



Nanopore strand-specific mismatch enables de novo detection of bacterial DNA modifications

Xudong Liu, Ying Ni, Lianwei Ye, et al.

Genome Res. published online October 2, 2024

Access the most recent version at doi:[10.1101/gr.279012.124](https://doi.org/10.1101/gr.279012.124)

P<P	Published online October 2, 2024 in advance of the print journal.
Accepted Manuscript	Peer-reviewed and accepted for publication but not copyedited or typeset; accepted manuscript is likely to differ from the final, published version.
Creative Commons License	This article is distributed exclusively by Cold Spring Harbor Laboratory Press for the first six months after the full-issue publication date (see https://genome.cshlp.org/site/misc/terms.xhtml). After six months, it is available under a Creative Commons License (Attribution-NonCommercial 4.0 International), as described at http://creativecommons.org/licenses/by-nc/4.0/ .
Email Alerting Service	Receive free email alerts when new articles cite this article - sign up in the box at the top right corner of the article or click here .

Advance online articles have been peer reviewed and accepted for publication but have not yet appeared in the paper journal (edited, typeset versions may be posted when available prior to final publication). Advance online articles are citable and establish publication priority; they are indexed by PubMed from initial publication. Citations to Advance online articles must include the digital object identifier (DOIs) and date of initial publication.

To subscribe to *Genome Research* go to:
<https://genome.cshlp.org/subscriptions>

Published by Cold Spring Harbor Laboratory Press

25 **Abstract**

26 DNA modifications in bacteria present diverse types and distributions, playing crucial functional roles.
27 Current methods for detecting bacterial DNA modifications via nanopore sequencing typically involve
28 comparing raw current signals to a methylation-free control. In this study, we found that bacterial DNA
29 modification induces errors in nanopore reads. And these errors are found only in one strand but not
30 the other, showing a strand-specific bias. Leveraging this discovery, we developed Hammerhead, a
31 pioneering pipeline designed for *de novo* methylation discovery that circumvents the necessity of raw
32 signal inference and a methylation-free control. The majority (14 out of 16) of the identified motifs
33 can be validated by raw signal comparison methods or by identifying corresponding
34 methyltransferases in bacteria. Additionally, we included a novel polishing strategy employing duplex
35 reads to correct modification-induced errors in bacterial genome assemblies, achieving a reduction of
36 over 85% in such errors. In summary, Hammerhead enables users to effectively locate bacterial DNA
37 methylation sites from nanopore FASTQ/FASTA reads, thus holds promise as a routine pipeline for a
38 wide range of nanopore sequencing applications, such as genome assembly, metagenomic binning,
39 decontaminating eukaryotic genome assembly, and functional analysis for DNA modifications.

40

41 **Keywords**

42 Long-read sequencing, nanopore sequencing, bacterial DNA modification

43 **Introduction**

44 DNA base modifications are crucial components of epigenetic changes and significantly influence
45 various biological functions. The most extensively studied and understood DNA base modification in
46 vertebrates is 5-methylcytosine (5mC). This type of modification is abundant and can mostly be found
47 in the CpG motif. In prokaryotes, there are three main types of DNA molecule modifications, all of
48 which involve methylation: N^6 -methyladenine (6mA), N^4 -methylcytosine (4mC), and 5mC
49 (Beaulaurier et al. 2019). These methylation forms differ in distribution and function in bacteria, but
50 all play pivotal roles in bacterial life processes. Long-read sequencing platforms, including PacBio
51 single-molecule real-time (SMRT) sequencing and Oxford Nanopore Technology (ONT) sequencing,
52 have enabled the direct detection of DNA modifications (Gouil and Keniry 2019; Amarasinghe et al.
53 2020; Ni et al. 2023a). This is due to the distinct differences in the raw signals produced by modified
54 DNA compared to those produced by canonical DNA sequences. Notably, the detection of 5mC in
55 vertebrate genomes using nanopore sequencing has led to significant advancements over traditional
56 bisulfite sequencing (Liu et al. 2023). However, the challenge of detecting bacterial DNA
57 modifications persists, primarily due to their unique and diverse motifs (Roberts et al. 2022).

58

59 Significant advancements have also been made in developing methods for detecting bacterial DNA
60 modifications through nanopore sequencing. The available methods, including Tombo (Stoiber et al.
61 2017), Dorado (<https://github.com/nanoporetech/dorado>), Snapper (Konanov et al. 2023), and
62 Nanodisco (Tourancheau et al. 2021), rely on interpreting nanopore raw current signals. This process,
63 however, is notably time-consuming and often restricted, as the raw data are typically accessible only
64 to the data producers. Additionally, *de novo* identification of DNA modification sites necessitates a
65 control sample without modifications, which is generally achieved through whole-genome
66 amplification (WGA). This requirement adds complexity and increases the cost associated with
67 bacterial DNA modification detection in nanopore sequencing. Consequently, *de novo* detection of
68 bacterial DNA modifications remains a challenging task.

69

70 The read accuracy of nanopores may be compromised by basecalling errors, which can be caused by
71 bacterial DNA modifications (Rand et al. 2017; Wick et al. 2019). Bacterial DNA methylation can
72 occur frequently in diverse sequence contexts. In nanopore direct RNA sequencing, mismatches in
73 modified RNA molecules can be used to detect modifications (Liu et al. 2019). However, this feature
74 has not been implemented in bacterial DNA sequencing. A more efficient method for identifying

75 bacterial DNA methylation via nanopore sequencing platforms is needed. A pipeline based on sequence
76 mismatches could serve as a more straightforward and time-efficient alternative.

77

78 In this study, we sequenced eight bacterial species, namely, *Acinetobacter pittii* (Liu et al. 2022),
79 *Bacillus cereus* (Cui et al. 2016), *Enterococcus faecium* (Zheng et al. 2007), *Escherichia coli* (Chen et
80 al. 2017), *Klebsiella pneumoniae* (Yang et al. 2019), *Pseudomonas aeruginosa* (Zhao et al. 2023),
81 *Salmonella enterica* (Chen et al. 2020), and *Staphylococcus aureus* (Wang et al. 2017). We identified
82 patterns of unbalanced mismatches between forward and reverse strands on these genomes. The
83 sequences of these gram-negative and gram-positive strains differ in terms of complexity in
84 chromosomes and plasmids, and these strains can be used as representatives for developing a method
85 for detecting DNA modifications *de novo*.

86

87 We developed the Hammerhead to detect bacterial DNA modifications using unbalanced mismatches
88 between forward and reverse strands. Most motifs and sites detected by “Hammerhead” can be further
89 validated using Nanodisco (Tourancheau et al. 2021) and the presence of methyltransferase in different
90 bacterial species. In summary, our results provide insights and solutions for bacterial genome
91 assemblies and identifying modification sites using Nanopore whole-genome reads.

92

93 **Results**

94 **The updated Nanopore sequencing platform observed improvements in read accuracy and** 95 **assembly quality.**

96 Read accuracy is key to successful genome assembly. To evaluate the performance of the R9.4.1 and
97 R10.4.1 flow cells in terms of read quality and genome assembly, we sequenced DNA samples from
98 eight distinct bacterial species under Nanopore R9.4.1, R10.4.1, and high-throughput short-read
99 sequencing (SRS) platforms (see **Supplementary methods**). The R9.4.1 platform produced simplex
100 reads only, while the R10.4.1 platform produced both simplex and duplex reads. A total of 9.64 Gb,
101 8.12 Gb, 15.06 Gb, and 442.4 Mb whole-genome shotgun (WGS) sequencing data were generated for
102 short reads and R9.4.1, R10.4.1 simplex, and duplex long reads, respectively (**Table S1**) (Ye et al.
103 2024). We first assembled the high-quality reference genomes for the eight samples using all R10.4.1
104 simplex and duplex reads, followed by a polishing phase using both long and short reads (see
105 **Methods**). Eight circular bacterial chromosomes and eighteen circular plasmids were obtained (**Table**
106 **S2**).

107

108 The new R10.4.1 reads had a 99% estimated modal read accuracy, outperforming the 97% in R9.4.1
109 reads (**Figs. S1**). The median mapping accuracy for R10.4.1 reads was 98%, 2% higher than the R9.4.1
110 reads (**Fig. S2**). This enhancement was particularly evident in homopolymer regions, where R10.4.1
111 reads achieved 85% accuracy, surpassing the 74% accuracy of R9.4.1 reads (**Figs. S2 and S3**).
112 Nanopore R10.4.1 can produce a small amount (2% to 7%) of duplex reads (**Table S1**). These reads
113 were self-corrected between forward and reverse strands and had even higher accuracy than normal
114 R10.4.1 reads, nearly 99.9% (**Fig. S2**).

115 For quality in bacterial genome assembly, we evaluated the efficacy of using solely R10.4.1 or R9.4.1
116 reads across a range of coverage from 10- to 120-fold (see **Methods**). Assemblies from sole R10.4.1
117 reads outperformed those from R9.4.1 in terms of genome completeness and indels proportion, but
118 these advantages can also be achieved through short-read polishing (**Fig. S4**).

119

120 **Species-dependent efficacy in mitigating single nucleotide substitutions.**

121 Our analyses highlighted distinct patterns of single nucleotide substitutions (SNSs) in bacterial species.
122 For some bacteria, such as *B. cereus* and *S. enterica*, additional short-read polishing did not provide
123 further benefits in decreasing the SNS rate. Conversely, in the case of *A. pittii*, *E. faecium*, *E. coli*, and
124 *K. pneumoniae*, the assemblies based solely on R10.4.1 reads consistently exhibited elevated SNS rates
125 compared to the other assemblies (**Fig. 1A**). Although one might anticipate that random SNS errors
126 would be rectified with increased long-read coverage, this was not observed for these bacterial
127 assemblies. This trend suggested the presence of systematic bias within the reads. We postulate that
128 the elevated SNS rates observed in these assemblies might be associated with unique genomic features,
129 potentially species-specific *k*-mer compositions, or DNA modifications.

130

131 **Substitution types within assemblies exhibit consistent and systematic patterns.**

132 To confirm whether SNSs observed in the *A. pittii*, *E. faecium*, *E. coli*, and *K. pneumoniae* R10.4.1
133 assemblies were technical errors, we quantified all twelve SNS types to determine whether the errors
134 were randomly distributed within the four bacterial assemblies (see **Methods**). In particular, the
135 substitutions cytosine to thymine (C2T) and guanine to adenine (G2A) dominated the SNS landscape
136 (**Fig. 1B**). The frequencies of both C2T and G2A were significantly larger than the simulated
137 distribution (**Fig. S5**). The pronounced prevalence of these specific substitutions, C2T and G2A,
138 suggested that such SNS occurrences are not merely due to random basecalling discrepancies.

139

140 To rule out potential biases from the Medaka (V1.8.0) polishing process

141 (<https://github.com/nanoporetech/medaka>), we additionally polished the assemblies utilizing a
142 different long-read polisher, Racon (V1.5.0) (Vaser et al. 2017). The polishing outcomes still
143 predominantly presented C2T and G2A substitution patterns (**Fig. S6**), indicating that the errors
144 observed cannot be attributed to the polishing process. We can conclude that the bias of substitution
145 type enriched at C2T and G2A in the assemblies originated from technical errors in the R10.4.1 reads.
146 From the SNS frequencies in the 35 assemblies derived from subsampled R10.4.1 reads, we identified
147 the error-prone sites. Specifically, C2T or G2A with more than two occurrences was identified as an
148 error-prone site for C or G, respectively (**Fig. 1C** and **Table S4**).

149

150 **Error-prone sites arise from bacterial DNA modifications.**

151 Given the distinctive SNS patterns observed in assemblies, we aimed to ascertain whether DNA
152 modifications were the primary cause of the identified SNSs. To this end, we utilized the whole-
153 genome amplification (WGA) sequencing approach to produce reads devoid of possible DNA
154 modifications, serving as a control (Ni et al. 2023b).

155

156 Analysis of the correct and incorrect mappings at the error-prone C and G sites revealed that the R9.4.1
157 WGS and WGA reads exhibited similar accuracies, approximately 98% (**Fig. 1D**). However, for the
158 R10.4.1 WGS reads, a significant fraction of the C and G bases were misinterpreted as T (C2T) and A
159 (G2A) (**Figs. 1D** and **S7**). In contrast, very few substitutions were detected at the error-prone sites in
160 R10.4.1 WGA sequencing, with 99% of the C and G being accurately basecalled (**Fig. 1D**). These
161 observations strongly point toward base modifications in these four genomes as the culprits behind the
162 SNS inconsistencies observed in R10.4.1 WGS and WGA reads.

163

164 **A strand-specific mismatch pattern was observed at error-prone sites.**

165 In a deep dive into the distinctions between forward and reverse strands at error-prone sites in the four
166 bacterial species, an intriguing pattern emerged. All the error-prone C sites across the genomes of the
167 two bacteria (*E. faecium* and *K. pneumoniae*) exhibited sharp contrasts: forward reads exhibited an
168 impressive mapping accuracy exceeding 99.9%, whereas the reverse reads plummeted below 20%
169 accuracy (**Fig. 2A** and **2B**). The scenario reversed for error-prone G sites, with low mapping accuracy
170 for forward reads and high mapping accuracy for reverse reads (**Fig. 2A** and **2B**). A similar result was
171 also observed at the error-prone C and G sites within the *A. pittii* and *E. coli* genomes (**Fig. S8**).
172 Considering that forward C is equivalent to reverse G in the genome, as are the single nucleotide
173 substitution patterns of C2T and G2A, error-prone C and G sites may arise from the same type of
174 systemic error.

175

176 To further confirm whether the strand-specific mismatch pattern originated from DNA modification,
177 we investigated the difference between the raw WGS and WGA current signals at the possible modified
178 sites. By checking the signal around a forward error-prone G site in *E. faecium*, we found that the
179 current signals in the WGS reads were significantly different from those in the WGA reads, while the
180 reverse strand was less different (**Fig. 2C**). The same example can be found in *K. pneumoniae* (**Fig.**
181 **2D**). These two sites had a significant proportion of G2A SNSs in the forward reads but not in the
182 reverse ones (**Fig. S9**). From these examples, we believe that the unevenness of the raw signal
183 disturbance between two DNA strands on a modified site is the cause of the strand-specific mismatch
184 pattern. Moreover, the strand-specific mismatch pattern can be used to indicate DNA modification loci
185 in the bacterial genome.

