

Correcting Modification-Mediated Errors in Nanopore Sequencing by Nucleotide Demodification and in silico Correction

Chien-Shun Chiou¹, Bo-Han Chen¹, You-Wun Wang¹, Nang-Ting Kuo²,
Chih-Hsiang Chang², Yao-Ting Huang^{2*}

¹ Centers for Disease Control, Taichung, Taiwan.

² Department of Computer Science and Information Engineering,

National Chung Cheng University, Chiayi, Taiwan.

*Corresponding email: ythuang@cs.ccu.edu.tw

1 **Abstract**

2 The accuracy of Oxford Nanopore Technology (ONT) sequencing has significantly
3 improved thanks to new flowcells, sequencing kits, and basecalling algorithms.
4 However, novel modifications untrained in the basecalling models can seriously reduce
5 the quality. This paper reports a set of ONT-sequenced genomes with unexpected low
6 quality (~Q30) due to extensive new modifications. Demodification by whole-genome
7 amplification (WGA) significantly improved the quality of all genomes (~Q50-60)
8 while losing the epigenome. We developed a computational method, Modpolish, for
9 correcting modification-mediated errors without WGA. Modpolish produced high-
10 quality genomes and uncovered the underlying modification motifs without loss of
11 epigenome. Our results suggested that novel modifications are prone to ONT errors,
12 which are correctable by WGA or Modpolish without additional short-read sequencing.

13 **Keywords:** DNA modifications, Nanopore Sequencing, Whole-Genome Amplification.

14 **Background**

15 The Oxford Nanopore Technology (ONT) is a popular long-read sequencing
16 platform that enables real-time sequencing for point-of-care medical applications, such
17 as the diagnosis of infectious and newborn diseases within hospitals [1, 2]. Despite its
18 great potential and popularity, the accuracy of ONT was inferior to those of other
19 platforms (e.g., Illumina and PacBio HiFi). Recently, the quality of ONT sequencing
20 has significantly improved thanks to new flowcells (e.g., R10.4), sequencing kits (e.g.,
21 Kit 14), and basecalling algorithms (e.g., Bonito). For example, by using the R10.4
22 flowcells, near-perfect microbial genomes from isolates or metagenomes can be
23 reconstructed by ONT-only sequencing without short-read polishing [3].

24 However, because the throughput of R10.4 is much lower due to slower
25 sequencing speed, most projects rely on the R9.4 flowcells for higher yield. Although
26 the upcoming sequencing kit will further enhance the accuracy (e.g., Kit 14),
27 postassembly genome polishing is still compulsory for removing ONT systematic
28 errors regardless of the flowcell or kit versions. Systematic errors are recurrent
29 basecalling errors at the same locus, which are not correctable by the consensus of read
30 pileups (e.g., Racon) [4]. Homopolymer errors (i.e., indels) were the primary source of
31 ONT systematic errors. Thanks to several machine-learning algorithms, these errors

32 have been significantly reduced by read-based (e.g., Medaka) or reference-based (e.g.,
33 Homopolish) polishing methods [5]. These algorithmic advances have produced high-
34 quality ONT genomes sufficient for downstream analysis (e.g., >Q50) [3, 6].

35 Unfortunately, the ONT signals are ultra-sensitive to various modifications (e.g.,
36 5mC, 6mA). More than 17 and 160 modification types have been found in DNA and
37 RNA, respectively, and the number is still growing (e.g., DNA adducts, N4-
38 acetyldeoxycytosine) [7, 8]. These modifications disturb the electrical current and
39 result in unfixable systematic errors [9]. Note that these modification-mediated errors
40 cannot be eliminated by new flowcells and sequencing kits (e.g., R10.4 and Kit 14)
41 which aim to reduce homopolymer errors. Furthermore, existing basecalling and
42 polishing algorithms (e.g., Guppy and Medaka) were trained for capturing only a few
43 modifications (e.g., 5mC, 5hmc, 6mA). Consequently, the quality of ONT sequencing
44 is unreliable when novel modifications extensively edit the genome.

45 This paper presents a set of unexpected low-quality genomes due to extensive
46 novel modifications. We show that the removal of modifications by whole-genome
47 amplification (WGA) significantly improves the quality of all genomes. A novel
48 computational method is developed for correcting these modification errors without
49 WGA.

50 **Results**

51 **Unusual low-quality of ONT genomes due to extensive modifications**

52 We sequenced 12 microbial strains of *Listeria monocytogenes* using Illumina
53 and ONT (~200-990Mbp) (Figure 1(a), Supplementary Tables S1 and S2). The ONT
54 reads were assembled into genomes with sequencing errors further polished by the-
55 state-of-the-art tools (Supplementary Table S3, see Methods). The Illumina and ONT
56 reads were hybrid assembled for evaluation purposes (Supplementary Table S4). When
57 compared with the Illumina/ONT hybrid assemblies (Figure 1(b)), seven ONT-only
58 genomes exhibited high quality (HQ) ranging from Q47 to Q60 (e.g., R19-2905 and
59 R20-0088). However, five isolates (R20-0026, R20-0030, R20-0127, R20-0148, and
60 R20-0150) showed unexpectedly low quality (LQ) varying from Q27 to Q34. The
61 accuracy of these five LQ genomes remained unimproved after replicated ONT
62 sequencing (data not shown). Further investigation of the five LQ genomes revealed
63 excessive amounts of mismatch errors (1,228-5,780) compared with the seven HQ ones
64 (3-36 mismatches) (Figure 1(c)). Homopolymer errors (i.e., indels) were not the source
65 of inferior quality (7-306, Supplementary Table S5).

66 Manual inspection revealed that these mismatches were ONT basecalling errors
67 uncorrected after genome polishing (Figure 1(d) and Supplementary Figure S1). As

68 mismatch errors in ONT are mainly due to epigenetic modifications, we computed the
69 frequency of well-known methylation in these isolates (see Method and Supplementary
70 Table S6). In terms of 5-methylcytosine (5mC), the numbers of modified loci in the
71 five LQ genomes (~240-340k) were not significantly higher than those in the HQ ones
72 (210-345k, $P=0.89$, Figure 1(e)). Similarly, the numbers of N⁶-methyladenine (6mA)
73 modifications also showed no significant difference between the LQ and HQ groups
74 (98-218k v.s. 126-223k, $P=0.34$, Figure 1(f)). Because the numbers of mismatch errors
75 in LQ genomes are significantly higher than those of HQ ones ($P=0.005$, Figure 1(g)),
76 we suspected ONT basecalling algorithms failed to distinguish the novel modifications
77 in the LQ isolates.

78 **High-quality ONT genomes by WGA demodification**

79 We removed the modifications in all microbial samples by WGA (Figure 2(a)),
80 which randomly amplifies the genome fragments without retaining any epigenetic
81 modification (see Methods). The WGA-demodified samples were sequenced by ONT,
82 assembled into chromosomes, and compared with the Illumina/ONT hybrid genomes
83 (Figure 2(a), Supplementary Tables S7 and S8). The five LQ genomes after WGA
84 exhibited significantly higher quality than those without demodifications (e.g., Q27 to
85 Q53 in R20-0026) (Figure 2(b), Supplementary Table S9). In particular, the amounts

86 of mismatch errors significantly reduced after demodification (e.g., 5,780 to 16 in R20-
87 0026) (Figure 2(c)). Consequently, the unexpected low quality of ONT was due to
88 excessive novel modifications untrained in their basecalling model. The demodification
89 by WGA can produce high-quality ONT genomes without the need for Illumina short
90 reads.

91 However, while WGA successfully erased these modifications, the sequencing
92 cost increased by two factors. First, WGA required a higher sequencing depth (~100x)
93 for assembling a complete genome when compared with ordinary ONT sequencing
94 (~30x) (Figure 2(d) and Supplementary Figures S2-3). It was due to the uneven
95 amplification of WGA, which led to non-uniform sequencing depth and a fragmented
96 assembly at moderate coverage. Second, the WGA-demodified samples may reduce the
97 ONT yields. We observed the numbers of available/active pores could sometimes
98 decrease quickly (e.g., less than 100 pores after 12h) (Figure 2(e)), which was possibly
99 owing to the hyperbranched structure unresolved after WGA. Consequently, the
100 sequencing cost of WGA-demodified samples using ONT is much higher than ordinary
101 sequencing.

