1     # Strain and lineage-level methylome heterogeneity in the multi-drug

2     # resistant pathogenic *Escherichia coli* ST101 clone.

3

4     Melinda M. Ashcroft[1,2,3], Brian M. Forde[1,2,3], Minh-Duy Phan[1,2], Kate M. Peters[1,2], Leah W.

5     Roberts[1,2,3], Kok-Gan Chan[4,5], Teik Min Chong[4], Wai-Fong Yin[4], David L. Paterson[2,6], Timothy R.

6     Walsh[7], Mark A. Schembri[1,2*], Scott A. Beatson[1,2,3*]

7

8     [1]School of Chemistry and Molecular Biosciences, The University of Queensland, St Lucia, QLD,

9     Australia.

10     [2]Australian Infectious Diseases Research Centre, The University of Queensland, St Lucia, QLD,

11     Australia.

12     [3]Australian Centre for Ecogenomics, The University of Queensland, Brisbane, QLD, Australia.

13     [4]Division of Genetics and Molecular Biology, Institute of Biological Sciences, Faculty of Science,

14     University of Malaya, Kuala Lumpur, Malaysia.

15     [5]International Genome Centre, Jiangsu University, Zhenjiang, China.

16     [6]UQ Centre for Clinical Research, The University of Queensland, Herston, QLD, Australia.

17     [7]Department of Medical Microbiology and Infectious Disease, Cardiff University, Cardiff, United

18     Kingdom.

19     *Contributed equally

20

21     **\*Corresponding authors:**

22     Scott A. Beatson, School of Chemistry and Molecular Biosciences, The University of Queensland, St

23     Lucia 4072, QLD, Australia; Telephone: +61-7-33654863; Email: s.beatson@uq.edu.au

24     Mark A. Schembri, School of Chemistry and Molecular Biosciences, The University of Queensland,

25     St Lucia 4072, QLD, Australia; Telephone: +61 7 336 53306; Email: m.schembri@uq.edu.au

26

27    **Short title:** *E. coli* ST101 methylome heterogeneity

28    **Keywords:**

29    DNA methylation; Restriction-Modification Systems; Pacific Biosciences; Mobile Genetic Elements;

30    epigenome

31

32

33

34

35

36

37

38

39

40

41

42

43

44

45

46

47

48

49

50

51 **Abstract**

52 *Escherichia coli* Sequence Type (ST)101 is an emerging, multi-drug resistant lineage associated

53 with carbapenem resistance. We recently completed a comprehensive genomics study on mobile

54 genetic elements (MGEs) and their role in $bla_{NDM-1}$ dissemination within the ST101 lineage. DNA

55 methyltransferases (MTases) are also frequently associated with MGEs, with DNA methylation

56 guiding numerous biological processes including genomic defence against foreign DNA and

57 regulation of gene expression. The availability of Pacific Biosciences Single Molecule Real Time

58 Sequencing data for seven ST101 strains enabled us to investigate the role of DNA methylation on

59 a genome-wide scale (methylome). We defined the methylome of two complete (MS6192 and

60 MS6193) and five draft (MS6194, MS6201, MS6203, MS6204, MS6207) ST101 genomes. Our

61 analysis identified 14 putative MTases and eight N6-methyladenine DNA recognition sites, with

62 one site that has not been described previously. Furthermore, we identified a Type I MTase

63 encoded within a Transposon 7-like Transposon and show its acquisition leads to differences in the

64 methylome between two almost identical isolates. Genomic comparisons with 13 previously

65 published ST101 draft genomes identified variations in MTase distribution, consistent with MGE

66 differences between genomes, highlighting the diversity of active MTases within strains of a single

67 *E. coli* lineage. It is well established that MGEs can contribute to the evolution of *E. coli* due to

68 their virulence and resistance gene repertoires. This study emphasises the potential for mobile

69 genetic elements to also enable highly similar bacterial strains to rapidly acquire genome-wide

70 functional differences via changes to the methylome.

71

72

73

74

75 **Impact Statement**

76 *Escherichia coli* ST101 is an emerging human pathogen frequently associated with carbapenem

77 resistance. *E. coli* ST101 strains carry numerous mobile genetic elements that encode virulence

78 determinants, antimicrobial resistance, and DNA methyltransferases (MTases). In this study we

79 provide the first comprehensive analysis of the genome-wide complement of DNA methylation

80 (methylome) in seven *E. coli* ST101 genomes. We identified a Transposon carrying a Type I

81 restriction modification system that may lead to functional differences between two almost

82 identical genomes and showed how small recombination events at a single genomic region can

83 lead to global methylome changes across the lineage. We also showed that the distribution of

84 MTases throughout the ST101 lineage was consistent with the presence or absence of mobile

85 genetic elements on which they are encoded. This study shows the diversity of MTases within a

86 single bacterial lineage and shows how strain and lineage-specific methylomes may drive host

87 adaptation.