186

187 **Bacterial DNA modifications can be identified through strand-specific mismatch patterns.**

188 With this idea, we developed a pipeline named “Hammerhead” to locate possible bacterial DNA
189 modifications using the mapping accuracy between forward and reverse strands in Nanopore R10.4.1
190 reads. Briefly, after mapping the reads to the genome, the nucleotide difference indices between
191 forward and reverse reads for each genomic site were calculated (**Fig. 3A** and **3B**). Theoretically, the
192 difference index for modification-free WGA reads should be zero. However, random errors still
193 occurred in the WGA reads, creating a background-level difference index. To achieve a false discovery
194 rate (FDR) less than $1e-06$ in all four WGA datasets, we obtained an empirical cutoff of 0.35 (**Figs.**
195 **3C, 3D, and S10**). By counting genomic sites with a difference index larger than 0.35 in WGS reads,
196 we obtained 1820 sites in *A. pittii*, 563 sites in *E. faecium*, 1860 sites in *E. coli*, and 1758 sites in *K.*
197 *pneumoniae* as potential modification sites. These sites included almost all the error-prone sites
198 identified in the genome assemblies (**Fig. S11**). The motifs linked with these sites in *E. faecium* and *K.*
199 *pneumoniae* were mostly enriched as GATC, which is likely to be a 6mA or 5mC motif in bacterial
200 genomes, confirming that the Hammerhead pipeline can locate DNA modifications using R10.4.1
201 reads (**Fig. 4E** and **4F**). Moreover, the motifs of the plasmids inside the two bacteria were similar to
202 those of the chromosome (**Fig. 4E** and **4F**).

203

204 To specify that strand-specific errors exist only in the bacterial sequencing reads, we applied
205 Hammerhead to our human R10.4 data (Ni et al. 2023b). The difference index distributions between
206 the R10.4 WGS and R10.4 WGA reads were similar (**Fig. S12**), confirming that this pattern is bacterial-
207 specific.

208

209 We then applied the Hammerhead pipeline to the other four bacteria to evaluate its performance.
210 Potential modification sites were also identified in the four other bacteria (**Fig. S13**). These sites in all
211 eight bacterial species were evenly distributed inside the chromosomes or the plasmid (**Fig. S14**).
212 Moreover, the motifs were enriched in sequences near these potential modification sites (-10 bp to +9
213 bp) with MEME E-values less than $1e-50$ (**Figs. 3E, 3F, S15, and S16**). The plasmids in bacteria
214 usually share a similar pattern as bacterial chromosomes (**Fig. S16**).

215

216 **Nanodisco validates and supports the identified bacterial DNA modifications from Hammerhead.**

217 Nanodisco (Tourancheau et al. 2021) is a well-recognized tool for detecting bacterial DNA
218 modifications using the current signal difference between R9.4.1 WGA and WGS reads. Fortunately,
219 we had access to the corresponding R9.4.1 datasets for all the bacterial samples. To validate these
220 potential modification motifs and sites identified from Hammerhead, we utilized Nanodisco to perform
221 *de novo* identification of methylation.

222

223 Modifications in bacteria are consistently linked to specific motifs, with more than 95% of
224 modification sequence motifs undergoing methylation (Casadesús and Low 2006; Wion and Casadesús
225 2006; Beaulaurier et al. 2019). Based on this rationale, the enriched motifs obtained from Hammerhead
226 were compared with those derived from Nanodisco to assess whether the identified motifs are indeed
227 indicative of true DNA modifications. A total of 14 motifs were identified by Hammerhead, while 10
228 were detected by Nanodisco. Six motifs were discovered by both approaches (**Fig. 4A**).

229

230 All the shared motifs (n=6) or their reverse complements (RCs) exhibited current differences (**Fig.**
231 **S17**). Moreover, four of these motifs were predicted to be involved in 6mA methylation, including
232 C6mATCTC in *A. pittii*, R6mAYCNNNNNTTRG and CYA6mANNNNNNGRTY in *E. faecium*, and
233 G6mATC in *K. pneumoniae* (**Figs. 4B-4D and S18**). Additionally, one motif, RTAGACGC (the RC of
234 GCGTCTAY), in *P. aeruginosa* was identified as the 4mC methylation type (**Figs. 4E and S18**).
235 Although YGAAGC in *A. pittii* was not characterized as a specific methylation type (4mC, 5mC, or
236 6mA), we believe that it belongs to the A-base methylation category based on the significant disparity
237 observed in the current signal profile of the first A base (**Figs. S17 and S18**).

238

239 There are eight motifs uniquely identified by Hammerhead. Among these motifs, GATC in *E. coli* is
240 a well-known 6mA motif. However, this motif was missed in the Nanodisco result. The reason might
241 be attributed to the high frequency of CCWGG motif in *E. coli* genomes (Breckell and Silander 2022).

242 Nanodisco calculates p-values for each base using the Mann–Whitney U test to indicate the
243 significance of the current signal difference between native DNA reads and the negative control. Only
244 the top 2000 peaks of 5 bp smoothed p-values were selected for final motif enrichment analysis
245 (Tourancheau et al. 2021). In the *E.coli* genome, the 5mC motif CCWGG had smaller p-values than
246 other motifs (Breckell and Silander 2022) thus occupying most of the top 2000 peaks. As a result, the
247 GATC motif could not be enriched using Nanodisco using default parameters.

248

249 To further confirm whether Nanodisco could validate our novel motif finding, we compared the
250 differences in the current signals between WGA and WGS R9.4.1 reads using Nanodisco (**Fig. S19**).
251 The differences among all the motifs were confirmed, indicating that these motifs could be identified
252 by Nanodisco by fine-tuning its cutoff. Additionally, the detection of six motif-related
253 methyltransferases in their corresponding species strongly supported the occurrence of modifications
254 within the motifs (**Table 1** and **Table S5**).

255

256 To further validate all potential modification sites identified from Hammerhead, a comparison was
257 conducted at the site level with Nanodisco. Specifically, potential methylation sites identified by
258 Hammerhead with a difference index ≥ 0.35 were selected. For the Nanodisco sites, the RDS files
259 generated from the function of “difference” were used to obtain the site position and mean current
260 differences. Sites with an absolute current difference greater than 2 pA were selected for downstream
261 comparison (**Fig. S20**).

262

263 Considering that methylation can impact the current of bases near the modified sites, a methylation
264 unit was employed to assess the overlaps between Hammerhead and Nanodisco. The methylation unit
265 consists of a 9-mer base unit centered on the potential methylation sites from Hammerhead or
266 Nanodisco. Based on the number of overlaps of these methylation units, the majority of modified sites
267 identified by Hammerhead were confirmed by Nanodisco in most bacterial species ($>75\%$ intersection),
268 except for *S. enterica* (about 30% intersection) (**Fig. 4F**).

269

270 **The 5mC motif GATC in *S. aureus* can be detected by decreasing the cutoff.**

271 Four motifs were identified by Nanodisco but not Hammerhead (**Figs. 4A** and **S21**). Among them, the
272 SATSNNSNNSNNS motif seemed to be a false positive result. The bit scores are low for all bases in
273 this motif, indicating low representation and high uncertainty during motif enrichment and typing
274 analysis (**Fig. S21**). Moreover, the difference in the current signal density between the WGA and WGS
275 reads was mild for this motif (**Fig. S22**). The other three motifs, RCCWGGHND, RCCWGGY, and

276 KGATCADYHDNHWR, were predicted to be the 5mC methylation type based on the Nanodisco
277 results and the presence of motif-related methyltransferases (**Fig. S21** and **Table S5**).

278

279 There are two possible reasons for the false negative results for detecting DNA modification motifs
280 using Hammerhead. The first issue is the cutoff issue. A cutoff difference index > 0.35 was used to
281 avoid false positives, with an empirical FDR less than $1e-06$. The cutoff may be too harsh for low-
282 frequency DNA modifications to be enriched. The other reason is that the motif and modification have
283 been incorporated into training datasets for the current basecalling models. If so, the read accuracy
284 around these motifs should be as high as that of other genomic regions in the WGS reads.

285

286 To confirm which of these mutations is the cause of the missing CCWGG and GATC motifs in
287 Hammerhead, we performed a search for the two motifs within the corresponding genomes and
288 examined whether there were any differences in the observed read accuracy. A noticeable decrease in
289 read accuracy was observed for the GATC motif at the G-base and C-base (**Fig. S23**), indicating that
290 this motif could be identified by Hammerhead using the raw basecalling model. The strict cutoff value
291 for the difference index (0.35) might explain the lack of identification of this motif in *S. aureus*. To
292 validate whether the cutoff made a difference, 3821 sites with a difference index greater than 0.1 (FDR
293 $< 1e-03$) were selected for downstream motif enrichment analysis. At this time, we identified three
294 motifs, GAT5mC and two other motifs identified earlier, and corresponding methyltransferases were
295 also detected (**Figs. S15C, S24, and Table S5**).

296

297 **The 5mC motif CCWGG can be retrieved by fine-tuning the basecalling model.**

298 In contrast, the base accuracy within the CCWGG motif in both *E. coli* and *K. pneumoniae* was
299 approximately 99%, which closely aligns with the expected accuracy for R10.4.1 reads (**Fig. S23**). The
300 basecalling model seems to have incorporated this modified motif into the training dataset. To detect
301 modifications within the CCWGG motif using Hammerhead, we needed to revert the model to increase
302 strand-specific errors in the modified CCWGG motif.

303

304 The modification system in *E. coli* has been well studied. C5mCWGG has been associated with the
305 *dcm* gene, while G6mATC has been associated with the *dam* gene (Marinus and Løbner-Olesen 2014).
306 We searched for known *dam* (NCBI Gene ID: 947893) and *dcm* (NCBI Gene ID: 946479) genes in
307 our *E. coli* assembly. The presence of methylation within the CCWGG and GATC motifs was
308 confirmed by detection of both genes (**Fig. S25** and **Tables S6-S8**). Therefore, we chose *E. coli* datasets
309 to test whether a new basecalling model could help Hammerhead identify methylation in the CCWGG

310 motif.

311

312 We trained a “modification aware” basecalling model for *E. coli*, starting from the super accuracy
313 basecalling (SUP) model ([dna_r10.4.1_e8.2_400bps_sup@v4.2.0](#)), fine-tuned using methylation-free
314 WGA reads (see **Methods**). To ensure that the “modification aware” model does not have an extensive
315 impact on WGS read quality, we first calculated the observed accuracy and mismatch proportion using
316 rebasecalled *E. coli* WGS reads. The overall read accuracy, mismatch ratio, and homopolymer
317 identification accuracy were consistent with those obtained using the SUP model (**Figs. S2 and S26**).

318

319 To assess the ability to detect modification sites using the new model, we applied the Hammerhead
320 pipeline to the rebasecalled *E. coli* WGA and WGS reads. A similar distribution of difference indices
321 was observed for the WGA reads as for the SUP model (**Fig. S27A**). With respect to the WGS reads,
322 a greater number (36503 vs. 1860) of potential modification sites were identified with the same
323 difference index cutoff of 0.35 (**Fig. S27A**). These sites were selected for downstream motif
324 enrichment analysis. Only two motifs were significantly enriched. One was GATC, and the other was
325 CCWGG (**Fig. S27B**). The significant increase in strand-specific errors around the modification sites
326 enabled the detection of the 5mC modification in CCWGG with Hammerhead, which was missed
327 using the SUP model.

328

329 **Hammerhead exhibited comparable precision and reliability to ONT official software Dorado in**
330 **the detection of 6mA methylation.**

331 After validating the ability for *de novo* bacterial methylation finding through comparison with
332 Nanodisco using R9.4.1 reads, we further benchmarked Hammerhead against a tool based on R10.4.1
333 reads to better showcase its performance. Dorado (Oxford Nanopore 2023) is the official software by
334 ONT for identifying modifications using machine learning. It is important to note that Dorado is not a
335 *de novo* bacterial methylation detection tool but specifically focuses on identifying 6mA or 5mC
336 methylation.

337

338 We needed highly confident DNA methylation sites for benchmarking methylation detection methods
339 using R10.4.1 reads. We used the sites from two 6mA motifs, which have been identified by Nanodisco
340 methylation motif typing with R9.4.1 reads. Only the motifs with the modified site prediction
341 percentage higher than 85% were selected. These two motifs were expected to be nearly 100%
342 modified in WGS reads, and 0% modified in WGA reads. The first motif is a type I motif,
343 R6mAYCNNNNNTTRG, in *E. faecium* (n=1063), and the second motif is a type II motif,

344 C6mATCTC, in *A. pittii* (n=1318). The genomics sites containing the two motifs could be considered
345 as true positives in WGS reads and true negatives in WGA methylation-free reads.

346

347 The distribution of the difference index from Hammerhead and methylation proportion from Dorado
348 demonstrated that both methods performed well in identifying positive and negative sites (**Fig. S28**).
349 Furthermore, the performance of the two methods was illustrated using receiver operating
350 characteristic (ROC) curves and precision-recall (PR) curves (**Figs. S29** and **S30**). Specifically, the
351 areas under the curve (AUC) of the two methods for both motifs were greater than 0.99 in the ROC
352 curve (**Fig. S29**). When it comes to the PR curve, Dorado achieved a higher AUCPR value (0.993)
353 compared to Hammerhead (0.980) in the type I motif, while Hammerhead (0.999) outperformed
354 Dorado (0.987) in the type II motif (**Fig. S30**).

355

356 In summary, our Hammerhead successfully identified a total of 16 methylation-related motifs across 8
357 representative bacteria using R10.4.1 reads (**Table 1**). Among them, 14 motifs were further validated
358 either by the results obtained from Nanodisco or through the detection of methyltransferases.
359 Meanwhile, Hammerhead exhibited comparable precision and reliability to Dorado in the detection of
360 6mA methylation. These validations have provided evidence that Hammerhead is effective in
361 accurately locating bacterial DNA modifications.

362

363 **A high-accuracy model can be used for modification detection but not a fast model.**

364 In the context of converting the current signal into base information, ONT provided three types of
365 basecalling models, namely, fast basecalling (FAST), high accuracy basecalling (HAC), and super
366 accuracy basecalling (SUP), for R10.4.1 reads. We developed the Hammerhead pipeline based on reads
367 basecalled on the SUP model. To test whether the strand-specific error pattern was also present in the
368 FAST and HAC models, the WGA and WGS reads of *A. pittii* were re-basecalled using the FAST and
369 HAC models, respectively. We first compared the read quality using the basecalled results obtained
370 from three different models. The results revealed that the reads from the SUP model exhibited the
371 highest quality performance compared to those from the other models. On the other hand, reads from
372 the FAST model showed the highest rate of inaccuracy (**Fig. S31**).

373

374 We applied the Hammerhead pipeline to both FAST and HAC reads. In WGA reads, the FAST model
375 identified more sites (3,561) with a difference index over 0.35 compared to the HAC model, which
376 identified only 57 sites (**Fig. S32A**). The 3,561 identified sites were false positives, which could impact
377 the further motif enrichment from the Hammerhead result. The HAC model reads posed a similar

378 WGA/WGS difference index distribution pattern, compared to the reads from the SUP model (**Figs.**
379 **S10** and **S32B**). Furthermore, from the potential modification sites detected using HAC model reads,
380 we were able to enrich the same modification motifs as those identified using SUP model reads (**Fig.**
381 **S32C**). As for FAST model reads, only one motif could be enriched (**Fig. S33**). To quantify the
382 performance of Hammerhead across different models, ROC and PR curves were generated. The results
383 of ROC and PR curves led to the same conclusion that the SUP model had the highest performance,
384 with an AUC of 1.000 and AUCPR of 0.999 (**Figs. S29B** and **S30B**). On the other hand, the FAST
385 model exhibited the lowest performance, with an AUC of 0.968 and AUCPR of 0.969 (**Figs. S34** and
386 **S35**). Additionally, the HAC model demonstrated a performance close to the SUP model, with an AUC
387 of 0.999 and AUCPR of 0.996 (**Figs. S34** and **S35**). Based on these findings, we recommended using
388 HAC reads and SUP reads to detect modifications using Hammerhead.

