pp. 170–175.
608 doi: 10.1038/s41592-020-01056-5.
609 Chin, C.-S. *et al.* (2013) 'Nonhybrid, finished microbial genome assemblies from
610 long-read SMRT sequencing data.', *Nature methods*, 10(6), pp. 563–569. doi:
611 10.1038/nmeth.2474.

612 Chin, C.-S. *et al.* (2016) 'Phased diploid genome assembly with single-molecule
613 real-time sequencing.', *Nature methods*. Nature Publishing Group, (October), pp.
614 1–7.

615 Derrington, I. M. *et al.* (2010) 'Nanopore DNA sequencing with MspA', *PNAS*,
616 107(37), pp. 16060–16065. doi: 10.1073/pnas.1001831107.

617 Ellis, E. A., Storer, C. G. and Kawahara, A. Y. (2021) 'De novo genome assemblies of
618 butterflies', *GigaScience*. Oxford University Press, 10(6), pp. 1–8. doi:
619 10.1093/gigascience/giab041.

620 Ghurye, J. *et al.* (2018) 'Integrating Hi-C links with assembly graphs for
621 chromosome-scale assembly', pp. 1–9.

622 Gilbert, K. J. *et al.* (2020) 'Transition from Background Selection to Associative
623 Overdominance Promotes Diversity in Regions of Low Recombination', *Current*
624 *Biology*. Elsevier Ltd., 30(1), pp. 101-107.e3. doi: 10.1016/j.cub.2019.11.063.

625 Guan, D. *et al.* (2020) 'Identifying and removing haplotypic duplication in primary
626 genome assemblies', *Bioinformatics*, 36(9), pp. 2896–2898. doi:
627 10.1093/bioinformatics/btaa025.

628 Hickey, G. *et al.* (2020) 'Genotyping structural variants in pangenome graphs using
629 the vg toolkit', *Genome Biology*. Genome Biology, 21(35). doi: 10.1101/654566.

630 Huang, S., Kang, M. and Xu, A. (2017) 'HaploMerger2 : rebuilding both haploid
631 sub-assemblies from high-heterozygosity diploid genome assembly',
632 *Bioinformatics*, 33(16), pp. 2577–2579. doi: 10.1093/bioinformatics/btx220.

633 Kajitani, R. *et al.* (2019) 'Platanus-allee is a de novo haplotype assembler enabling
634 a comprehensive access to divergent heterozygous regions', *Nature*
635 *Communications*. Springer US, 10, p. 1702. doi: 10.1038/s41467-019-09575-2.

636 Kalhor, R. *et al.* (2012) 'Genome architectures revealed by tethered chromosome
637 conformation capture and population-based modeling', *Nature Biotechnology*.
638 Nature Publishing Group, 30(1), pp. 90–98. doi: 10.1038/nbt.2057.

639 Karst, S. M. *et al.* (2021) 'High-accuracy long-read amplicon sequences using
640 unique molecular identifiers with Nanopore or PacBio sequencing', *Nature*
641 *Methods*, 18(2), pp. 165–169. doi: 10.1038/s41592-020-01041-y.

642 Kingan, S. B. *et al.* (2019) 'A High-Quality De novo Genome Assembly from a Single
643 Mosquito Using PacBio Sequencing', *Genes*, 10, p. 62. doi:
644 10.3390/genes10010062.

645 Kolmogorov, M. *et al.* (2019) 'Assembly of long, error-prone reads using repeat
646 graphs', *Nature Biotechnology*. Springer US, 37(5), pp. 540–546. doi:
647 10.1038/s41587-019-0072-8.

648 Koren, S. *et al.* (2017) 'Canu : scalable and accurate long-read assembly via
649 adaptive k -mer weighting and repeat separation', *Genome research*, 27, pp. 722–
650 736. doi: 10.1101/gr.215087.116.Freely.

651 Koren, S. *et al.* (2018) 'De novo assembly of haplotype-resolved genomes with trio
652 binning', *Nature Biotechnology*, 36, pp. 1174–1182.

653 Martins, S. *et al.* (2012) 'Germline transformation of the diamondback moth,
654 *Plutella xylostella* L., using the piggyBac transposable element.', *Insect molecular*
655 *biology*, 21(4), pp. 414–421.

656 Nowell, R. W. *et al.* (2017) 'A high-coverage draft genome of the mycalesine
657 butterfly *Bicyclus anynana*', *Giga Science*, 6, pp. 1–7.

658 Nurk, S. *et al.* (2020) 'HiCanu: Accurate assembly of segmental duplications,

659 satellites, and allelic variants from high-fidelity long reads', *Genome Research*,
660 30(9), pp. 1291–1305. doi: 10.1101/GR.263566.120.

661 Ohta, T. and Kimura, M. (1970) 'Development of associative overdominance
662 through linkage disequilibrium in finite populations', *Genetical Research*.
663 University of Liverpool Library, 16(2), pp. 165–177. doi:
664 10.1017/S0016672300002391.

665 Rhie, A. *et al.* (2020) 'Mercury: Reference-free quality, completeness, and phasing
666 assessment for genome assemblies', *Genome Biology*. *Genome Biology*, 21(1), pp.
667 1–27. doi: 10.1186/s13059-020-02134-9.