88

89 **Data Summary**

90 Sequence data including reads, assemblies and motif summaries have previously been submitted

91 to the National Center for Biotechnology Information (https://www.ncbi.nlm.nih.gov)  under the

92 BioProject Accessions: PRJNA580334, PRJNA580336, PRJNA580337, PRJNA580338, PRJNA580339,

93 PRJNA580341 and PRJNA580340 for MS6192, MS6193, MS6194, MS6201, MS6203, MS6204 and

94 MS6207 respectively. All supporting data, code, accessions, and protocols have been provided

95 within the article or through supplementary data files.

96

97

98

99

4

100

## Introduction

*Escherichia coli* sequence type (ST)101 is a pathogenic clone that has recently been associated with urinary tract and bloodstream infections in humans [1-4]. ST101 represents one of the major, emerging *E. coli* clones associated with the carriage of the $bla_{NDM-1}$ gene, causing carbapenem resistance [1, 5-9]. Recently, we undertook the most comprehensive genomics study on mobile genetic elements (MGEs) and their role in $bla_{NDM-1}$ dissemination within the ST101 lineage to date [10]. We sequenced the genomes of seven $bla_{NDM-1}$-positive ST101 isolates using Pacific Biosciences (PacBio) Single Molecule Real Time (SMRT) sequencing, generating two complete (MS6192 and MS6193) and five high-quality draft genomes (MS6194, MS6201, MS6203, MS6204, MS6207) [10]. Using an additional thirteen previously published and publicly available draft ST101 genomes, we showed that ST101 strains formed two distinct clades (Clades 1 and 2) with clustering based on infection site, *fimH* profile and antimicrobial resistance gene repertoire. Notably multidrug resistance and the carriage of the $bla_{NDM-1}$ gene were restricted to a subset of Clade 1 isolates. ST101 strains have a variable mobile genetic element (MGE) complement including prophages, genomic islands, transposons, and plasmids that encode genes for virulence, fitness, and antimicrobial resistance. Many MGEs also contained DNA methyltransferase (MTase) genes, which may result in differential methylation patterns.


In bacteria, DNA methylation is catalysed by MTases, where it guides many biological processes including defence mechanisms against foreign DNA, DNA replication and repair, timing of transposition and regulation of gene expression [11]. Three methylated nucleotides are known to occur in bacteria: N6-methyladenine, ($^{6m}$A), N4-methylcytosine ($^{4m}$C) and C5-methylcytosine ($^{5m}$C) [12]. MTases are often encoded alongside, or as part of, restriction endonucleases (REases), which have the same DNA recognition site, forming restriction-modification (RM) systems that play a

5

125    central role in defence against foreign, invading DNA [13]. Additionally, MTases can act

126    independently of REases and such DNA-modifying enzymes are known as orphan MTases. MTases

127    and RM systems are ubiquitous and extremely diverse in prokaryotes, and are classified into four

128    major groups: Type I, II, III and IV based on subunit composition, DNA recognition site specificity,

129    site of cleavage and reaction substrates (for a comprehensive review: [14]). In *E. coli*, RM systems

130    and orphan MTases are most commonly Type I or Type II [14]. Type I systems are comprised of

131    three subunits: restriction (R), modification (M) and specificity (S) [15]. The S subunit contains the

132    DNA target recognition domain (TRD) and recognises bipartite asymmetric recognition sequences

133    separated by 4-9 degenerate bases [15]. Type II systems are the most widespread, and in their

134    simplest form comprise separate R and M genes with identical DNA binding specificity, and often

135    recognise 4-6 bp palindromic sequences [16]. Exceptions include Type IIG, where the R and M

136    subunits are contained in one polypeptide and in general, bind to short, non-palindromic

137    sequences, resulting in hemi-methylation [17, 18]. Knowledge of MTase binding specificities is

138    critical for pairing motifs with their cognate MTase.

139

140    Genes encoding DNA MTases have been identified in most prokaryote genomes available to date

141    [13, 19]. However, despite the rapid growth of genomic information in public databases,

142    epigenomic information such as methylation has lagged due to methodological limitations of

143    previous technologies [20]. PacBio SMRT sequencing produces long reads, enabling the resolution

144    of complex genetic structures such as MGEs and *de novo* assembly of complete bacterial

145    chromosomes and plasmids [21]. Additionally, SMRT sequencing can be used to identify DNA

146    modifications such as methylation at a single base resolution, based on the kinetics of the

147    sequencing reaction [20]. PacBio SMRT sequencing can directly detect $^{6m}A$ and $^{4m}C$ modifications

148    due to their robust kinetic signatures, however it is only moderately sensitive for $^{5m}C$

149    modifications [22]. The impact of SMRT sequencing on cataloguing genome-wide methylation in

150    bacteria has been demonstrated recently, with the complete methylome of hundreds of bacterial

151    pathogens and environmental species now characterised (for example [23]). MTases and RM

152    systems that have been characterised in bacteria are often encoded on MGEs [13, 19] and have

153    additional biological roles including the generation of genomic diversity required for host fitness

154    [13, 24]. However, there are relatively few studies on the genomic context and functional and

155    evolutionary consequences of most identified MTases.