389

390 **Duplex polishing resolves nucleotide substitution errors in assemblies.**

391 Although Hammerhead can be used to effectively identify modification sites in bacterial assemblies
392 derived from R10.4.1 reads, the issue of substitution errors (C2T and G2A) in these assemblies remains
393 a concern when these assemblies are used as references. The “duplex basecalling” technique utilizes
394 both the forward strand and reverse strand of a DNA fragment for sequencing and basecalling
395 (Silvestre-Ryan and Holmes 2021). In such a situation, the Nanopore duplex reads can achieve even
396 greater quality than normal ones (**Figs. S1-S3**). We wanted to investigate whether error-prone sites in
397 assemblies can be corrected by employing duplex reads for polishing.

398

399 To this end, we calculated the observed accuracy for each error-prone C and G site. For the R10.4.1
400 reads, the average observed accuracy at error-prone sites ranged from 54 to 65% (**Figs. 5A** and **S36**).
401 In contrast, the average observed accuracies for duplex reads exceeded 88%, except for that for *E.*
402 *faecium*, which was approximately 78%, suggesting the potential to correct substitution errors in
403 assemblies by polishing with highly accurate duplex reads (**Figs. 5A** and **S36**). We then polished the
404 assemblies with duplex reads for *E. faecium* and *K. pneumoniae*. The SNS percentages in both bacterial
405 genome assemblies decreased after polishing with duplex reads and were comparable to those polished
406 with short reads (**Fig. 5B** and **5C**). This result suggested that the SNSs caused by modification could
407 be effectively corrected using duplex reads.

408

409 Considering that most error-prone sites are potential modification sites (**Fig. S11**), we wanted to
410 determine whether focusing solely on polishing those potential modification sites could effectively
411 minimize SNS errors via duplex reads with limited yields. To validate this idea, the Hammerhead

412 method was used to identify potential modification sites. Subsequently, the alignment files of the
413 duplex reads were utilized to polish these sites (**Fig. 5D**). The duplex read-polished assemblies
414 exhibited SNS accuracy comparable to that of the short-read (50-fold) polishing, with a coverage of
415 approximately 20-fold (**Fig. 5E and 5F**). The duplex-read polishing process works well even for
416 bacteria with limited duplex reads, such as *A. pittii* (10-fold) and *E. coli* (4-fold), in our study (**Fig.**
417 **S37 and Table S3**). The polishing pipeline showed a reliable reduction in SNS density and could be
418 incorporated into the bacterial assembly pipeline using only R10.4.1 reads.

419

420 **Discussion**

421 In this work, we conducted a benchmark analysis to evaluate the quality of reads and assemblies
422 generated from the most recent ONT R10.4.1 flow cell. Our results demonstrated that R10.4.1 reads
423 were superior in terms of read accuracy and homopolymer detection compared with the R9.4.1 reads.
424 The genome assembly is the consensus of sequenced reads. The quality of raw reads directly impacts
425 genome assembly when using only nanopore reads. In theory, random errors in reads, with a proportion
426 less than 50%, can be corrected during genome assembly or the self-polishing stage. Consequently, the
427 accuracy of the assembly is usually higher than that of the raw reads. However, some error sites,
428 particularly those with error proportions greater than 50%, which were located within complex
429 repetitive regions or close to other error sites, are difficult to be corrected using only nanopore reads.
430 These errors persist in the final genome assembly. Our analysis revealed that bacterial assemblies
431 obtained using only R10.4.1 reads were comparable to those obtained with short-read polishing in
432 terms of completeness and number of indels. However, in certain R10.4.1 assemblies, the occurrence
433 of base modifications results in a considerable number of error-prone sites with single nucleotide
434 substitutions, particularly C2T and G2A substitutions. We hypothesized that the base methylation is
435 the root cause of the enriched error substitution types. We confirmed this idea by comparing native
436 reads with negative controls at error-prone C and G sites. Based on the novel finding of strand-specific
437 mismatch patterns, we developed a method named Hammerhead to identify bacterial DNA
438 modification. This method was further validated to work effectively and accurately.

439

440 The R10.4.1 reads have better sensitivity in bacterial DNA methylation detection than the R9.4.1 reads.
441 The specific strand mismatch pattern was clearly evident in R10.4.1 reads, while not easily observed
442 in R9.4.1 reads at error-prone C and G sites identified from R10.4.1 assemblies (**Fig. 1D**). The
443 assembly errors from the sole R10.4.1 read have also been identified in different *K. pneumoniae* strains
444 (Lohde et al. 2024). The deviation of Hammerhead's difference index between WGS and WGA from

445 R9.4.1 reads is less pronounced than that of R10.4.1 (**Figs. S28B** and **S38**). In general, the
446 improvement of read accuracy in R10.4.1 highlighted the technical errors, which are utilized in
447 Hammerhead. The read quality of R10.4.1 is higher than that of R9.4.1 (**Fig. S2**), resulting in lower
448 random errors and reduced background noise, thereby making it easier to detect mismatches caused
449 by methylation. Similarly, we have observed that the false-discovery rate (FDR) for CpG methylation
450 in the human genome is lower in R10.4 compared to R9.4.1, due to the higher read accuracy (Ni et al.
451 2023b). We could conclude that the difference between R10.4.1 and R9.4.1 reads is the combination
452 effect from the differences of pore protein, motor protein, and basecaller. The library preparation
453 protocols we used for R9.4.1 and R10.4.1 are both “ligation” kits (LSK110 and LSK114). The two kits
454 use different types of motor proteins, which were ligated to DNA molecules during library preparation.
455 The raw current signal is generated when DNA passes through the pore protein, and different protein
456 structures have distinctive signal patterns (Deamer et al. 2016; Peraro and van der Goot 2016; Bhatti
457 et al. 2021; Mayer et al. 2022). The motor protein, however, binds with the pore protein and controls
458 the sequencing speed of the DNA, affecting the sampling frequency and the current pattern. The
459 basecall process, which translates current signal to nucleotide, is now handled by deep learning based
460 basecalling models. The model architecture for R9.4.1 and R10.4.1 is similar and embedded in
461 basecallers. However, the models were trained with different datasets, which may also affect the error
462 pattern (Seymour 2019; Xu et al. 2021).

463

464 Bacterial DNA methylation differs from that in mammals in both motifs and functions. In mammals,
465 the most abundant form of DNA methylation is 5mC, and cytosine methylation mainly occurs within
466 CpG dinucleotides (Petryk et al. 2020). DNA methylation in mammals is essential for embryonic
467 development and cellular function (Greenberg and Bourc’his 2019; Dahlet et al. 2020; Grosswendt et
468 al. 2020). In bacteria, there are three primary forms of methylation typing: 6mA, 4mC, and 5mC
469 (Casadesús and Low 2006; Beaulaurier et al. 2019). The motifs of certain modified sites vary; for
470 example, 6mA can be detected at GATC and CATCTC sequences. The primary function of bacterial
471 DNA methylation is associated with restriction-modification (RM) systems (Casadesús and Low 2006;
472 Loenen et al. 2013; Roberts et al. 2022). After validating Hammerhead’s ability to detect bacterial
473 DNA methylation, we applied our method to human cell line reads to investigate whether technical
474 issues were present in human samples. The similar distribution of the difference index between WGS
475 and WGA reads indicated that this specific-stand error pattern was not observed at human reads (**Fig.**
476 **S12**). The absence of this pattern may be due to the inclusion of human native DNA reads in the
477 basecall training datasets.

478

479 Hammerhead is easy to implement for single-bacterial sequencing or metagenomic sequencing data
480 because it requires only WGS reads, without the need to infer the raw current signal. We set a cutoff
481 of 0.35 for the reference index, which is a $1e-06$ false discovery rate based on our R10.4.1 WGA
482 datasets. This cutoff may be too harsh for some species/strains that have a distinct k -mer composition
483 compared to the four strains in our study. The harsh cutoff may result in the absence of some motifs,
484 such as the GATC in *S. aureus* (**Fig. S24**). We have provided the “cutoff” parameter in Hammerhead,
485 with the default value of 0.35. Users can adjust the cutoff for their own data. Meanwhile, we have
486 added an option to select the top N sites (with a default number of 2000) from Hammerhead. The top
487 sites can be used for downstream motif enrichment analysis, mimicking the behavior of Nanodisco, to
488 avoid further cutoff adjustment.

489

490 Although the negative controls such as WGA reads are not required by Hammerhead, we have
491 showcased the three usages of negative controls as follows. First, the additional negative controls can
492 help to further validate the results. After the user gets the enriched motifs, a comparison of the
493 difference index distribution of resulting motifs between native reads and methylation-free reads can
494 be conducted to check whether there is a significant difference (**Fig S28**). Second, a negative control
495 sample could be used to calculate the background noise for the difference index in a specific species.
496 Based on the difference index distribution of control reads, users can adjust the cutoff based on the
497 required FDR. Third, the methylation-free sequencing data can be used to retrain a “modification-
498 aware” model, which can also increase the sensitivity of modification detection. We have demonstrated
499 the process and effect of retraining in the detection of the CCWGG motif in our *E. coli* strain (**Figs**
500 **S10, S16B, and S27**). However, the retraining process takes a longer time and requires GPU resources.
501 For users who do not have WGA reads or GPU resources but still want the high sensitivity from the
502 “modification-aware” basecall model, we have provided our “modification-aware” model trained from
503 our four methylation-free bacterial datasets (<https://doi.org/10.6084/m9.figshare.25858072>).

504

505 Hammerhead was developed based on the stand-specific mismatch pattern. Hammerhead located the
506 potential modifications by selecting genomic sites with a difference index higher than the cutoff. Based
507 on those potential modification sites, the methylation motif could be enriched. However, one of the
508 limitations is the quantitative measurement for individual sites. Hammerhead’s difference index cannot
509 directly reflect the proportion of DNA methylation at specific sites. Although the difference index can
510 be used to compare the change of methylation between different samples (**Fig. S28**), it cannot be used
511 to measure the methylation proportion between different motifs. Bacterial DNA methylation is highly
512 motif-driven, with over 95 percent of nucleotides being modified at methylation motifs, indicating an

513 all-or-none case for a given methylation motif (Casadesús and Low 2006; Beaulaurier et al. 2019). We
514 have demonstrated the all-or-none principle in two 6mA motifs, namely C6mATCTC and
515 R6mAYCNNNNNTTRG (**Fig. S28**). Both the difference index from Hammerhead and the
516 methylation proportion from Dorado exhibited a near “one to zero” difference between the native reads
517 and WGA reads (**Fig. S28**), giving an insight that counting methylated or non-methylated sites in one
518 motif could be an alternative quantitative result from Hammerhead. Another limitation of
519 Hammerhead is the lack of ability to judge the methylation type and the position of the modified site
520 within the motif. To fill this gap, a comprehensive database of bacterial methylation like Rebase
521 (Roberts et al. 2022) is considered to be included in the Hammerhead pipeline in the future.

522

523 Hammerhead, as an easy way to detect bacterial DNA methylation, can bring new biological insights
524 and downstream technical applications. Bacterial DNA methylation, primarily mediated through the
525 RM system, governs various cellular functions such as crucial aspects of virulence and metabolism.
526 For instance, the newly identified Type I RM system in *Pseudomonas syringae* has been shown to play
527 pivotal roles in virulence and metabolic pathways, influencing critical processes like the secretion
528 system, biofilm formation, and translational efficiency (Huang et al. 2024). The Hammerhead could
529 quickly detect modification motifs in bacterial genomes, which adds a new layer of information to
530 studying bacterial virulence and horizontal gene transfer (Tisza et al. 2023). Furthermore, there are
531 vast technical applications from the downstream analysis of the Hammerhead result. For instance, in
532 metagenomic sequencing, the contig binning now relies on GC content, *k*-mer composition, and length.
533 In some complex environment samples, the traditional classification does not work well. We showed
534 that the plasmid generally follows the methylation pattern of chromosomes for certain species (**Figs**
535 **3E, 3F, and S16**), thus the methylation motif identified from Hammerhead can be further incorporated
536 into the binning processing in metagenomic analysis (Tourancheau et al. 2021). Additionally,
537 according to the outputs of Hammerhead, the strand-specific error pattern is absent in human reads,
538 but present in bacteria reads (**Figs. 3C, 3D, S10, and S12**). This pattern could be used to distinguish
539 the host and bacterial reads, and help to remove DNA contamination from bacterial DNA, or vice versa.

540

541 Hammerhead can *de novo* identify all potential methylation sites with one command, from
542 FASTQ/FASTA input. The recent development of new deep learning models has enabled Dorado, the
543 basecaller software of ONT, to identify bacterial methylation sites using raw current signal files
544 (FAST5 or POD5). The raw signal files are roughly five to ten times larger than the FASTQ/FASTA
545 file, typically discarded by users soon after sequencing due to their size. Additionally, Dorado's
546 modification calling requires the use of different deep learning models, separate from the basecalling

547 model. For bacterial DNA modification, users must run three different models: one for basecalling,
548 one for 6mA, and another for 5mC/4mC. These deep learning-based processes also require GPU-based
549 computation, which significantly increases the time and resources needed.

550

551 In summary, our study offers valuable insights into methylation detection and solutions for bacterial
552 genome assemblies using only ONT reads. In the R10.4.1 reads, mismatches between the forward and
553 reverse strands were found to be linked with DNA modifications and thus can be used to identify
554 possible modification sites. Building upon the identification of a strand-specific error pattern caused
555 by the base methylation, we developed the "Hammerhead", the first tool enabling *de novo* DNA
556 modification calling from ONT basecalled reads, eliminating the need to infer raw ionic currents.
557 Hammerhead demonstrated strong performance in methylation detection, validated by the Nanodisco,
558 Dorado, and the presence of methyltransferases. Importantly, Hammerhead holds promise as a routine
559 pipeline for identifying and polishing bacterial DNA methylation sites for a wide range of nanopore
560 sequencing applications, such as genome assembly, metagenomic binning, decontaminating eukaryotic
561 genome assembly, and functional analysis for DNA modifications.

562

563 **Methods**

564 **Read processing.**

565 The raw Nanopore data in FAST5 format were subjected to basecalling using Guppy (V6.4.6)
566 (<https://community.nanoporetech.com>), which employed the basecalling model file
567 dna_r9.4.1_450bps_sup.cfg for R9.4.1 sequencing data and dna_r10.4.1_e8.2_400bps_sup.cfg for
568 R10.4.1 sequencing data. To mitigate the presence of informatic chimeras and concatemeric reads, the
569 split_on_adapter function from the duplex_tools (V0.3.2) ([https://github.com/nanoporetech/duplex-
571 tools](https://github.com/nanoporetech/duplex-
570 tools)) was employed with the following arguments: "--allow_multiple_splits --trim_start 50 --trim_end
572 50".

572

573 To acquire duplex reads, we utilized the duplex_tools (V0.3.2) package to initially extract potential
574 paired-read information. Subsequently, we rebasecalled the raw current signal data in FAST5 format
575 by employing the guppy (V6.4.6) software guppy_basecaller_duplex function, which incorporates
576 paired-read information.