102 ***in silico* correction of modification-mediated errors by Modpolish**

103 We developed a novel computational method (called Modpolish) for correcting
104 these modification-mediated errors without WGA and prior knowledge of the
105 modifications. Modpolish identifies and corrects the modification-mediated errors by
106 investigating basecalling quality, basecalling consistency, and evolutionary
107 conservation (Figure 3(a), see Method). Briefly, because the ONT signals are disturbed
108 by modifications, the basecalling quality is usually low, and the basecalled nucleotides
109 are often inconsistent at the modified loci. In conjunction with the conservation degree
110 measured by closely-related genomes, only the modified loci with ultra-high
111 conservation will be corrected by Modpolish, avoiding false corrections of strain
112 variations.

113 We assessed the accuracy of Modpolish by comparing the quality of the ONT-
114 only genomes (polished by Medaka/Homopolish) with those further polished by
115 Modpolish. The results indicated that Modpolish significantly improved the genome
116 quality of all LQ genomes (Figure 3(b), Supplementary Table S10). For instance, the
117 quality of R20-0030 improved from Q34 to Q60, and the number of mismatches
118 decreased from 1,228 to 33 (Figure 3(c)). We observed that the number of mismatches
119 in R20-0026 reduced dramatically (i.e., from 5,780 to 143). However, the quality
120 improvement (i.e., from Q27 to Q45) was slightly inferior to the others due to the 143

121 uncorrected mismatches. Note that no false corrections were made on the seven HQ
122 genomes, implying the correction specificity of Modpolish is high.

123 The multilocus sequencing typing (MLST) indicated that R20-0026 belonged to
124 the sequence type ST1081 and the remaining four LQ strains (i.e., R20-0030, R20-0127,
125 R20-0148, R20-150) were ST87. Hence, we investigated whether an identical
126 modification system extensively edited the genomes of these two lineages. Sequence
127 analysis of the modified loci revealed that the modifications of ST1081 were on the
128 GCTGG motif (Figure 3(d)). On the other hand, the modification sites of all ST87
129 strains centered on the GCAGC motif (Figure 3(e)). Therefore, two modification
130 systems seem specific to each of the two lineages. In addition, while both motifs are
131 not palindromic, their reverse complements (i.e., CCAGC, GCTGC) were also hotspots
132 of modifications (Supplementary Figure S4). Because the mismatches frequently
133 appeared on both strands at the same loci (Supplementary Figure S5), the unknown
134 modification may symmetrically edit both strands. Although their underlying
135 mechanisms remained unclear, the two systems extensively modified the genomes at
136 specific motifs with high conservation within each lineage, leading to excessive
137 amounts of basecalling yet correctable errors.

138 **Comparison of phylogeny reliability under extensive modifications**

139 Because sequencing errors alter the genetic distances between strains, we assessed
140 the reliability of phylogeny using ONT with or without modification-error removal. We
141 reconstructed the core-genome MLST (cgMLST) phylogeny of the five LQ strains
142 sequenced and assembled by four methods: ONT-only sequencing, WGA-demodified
143 ONT, ONT with Modpolish, and hybrid ONT/Illumina sequencing (Figure 4(a)). The
144 WGA-demodified genomes perfectly clustered with the ONT/Illumina hybrid for each
145 strain in both clades (ST87 and ST1081). The ONT genomes corrected by Modpolish
146 clustered with the hybrid and WGA-demodified genomes in both clades. But the
147 genetic distance slightly deviated from them, especially in the ST1081 clade. The ONT-
148 only genomes were phylogenetic distant from the others due to excessive amounts of
149 modification-mediated errors.

150 When comparing each method in the seven HQ isolates, ONT with WGA was
151 slightly worse than the original ONT and Modpolish in six strains (e.g., 47 v.s. 6
152 mismatches in R20-0088) (Figure 4(b)), except for the R19-2905 isolate (i.e., 12 v.s.
153 36 mismatches). These mismatches slightly increased the genetic distance to the others
154 (Supplementary Figure S6). Nevertheless, phylogenetic analysis indicated that the
155 genomes of all methods were perfectly clustered for each HQ strain (Supplementary
156 Figure S7), implying the number of mismatches is less than that of strain variations.
157 Consequently, all methods can produce reliable phylogeny when free of novel

158 modifications. But when new modifications extensively edit the genome, only ONT
159 with WGA or Modpolish can provide sufficient typing accuracy without additional
160 Illumina sequencing.

161 **Discussion**

162 This paper presented a set of unexpected low-quality ONT genomes due to extensive
163 modifications untrained in the basecalling models. Demodification by WGA
164 successfully improved the genome quality while losing the epigenome. The *in silico*
165 method, Modpolish, removed these modification-mediated errors without prior
166 knowledge of modifications and uncovered the modified motifs while retaining the
167 epigenome. When unknown modifications extensively shaped the genome, ONT with
168 WGA or Modpolish produced nearly identical cgMLST profiles as hybrid
169 ONT/Illumina did. On the other hand, the phylogeny of ONT-only genomes was
170 disturbed by modification-mediated errors. Therefore, ONT with WGA or Modpolish
171 is robust to modification-mediated errors without the need for additional Illumina
172 sequencing.

173 **Quality reduction of ONT on novel modifications**

174 Existing ONT basecalling algorithms only capture a few methylations (e.g., 5mC,
175 5hmc, 6mA) and ignore the vast amount of other modifications. Theoretically, species-

176 specific modifications can be distinguished by training bespoke models for one
177 organism (e.g., *Taiyaki*). But practically, it is infeasible to train models for hundreds of
178 modifications in the biosphere. Especially in metagenomic sequencing, the usage of
179 any particular model is biased against other modifications. For instance, a meta-
180 epigenomic sequencing uncovered 22 methylation systems in a single microbial
181 community [10]. Hence, if WGA is not an option, modification-mediated errors are
182 better removed at the postassembly stage as each assembled contig can be polished
183 independently.

184 **Limitations of ONT with WGA**

185 The cost of WGA ONT is higher than ordinary sequencing due to several side
186 effects of the amplification protocol. First, the amplified DNA may still contain a
187 hyperbranched structure after Flap endonuclease (e.g., T7) cleavage. The
188 hyperbranched DNA may block the pores during ONT sequencing and reduce the
189 available pores and yields. In addition, the usage of endonuclease cleavage also
190 decreased the read lengths. In conjunction with the uneven amplification, WGA
191 requires higher coverage (~100x) for reconstructing a complete genome than ordinary
192 ONT sequencing (~30x). Notably, the usage of WGA discards the entire methylome.
193 The loss of modifications would prohibit any epigenetic study using ONT.

194 **Limitations of ONT with Modpolish**

195 While Modpolish eliminated most modification-mediated errors, the correction
196 power was lower in the ST1081 isolate. The lack of ST1081 genomes in NCBI RefSeq
197 decreased the sensitivity of Modpolish. As the algorithm only corrects the loci of high
198 evolutionary conservation, a sufficient number of closely-related genomes is necessary.
199 Therefore, Modpolish is more suitable for common instead of rare lineages.

200 Nevertheless, Modpolish retains all modifications after ONT sequencing while
201 WGA loses the epigenome. Epigenetic methylation has been thought to contribute to
202 the rapid adaptation of resistance [11]. For instance, phase-variable adenine DNA
203 methyltransferases (e.g., ModA11 and ModA12) increase susceptibility to cloxacillin
204 and ciprofloxacin in *Neisseria meningitidis* [12]. The resistance due to overexpression
205 of efflux pumps (e.g., sugE) has been linked to the lack of the Dcm-mediated 5mC
206 silencing [13]. Therefore, Modpolish should be used when the epigenome is the focus
207 of the study.

208 **Functional implications of the two modification systems**

209 We discovered two pentanucleotide motifs, GCTGG (CCAGC) and GCAGC
210 (GCTGC), specific to each of the two lineages (ST1081 and ST87). In ST1081, the
211 GCTGG (CCAGC) motif is part of *chi* sites, hotspots of homologous recombination
212 mediated by the RecBC enzyme [14, 15]. As phages cut by restriction enzymes are

213 further degraded by RecBC [16], modifications on the GCTGG motif may be part of
214 the defending system of ST1081, which protect itself against the RecBC cleavage.

215 In ST87 strains, the GCAGC/GCTGC (i.e., GCWGC) motif was the known target
216 of the orphan methyltransferase M.BatI [17]. M.BatI produced fully-methylation on 5'-
217 GCWGC-3' and hemimethylation on 5'-GCSGC-3'. Reinvestigation of the modified
218 sites in ST87 showed the existence of both GCWGC and GCSGC (Supplementary
219 Figure S4). Interestingly, M.BatI increased toxicity when expressed in *E coli* in their
220 study, which was concordant with the elevated virulence of ST87 strains.