668 Roach, M. J., Schmidt, S. A. and Borneman, A. R. (2018) 'Purge Haplotigs: allelic
669 contig reassignment for third-gen diploid genome assemblies', *BMC*
670 *Bioinformatics*. *BMC Bioinformatics*, 19, p. 460.

671 Ruan, J. and Li, H. (2020) 'Fast and accurate long-read assembly with wtdbg2',
672 *Nature Methods*. Springer US, 17(2), pp. 155–158. doi: 10.1038/s41592-019-0669-
673 3.

674 Saccheri, I. J. and Bruford, M. W. (1993) 'DNA fingerprinting in a butterfly, *Bicyclus*
675 *anyana* (Satyridae)', *Journal of Heredity*, 84(3), pp. 195–200. doi:
676 10.1093/oxfordjournals.jhered.a111316.

677 Schneider, C. *et al.* (2021) 'Two high-quality de novo genomes from single
678 ethanol-preserved specimens of tiny metazoans (*Collembola*)', *GigaScience*, 10(5),
679 pp. 1–12. doi: 10.1093/gigascience/giab035.

680 Simão, F. A. *et al.* (2015) 'BUSCO : assessing genome assembly and annotation
681 completeness with single-copy orthologs', *Bioinformatics*, 31(19), pp. 3210–3212.
682 doi: 10.1093/bioinformatics/btv351.

683 The Heliconius Genome Consortium (2012) 'Butterfly genome reveals
684 promiscuous exchange of mimicry adaptations among species.', *Nature*. Nature
685 Publishing Group, 487(7405), pp. 94–98. doi: 10.1038/nature11041.

686 Traut, W. *et al.* (2013) 'High-throughput sequencing of a single chromosome: a
687 moth *W* chromosome.', *Chromosome research*, 21(5), pp. 491–505. doi:
688 10.1007/s10577-013-9376-6.

689 Vurture, G. W. *et al.* (2017) 'GenomeScope: Fast reference-free genome profiling
690 from short reads', *Bioinformatics*, 33(14), pp. 2202–2204. doi:
691 10.1093/bioinformatics/btx153.

692 Waller, D. M. (2021) 'Addressing Darwin's dilemma: Can pseudo-overdominance
693 explain persistent inbreeding depression and load?', *Evolution*, 75(4), pp. 779–
694 793. doi: 10.1111/evo.14189.

695 Ye, C. *et al.* (2016) 'DBG2OLC: Efficient assembly of large genomes using long
696 erroneous reads of the third generation sequencing technologies', *Scientific*
697 *Reports*. Nature Publishing Group, 6(7), pp. 1–9. doi: 10.1038/srep31900.

698 You, M. *et al.* (2013) 'A heterozygous moth genome provides insights into
699 herbivory and detoxification.', *Nature genetics*. Nature Publishing Group, 45(2),
700 pp. 220–225. doi: 10.1038/ng.2524.

701 You, M. *et al.* (2020) 'Variation among 532 genomes unveils the origin and
702 evolutionary history of a global insect herbivore', *Nature Communications*.
703 Springer US, 11(1). doi: 10.1038/s41467-020-16178-9.

704 Zhao, L. and Charlesworth, B. (2016) *Resolving the conflict between associative*
705 *overdominance and background selection*, *Genetics*. doi:

706 10.1534/genetics.116.188912.
707 Zheng, G. X. Y. *et al.* (2016) 'Haplotyping germline and cancer genomes with high-
708 throughput linked-read sequencing', *Nature Biotechnology*. Nature Publishing
709 Group, 34(3), pp. 303–311. doi: 10.1038/nbt.3432.
710 Zieleszinski, A. *et al.* (2019) 'Benchmarking of alignment-free sequence comparison
711 methods', *Genome biology*. Genome Biology, 20, p. 144. doi: 10.1101/611137.
712 Zimin, A. V *et al.* (2017) 'Hybrid assembly of the large and highly repetitive
713 genome of *Aegilops tauschii*, a progenitor of bread wheat, with the MaSuRCA
714 mega-reads algorithm', *Genome research*, 27, pp. 1–6. doi:
715 10.1101/gr.213405.116.27.

716
717

718 **Figure 1. Contiguity and BUSCO content and of alternative genome assembly**
719 **methods and the effects of removing putative allelic redundancy.** In each panel,
720 "canu", "flye" and "wtdbg" refer to the preliminary assemblies produced by each
721 algorithm. "+ purge_dups + HiC" refers to these same assemblies with the
722 additional application of the purge_dups program followed by HiC scaffolding or,
723 Haplomerger2 followed by HiC scaffolding (A) depicts the differences in overall
724 contig size and contiguity between the different methods. The dotted curve
725 describes a previously published reference genome (accession:
726 GCA_000330985.1). The dashed straight line indicates the estimated genome size
727 from an independent flow cytometry estimate (Baxter *et al.*, 2011). (B) shows
728 overall BUSCO scores from a database of 5286 genes. BUSCO Scores from the
729 aforementioned accession are also included. (C) details the relationships of genes
730 within these sets. Groups of genes are coloured by BUSCO score in the initial
731 assembly. BUSCO genes that are single-copy and complete in all assemblies are
732 omitted to emphasise differences between assemblies.