577

578 To demultiplex the reads based on the barcoding information (SQK-NBD114-24 for R10.4.1 and EXP-
579 NBD104 for R9.4.1), the guppy_barcode function of Guppy (V6.4.6) was used with the following

580 arguments: "--enable_trim_barcodes --num_extra_bases_trim 3".

581

582 To distinguish between R10.4.1 and duplex sequencing reads, we utilized SeqKit (V2.3.0) (Shen et al.
583 2016) and applied the grep command to filter reads based on their read ID. Finally, any reads with a
584 length less than 200 bp were removed using SeqKit (V2.3.0) (Shen et al. 2016).

585

586 **Read quality analysis and homopolymer identification.**

587 The reads for each bacterial species were aligned to the high-quality reference using minimap2 (V2.22)
588 (Li 2021) with default arguments. To evaluate the accuracy of the different sequencing libraries, we
589 computed several metrics for each primary aligned read, including the substitution rate, insertion rate,
590 and deletion rate, and observed read accuracy using the following equations:

591

$$592 \quad N(\text{total}) = N(\text{sub}) + N(\text{mat}) + N(\text{ins}) + N(\text{del})$$

$$593 \quad \text{Substitution rate} = N(\text{sub})/N(\text{total})$$

$$594 \quad \text{Insetion rate} = N(\text{ins})/N(\text{total})$$

$$595 \quad \text{Deletion rate} = N(\text{del})/N(\text{total})$$

$$596 \quad \text{Observed read accuracy} = N(\text{mat})/N(\text{total})$$

597

598 Here, $N(\text{sub})$, $N(\text{mat})$, $N(\text{ins})$, and $N(\text{del})$ are the number of substitutions, matches, insertions, and
599 deletions, respectively, in each read. All the functions could be achieved by the Giraffe (V0.1.0.14)
600 (Liu et al. 2024).

601

602 To assess the accuracy of homopolymer identification across multiple flow cells, we filtered out
603 homopolymers with a length less than 3 bp and coverage lower than 3. The remaining homopolymers
604 were analyzed using custom scripts, which calculated the proportion of homopolymers that achieved
605 100% match accuracy. All the scripts used for read quality benchmarking were packaged as the
606 “observe” function of Giraffe (https://github.com/lrslab/Giraffe_View), and the plotting steps are
607 available at https://github.com/lrslab/Bacteria-Multisequencing/tree/main/code/read_quality.

608

609 **Genome assembly and read subsampling.**

610 To obtain a high-quality reference genome assembly for eight bacterial strains, we employed long-read
611 sequencing data from both the R10.4.1 and duplex datasets as inputs for the different assembly
612 programs: Flye (V2.9.2) (Kolmogorov et al. 2019), Canu (V2.2) (Koren et al. 2017), and Unicycler

613 (V0.5.0) (Wick et al. 2017) (for only *S. enterica*). The default parameters were used in all the cases.
614 Contigs were selected based on their length, circularity, and synteny across different assembly
615 programs, yielding draft assemblies for each bacterial chromosome and plasmid(s). To correct
616 substitution errors and indels in the draft assemblies, we performed polishing using both long-read and
617 short-read data. Specifically, we used a two-step polishing approach, beginning with Racon (V1.5.0)
618 (Vaser et al. 2017), to polish the drafts in three rounds using ONT R10.4.1 reads, followed by three
619 rounds of polishing using Pilon (V1.24) (Walker et al. 2014) with short reads.

620

621 To evaluate genome assembly performance using long-read data from R10.4.1 and R9.4.1, we used
622 50-fold short-read data, and each long-read dataset was subsampled at 10- to 120-fold intervals using
623 Rasusa (V0.3.0) (Hall 2022). Given the potential variability in random subsampling, we performed the
624 subsampling process five times for each long-read dataset and coverage level, using seeds ranging
625 from 1 to 5. In brief, we obtained a total of 960 assemblies by performing five repetitions for each
626 long-read dataset and coverage level (8 bacterial species * 2 datasets, R10.4.1 and R9.4.1, * 5
627 repetitions * 12 read coverages). We used Flye (V2.9.2) (Kolmogorov et al. 2019) for genome
628 assembly and polished the resulting drafts with two rounds of Medaka (V1.8.0)
629 (<https://github.com/nanoporetech/medaka>) using the corresponding long-read dataset, with or without
630 three rounds of Racon (V1.5.0), with 50-fold short-read data for each bacterium at each read coverage.
631 We compared the assemblies to the high-quality reference, which included the detection of indels and
632 substitutions, using Quast (V5.2.0) (Gurevich et al. 2013) with arguments of “--min-alignment 1000 -
633 -min-identity 99”. To assess the genome completeness of each bacterium, we used benchmarking
634 universal single-copy orthologs (BUSCO) (V5.4.3) (Manni et al. 2021) with the following databases:
635 pseudomonadales_odb10 for *A. pittii* and *P. aeruginosa*; bacillales_odb10 for *B. cereus* and *S. aureus*;
636 lactobacillales_odb10 for *E. faecium*; and enterobacterales_odb10 for *E. coli*, *S. enterica*, and *K.*
637 *pneumoniae*.

638

639 To ensure the accuracy and reliability of our results, we implemented a stringent quality control
640 approach that involved removing the highest and lowest values for each metric at each assembly
641 coverage for each bacterium. We then calculated the average value from the remaining three values to
642 represent the performance of each assembly method for each bacterium. All the codes for processing
643 temperature files are available at [https://github.com/lrslab/Bacteria-](https://github.com/lrslab/Bacteria-Multisequencing/tree/main/code/genome_assembly)
644 [Multisequencing/tree/main/code/genome_assembly](https://github.com/lrslab/Bacteria-Multisequencing/tree/main/code/genome_assembly).

645

646 **Single nucleotide substitution error statistics and current signal comparison.**

647 To determine the frequencies of the different substitution types, we compared the 35 drafts at the
 648 chromosome level assembled from the long-read data (R10.4.1 or R9.4.1) at read coverage ranging
 649 from 40- to 100-fold by Quast (V5.2.0) (Gurevich et al. 2013) for four bacteria, namely, *A. pittii*, *E.*
 650 *faecium*, *E. coli*, and *K. pneumoniae*. The codes are available at [https://github.com/lrslab/Bacteria-](https://github.com/lrslab/Bacteria-Multisequencing/tree/main/code/mismatch_fre)
 651 [Multisequencing/tree/main/code/mismatch_fre](https://github.com/lrslab/Bacteria-Multisequencing/tree/main/code/mismatch_fre). Only the sites with coverage over two were considered
 652 error-prone.

653

654 To count the number of matched base types at those error-prone sites, we utilized minimap2 (V2.22)
 655 (Li 2021) to align the four datasets, including R9.4.1, R9.4.1 whole-genome amplification (WGA),
 656 R10.4.1, and R10.4.1 WGA, against the reference for *K. pneumoniae*, *E. faecium*, *A. pittii*, and *E. coli*.
 657 The sorted BAM files were subsequently input to the mpileup function from SAMtools (V1.17)
 658 (Danecek et al. 2021) with arguments of “--no-output-ends --no-output-ins --no-output-ins --no-
 659 output-del --no-output-del” to calculate the mapped proportion of A, C, G, and T bases at these error-
 660 prone sites. The codes can be found at [https://github.com/lrslab/Bacteria-](https://github.com/lrslab/Bacteria-Multisequencing/tree/main/code/modification_test)
 661 [Multisequencing/tree/main/code/modification_test](https://github.com/lrslab/Bacteria-Multisequencing/tree/main/code/modification_test).

662

663 To directly compare the difference in the current signal between WGS and WGA reads, we used f5c
 664 (V1.2) (Gamaarachchi et al. 2020) to resquiggle the reads with the reference and visualized them via
 665 nanoCEM (Guo et al. 2024).

666

667 **Potential modification site identification and motif enrichment with Hammerhead.**

668 To determine the sites where the proportion of mapped bases differed between the forward strand and
 669 reverse strand, we used a self-defined index to reflect this difference at a single site and calculated the
 670 function via the software Hammerhead (<https://github.com/lrslab/Hammerhead>).

671

$$672 \quad Pf(a) = \frac{N(a)}{N(a) + N(t) + N(g) + N(c)}$$

$$673 \quad Pr(a) = \frac{N(t)}{N(a) + N(t) + N(g) + N(c)}$$

674

675 Here, $N(a)$, $N(t)$, $N(g)$, and $N(c)$ represent the number of mapped occurrences of the bases A, T, G, and
 676 C, respectively, in the forward strand at a site in the reference sequence. Given that the DNA strands
 677 are reverse complements, $N(a)$, $N(t)$, $N(g)$, and $N(c)$ represent the number of mapped occurrences of

678 the bases T, A, C, and G in the reverse strand at the same site in the reference. $Pf(a)$ and $Pr(a)$ represent
 679 the proportions of mapped occurrences of base A in the forward strand and reverse strand, respectively.
 680 Additionally, we calculated the proportions of three other bases (T, G, and C) via the same method.
 681 Notably, only the sites with a depth greater than 50 and a depth in forward and reverse strands greater
 682 than 25 were retained for downstream analysis. Only the reads with a mapping quality (mapQ) greater
 683 than 30 were used for counting.

$$684$$

$$685 \quad \text{Dif}(A) = \text{ABS}(Pf(a) - Pr(a))$$

$$686$$

687 Then, we can calculate the absolute difference in the A bases in the two strands by using the proportion
 688 in the forward strand minus that in the reverse strand. Moreover, we calculated the absolute differences
 689 in T, G, and C bases via the same method.

$$690$$

$$691 \quad \text{Difference index} = \frac{\text{Dif}(A) + \text{Dif}(T) + \text{Dif}(G) + \text{Dif}(C)}{2}$$

$$692$$

693 Finally, to calculate the difference index, we sum the absolute differences in A, T, G, and C bases and
 694 divide the sum by two.

695

696 Considering the presence of random errors during sequencing, a cutoff is needed to filter these
 697 background noises. Whole-genome amplification (WGA) reads, which contain no methylation
 698 information, were regarded as negative controls. The differences in expression between the four
 699 bacterial WGA datasets were calculated by Hammerhead. To estimate the false discovery rate (FDR)
 700 for the WGA reads, we inferred the distribution of difference indices for the four datasets. The FDR
 701 was calculated based on the number of sites over the cutoff divided by the total number of sites. The
 702 sites over the cutoff can be considered false positives caused by random errors rather than methylation
 703 in WGA reads. To limit the number of false positives to less than 10 in WGA reads, a cutoff of 0.35
 704 was selected to achieve an FDR of less than 1e-06. Therefore, the sites with a difference index greater
 705 than 0.35 in WGS reads were regarded as potential modification sites.

706

707 To identify regular motifs, we selected sequences near the potential modified sites (-10 bp to +9 bp)
 708 for motif enrichment. The MEME (V5.5.3) (Bailey et al. 2015) was utilized to identify the motif with
 709 the command “meme inputFile.fa -dna -oc . -nostatus -time 14400 -mod zoops -nmotifs 10 -minw 4 -
 710 maxw 16 -objfun classic -revcomp -markov_order 0”.

711

712 Methylation identification and motif enrichment using Nanodisco.

713 Nanodisco (V1.0.3) was utilized to identify the modifications in the R9.4.1 WGA and WGS datasets.
714 The functions "process," "chunk_info," "difference," "merge," "motif," and "characterize" from
715 Nanodisco were used for our eight R9.4.1 bacterial datasets, all with default parameter as the detailed
716 tutorial in the Nanodisco documentation
717 (https://nanodisco.readthedocs.io/en/latest/detailed_tutorial.html).

718

719 To select the sites with an absolute value of current signal difference greater than 2 pA, the information
720 in the RDS file was processed using the readRDS function of R. These raw RDS files are available at
721 <https://doi.org/10.6084/m9.figshare.24298774>.

722

723 Intersection comparison of methylation units between Hammerhead and Nanodisco.

724 Considering that the difference in the current signal and basecalling error may be caused by the nearby
725 modified base (Laszlo et al. 2013; Rand et al. 2017; Wick et al. 2019; Tourancheau et al. 2021), the
726 potential modification site, plus four 5' and 3' flanking sites, was considered a modification unit (9 bp).
727 We subsequently conducted an intersection analysis between the modification units identified by
728 Hammerhead and Nanodisco using BEDTools (V2.18) with the command "bedtools intersect -wa -a
729 A.bed -b B.bed | sort | uniq -c".

730

731 Methylation metrics comparison between WGS and WGA reads in motifs.

732 To assess the performance of Hammerhead and Dorado in 6mA methylation calling, the metrics
733 between WGS and WGA in two methylation motifs, C6mATCTC and R6mAYCNNNNNTTRG
734 were used to compare. For Hammerhead, the 15-mer methylation units were chosen, with the center
735 positioned on the 6mA bases and spanning 7 bp upstream and downstream. The maximum value of the
736 difference index was considered as the representative for each unit in WGS reads. Similarly,
737 representatives in WGA reads were selected based on the positions of sites with the maximum values,
738 ensuring that the prediction values were chosen from the same site. For Dorado, the prediction values
739 in both WGS and WGA data were selected based on the positions of the modified 6mA sites.

740

741 A methylation-aware basecalling model was retrained using *E. coli* WGA reads.

742 The new basecalling model was fine-tuned from the original DNA super accuracy model (SUP) using
743 Bonito (V0.7.2) with the following parameters: "--epochs 40 --lr 5e-4 --batch 32 --pretrained
744 [dna_r10.4.1_e8.2_400bps_sup@v4.2.0](#)". To prepare the input SAM file for training, the parameter

745 “save-ctc” must be enabled in the Bonito basecaller.

746

747 The fine-tuned methylation-aware model for *E. coli* can be downloaded from the Hammerhead GitHub
748 repository (<https://github.com/lrslab/Hammerhead>). The new model was subsequently used for
749 rebasecalling our *E. coli* WGS and WGA reads.

750

751 **Genome polishing at potential modification sites with duplex reads.**

752 To validate whether polishing only the potential modification sites with duplex reads can reduce the
753 substitution error caused by modification, we used 15 assemblies with 40-, 50-, and 60-fold assembly
754 coverage to compare substitution rates in *A. pittii*, *E. faecium*, *E. coli*, and *K. pneumonia*, respectively.
755 First, all the potential modification sites with a difference index greater than 0.35 were selected. The
756 BAM file generated from duplex reads mapped against the reference was subsequently used to call the
757 pileup file for these potential modification sites. Ultimately, the decision to rectify or preserve the site
758 would be made based on the proportion accurately depicted in the mapping. The pipeline can be found
759 in the documentation ([https://hammerhead-documentation.readthedocs.io/en/latest/#assemblies-](https://hammerhead-documentation.readthedocs.io/en/latest/#assemblies-polish)
760 [polish](https://hammerhead-documentation.readthedocs.io/en/latest/#assemblies-polish)).