221 Hence, the two lineages possessed two distinct modification systems for defensive
222 purposes and increasing virulence. Although further investigations are required to
223 assess their biological function, modifications that have acquired regulatory effects in
224 bacteria are usually conservative within a clade [18]. Consequently, our *in silico*
225 algorithm successfully utilize the conservation for correcting modification errors.

226 **Conclusion**

227 This paper reported a set of unexpectedly low-quality genomes due to novel
228 modifications untrained in the ONT basecalling model. The increasing number of new
229 modifications found by single-molecular sequencing or high-resolution mass
230 spectrometry will unavoidably reduce the ONT accuracy. New ONT flowcells,

231 sequencing kits, and basecalling algorithms aim to resolve the homopolymer issue but
232 not modification-mediated errors. Our study showed that these modification-mediated
233 errors can be effectively corrected by preassembly amplification or postassembly
234 polishing without additional short-read sequencing, producing high-quality genomes
235 reliable for downstream analysis.

236 **Materials and Methods**

237 **Bacterial isolates.** Twelve *Listeria monocytogenes* isolates used in this study were
238 obtained from hospitals recovered from listeriosis patients in Taiwan between 2019 and
239 2020. The isolates were submitted to the Taiwan Centers for Disease Control for further
240 identification and genotyping. The isolates belonged to serogroups IIa (5 isolates), IIb
241 (6 isolates), and IVb (1 isolate) and sequence type (ST) 1, ST5 (2 isolates), ST87 (4
242 isolates), ST101, ST155, ST378, ST1081, and ST1532.

243 **Whole genome sequencing.** WGS of bacterial isolates was conducted in the Central
244 Region laboratory of Taiwan CDC using the Illumina MiSeq sequencing platform
245 (Illumina Co., USA) and the Nanopore sequencing platform (Oxford Nanopore
246 Technologies, Inc., UK). DNA of bacterial isolates was extracted using the Qiagen
247 DNeasy blood and tissue kit (Qiagen Co., Germany). Illumina DNA library
248 construction was performed using the Illumina DNA Prep, (M) Tagmentation system

249 (Illumina Co.), and sequencing was run with the MiSeq reagent kit version 3 (2X 300
250 bp), manipulated according to the manufacturer's instructions. Nanopore DNA library
251 construction was performed using the Rapid Barcoding Kit and sequencing was run
252 using the MinION device and R9.4 chemistry.

253 **Removal of modifications of nucleotides using whole-genome amplification.** DNA
254 Bacterial Genomic DNA was amplified using the REPLI-g Advanced DNA Single Cell
255 Kit (Qiagen, Hilden, Germany), manipulated according to the manufacturer's
256 instructions. The amplified DNA was purified using the KAPA HyperPure Beads
257 (Roche, Basel, Switzerland) before subjecting to Nanopore sequencing.

258 **Assembly of sequence reads.** Illumina sequence reads for each isolate were assembled
259 using the SPAdes assembler version 3.12.0 (<http://cab.spbu.ru/software/spades/>) [19];
260 both Illumina sequence reads and Nanopore sequence reads for each isolates together
261 were assembled to complete the full genomic sequences using the Unicycler Assembler
262 [20]. The Nanopore reads for each isolate (in FAST5 file) were subjected to basecalling
263 using the Guppy basecaller (<https://nanoporetech.com/>). In the ONT-only assembly, the
264 sequences (in FASTQ file) were assembled using Flye
265 (<https://github.com/fenderglass/Flye>)[21], then polished using the Racon
266 (<https://github.com/lbcb-sci/racon>) [4], the Medaka

267 (<https://github.com/nanoporetech/medaka>), and the Homopolish
268 (<https://github.com/ythuang0522/homopolish>) [5]. Methylations (i.e., 5mC, 6mA) in
269 the ONT-only genomes were called by Megalodon
270 (<https://github.com/nanoporetech/megalodon>). The Integrative Genome Viewer (IGV)
271 was used for visualizing the ONT modification errors [22]. The genome quality was
272 assessed by fastmer (https://github.com/jts/assembly_accuracy).

273 **cgMLST analysis.** Assembled Illumina contigs, assembled and polished Nanopore
274 contigs, and assembled complete genomic sequence (obtained from assembling
275 Illumina sequences and Nanopore sequences) for each isolate were used to generate
276 core-gene multilocus sequence typing (cgMLST) profiles (based on 2,172 core genes)
277 using an in-house-developed cgMLST profiling tool available on the
278 cgMLST@Taiwan website (<http://rdvd.cdc.gov.tw/cgMLST>). Phylogenetic trees were
279 constructed with cgMLST profiles using the minimum spanning tree algorithm and the
280 tool provided on the cgMLST@Taiwan website.

281 **Overview of Modpolish.** The proposed computational method, Modpolish, aims to
282 remove modification-mediated errors by investigating the inconsistency of basecalled
283 nucleotides, qualities of basecalled alleles, and evolutionary conservation at the
284 modified loci. Modpolish is an extension from Homopolish, a polishing algorithm

285 designed for correcting ONT homopolymer errors [5]. Figure 5 depicts the workflow
286 of Modpolish. The closely-related genomes are first identified by screening against a
287 compressed representation of microbial genomes. The genome sequences are then
288 retrieved on the fly and compared with the draft genome. We only retain closely-related
289 genomes of high nucleotide and structural similarity. Given the alignment matrix of
290 reads, qualities, and homologs, Modpolish identifies potential-modified loci of
291 inconsistent basecalling and low quality and only corrects the mismatch errors highly
292 conserved in homologs. The details are described in the following sections.

293 **Collection of homologs by nucleotide and structural similarity.** The draft genome
294 (to be polished) is scanned against the virus, bacteria, or fungus genomes compressed
295 by Mash as (MinHash) sketches, which is a reduced representation of all microbial
296 genomes in NCBI RefSeq [23]. Subsequently, top t (default 20) closely related genomes
297 will be retrieved on the fly. Mash estimated the Jaccard similarity between the draft and
298 related genomes over a subset of k -mers. Though very fast, this method has low
299 resolution at distinguishing closely-related genomes because the small subset of k -mers
300 may not capture the few strain variations. Consequently, the genome similarity has to
301 be re-estimated using more sensitive approaches.

302 Subsequently, each downloaded genome is compared against the draft genome
303 using FastANI for computing the average nucleotide identity (ANI) at a higher

304 resolution than Mash [24]. FastANI chops the two genomes into pieces and aligns them
305 against each other for speedup. However, it only considers the aligned segments for
306 ANI estimation and ignores the unaligned portions (Supplementary Figure S8(a)). The
307 unaligned segments imply these two genomes differ by structural variations (i.e.,
308 vertically-/horizontally-transferred genes). As small and large variants are both genetic
309 footprints of strain variations during evolution, Modpolish also computes the structural
310 similarity (average-structural identity, ASI), defined as the percentage of aligned
311 segments. We only retain the related genomes with sufficient ANI (>99%) and ASI
312 (>90%) for subsequent error correction. These empirical cutoffs were determined by
313 investigating the distributions of ANI and ASI in real microbial genomes.

314 **Correction of modification-mediated errors by reads and homologs.** These closely-
315 related genomes with sufficient ANI and ASI are aligned against the draft genome via
316 minimap2 (with asm5 option) [25]. The raw ONT reads are also mapped against the
317 draft genome by minimap2 (with map-ont option). We extract the basecalled
318 nucleotides, basecalling qualities, and homologous alleles from the alignments. The
319 aligned homologs, reads, and qualities are converted into a table of several summary
320 statistics (Supplementary Figure S8(b)).

321 The summary statistics include the allele counts of A, T, C, and G separately for
322 homologs and ONT reads, ignoring the insertion and deletion gaps. We identify the

323 potentially modified sites according to the allele discordancy and average quality (see
324 also Supplementary Figure S8(b)). The allele discordancy is the frequency of
325 alternative alleles (i.e., non-major ones) at one locus. The average quality was
326 computed by averaging the qscores from all read bases at the same locus. A potentially-
327 modified locus is defined as the allele discordancy greater than 5% and the average
328 quality score below 15, which were empirically observed from the modification-
329 mediated errors.

330 For each potentially-modified locus, if all the homologous alleles are 100%
331 conserved, we will correct the erroneous nucleotide into the alternative allele
332 concordant with the homologs. These stringent criteria aimed for specificity instead of
333 sensitivity, ensuring little or no false corrections would be made. We also implemented
334 a motif-aware mode when the modification system is known in advance. If the user
335 specifies a known modification motif (e.g., CCGAC), the program will additionally
336 correct loci according to the provided pattern by lowering the homologous conservation
337 ratio from 100% to 80%.