733

734 **Figure 2. A k-mer based validation of the alternative genome assembly methods**
735 **and effects of removing putative allelic redundancy.** (A) shows an example of
736 stacked k-mer distributions subdivided by assembly representation (spectra-cn
737 plot) and an overlay of the modelled contributions of sequencing errors,
738 heterozygous content and homozygous content (dotted lines from left to right
739 respectively). (B) shows the spectra-cn plots for each of the assembly versions (C)
740 shows the number of k-mers present in the intersections between the modelled k-
741 mer content distributions and individual assembly coverage categories present in
742 the spectra-cn plots.

743

744 **Figure 3. Quantifying divergence between duplicated BUSCO genes.** (A) shows
745 the distribution of π_N/π_S scores for duplicated (N. copies = 2) BUSCO genes
746 remaining after the application of purge_dups or Haplomerger2. (B) shows and
747 alignment free quantification of the dissimilarity of intronic and exonic sequence
748 between the same duplicated BUSCO genes (see methods for details). Panels
749 labelled "Tandem" indicate that the BUSCO copies were found on the same
750 assembly contig, whereas "Unique" indicates that the copies were found on
751 different assembly contigs.

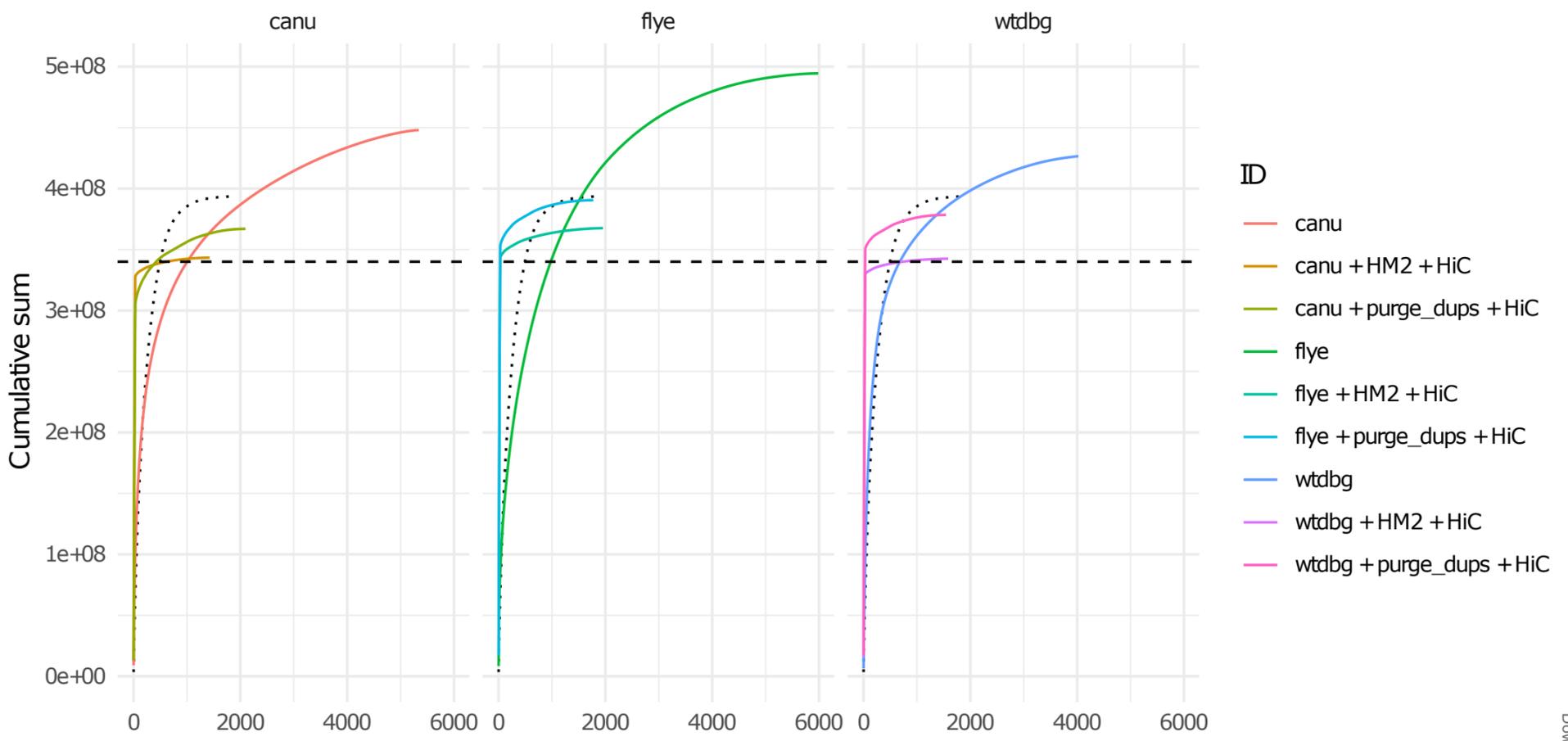
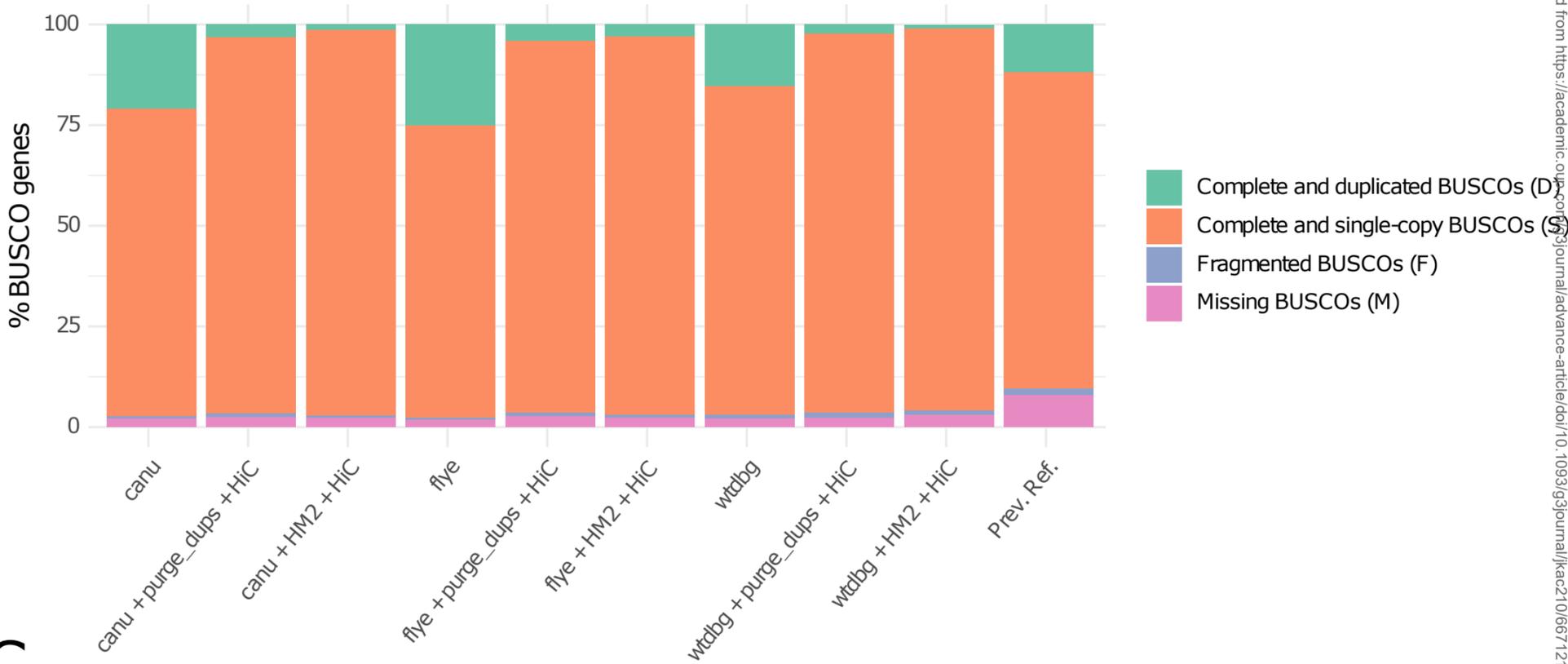
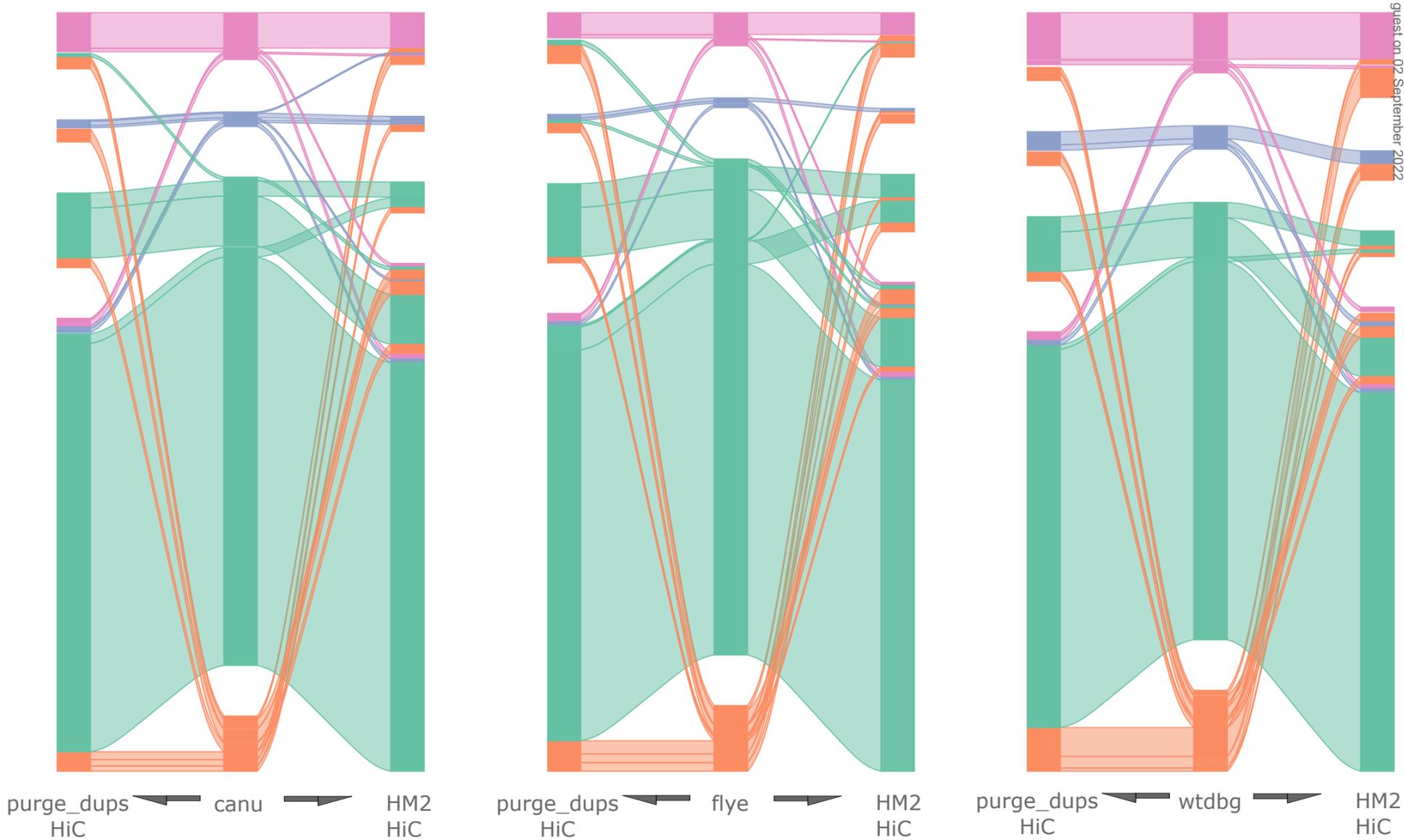
752