761

762 **Data access**

763 The whole-genome amplification (WGA) data generated in this study have been submitted to the NCBI
764 BioProject database (<https://www.ncbi.nlm.nih.gov/bioproject/>) under accession number
765 PRJNA980403. The Hammerhead software, including modification sites finding and duplex read
766 polish, are available at GitHub (<https://github.com/lrslab/Hammerhead>) and as Supplemental Code.
767 All the processed files and other analysis scripts used in this study are available at GitHub
768 (<https://github.com/lrslab/Bacteria-Multisequencing>) and as Supplemental Material.

769

770 **Acknowledgments**

771 We thank Yating Xu, Kaichao Chen, Qiao Hu, and Wai Chi Chan for providing us with valuable
772 bacterial samples. This work was supported by the Early Career Scheme from the Research Grants
773 Council of the Hong Kong Special Administrative Region, China (CityU 21100521); the Hong Kong
774 Health and Medical Research Fund (project number 08194126); the Guangdong General Research
775 Fund (project number 9240054) from the Natural Science Foundation of Guangdong Province; new
776 Research Initiatives support from City University of Hong Kong (project number 9610497) to R.L.;

777 the Theme-based Research Scheme (T11-104/22-R) to S.C.; the Hetao Shenzhen-Hong Kong Science
 778 and Technology Innovation Cooperation Zone Shenzhen Park Project (HZQB-KCZYZ-2021017); and
 779 the City University of Hong Kong Project (project number 9680217 and number 9678223) to M.Y.

780

781 **Contributions**

782 **Xudong Liu:** Conceptualization, methodology, software, validation, formal analysis, investigation,
 783 data curation, writing (original draft, review, and editing). **Ying Ni:** Methodology (whole-genome
 784 amplification, nanopore library construction, and sequencing) and writing (review and editing).
 785 **Lianwei Ye:** Methodology (DNA extraction, short-read library construction, and short-read
 786 sequencing), writing (review and editing). **Zhihao Guo:** Writing (review and editing). **Lu Tan:**
 787 Writing (review and editing). **Jun Li:** Writing (review and editing). **Mengsu Yang:** Supervision,
 788 funding acquisition, writing (review and editing). **Sheng Chen:** Supervision, funding acquisition,
 789 writing (review and editing). **Runsheng Li:** Conceptualization, methodology, validation, investigation,
 790 writing (original draft, review, and editing), supervision, and funding acquisition.

791

792 **Competing interests**

793 The authors declare that there are no conflicts of interest associated with this study.

794

795 **References**

- 796 Amarasinghe SL, Su S, Dong X, Zappia L, Ritchie ME, Gouil Q. 2020. Opportunities and challenges
 797 in long-read sequencing data analysis. *Genome Biology* **21**: 30.
- 798 Bailey TL, Johnson J, Grant CE, Noble WS. 2015. The MEME Suite. *Nucleic Acids Res* **43**: W39-49.
- 799 Beaulaurier J, Schadt EE, Fang G. 2019. Deciphering bacterial epigenomes using modern sequencing
 800 technologies. *Nature Reviews Genetics* **20**: 157-172.
- 801 Bhatti H, Jawed R, Ali I, Iqbal K, Han Y, Lu Z, Liu Q. 2021. Recent advances in biological nanopores
 802 for nanopore sequencing, sensing and comparison of functional variations in MspA mutants.
 803 *RSC Advances* **11**: 28996-29014.
- 804 Breckell GL, Silander OK. 2022. Growth condition-dependent differences in methylation imply
 805 transiently differentiated DNA methylation states in *Escherichia coli*. *G3*
 806 *Genes|Genomes|Genetics* **13**.
- 807 Casadesús J, Low D. 2006. Epigenetic Gene Regulation in the Bacterial World. *Microbiology and*
 808 *Molecular Biology Reviews* **70**: 830-856.
- 809 Chen K, Chan EW-C, Xie M, Ye L, Dong N, Chen S. 2017. Widespread distribution of mcr-1-bearing
 810 bacteria in the ecosystem, 2015 to 2016. *Eurosurveillance* **22**.
- 811 Chen K, Yang C, Dong N, Xie M, Ye L, Chan EWC, Chen S. 2020. Evolution of Ciprofloxacin
 812 Resistance-Encoding Genetic Elements in *Salmonella*. *mSystems* **5**: 10.1128/msystems.01234-
 813 01220.
- 814 Cui Y, Liu X, Dietrich R, Märklbauer E, Cao J, Ding S, Zhu K. 2016. Characterization of *Bacillus*
 815 *cereus* isolates from local dairy farms in China. *FEMS microbiology letters* **363**: fnw096.

- 816 Dahlet T, Argüeso Lleida A, Al Adhami H, Dumas M, Bender A, Ngondo RP, Tanguy M, Vallet J,
817 Auclair G, Bardet AF et al. 2020. Genome-wide analysis in the mouse embryo reveals the
818 importance of DNA methylation for transcription integrity. *Nature Communications* **11**: 3153.
- 819 Danecek P, Bonfield JK, Liddle J, Marshall J, Ohan V, Pollard MO, Whitwham A, Keane T, McCarthy
820 SA, Davies RM et al. 2021. Twelve years of SAMtools and BCFtools. *Gigascience* **10**.
- 821 Deamer D, Akeson M, Branton D. 2016. Three decades of nanopore sequencing. *Nature Biotechnology*
822 **34**: 518-524.
- 823 Gamaarachchi H, Lam CW, Jayatilaka G, Samarakoon H, Simpson JT, Smith MA, Parameswaran S.
824 2020. GPU accelerated adaptive banded event alignment for rapid comparative nanopore signal
825 analysis. *BMC Bioinformatics* **21**: 343.
- 826 Gouil Q, Keniry A. 2019. Latest techniques to study DNA methylation. *Essays in Biochemistry* **63**:
827 639-648.
- 828 Greenberg MVC, Bourc'his D. 2019. The diverse roles of DNA methylation in mammalian
829 development and disease. *Nature Reviews Molecular Cell Biology* **20**: 590-607.
- 830 Grosswendt S, Kretzmer H, Smith ZD, Kumar AS, Hetzel S, Wittler L, Klages S, Timmermann B,
831 Mukherji S, Meissner A. 2020. Epigenetic regulator function through mouse gastrulation.
832 *Nature* **584**: 102-108.
- 833 Guo Z, Ni Y, Tan L, Shao Y, Ye L, Chen S, Li R. 2024. Nanopore Current Events Magnifier (nanoCEM):
834 a novel tool for visualizing current events at modification sites of nanopore sequencing. *NAR*
835 *Genomics and Bioinformatics* **6**.
- 836 Gurevich A, Saveliev V, Vyahhi N, Tesler G. 2013. QUAST: quality assessment tool for genome
837 assemblies. *Bioinformatics* **29**: 1072-1075.
- 838 Hall M. 2022. Rasusa: Randomly subsample sequencing reads to a specified coverage. *Journal of*
839 *Open Source Software* **7**: 3941.
- 840 Huang J, Chen F, Lu B, Sun Y, Li Y, Hua C, Deng X. 2024. DNA Methylome Regulates Virulence and
841 Metabolism in *Pseudomonas syringae*. doi:10.1101/2024.02.12.579912. Cold Spring Harbor
842 Laboratory.
- 843 Kolmogorov M, Yuan J, Lin Y, Pevzner PA. 2019. Assembly of long, error-prone reads using repeat
844 graphs. *Nature Biotechnology* **37**: 540-546.
- 845 Konanov DN, Babenko VV, Belova AM, Madan AG, Boldyreva DI, Glushenko OE, Butenko IO,
846 Fedorov DE, Manolov AI, Krivonos DV et al. 2023. Snapper: high-sensitive detection of
847 methylation motifs based on Oxford Nanopore reads. *Bioinformatics* **39**.
- 848 Koren S, Walenz BP, Berlin K, Miller JR, Bergman NH, Phillippy AM. 2017. Canu: scalable and
849 accurate long-read assembly via adaptive *k*-mer weighting and repeat separation. *Genome*
850 *Research* **27**: 722-736.
- 851 Laszlo AH, Derrington IM, Brinkerhoff H, Langford KW, Nova IC, Samson JM, Bartlett JJ, Pavlenok
852 M, Gundlach JH. 2013. Detection and mapping of 5-methylcytosine and 5-
853 hydroxymethylcytosine with nanopore MspA. *Proceedings of the National Academy of*
854 *Sciences* **110**: 18904-18909.
- 855 Li H. 2021. New strategies to improve minimap2 alignment accuracy. *Bioinformatics* **37**: 4572-4574.
- 856 Liu C, Chen K, Wu Y, Huang L, Fang Y, Lu J, Zeng Y, Xie M, Chan EWC, Chen S et al. 2022.
857 Epidemiological and genetic characteristics of clinical carbapenem-resistant *Acinetobacter*
858 *baumannii* strains collected countrywide from hospital intensive care units (ICUs) in China.
859 *Emerging Microbes & Infections* **11**: 1730-1741.
- 860 Liu H, Begik O, Lucas MC, Ramirez JM, Mason CE, Wiener D, Schwartz S, Mattick JS, Smith MA,
861 Novoa EM. 2019. Accurate detection of m6A RNA modifications in native RNA sequences.
862 *Nature Communications* **10**.
- 863 Liu X, Ni Y, Wang D, Ye S, Yang M, Sun X, Leung AYH, Li R. 2023. Unraveling the whole genome
864 DNA methylation profile of zebrafish kidney marrow by Oxford Nanopore sequencing.
865 *Scientific Data* **10**: 532.

- 866 Liu X, Shao Y, Guo Z, Ni Y, Sun X, Hung Leung AY, Li R. 2024. Giraffe: a tool for comprehensive
867 processing and visualization of multiple long-read sequencing data. *bioRxiv*
868 doi:10.1101/2024.05.10.593289: 2024.2005.2010.593289.
- 869 Loenen WAM, Dryden DTF, Raleigh EA, Wilson GG, Murray NE. 2013. Highlights of the DNA
870 cutters: a short history of the restriction enzymes. *Nucleic Acids Research* **42**: 3-19.
- 871 Lohde M, Wagner GE, Dabernig-Heinz J, Viehweger A, Braun SD, Monecke S, Diezel C, Stein C,
872 Marquet M, Ehricht R et al. 2024. Accurate bacterial outbreak tracing with Oxford Nanopore
873 sequencing and reduction of methylation-induced errors. *bioRxiv*
874 doi:10.1101/2023.09.15.556300: 2023.2009.2015.556300.
- 875 Manni M, Berkeley MR, Seppey M, Simão FA, Zdobnov EM. 2021. BUSCO Update: Novel and
876 Streamlined Workflows along with Broader and Deeper Phylogenetic Coverage for Scoring of
877 Eukaryotic, Prokaryotic, and Viral Genomes. *Molecular Biology and Evolution* **38**: 4647-4654.
- 878 Marinus MG, Løbner-Olesen A. 2014. DNA Methylation. *EcoSal Plus* **6**: 10.1128/ecosalplus.ESP-
879 0003-2013.
- 880 Mayer SF, Cao C, Dal Peraro M. 2022. Biological nanopores for single-molecule sensing. *iScience* **25**.
- 881 Ni P, Nie F, Zhong Z, Xu J, Huang N, Zhang J, Zhao H, Zou Y, Huang Y, Li J et al. 2023a. DNA 5-
882 methylcytosine detection and methylation phasing using PacBio circular consensus sequencing.
883 *Nature Communications* **14**: 4054.
- 884 Ni Y, Liu X, Simeneh ZM, Yang M, Li R. 2023b. Benchmarking of Nanopore R10.4 and R9.4.1 flow
885 cells in single-cell whole-genome amplification and whole-genome shotgun sequencing.
886 *Computational and Structural Biotechnology Journal* **21**: 2352-2364.
- 887 Oxford Nanopore PLC. 2023. Dorado. GitHub.
- 888 Peraro MD, van der Goot FG. 2016. Pore-forming toxins: ancient, but never really out of fashion.
889 *Nature Reviews Microbiology* **14**: 77-92.
- 890 Petryk N, Bultmann S, Bartke T, Defossez P-A. 2020. Staying true to yourself: mechanisms of DNA
891 methylation maintenance in mammals. *Nucleic Acids Research* **49**: 3020-3032.
- 892 Rand AC, Jain M, Eizenga JM, Musselman-Brown A, Olsen HE, Akeson M, Paten B. 2017. Mapping
893 DNA methylation with high-throughput nanopore sequencing. *Nature Methods* **14**: 411-413.
- 894 Roberts RJ, Vincze T, Posfai J, Macelis D. 2022. REBASE: a database for DNA restriction and
895 modification: enzymes, genes and genomes. *Nucleic Acids Research* **51**: D629-D630.
- 896 Seymour C. 2019. Bonito: A PyTorch Basecaller for Oxford Nanopore Reads. Oxford Nanopore
897 Technologies Ltd.
- 898 Shen W, Le S, Li Y, Hu F. 2016. SeqKit: A Cross-Platform and Ultrafast Toolkit for FASTA/Q File
899 Manipulation. *PLOS ONE* **11**: e0163962.
- 900 Silvestre-Ryan J, Holmes I. 2021. Pair consensus decoding improves accuracy of neural network
901 basecallers for nanopore sequencing. *Genome Biology* **22**: 38.
- 902 Stoiber M, Quick J, Egan R, Lee JE, Celniker S, Neely RK, Loman N, Pennacchio LA, Brown J. 2017.
903 De novo Identification of DNA Modifications Enabled by Genome-Guided Nanopore Signal
904 Processing. *bioRxiv* doi:10.1101/094672: 094672.
- 905 Tisza MJ, Smith DDN, Clark AE, Youn J-H, Barnabas BB, Black S, Bouffard GG, Brooks SY,
906 Crawford J, Marfani H et al. 2023. Roving methyltransferases generate a mosaic epigenetic
907 landscape and influence evolution in *Bacteroides fragilis* group. *Nature Communications* **14**:
908 4082.
- 909 Tourancheau A, Mead EA, Zhang XS, Fang G. 2021. Discovering multiple types of DNA methylation
910 from bacteria and microbiome using nanopore sequencing. *Nat Methods* **18**: 491-498.
- 911 Vaser R, Sović I, Nagarajan N, Šikić M. 2017. Fast and accurate de novo genome assembly from long
912 uncorrected reads. *Genome Research* **27**: 737-746.
- 913 Walker BJ, Abeel T, Shea T, Priest M, Abouelliel A, Sakthikumar S, Cuomo CA, Zeng Q, Wortman J,
914 Young SK et al. 2014. Pilon: An Integrated Tool for Comprehensive Microbial Variant
915 Detection and Genome Assembly Improvement. *PLoS ONE* **9**: e112963.