156

157    Except for Ashcroft *et al.* [10], there have been limited genomics studies of the *E. coli* ST101

158    lineage and no methylome analyses to date. Here, we present the first methylome analysis of

159    ST101 using PacBio SMRT sequencing data for seven *E. coli* ST101 genomes. We defined the

160    patterns of DNA methylation across all seven ST101 genomes, pairing recognition sites with their

161    cognate MTase. Notably, we found a functional Type I RM system encoded within a Transposon 7-

162    like Transposon (Tn) was responsible for extensive methylome differences in otherwise identical

163    strains. By including an additional 13 previously published, draft ST101 genomes, we found that

164    the majority of MTases were encoded on variably distributed MGEs, giving the potential for an

165    unprecedented level of differential methylation within a single *E. coli* lineage.

166

167    **Methods**

168    **SMRT sequencing and whole-genome detection of methylated bases.** Genomic DNA (gDNA) of

169    seven *E. coli* ST101 strains; MS6192, MS6193, MS6194, MS6201, MS6203, MS6204 and MS6207

170    was extracted from overnight cultures and sequenced on either a PacBio RSI or RSII sequencer as

171    previously described [10]. Detection of methylated bases and the identification of associated

172    methyltransferase recognition sites across the seven genomes (2 complete, 5 draft), was

173    performed using the RS_Modification_and_Motif_Analysis protocol within the SMRT Analysis suite

174    v2.3.0. Interpulse durations (IPDs) were calculated based on the kinetics of the nucleotide

7

175     incorporation and were processed as previously described [20]. Sequence motifs were identified

176     using Motif Finder v1, implemented in the SMRT Portal v2.3.0. Quality value cut-offs of 20 and 30

177     were applied for the draft and complete genomes, respectively. Here we report only $^{6m}A$

178     methylation. As the DNA was not Ten-eleven translocation (Tet) treated prior to sequencing, $^{5m}C$

179     modifications were not quantitated and $^{4m}C$ modifications were not identified in any genome.

180

181     **Analysis of methyltransferase target site enrichment in gene regulatory regions.** Putative gene

182     regulatory (promoter) regions were defined as up to 300 bp upstream of the start codon of each

183     CDS. To identify RM.EcoST101V recognition sites that were within intergenic regions we used the

184     Bedtools v2.23.0 [25] closest flag, which reports the nearest genomic distance between

185     recognition sites and CDSs. Sites that were within or overlapped the ends of CDSs were removed.

186     A list of all protein-coding genes that contained a 5`-$^{6m}$**A**CG$N_5$G$\underline{T}$TG-3` site within 300 bp upstream

187     of a start codon in MS6193 was generated (Supplementary Dataset, Table S1). This was used as a

188     target gene list to compare with a background gene list formed by all genes within the *E. coli* K12

189     genome. If a gene was within an operon, all members of the operon were included. This target

190     and background gene list comparison was performed using the functional enrichment analysis

191     within the Database for Annotation, Visualization and Integrated Discovery (DAVID) v6.7 [26, 27].

192     Genes were annotated based on known function, gene ontologies and pathways. Results were

193     determined as significant if post-hoc Benjamini-Hochberg correction for multiple testing reported

194     a P value of <0.05.

195

196     **Methyltransferase diversity.** To further analyse the distribution of MTases across the *E. coli* ST101

197     lineage, 13 additional published and publicly available ST101 draft genomes were downloaded

198     from Genbank or the SRA as previously described [10]. An additional eight ExPEC complete

199     genomes (accession details available in Supplementary Dataset, Table S2) were also included to

200    emphasise ST101 lineage specific MTases. Active MTase genes identified in the *E. coli* ST101 draft

201    genomes, all MTase genes from MS6192 and MS6193 and MTases from the REBASE Gold Standard

202    database were searched against the 20 ST101 genomes and eight ExPEC complete genomes

203    (Nucleotide Blast, ≥90% nucleotide identity and >95% sequence coverage) with redundancy