338 **Data and software availability**

339 The genomes sequenced and assembled by Illumina, ONT, and WGA ONT are
340 deposited in the NCBI with BioProject (xxxxxx). Modpolish was implemented as a

341 subcommand in the Homopolish package, which is freely available at
342 (<https://github.com/ythuang0522/homopolish/>).

343 **Conflict of interests**

344 The authors declare no conflict of interests.

345

346 **Reference**

- 347 1. Gorzynski JE, Goenka SD, Shafin K, Jensen TD, Fisk DG, Grove ME, Spiteri E,
348 Pesout T, Monlong J, Baid G *et al*: **Ultrarapid Nanopore Genome**
349 **Sequencing in a Critical Care Setting**. *New England Journal of Medicine*
350 2022, **386**(7):700-702.
- 351 2. Gu W, Deng X, Lee M, Sucu YD, Arevalo S, Stryke D, Federman S, Gopez A,
352 Reyes K, Zorn K *et al*: **Rapid pathogen detection by metagenomic next-**
353 **generation sequencing of infected body fluids**. *Nat Med* 2021,
354 **27**(1):115-124.
- 355 3. Sereika M, Kirkegaard RH, Karst SM, Michaelsen TY, Sørensen EA,
356 Wollenberg RD, Albertsen M: **Oxford Nanopore R10.4 long-read**
357 **sequencing enables near-perfect bacterial genomes from pure**
358 **cultures and metagenomes without short-read or reference polishing**.
359 *bioRxiv* 2021:2021.2010.2027.466057.
- 360 4. Vaser R, Sović I, Nagarajan N, Šikić M: **Fast and accurate de novo genome**
361 **assembly from long uncorrected reads**. *Genome Res* 2017, **27**(5):737-
362 746.
- 363 5. Huang Y-T, Liu P-Y, Shih P-W: **Homopolish: a method for the removal of**
364 **systematic errors in nanopore sequencing by homologous polishing**.
365 *Genome biology* 2021, **22**(1):1-17.
- 366 6. Zhang P, Jiang D, Wang Y, Yao X, Luo Y, Yang Z: **Comparison of De Novo**
367 **Assembly Strategies for Bacterial Genomes**. *International Journal of*
368 *Molecular Sciences* 2021, **22**(14):7668.

- 369 7. Wang S, Xie H, Mao F, Wang H, Wang S, Chen Z, Zhang Y, Xu Z, Xing J, Cui Z
370 *et al*: **N4-acetyldeoxycytosine DNA modification marks euchromatin**
371 **regions in *Arabidopsis thaliana***. *Genome Biology* 2022, **23**(1):5.
- 372 8. Xu L, Seki M: **Recent advances in the detection of base modifications**
373 **using the Nanopore sequencer**. *Journal of Human Genetics* 2020,
374 **65**(1):25-33.
- 375 9. Schatz MC: **Nanopore sequencing meets epigenetics**. *Nature Methods*
376 2017, **14**(4):347-348.
- 377 10. Hiraoka S, Okazaki Y, Anda M, Toyoda A, Nakano S-i, Iwasaki W:
378 **Metaepigenomic analysis reveals the unexplored diversity of DNA**
379 **methylation in an environmental prokaryotic community**. *Nature*
380 *Communications* 2019, **10**(1):159.
- 381 11. Ghosh D, Veeraraghavan B, Elangovan R, Vivekanandan P: **Antibiotic**
382 **Resistance and Epigenetics: More to It than Meets the Eye**.
383 *Antimicrobial Agents and Chemotherapy* 2020, **64**(2):e02225-02219.
- 384 12. Jen FE-C, Seib KL, Jennings MP: **Phasevarions Mediate Epigenetic**
385 **Regulation of Antimicrobial Susceptibility in *Neisseria meningitidis***.
386 *Antimicrobial Agents and Chemotherapy* 2014, **58**(7):4219-4221.
- 387 13. He G-X, Zhang C, Crow RR, Thorpe C, Chen H, Kumar S, Tsuchiya T, Varela
388 MF: **SugE, a New Member of the SMR Family of Transporters,**
389 **Contributes to Antimicrobial Resistance in *Enterobacter***
390 ***cloacae***. *Antimicrobial Agents and Chemotherapy* 2011, **55**(8):3954-
391 3957.
- 392 14. Aoki H, Kajino K, Arakawa Y, Hino O: **Molecular cloning of a rat**
393 **chromosome putative recombinogenic sequence homologous to the**
394 **hepatitis B virus encapsidation signal**. *Proceedings of the National*
395 *Academy of Sciences* 1996, **93**(14):7300-7304.
- 396 15. Ponticelli AS, Schultz DW, Taylor AF, Smith GR: **Chi-dependent DNA**
397 **strand cleavage by RecBC enzyme**. *Cell* 1985, **41**(1):145-151.
- 398 16. Simmon VF, Lederberg S: **Degradation of bacteriophage lambda**
399 **deoxyribonucleic acid after restriction by *Escherichia coli* K-12**. *J*
400 *Bacteriol* 1972, **112**(1):161-169.
- 401 17. Furuta Y, Miura F, Ichise T, Nakayama SMM, Ikenaka Y, Zorigt T,
402 Tsujinouchi M, Ishizuka M, Ito T, Higashi H: **A GCDGC-specific DNA**
403 **(cytosine-5) methyltransferase that methylates the GCWGC sequence**
404 **on both strands and the GCSGC sequence on one strand**. *PLOS ONE*
405 2022, **17**(3):e0265225.

- 406 18. Seshasayee ASN, Singh P, Krishna S: **Context-dependent conservation of**
407 **DNA methyltransferases in bacteria.** *Nucleic acids research* 2012,
408 **40(15):7066-7073.**
- 409 19. Bankevich A, Nurk S, Antipov D, Gurevich AA, Dvorkin M, Kulikov AS, Lesin
410 VM, Nikolenko SI, Pham S, Prjibelski AD *et al*: **SPAdes: a new genome**
411 **assembly algorithm and its applications to single-cell sequencing.** *J*
412 *Comput Biol* 2012, **19(5):455-477.**
- 413 20. Wick RR, Judd LM, Gorrie CL, Holt KE: **Unicycler: Resolving bacterial**
414 **genome assemblies from short and long sequencing reads.** *PLoS*
415 *Comput Biol* 2017, **13(6):e1005595.**
- 416 21. Kolmogorov M, Bickhart DM, Behsaz B, Gurevich A, Rayko M, Shin SB, Kuhn
417 K, Yuan J, Polevikov E, Smith TPL *et al*: **metaFlye: scalable long-read**
418 **metagenome assembly using repeat graphs.** *Nature Methods* 2020,
419 **17(11):1103-1110.**
- 420 22. Robinson JT, Thorvaldsdóttir H, Winckler W, Guttman M, Lander ES, Getz
421 G, Mesirov JP: **Integrative genomics viewer.** *Nature biotechnology* 2011,
422 **29(1):24-26.**
- 423 23. Ondov BD, Treangen TJ, Melsted P, Mallonee AB, Bergman NH, Koren S,
424 Phillippy AM: **Mash: fast genome and metagenome distance estimation**
425 **using MinHash.** *Genome Biology* 2016, **17(1):132.**
- 426 24. Jain C, Rodriguez-R LM, Phillippy AM, Konstantinidis KT, Aluru S: **High**
427 **throughput ANI analysis of 90K prokaryotic genomes reveals clear**
428 **species boundaries.** *Nature Communications* 2018, **9(1):5114.**
- 429 25. Li H: **Minimap2: pairwise alignment for nucleotide sequences.**
430 *Bioinformatics* 2018, **34(18):3094-3100.**

431

432 **Figure Legends**

433 **Figure 1.** Quality comparison of 12 microbial strains using ONT-only and
434 ONT/Illumina hybrid sequencing. (a) Workflow of ONT-only and ONT/Illumina
435 hybrid assembly; (b) Q scores; (c) number of mismatches; (d) comparison of ONT and
436 Illumina reads by IGV; (e) numbers of 5mC, 6mA, and mismatches between HQ/LQ
437 strains.

438 **Figure 2.** Quality improvement of ONT by WGA demodification. (a) Workflow of
439 WGA-demodified ONT; (b) Q scores of the WGA-demodified and ONT-only genomes;
440 (c) numbers of mismatches of the WGA-demodified and ONT-only genomes; (d) WGA
441 genome quality with respect to sequencing depth; (e) numbers of active/available pores
442 during WGA-demodified and ordinary ONT sequencing.

443 **Figure 3.** Correction of modification-mediated errors by Modpolish. (a) Workflow of
444 Modpolish; (b) Q scores before and after Modpolish; (c) numbers of mismatches before
445 and after Modpolish; (d) the sequence motif of modification on ST1081; (e) the
446 sequence motif of modifications on ST87.