- 916 Wang W, Baloch Z, Jiang T, Zhang C, Peng Z, Li F, Fanning S, Ma A, Xu J. 2017. Enterotoxigenicity
917 and Antimicrobial Resistance of *Staphylococcus aureus* Isolated from Retail Food in China.
918 *Front Microbiol* **8**: 2256.
- 919 Wick RR, Judd LM, Gorrie CL, Holt KE. 2017. Unicycler: Resolving bacterial genome assemblies
920 from short and long sequencing reads. *PLOS Computational Biology* **13**: e1005595.
- 921 Wick RR, Judd LM, Holt KE. 2019. Performance of neural network basecalling tools for Oxford
922 Nanopore sequencing. *Genome Biology* **20**: 129.
- 923 Wion D, Casadesús J. 2006. N6-methyl-adenine: an epigenetic signal for DNA–protein interactions.
924 *Nature Reviews Microbiology* **4**: 183-192.
- 925 Xu Z, Mai Y, Liu D, He W, Lin X, Xu C, Zhang L, Meng X, Mafofo J, Zaher WA et al. 2021. Fast-
926 bonito: A faster deep learning based basecaller for nanopore sequencing. *Artificial Intelligence*
927 *in the Life Sciences* **1**: 100011.
- 928 Yang X, Wai-Chi Chan E, Zhang R, Chen S. 2019. A conjugative plasmid that augments virulence in
929 *Klebsiella pneumoniae*. *Nature Microbiology* **4**: 2039-2043.
- 930 Ye L, Liu X, Ni Y, Xu Y, Zheng Z, Chen K, Hu Q, Tan L, Guo Z, Wai CK et al. 2024. Comprehensive
931 genomic and plasmid characterization of multidrug-resistant bacterial strains by R10.4.1
932 nanopore sequencing. *Microbiological Research* **283**: 127666.
- 933 Zhao Y, Chen D, Chen K, Xie M, Guo J, Chan EWC, Xie L, Wang J, Chen E, Chen S et al. 2023.
934 Epidemiological and Genetic Characteristics of Clinical Carbapenem-Resistant *Pseudomonas*
935 *aeruginosa* Strains in Guangdong Province, China. *Microbiology Spectrum* **11**: e04261-04222.
- 936 Zheng B, Tomita H, Xiao YH, Wang S, Li Y, Ike Y. 2007. Molecular Characterization of Vancomycin-
937 Resistant *Enterococcus faecium* Isolates from Mainland China. *Journal of Clinical*
938 *Microbiology* **45**: 2813-2818.
939

940 **Figure 1. Bacterial DNA modifications influence substitutions in R10.4.1 read-based assemblies.**

941 (A) Substitution per 100 kbps of assemblies generated using different coverage of R10.4.1 and R9.4.1
 942 reads, with or without high-quality short read polishing. The X-axis shows the subsampled read
 943 coverage for ONT reads. (B) Normalized per-assembly counts of all 12 single nucleotide substitution
 944 (SNS) types in *A. pittii*, *E. faecium*, *E. coli*, and *K. pneumoniae* assemblies generated from R10.4.1
 945 (n=35) or R9.4.1 (n=35) reads. The R10.4.1 and R9.4.1 assemblies are based on subsampled read
 946 coverage from 40 to 100-fold. * p-value < 0.01, Student's t-test. (C) Identification of error-prone sites
 947 in four bacteria defined by C2T or G2A substitution frequencies exceeding two. (D) Proportions and
 948 counts of accurately and inaccurately mapped bases at error-prone C and G sites in four bacterial
 949 genomes, respectively. Mapping reads were obtained from R9.4.1, R9.4.1 whole-genome amplification
 950 (WGA), R10.4.1, and R10.4.1 WGA. Notably, the C2T and G2A substitutions are more prevalent in
 951 R10.4.1 reads. SRS: short-read sequencing; *Api*: *Acinetobacter pittii*; *Eco*: *Escherichia coli*; *Efa*:
 952 *Enterococcus faecium*; *Kpn*: *Klebsiella pneumoniae*.

953

954 **Figure 2. The error-prone sites arising from bacterial DNA modifications show strand bias in**

955 **R10.4.1 reads. (A) and (B).** Proportions and counts of accurately and inaccurately mapped forward
 956 and reverse read at error-prone C and G sites in *E. faecium* and *K. pneumoniae*, respectively. *** means
 957 the p-values of the Fisher test less than 1e-30. (C) and (D). Illustration of raw current signals at an
 958 error-prone G site in *E. faecium* and *K. pneumoniae* genome, respectively. Notably, significant
 959 differences highlighted between WGS and WGA reads are observed in the forward strand (upper panel),
 960 but not in the reverse strand (lower panel). wgs: whole genome shotgun; wga: whole genome
 961 amplification; *Efa*: *Enterococcus faecium*; *Kpn*: *Klebsiella pneumoniae*. *, **, and *** p-value < 0.05,
 962 0.01, and 0.001, Student's t-test.

963

964 **Figure 3. Bacterial DNA modifications can be identified by comparing the mapping accuracy**

965 **between forward and reverse strands in R10.4.1 reads. (A) and (B).** The workflow and demo of
 966 Hammerhead, which is designed to identify the potential modification sites based on the degree of
 967 nucleotide difference between forward and reverse reads. (C) and (D). The distribution of site
 968 difference index in WGA and WGS sequencing reads for *E. faecium* and *K. pneumoniae*, respectively.
 969 The WGA sequencing, representing random read errors, serves as a background filter. High
 970 discrepancies between forward and reverse strands in WGS reads suggest potential DNA modifications.
 971 A cutoff of 0.35 (FDR < 1e-06 in WGA reads) is used here to identify possible DNA modification sites
 972 in WGS reads. (E) and (F). The motif was enriched from possible DNA modification sites identified
 973 by strand accuracy comparison for chromosome and plasmid sequences in *E. faecium* and *K.*

974 *pneumoniae*, respectively. “GATC” is the dominating motif. Note: Only one motif was enriched from
975 *E. faecium* plasmid (MEME E-value = 1.1e-07). FDR: false discovery rate.

976

977 **Figure 4. Validation of the potential modification sites from Hammerhead results.** (A) The number
978 of shared and unique motifs identified by Hammerhead and Nanodisco. (B)-(E). The comparison of
979 six shared motifs from two pipelines. (F) The proportion of overlap between two groups of
980 modification units identified by Hammerhead and Nanodisco. Each modification unit consists of a
981 potential modification site and the surrounding bases (-4bp to +4bp). Note: Hammerhead was designed
982 for R10.4.1 reads, while Nanodisco was for R9.4.1 reads. *Api*: *Acinetobacter pittii*; *Bce*: *Bacillus*
983 *cereus*; *Eco*: *Escherichia coli*; *Efa*: *Enterococcus faecium*; *Kpn*: *Klebsiella pneumoniae*; *Pae*:
984 *Pseudomonas aeruginosa*; *Sau*: *Staphylococcus aureus*; *Sen*: *Salmonella enterica*.

985

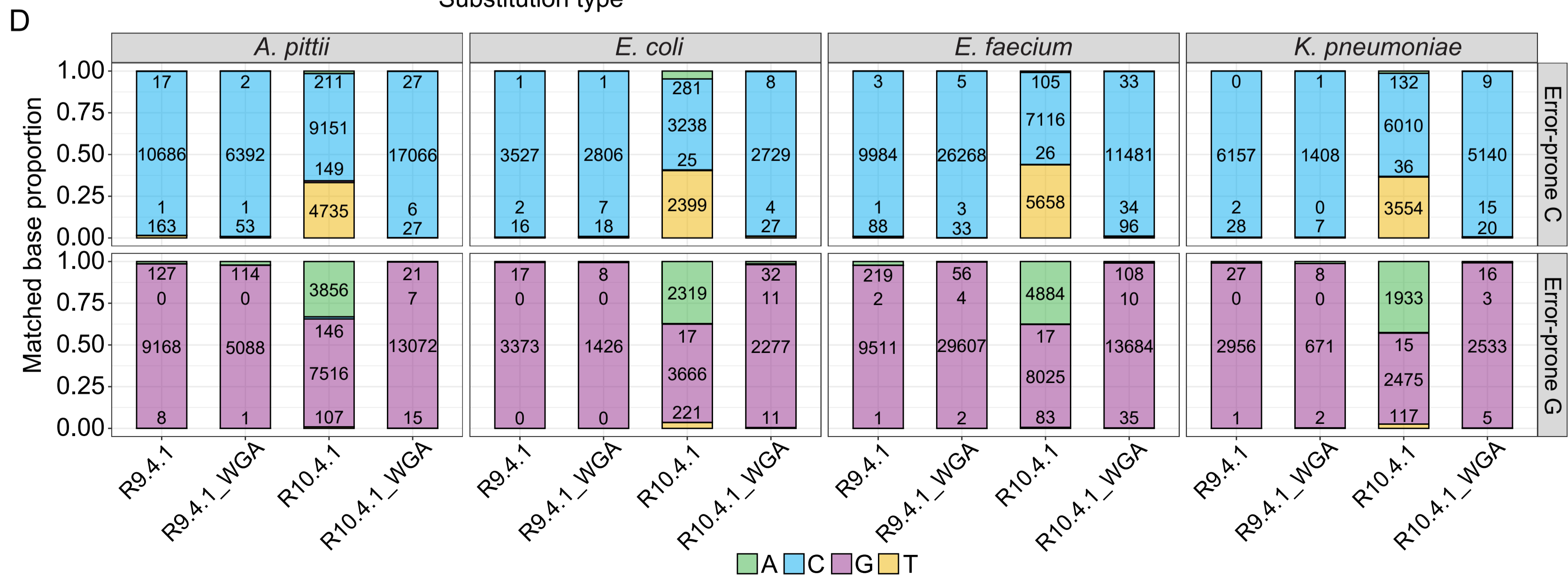
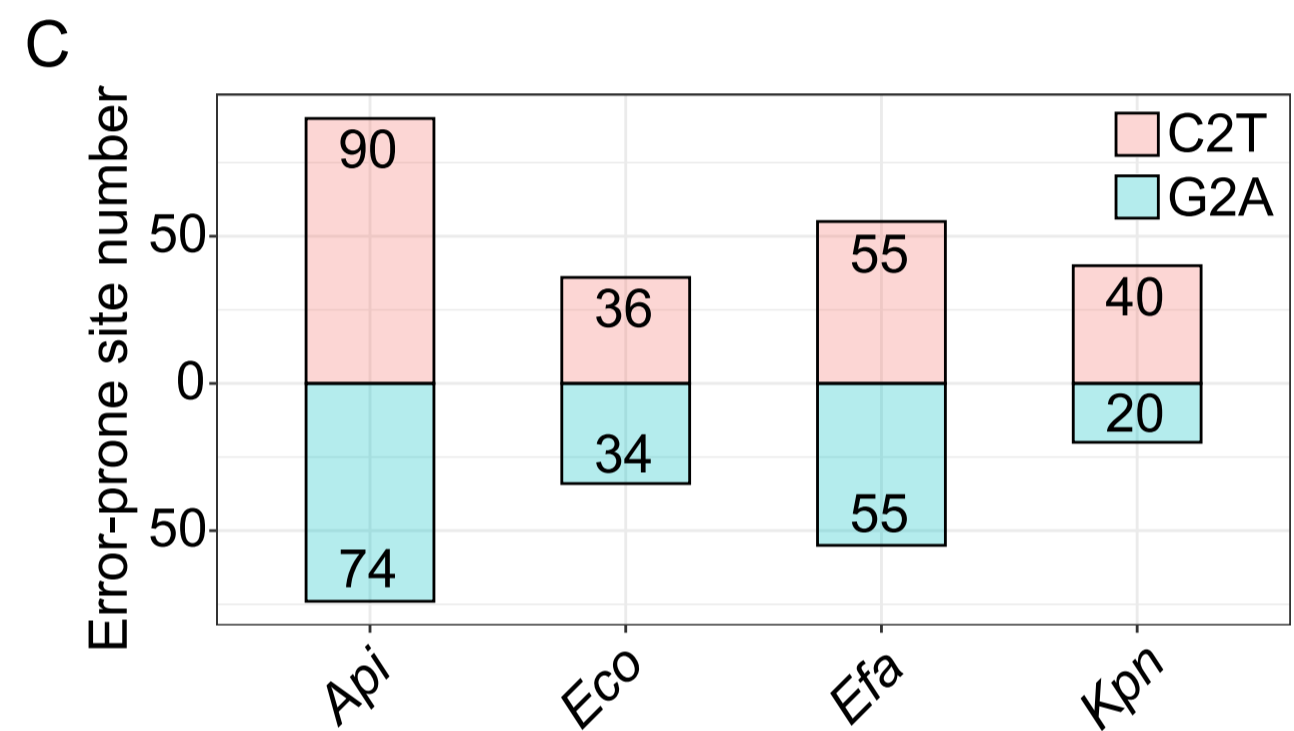
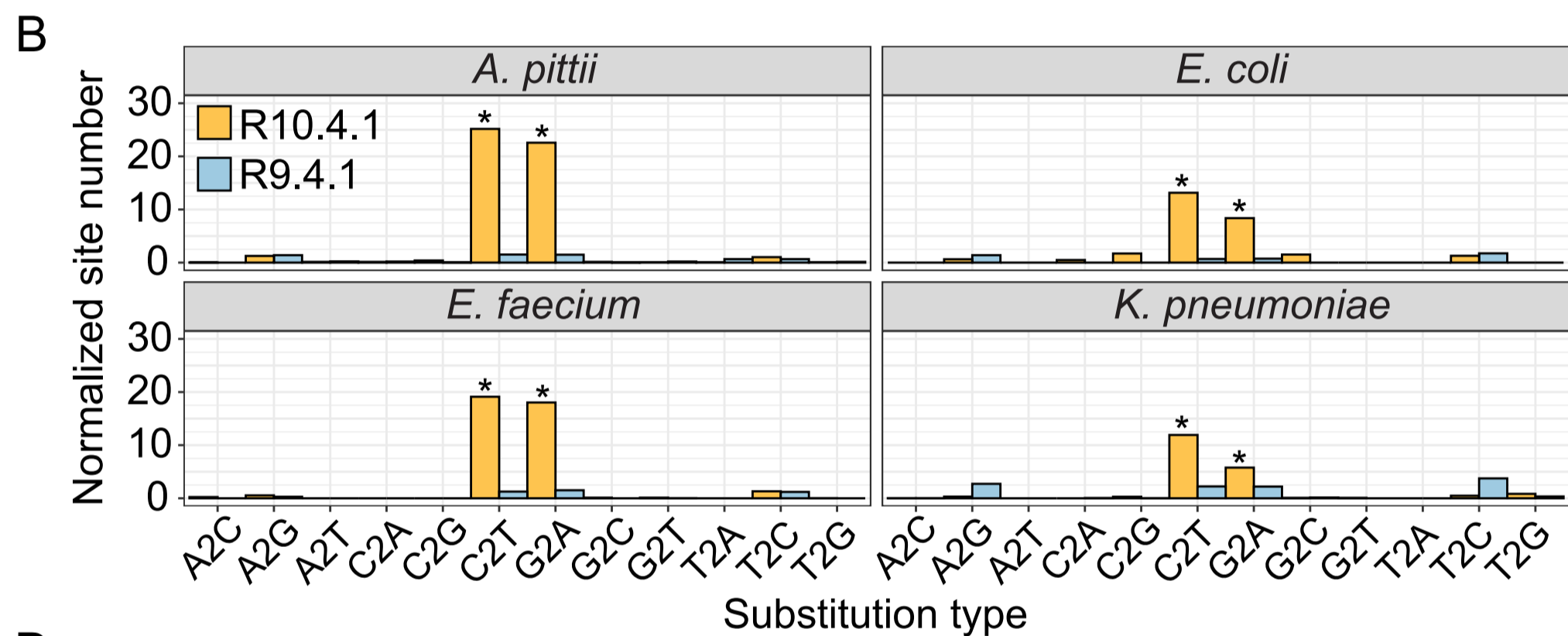
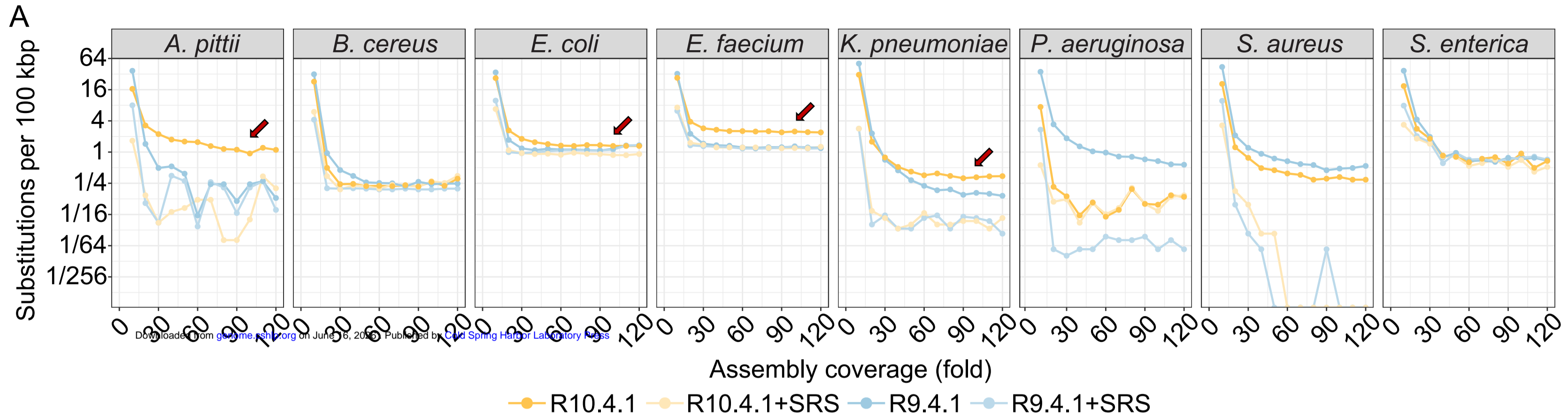
986 **Figure 5. Polishing with duplex reads resolves DNA modification-induced errors in R10.4.1**
987 **assembly.** (A) Read accuracy comparison for R10.4.1 and duplex reads at error-prone C and G sites.
988 (B) and (C). Rates of single nucleotide substitution (SNS) in R10.4.1 assemblies following short-read
989 or duplex polishing in *E. faecium* and *K. pneumoniae*, respectively. The assembly only with duplex
990 reads is also shown as a reference. (D) The workflow exclusively focuses on polishing the potential
991 modification sites using duplex reads. (E) and (F). Comparison of SNS rates in R10.4.1 assemblies
992 following targeted duplex and short-read polishing, focusing on 601 and 369 possible modification
993 sites identified by Hammerhead in *E. faecium* and *K. pneumoniae*, respectively. *** p-value < 1e-08,
994 Student’s t-test. SRS: short-read sequencing.