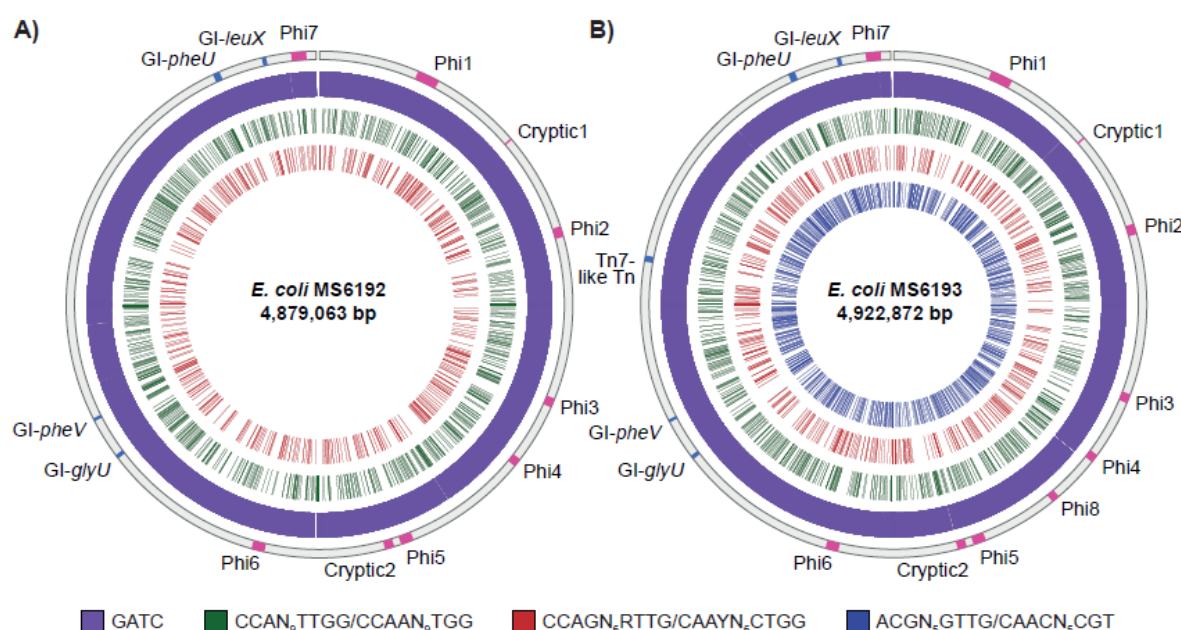204    removed. The presence or absence of MTase genes were visualised using Seqfindr

205    (http://github.com/mscook/seqfindr).

206

207    **Results**

208    ***E. coli* ST101 complete genomes MS6192 and MS6193 encode an almost identical complement**

209    **of DNA methyltransferases.** To characterise the role of methylation in shaping the *E. coli* ST101

210    lineage, we first defined the MTase complement of two near-identical Clade 1 *E. coli* ST101 strains

211    (MS6192 and MS6193) for which we had previously determined the complete genomes [10]. The

212    MS6192 genome encodes 12 putative MTases, with 8 on the chromosome, two on the large

213    *bla*$_{NDM-1}$-positive plasmid pMS6192A-NDM and one on each of the other large plasmids

214    (pMS6192B and pMS6192C) (Table 1). Three chromosomal MTases correspond to enzymes that

215    have been characterised in other *E. coli* strains, including Dam (5'-G$^{6m}$A$\underline{T}$C-3'), Dcm

216    (5'-C$^{5m}$CW$\underline{GG}$-3') and a homolog of the orphan MTase gene *yhdJ* encoding M.EcoST101III

217    (5'-A$\underline{T}$GC$^{6m}$AT-3') (by convention, underlined bases indicate methylation on the opposite strand),

218    which has previously been reported to be inactive in other *E. coli* [28]. Additionally, we identified

219    two Dam-like, orphan, Type II MTases (M.EcoST101I and M.EcoST101II) of unknown specificity

220    located on the prophages Phi2 and Phi6, respectively. Three orphan, Type II MTases with unknown

221    recognition sites also exist, with M.EcoST101VI and M.EcoST101VII encoded on Phi7 and

222    M.EcoST101VIII encoded on the *bla*$_{NDM-1}$-positive F-type plasmid pMS6192A-NDM. Also present

223    are two orphan, Type II MTases encoded on each of the plasmids pMS6192B (M.EcoST101X) and

224    pMS6192C (M.EcoST101XI); the recognition sites of these two MTases remains unknown. The

225    remaining two MTases correspond to Type I RM systems. RM.EcoST101V is carried on the

226    chromosome in an ST101 region of difference (RD12), with RM.EcoST101IX encoded on

227    pMS6192A-NDM.

228

229    Consistent with their close evolutionary relationship, MS6193 encodes all MTases found in

230    MS6192 except for the Type II MTase M.EcoST101XI, as there is no plasmid corresponding to

231    pMS6192C in the MS6193 genome. MS6193 also encodes an additional Type I RM system

232    RM.EcoST101IV, carried on the Tn7-like Transposon that is not found in the MS6192 genome.

233    Despite MS6193 also encoding an additional prophage (Phi8), no MTases were identified on this

234    MGE.

235

236    **_E. coli_ ST101 MS6192 and MS6193 genomes are differentially methylated.** The genome-wide

237    distribution of methylated bases in the complete genomes of MS6192 and MS6193 was

238    determined using PacBio SMRT sequencing. Three distinct MTase recognition sites were detected

239    as $^{6m}$A methylated in both genomes: 5'-G$^{6m}$A$\underline{T}$C-3', 5'-CC$^{6m}$AN$_9$$\underline{T}$TGG-3' and 5'-CC$^{6m}$AGN$_6$R$\underline{T}$TG-3'.

240    The recognition site 5'-$^{6m}$ACGN$_5$G$\underline{T}$TG-3' was also detected in MS6193, but not MS6192 (Figure 1).