447 **Figure 4.** Comparison of phylogeny reliability of four methods. (a) The cgMLST
448 phylogeny of the five LQ strains sequenced and assembled by four methods: ONT-only
449 sequencing (ONT), WGA-demodified ONT (ONT_WGA), ONT with Modpolish

450 (ONT_Modpolish), and hybrid ONT/Illumina sequencing (Hybrid_WGS); (b) the
451 cgMLST distances of ONT, ONT_WGA, and ONT_Modpolish to the Hybrid_WGS
452 assembled genomes.

453 **Figure 5.** Illustration of Modpolish workflow. A set of closely-related genomes are first
454 retrieved by screening the compressed sketches of RefSeq genomes. We retain the
455 genomes with sufficient nucleotide and structural similarity. The selected genomes and
456 ONT reads are aligned onto the draft genome, generating a pileup matrix of homologs,
457 reads, and qualites. Modpolish only corrects modification-mediated errors with
458 inconsistent read alleles, low quality, and high conservation in homologs.

459

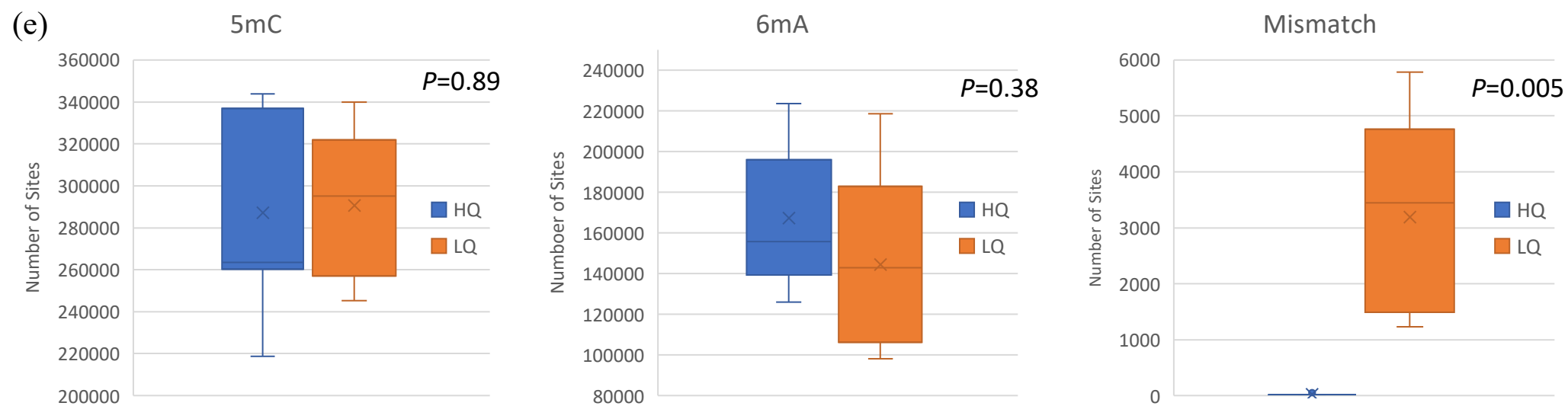
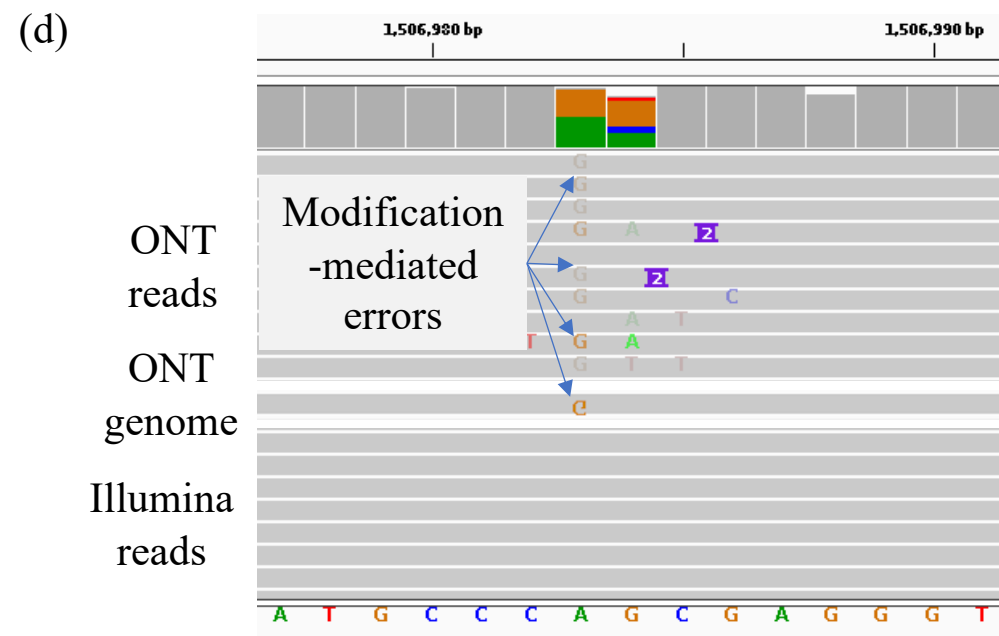
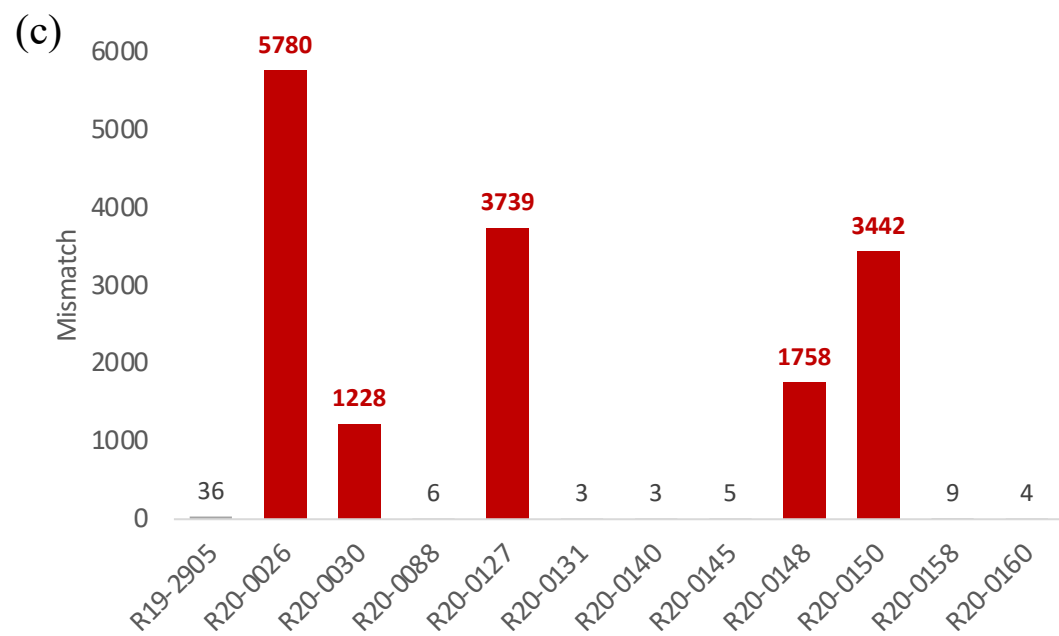
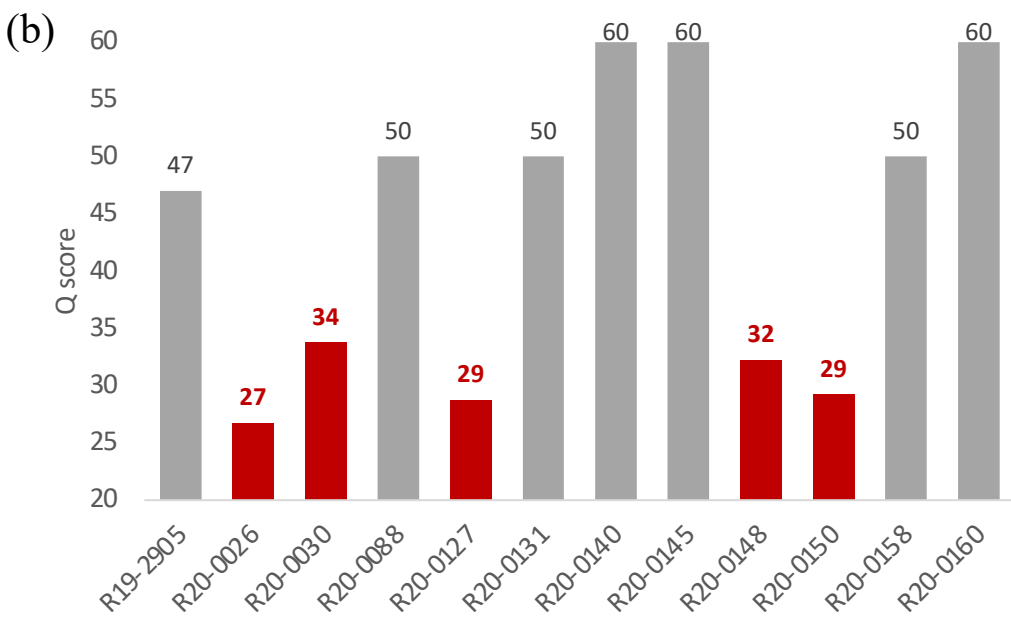
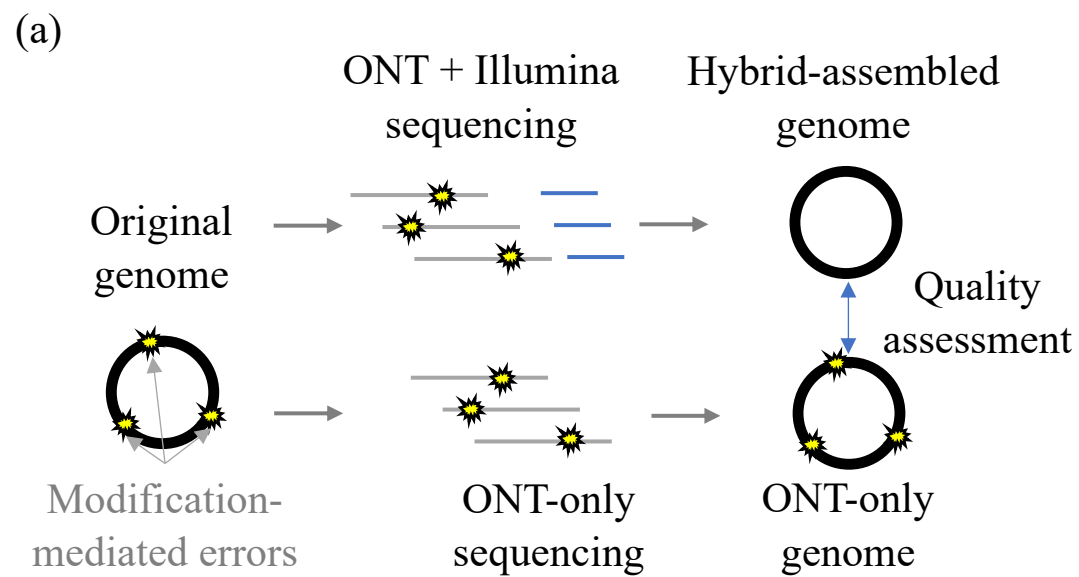
460 **Supplementary Information**