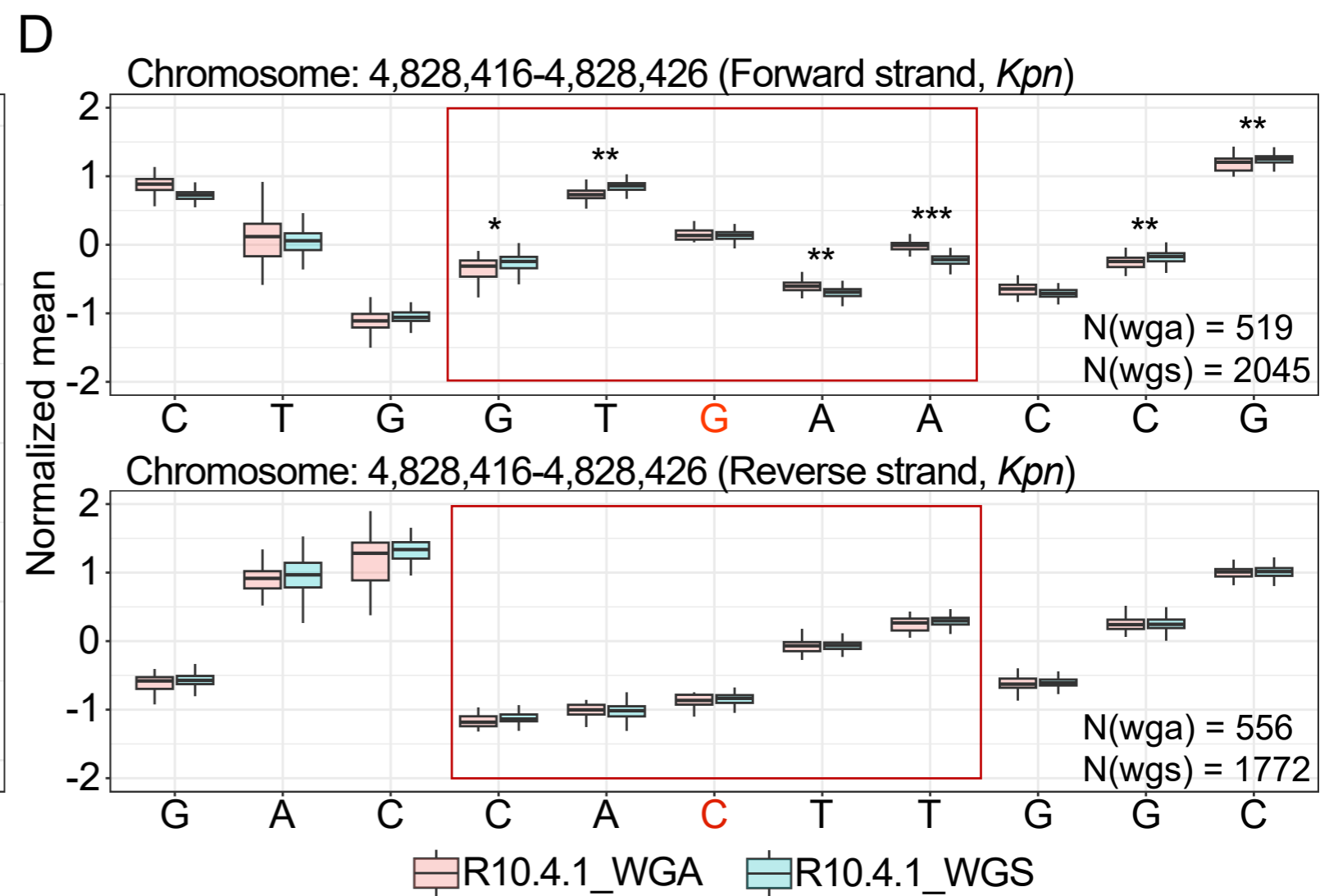
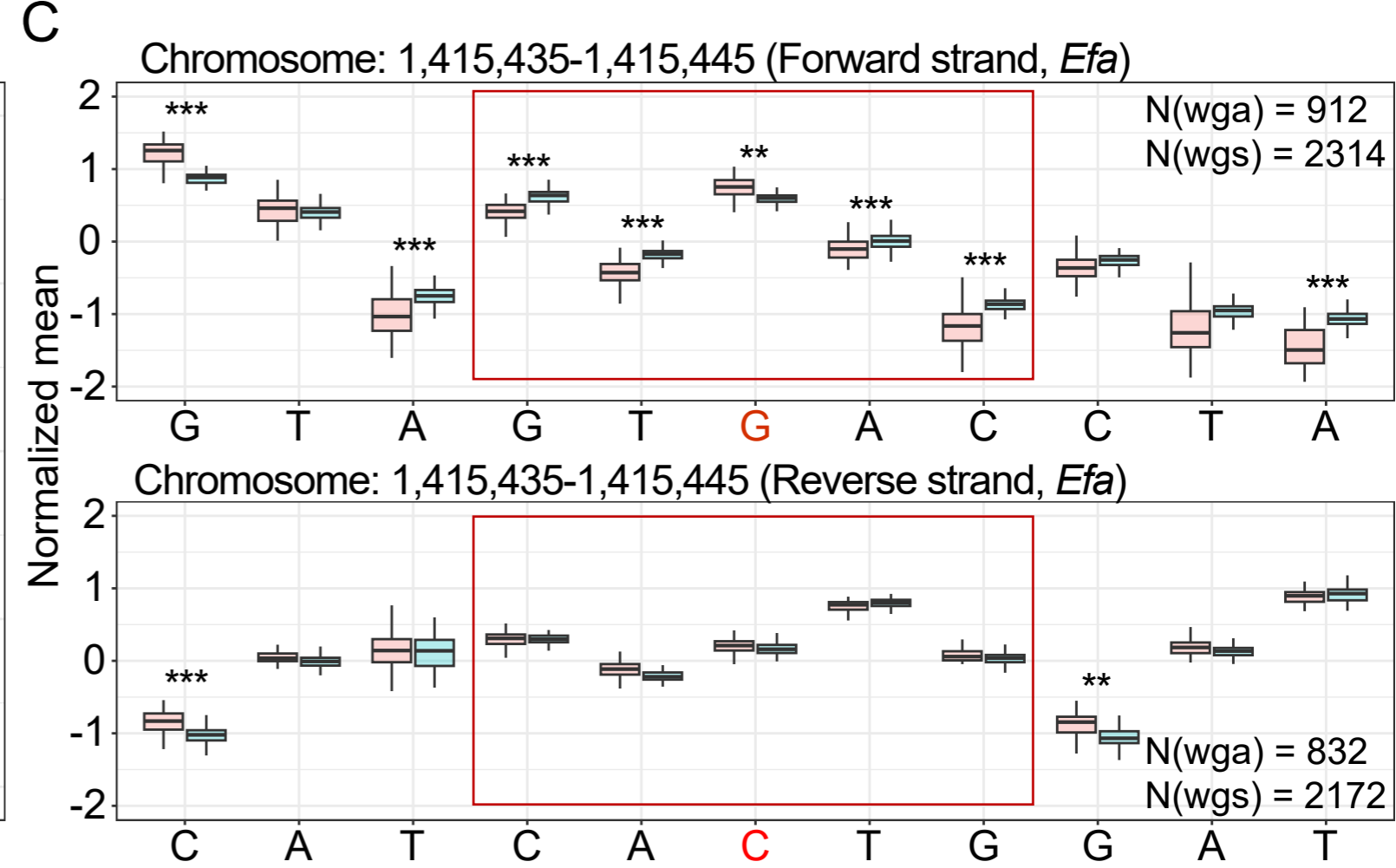
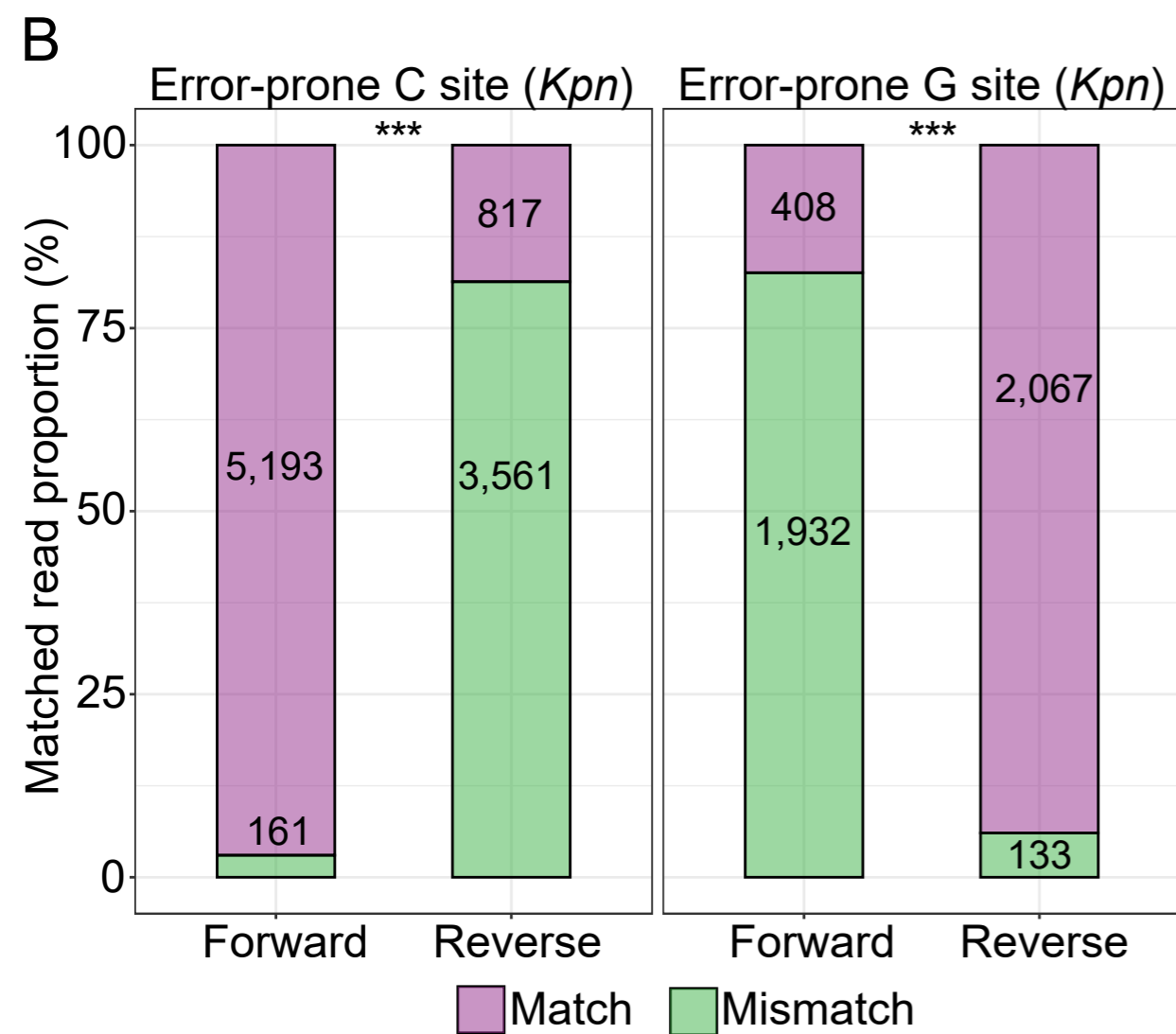
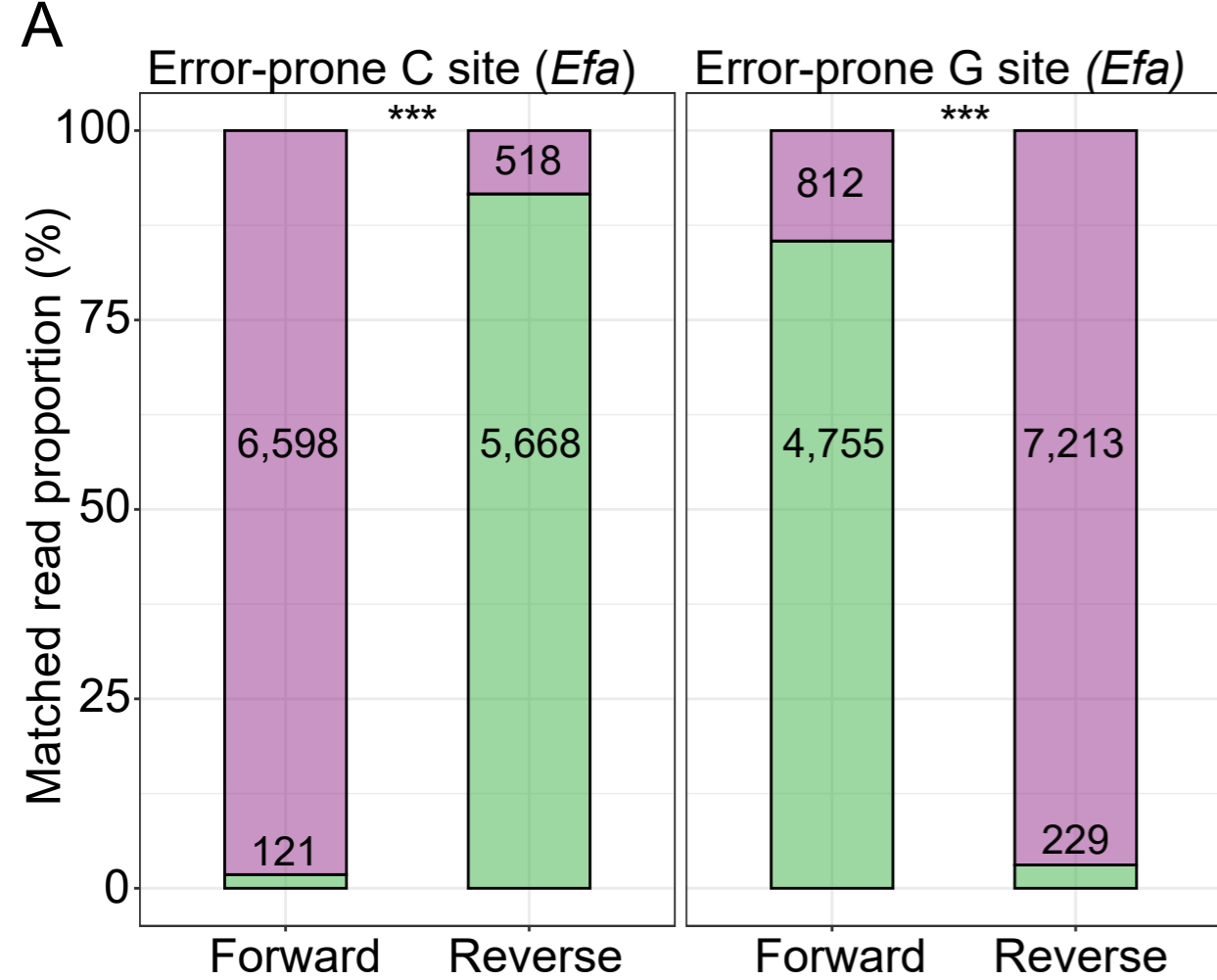
995