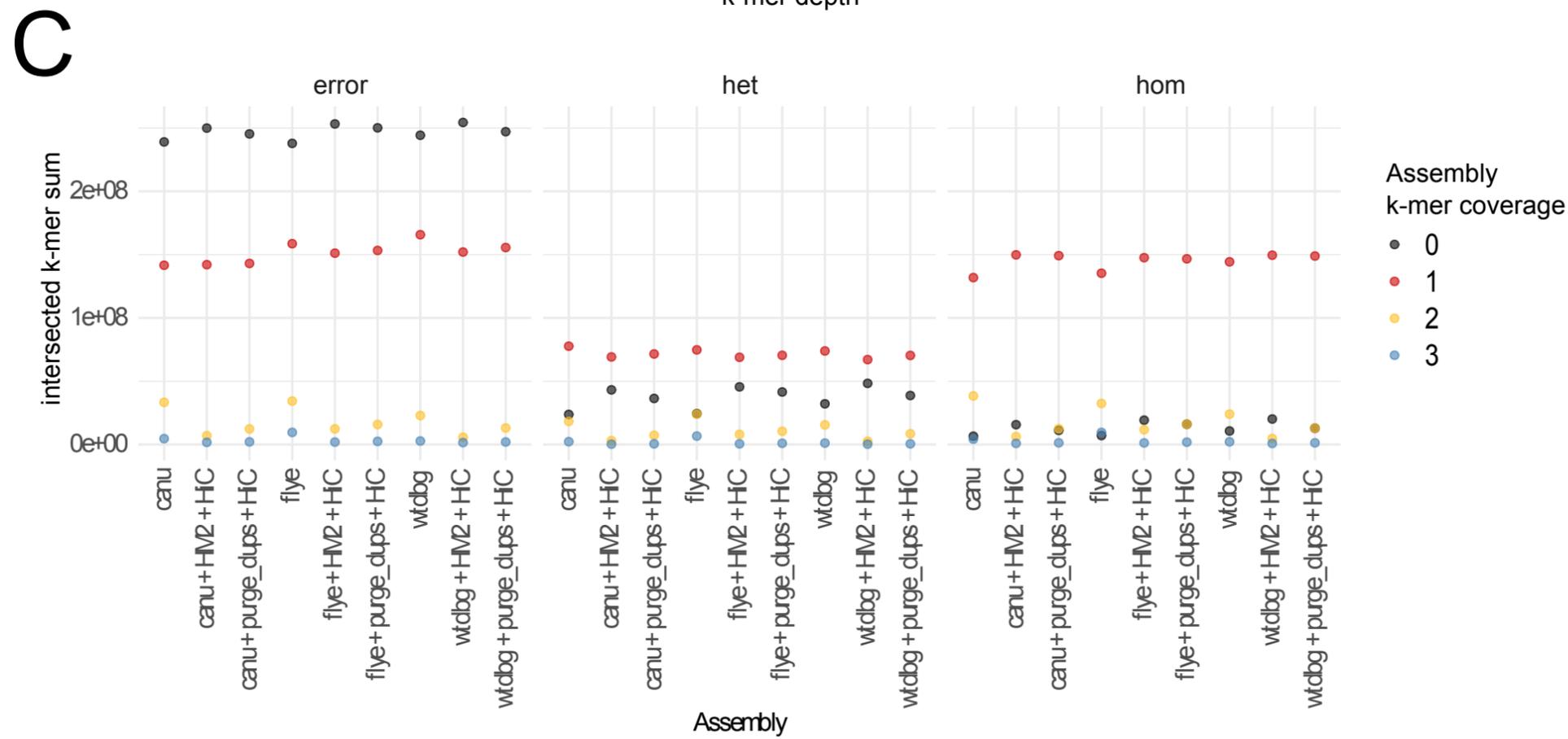
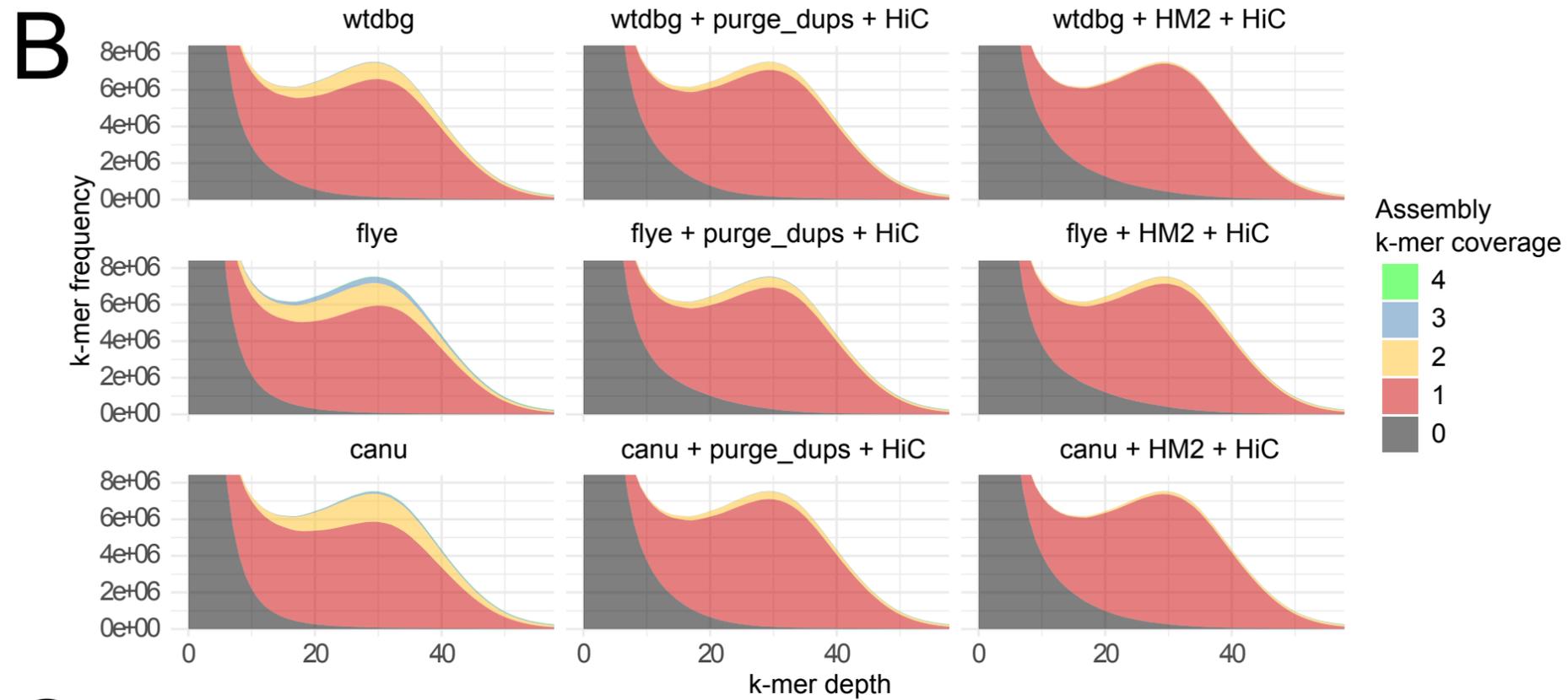
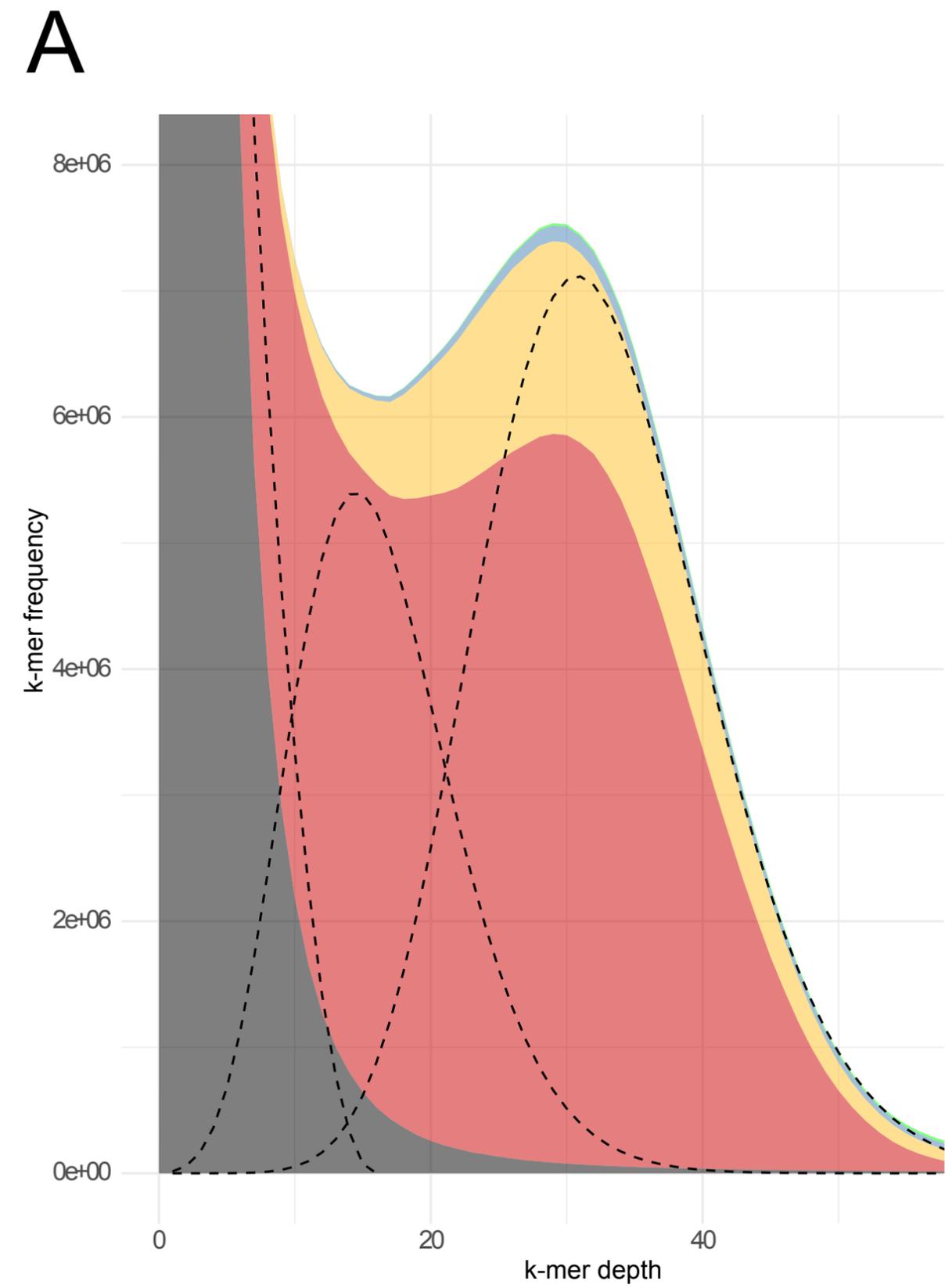
753
 754
 755
 756
 757
 758
 759
 760
 761
 762
 763
 764
 765
 766
 767
 768
 769
 770
 771
 772
 773