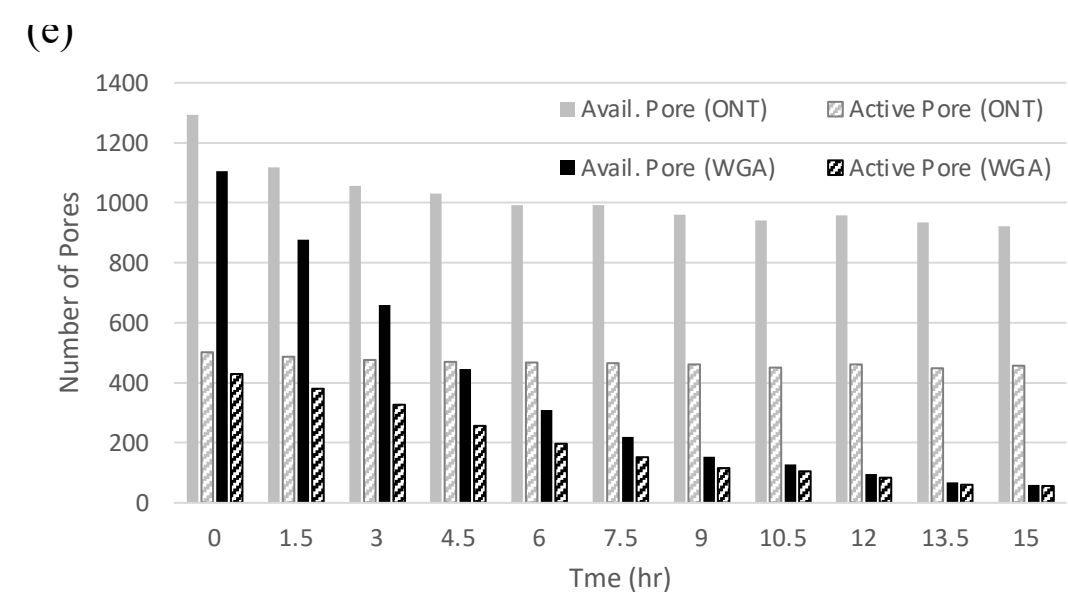
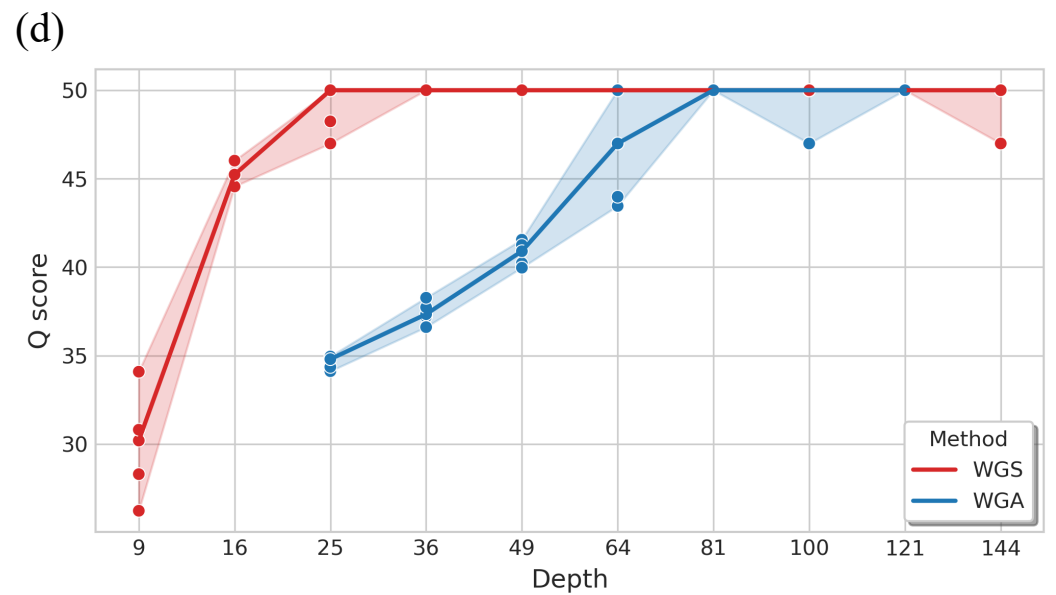
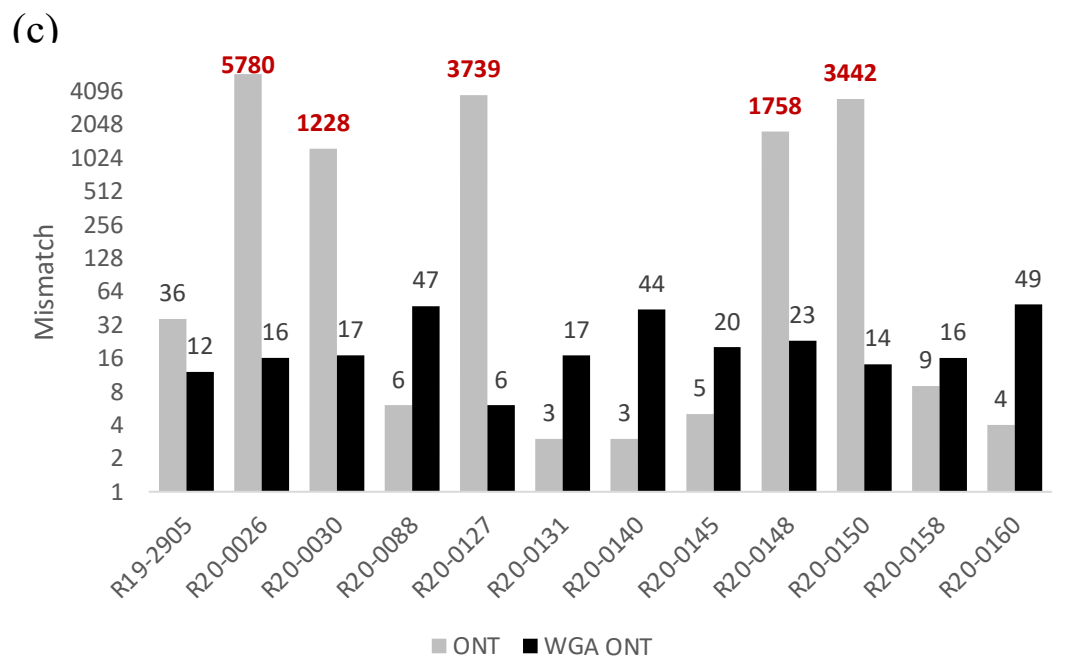
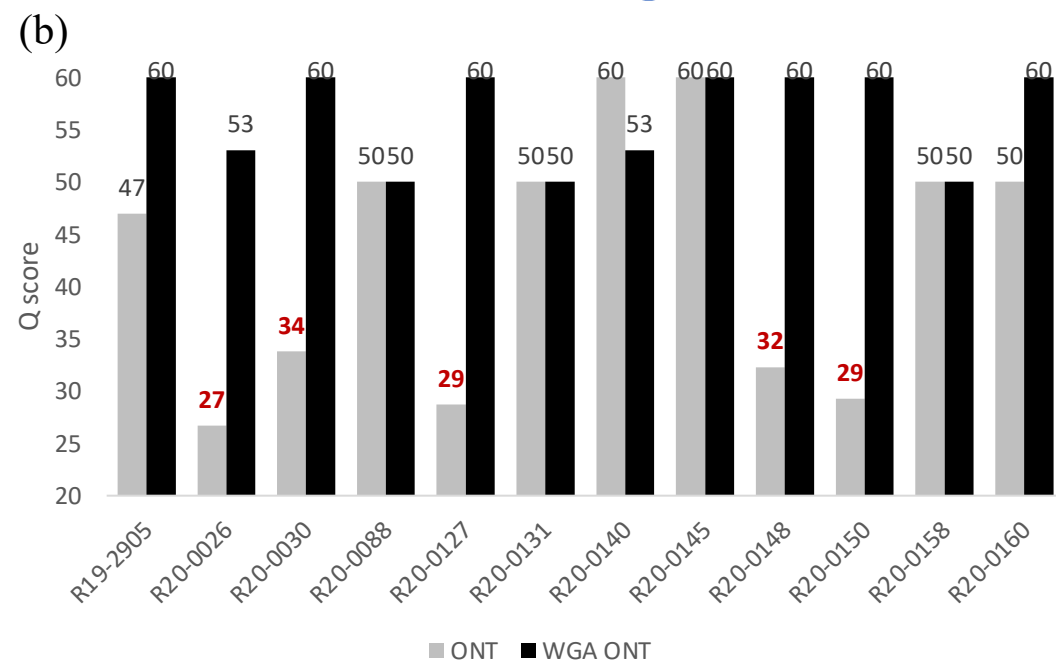
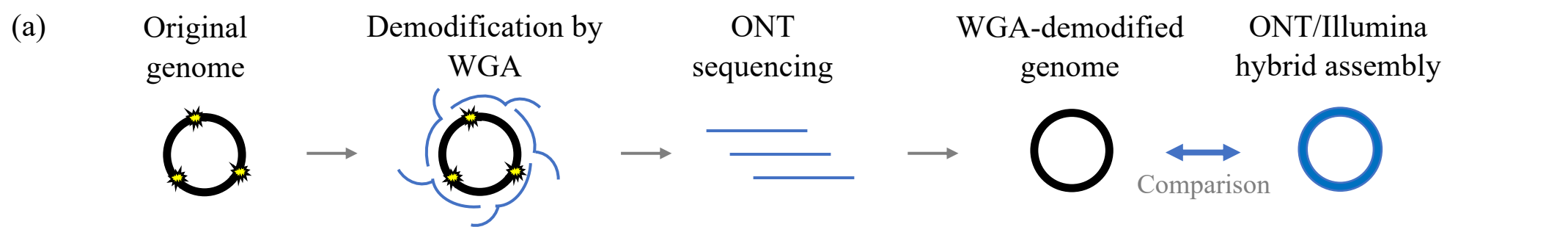
461 **Additional file 1**