241    To assign methylated sites to their cognate enzyme we used a process of elimination. As expected,

242    one of the four recognition sites matched the well-characterised orphan, Type II MTase Dam

243    (M.EcoST101Dam), with known specificity: 5'-G$^{6m}$A$\underline{T}$C-3'. Recently, the Type I RM recognition site

244    5'-CC$^{6m}$AN$_9$$\underline{T}$TGG-3' was identified in the _E. coli_ strain GB089, however a cognate MTase was not

245    assigned in REBASE [14]. Nucleotide comparisons of all Type I RM systems in GB089 and MS6192

246    revealed a single match between RM.EcoG089ORF25920P and RM.EcoST101V (100% identity),

247    thus we deduce that the 5'-CC$^{6m}$AN$_9$$\underline{T}$TGG-3' recognition site is methylated by RM.EcoST101V. To

248    investigate the third recognition site shared by both MS6192 and MS6193 (5'-CC$^{6m}$AGN$_6$R$\underline{T}$TG-3'),

249    we searched the motif against REBASE [14] and confirmed that it matches the recognition site of

250    RM.Eco067II, identified in the *E. coli* strain AR_0067 (Genbank accession: CP032258). This motif is

251    characteristic of a Type I RM system and with only one other Type I RM system identified in

252    MS6192, we deduce that RM.EcoST101IX is responsible for methylation of the 5'-CC$^{6m}$AGN$_6$R$\underline{T}$TG-3'

253    site. Amino acid comparisons of the specificities subunits (HsdS) S.Eco067II and S.EcoST101XI

254    confirm this match (99.77% identity, single amino acid substitution Y204H). The final $^{6m}$A

255    recognition site 5'-$^{6m}$ACGN$_5$G$\underline{T}$TG-3', detected only in MS6193, has previously been identified in

256    the *Klebsiella pneumoniae* strain AATZP [29], and was assigned to the Type I RM system

257    RM.KpnAATIII in REBASE. A nucleotide comparison showed that RM.KpnAATIII and RM.EcoST101IV

258    share 100% nucleotide identity. Thus the Type I RM system RM.EcoST101IV in MS6193 must be

259    responsible for methylation at 5'-$^{6m}$ACGN$_5$G$\underline{T}$TG-3'. Also observed in both genomes were

260    variations of the 5'-C$^{5m}$CW$\underline{GG}$-3' motif, which is characteristic of Dcm methylation. Despite the

261    presence of M.EcoST101Dcm in both genomes, the DNA was not Ten-eleven translocation (Tet)-

262    treated and the SMRT sequencing coverage is lower than 250X, therefore accurate detection and

263    quantification of $^{5m}$C in these genomes was limited (Supplementary Dataset, Table S3).

264

265



11

266

267 **Figure 1. Circos plot displaying the distribution of $^{m6}$A methylated bases in the *E. coli* MS6192**

268 **and MS6193 chromosomes.** The location of MGEs on the chromosome is indicated on the

269 outermost track; Prophages (Pink), Genomic Islands and Transposons (Blue). The remaining tracks

270 display the location of methylated recognition sites for each motif. A) For *E. coli* MS6192, from

271 outer to inner: GATC, purple (M.EcoST101Dam); CCAN$_9$TTGG/CCAAN$_9$TGG, green (RM.EcoST101V);

272 CCAGN$_6$RTTG/CAAYN$_6$CTGG, red (RM.EcoST101VIII). B) For *E. coli* MS6193, from outer to inner:

273 GATC, purple (M.EcoST101Dam); CCAN$_9$TTGG/CCAAN$_9$TGG, green (RM.EcoST101V);

274 CCAGN$_6$RTTG/CAAYN$_6$CTGG, red (RM.EcoST101VIII) and ACGN$_5$GTTG/CAACN$_5$CGT, blue

275 (RM.EcoST101V).

276

277 ***RM.EcoST101IV may have acquired a secondary role in gene regulation.*** The additional 18.9 Kb

278 Tn7-like Tn in MS6193 encoding RM.EcoST101IV (Supplementary Figure S1) is one of the major

279 differences between the two complete genomes MS6192 and MS6193. We hypothesised that the

280 acquisition of this additional RM system may lead to functional differences between the MS6192

281 and MS6193 strains. While the functional role of M.EcoST101IV is not currently known, the

282 majority of the 788 5'-$^{6m}$ACGN$_5$G$\underline{T}$TG-3' sites (96%) are found in coding regions of the MS6193

283 genome. As methylation sites in intragenic regions are more likely to be associated with gene

284 regulation [23], this suggests a primary role for RM.EcoST101IV in defence against foreign DNA.

285 We also identified the presence of two methylated 5'-$^{6m}$ACGN$_5$G$\underline{T}$TG-3' sites within the Tn7-like Tn

286 itself, found in MS6193_03822 encoding a putative DNA repair ATPase (UniProt), immediately

287 upstream of the *hsdS* gene of RM.EcoST101IV. Although the functional consequence of these

288 methylated sites is unknown, this may protect the Tn7-like Tn itself from degradation.