996 **Table1.** The summary of validated motifs.

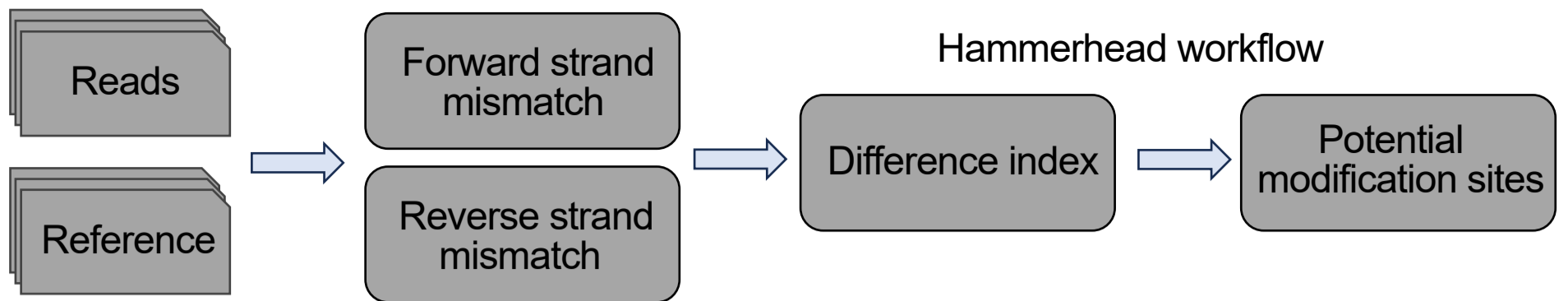
Species	Motif	Hammerhead	Nanodisco*	Methyltransferases
<i>A. pittii</i>	CATCTC	√	√	
<i>A. pittii</i>	CGAAG	√	√	
<i>B. cereus</i>	CTTCTG	√		
<i>E. coli</i>	GATC	√		√
<i>E. coli</i>	CCWGG	√	√	√
<i>E. coli</i>	GCYNNNNNCT	√		
<i>E. faecium</i>	RAYCNNNNNNNTTRG	√	√	
<i>E. faecium</i>	CYAANNNNNNNGRTY	√	√	√
<i>K. pneumoniae</i>	GATC	√	√	√
<i>K. pneumoniae</i>	GTGANNNNNNTGG	√		√
<i>K. pneumoniae</i>	CCWGG		√	
<i>P. aeruginosa</i>	TAGACGC	√	√	
<i>P. aeruginosa</i>	SATSNNNSNNSNNS		√	
<i>S. aureus</i>	GWAGNNNNNNNTAAA	√		√
<i>S. aureus</i>	GGANNNNNNNTGG	√		√
<i>S. aureus</i>	GATC	√	√	√
<i>S. enterica</i>	GATC	√		√
<i>S. enterica</i>	CAGAG	√		√

997 * **Nanodisco was run with default parameters.**

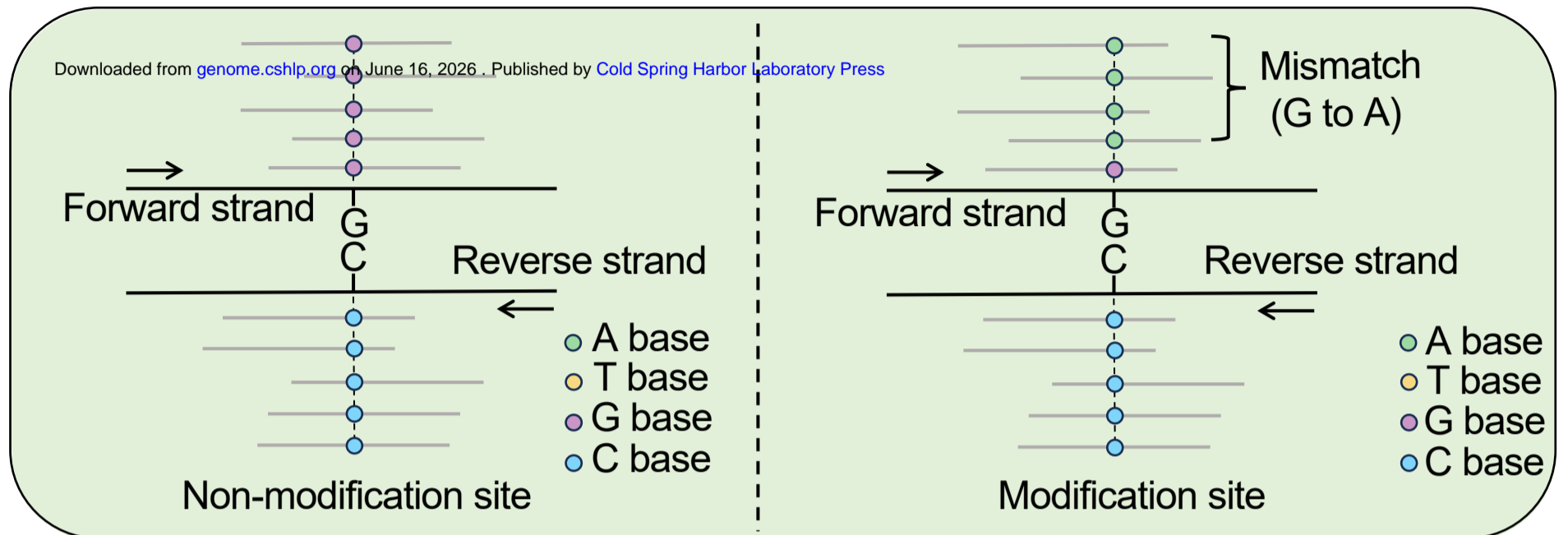




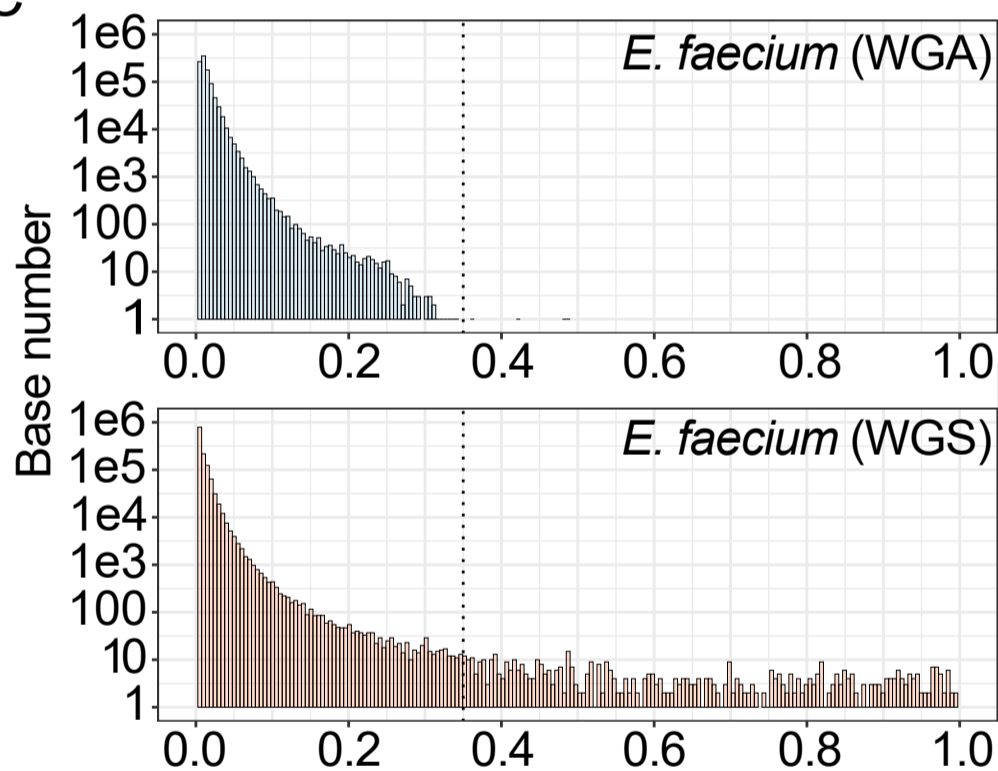
A



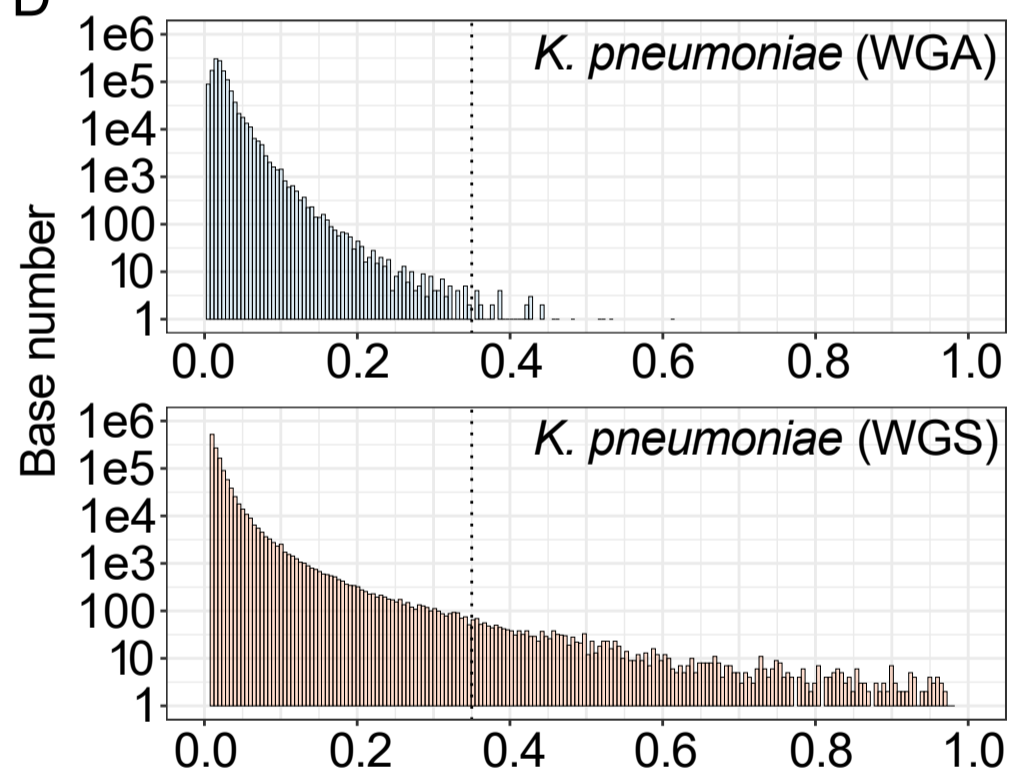
B



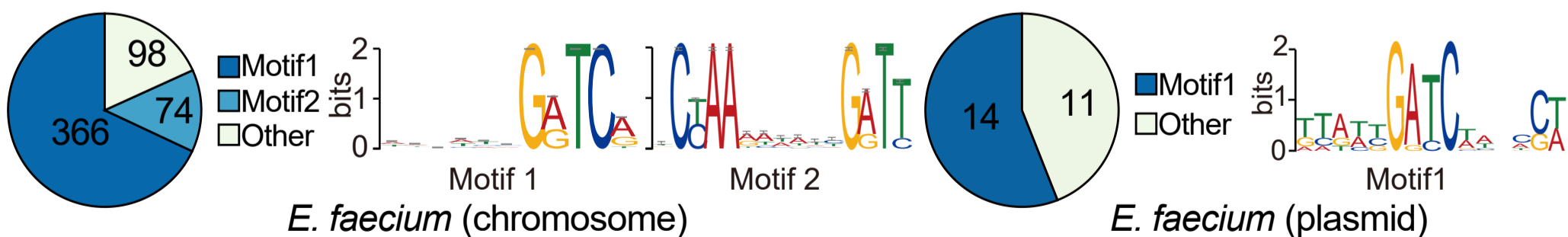
C



D



E



F