462 Additional file 1 includes Supplementary Figures S1-8.

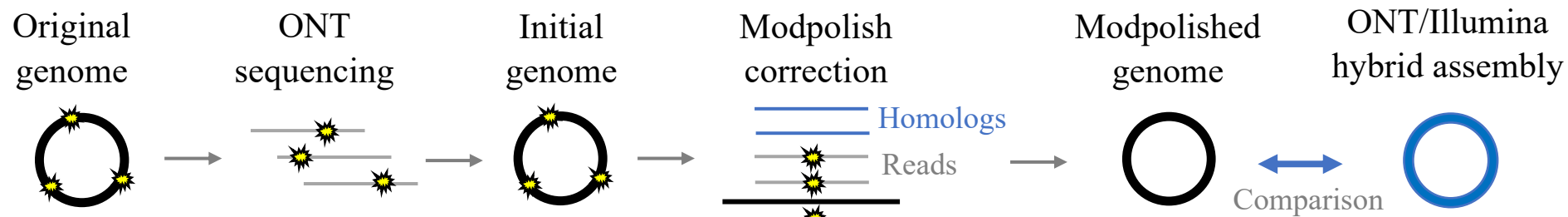
463 **Additional file 2**

464 Additional file 2 includes Supplementary Tables S1-10.

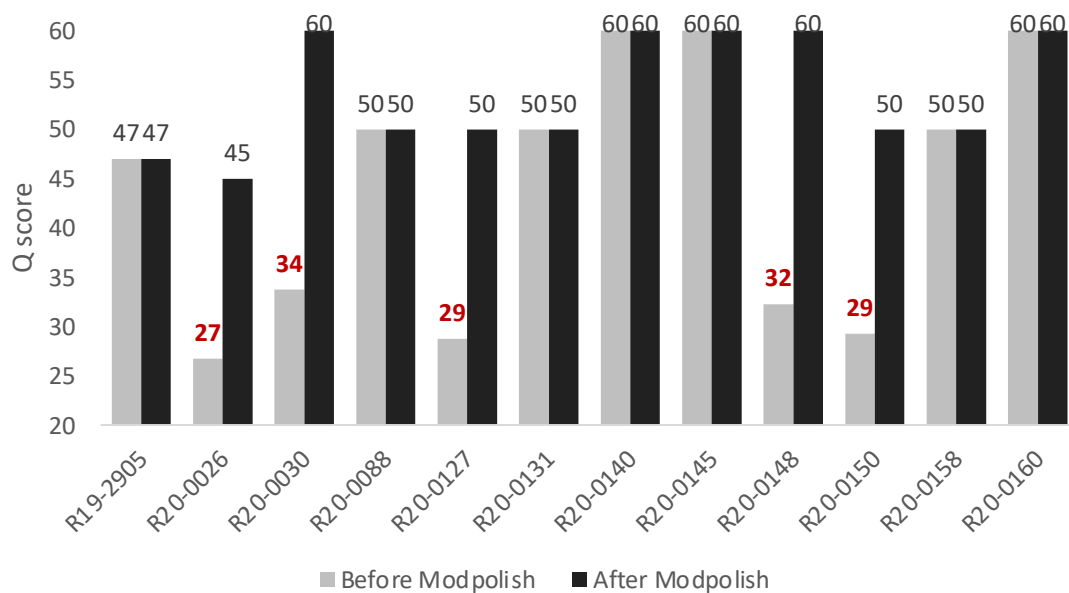




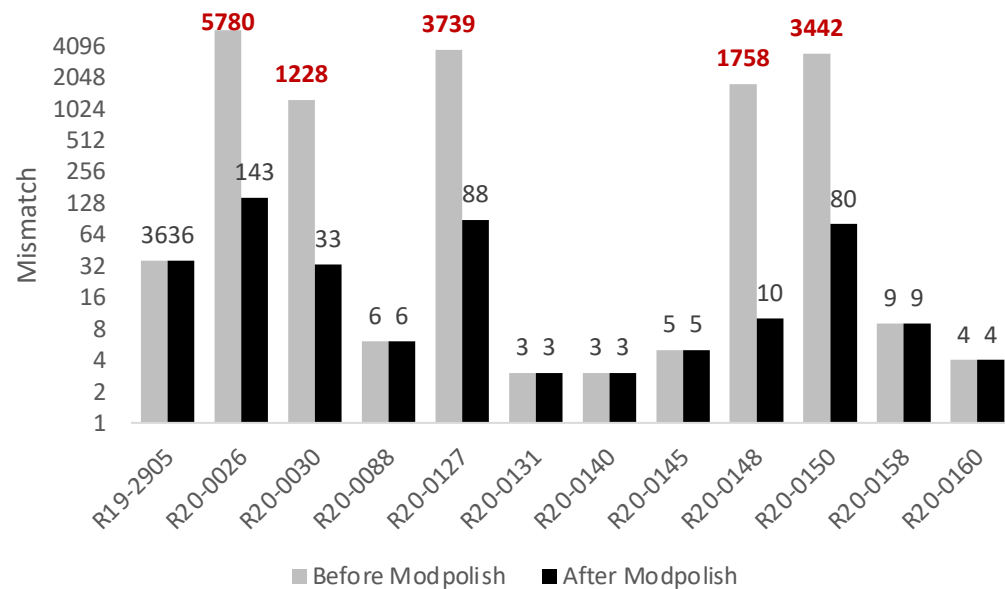
(a)