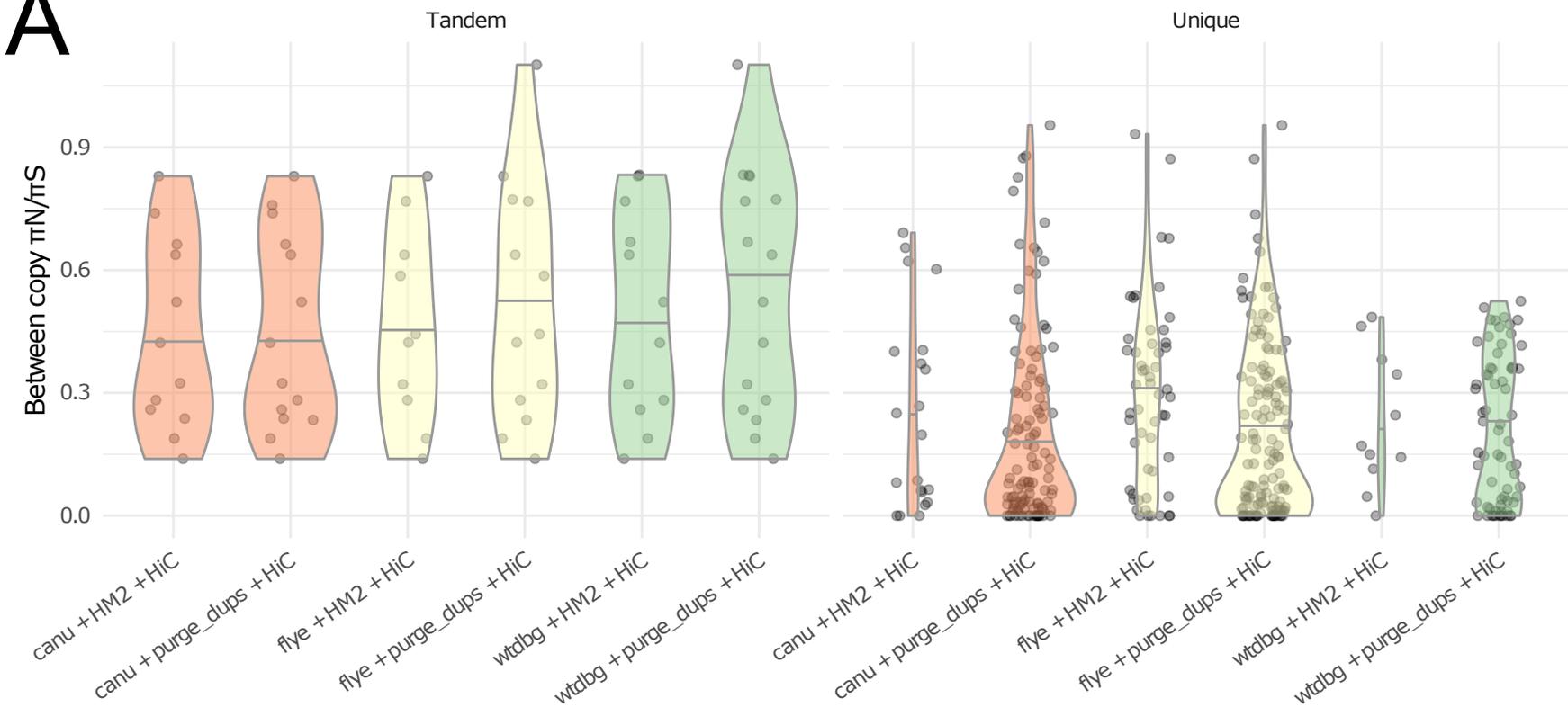
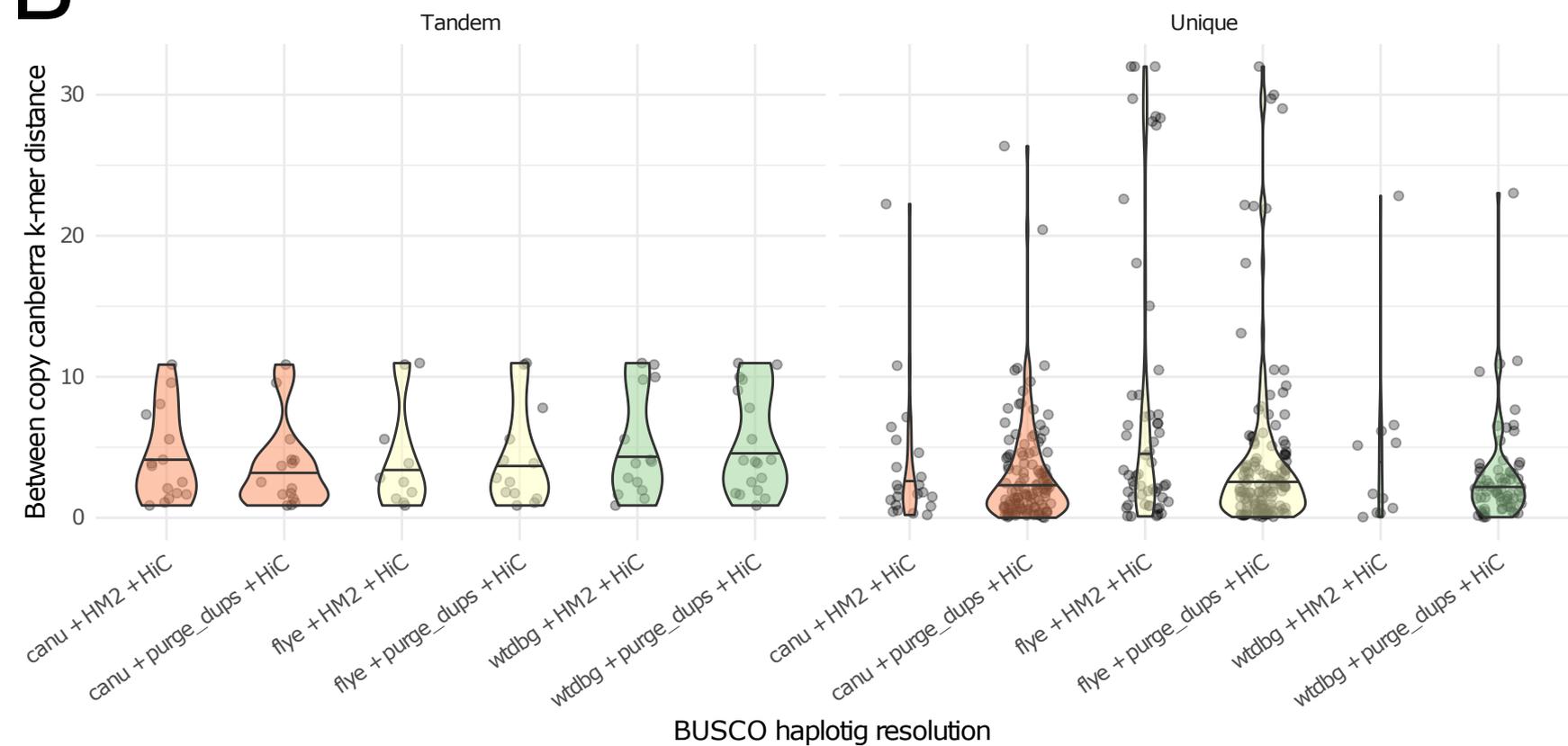
Table 1: Number of misassembly region candidates detected within the 3D-DNA pipeline. We used the default resolution parameters “wide_res = 25000bp” and “narrow_res = 1000bp”

	Total size (initial size) Mbp	Pre-HiC N50 (Mbp)	Post-HiC N50 (Mbp)	Narrow misassemblies		Wide misassemblies		Total iteration 2 misassemblies	Total disparity
				Iteration 1	Iteration 2	Iteration 1	Iteration 2		
canu + HM2	343 (448)	1.14	11.24	719	806	134	107	913	619
canu + purge_dups	367 (448)	0.62	9.49	794	1247	164	285	1532	
flye + HM2	368 (494)	0.41	11.42	1065	1225	185	115	1340	55
flye + purge_dups	391 (494)	0.32	11.44	835	1187	225	208	1395	
wtdbg + HM2	343 (427)	1.83	11.1	673	721	127	92	813	382
wtdbg + purge_dups	378 (427)	0.94	11.49	827	993	204	202	1195	

774

A**B****C**



A**B**

Program

- canu
- flye
- wtdbg