289

290 We identified 31 5'-$^{6m}$ACGN$_5$G$\underline{T}$TG-3' sites on the MS6193 chromosome and four on the plasmid

291 pMS6193A-NDM that were in intergenic regions. Of these sites, all but one were within 300 bp of

292 a start codon, which highlights the potential for RM.EcoST101IV to have acquired a secondary role

293 in gene regulation (Supplementary Dataset, Table S1). From this, we generated a target gene list

294 of 36 genes (including all genes within an operon if the RM.EcoST101IV site was within the

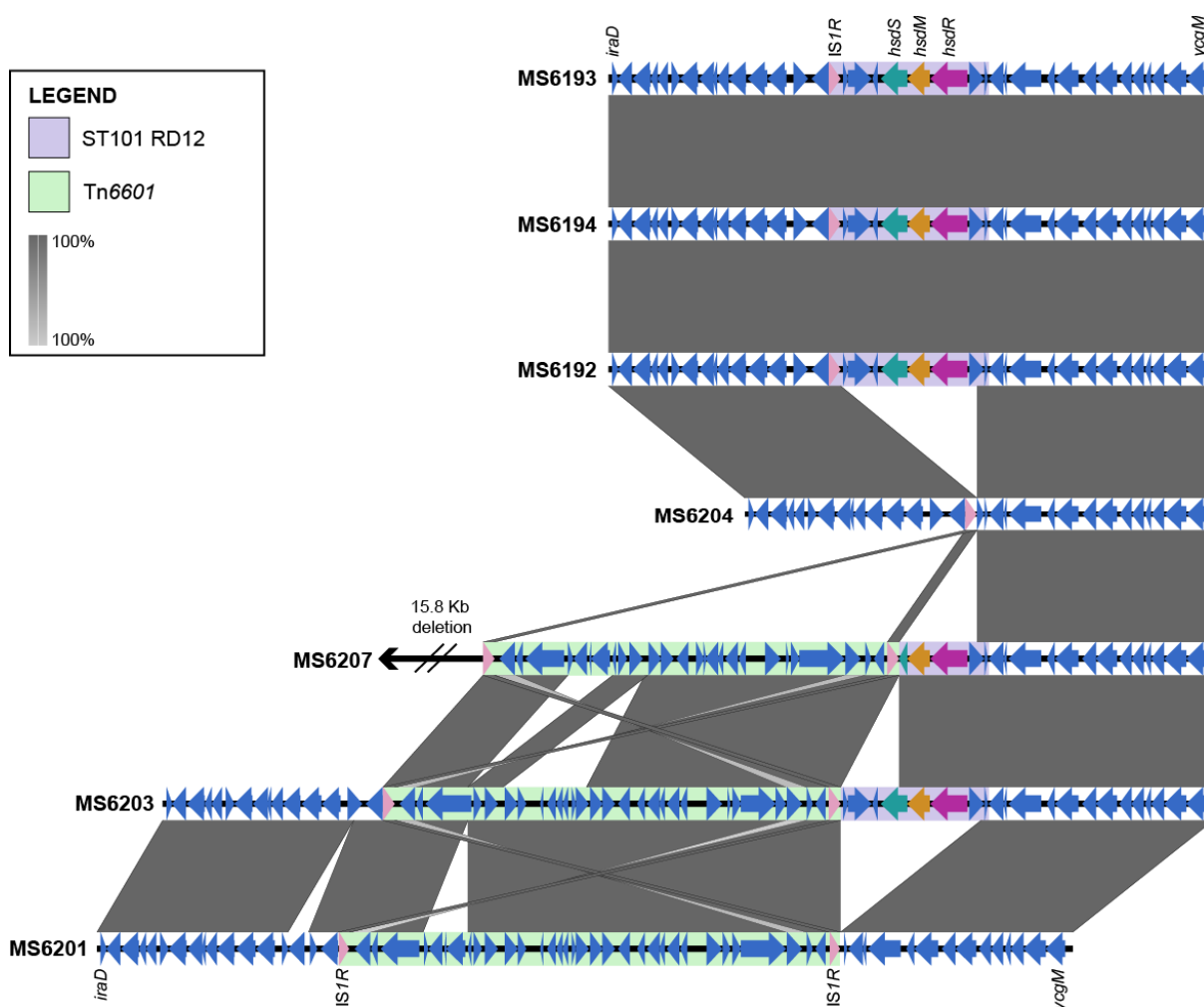295 putative promoter region for that operon). These genes include the transcriptional regulators

296  *mcbR* and *fimZ*, *mntP* (putative manganese efflux pump), *yejO* (predicted autotransporter outer

297  membrane protein), *ydcM* (putative transposase and virulence-associated protein) and *pagN*

298  (outer membrane protein and virulence-associated protein). Notably, a single RM.EcoST101IV site

299  was overlapping the start of an IS*26* element, which is 124 bp upstream of the truncated IS*Aba125*

300  element that provides the -35 promoter region for the $bla_{NDM-1}$ gene [30], responsible for

301  carbapenem resistance in this lineage. Carbapenem Minimum Inhibitory Concentration (MIC)

302  values were however identical between MS6192 and MS6193 except for Doripenem, which was

303  lower in MS6192 by one serial dilution, yet still above the resistance breakpoint [10]. Despite no

304  significant enrichment of functional pathways, these genes were primarily associated with

305  cofactor binding, cell walls and membranes, ATP binding, nucleotide binding and metal ion binding

306  (Supplementary Dataset, Table S4).