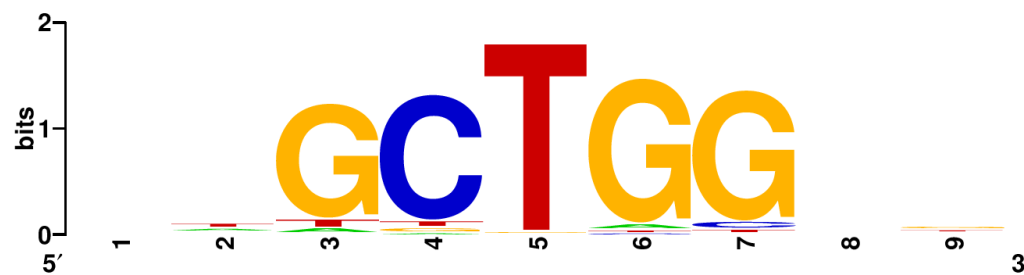
(b)



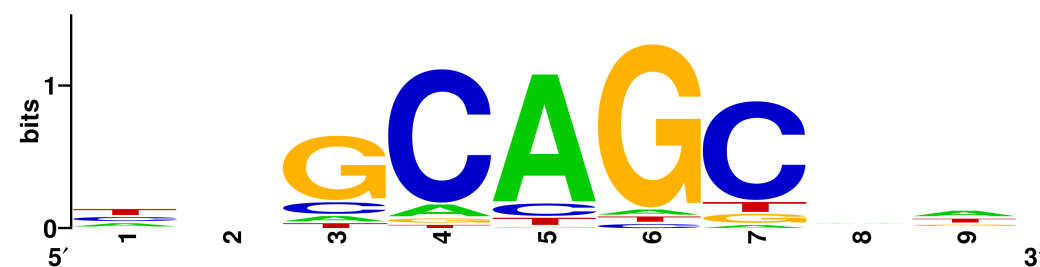
(c)



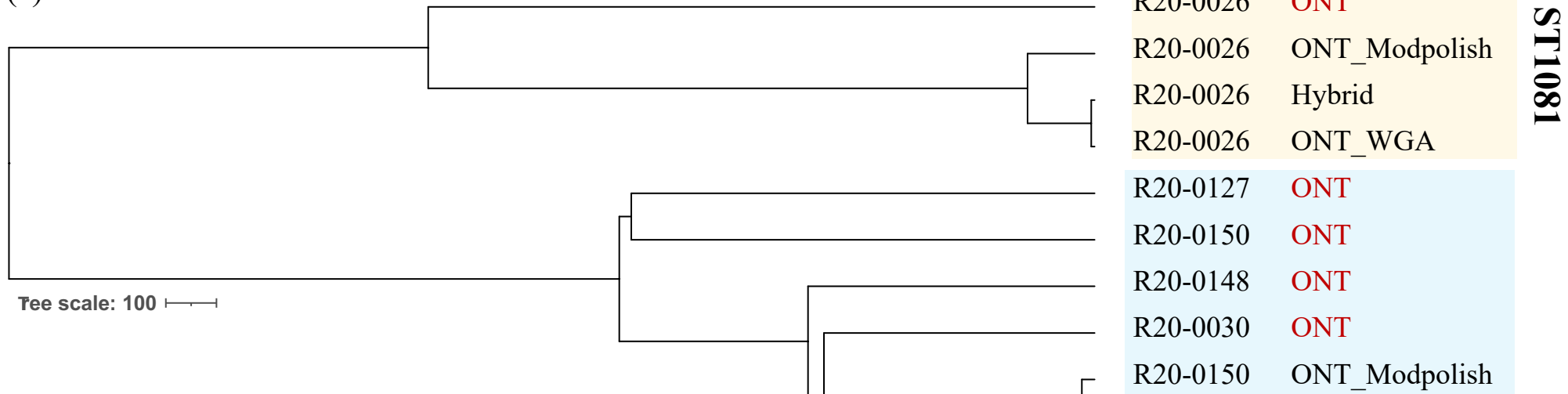
(d)



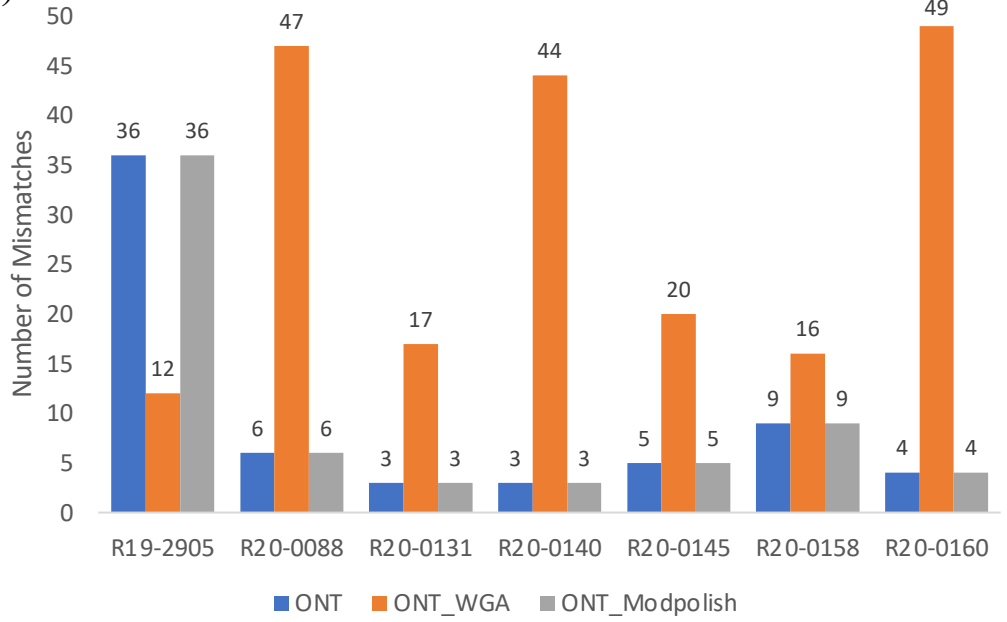
(e)



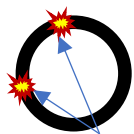
(a)



(b)



Draft genome



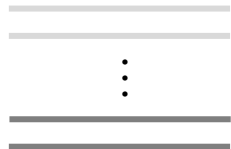
Modification-mediated errors



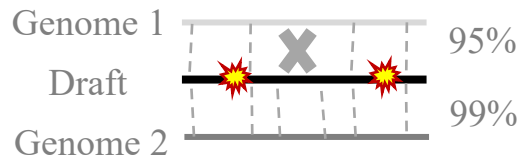
Screen related genomes



Retrieval of related genomes



Estimate structural and nucleotide similarity

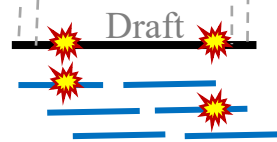


Select closely-related genomes

Whole-genome alignment

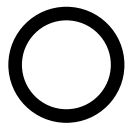


Read alignment



Draft	T	C	A	C	G	A	T	C	C
Homologs	T	C	A	C	C	A	T	C	C
Homologs	T	G	A	C	C	A	T	C	C
Reads	T	G	A	C	G	A	C	C	C
Reads	T	G	A	C	C	A	C	C	C
Qualities	4	4	3	2	1	1	3	3	4
Qualities	5	3	2	2	1	1	2	4	4

Pileups of homologs, reads, qualities



Polished genome