307

308  **Variation in *E. coli* ST101 Clade 1 methylomes is associated with variability in the accessory**

309  **genome.** To further investigate the ST101 methylome diversity, we included in our analyses five

310  draft genome assemblies (MS6194, MS6201, MS6203, MS6204 and MS6207) that we have

311  previously described [10]. In total, we identified four active $^{6m}$A MTases described above in

312  MS6192 and MS6193, plus an active Type I RM system (RM.EcoST101XII), found only in MS6201

313  and MS6203 (Table 2). We also identified a novel $^{6m}$A Type II-like motif (5'-AGG$^{6m}$ANTT-3') in

314  MS6203, resulting in hemi-methylation, however we could not definitively match it to its cognate

315  MTase. The following cases illustrate several different scenarios that lead to differential

316  methylation due to MGE-borne MTases in the ST101 lineage.

317

318  ***Differential methylation due to truncation of the DNA specificity gene in RM.EcoST101V.***

319  Methylation at the 5'-CC$^{6m}$AN₉TTGG-3' site, which we have assigned to the Type I RM system

320  RM.EcoST101V, is also present in two of the five draft genomes: MS6194 and MS6203.

13

321    RM.EcoST101V is encoded within the 7.2 Kb ST101 region of difference 12 (RD12) of the complete

322    genomes MS6192 and MS6193 [10] and appears to be in the same location in MS6194 and

323    MS6203. We also identified a homolog of the MTase M.EcoST101V in MS6207 (100% nucleotide

324    ID), however in MS6207 no methylation was observed at the corresponding 5'-CC$^{6m}$AN$_9$T$\underline{T}$GG-3'

325    recognition site. Further investigation revealed that the specificity gene (*hsdS*) of RM.EcoST101V in

326    the MS6207 genome was truncated at the 3' end by the upstream insertion of a large ~28 Kb

327    composite transposon: Tn*6601* . This truncation resulted in the loss of DNA specificity and

328    therefore loss of methylation and restriction activity. Additionally, in MS6207, the acquisition of

329    this Tn*6601* resulted in the deletion of a 15.8 Kb genomic region immediately upstream of the IS*1R*

330    (Figure 2). Tn*6601* is also present in both MS6203 and MS6201. In MS6203, Tn*6601* has also

331    inserted upstream of the original IS*1R*, leaving the RD12 locus intact, however in MS6201, Tn*6601*

332    has completely replaced the RD12 locus. In MS6204 however, Tn*6601* is not present and the

333    absence of RM.EcoST101V is due to the absence of the RD12 locus, leaving only the conserved

334    IS*1R* element remaining.

335

**Figure 2. Conservation of RM.EcoST101V and the ST101 Region of Difference 12 (RD12) locus.** Grey shading indicates nucleotide identity between sequences according to BLASTn (100%). MS6201 was reverse complemented for easier visualisation. Strains are in order as they appear in the ST101 phylogenetic tree [10]. ST101 Region of Difference (RD)12 (purple), IS*1R* flanked Tn*6601* (pale green), CDSs (blue), IS*1R* (pale-pink), *hsdS* specificity gene (teal), *hsdM* methyltransferase gene (orange), *hsdR* restriction gene (dark-pink). Image prepared using EasyFig.

To determine the distribution of Tn*6601* in ST101 strains, we included an additional seven published and publicly available Clade 1 ST101 draft genomes. All or most of this transposon is present in subclade 1.2 strains (MS6207, PI7, MS6203, NA086, NA084, MS6201 and NA099) (Supplementary Figure S2). This sequence of events is consistent with the acquisition of the IS*1R* flanked Tn*6601* into the same genomic locus as RD12 and then subsequent, independent transposition events leading to a) no change to the RD12 locus, b) truncation of the RD12 locus or c) loss of the RD12 locus. This result highlights how small recombination events at a single

15

351 genomic locus between otherwise highly similar bacterial strains can result in global methylome

352 changes within a single lineage.

353

354 ***Differential methylation due to acquisition of plasmid-borne RM.EcoST101IX.*** The

355 5'-CC$^{6m}$AGN$_6$R$\underline{T}$TG-3' recognition site encoded by the Type I RM system RM.EcoST101IX was also

356 identified in the draft genomes of MS6194 and MS6204. In all four genomes that carry this Type I

357 RM system (MS6192, MS6193, MS6194 and MS6204), this system is encoded on a pIP1206-like

358 plasmid within the F-type plasmid backbone and is carried on the IS*1*-family flanked composite

359 transposon Tn*6602* (Figure 3a). A BLAST comparison of RM.EcoST101IX confirmed that this Type I

360 RM system is also present (100% nucleotide ID) in the original plasmid pIP1206 (Genbank

361 accession: AM886293), with numerous homologs (>99% nucleotide ID) in other *E. coli*,

362 *K. pneumoniae* and *Salmonella enterica* plasmid sequences (Supplementary Dataset, Table S5),

363 highlighting the promiscuity of this RM system.

364



365

366

**Figure 3. Genomic context of the ST101 Type I RM systems RM.EcoST101IX and RM.EcoST101XII.**
Schematic diagrams illustrating the genetic organisation and conservation of active RM systems.
A) Tn*6602* encoded RM.EcoST101IX. B) Genomic Island GI-*pheU* encoded RM.EcoST101XII. Grey
shading indicates nucleotide identity between sequences according to BLASTn. Key genomic
features are indicated including integrase gene (red), Insertion sequences (pale-pink), *hsdS*
specificity gene (teal), *hsdR* restriction gene (orange), *hsdM* methyltransferase gene (dark-pink),
CDSs (pale blue). tRNA-*pheU* position labelled. GIs are indicated by the bright-green lines. Image
prepared using EasyFig.

***Differential methylation due to the acquisition of the chromosomal MGE-encoded***

***RM.EcoST101XII.*** The methylated recognition site 5'-GC$^{6m}$AN$_5$G$\underline{T}$TC-3' is also characteristic of a

Type I RM system and is present only in MS6201 and MS6203; it was not identified in MS6192 and

MS6193. We searched the motif against REBASE and confirmed that it matches the recognition

site of RM.Dso4321II, identified in the plant pathogen *Dickeya solani* D strain s0432-1 [31], a

member of the Enterobacterales. Only a single Type I RM system was identified in MS6201

(designated RM.EcoST101XII), thus this is the probable cause of methylation at the

5'-GC$^{6m}$AN$_5$G$\underline{T}$TC-3' site. Comparisons of the specificity subunits S.Dso4321II and S.EcoST101XII

however, indicate that they share only 81.22% amino acid identity, with several substitutions in

the specificity domains (amino acids 4-183 and 244-366). While this recognition site has not

previously been characterised in *E. coli*, the specificity gene is present in several *E. coli* genomes in

REBASE (>99% amino acid identity) including one genome (*E. coli* O118:H16 str. 07-4255, Genbank

accession: JASP01000001) that has associated PacBio sequence data. Despite the presence of this

specificity gene in strain 07-4255, the recognition site was not identified in this genome. Further

investigation revealed that the S and M subunits were present, however the R subunit was missing.

It is currently unknown if the missing R subunit is the cause of the inactivity of this RM system in

strain 07-4255.

To determine the genomic context of RM.EcoST101XII in MS6201, we characterised the

surrounding genes. The presence of several IS elements, phage-like genes and hypothetical genes

17

396   in close proximity to RM.EcoST101XII suggested carriage on a MGE. Comparative genomic analysis

397   between our seven ST101 genomes characterised this region as a tRNA-*pheU* integrated GI

398   (GI-*pheU*), different to the GI-*pheU* encoded in MS6192 and MS6193. MS6203 also contained a

399   tRNA-*pheU* integrated GI, highly similar to that of the MS6201_GI-*pheU*, encoding the same Type I

400   RM system RM.EcoST101XII as MS6201 (Figure 3b).

401

402   **DNA MTase distribution is reflected by differences in the accessory genome.** To analyse the

403   distribution of MTases across the ST101 lineage, we supplemented the seven PacBio sequenced

404   genomes with 13 publicly available and published draft ST101 genomes and eight publicly

405   available extraintestinal pathogenic *E. coli* (ExPEC) complete genomes (Supplementary Dataset,

406   Table S2). *E. coli* ST101 genomes contain many MTases that are both conserved and variable

407   across the lineage and other complete, reference ExPEC strains (Figure 4 and Supplementary

408   Dataset, Table S6). M.EcoST101Dam, M.EcoST101Dcm and the YhdJ homolog M.EcoST101III are

409   encoded in syntenic positions in all genomes, with all seven ST101 PacBio sequenced genomes

410   showing methylation at the 5'-G$^{6m}$**A**$\underline{T}$C-3' Dam recognition site. Aside from these core-genome

411   conserved MTases, the distribution of all other ST101 MTases is consistent with the presence or

412   absence of MGEs on which they are encoded (Supplementary Figure S3). For example, the Tn*7*-like

413   Tn-encoded Type I RM system RM.EcoST101IV in MS6193 is present in only two other ST101

414   genomes (NA086 and NA084) that also carry the Tn*7*-like Tn and is completely absent in all other

415   genomes surveyed. Likewise, the Type I RM system RM.EcoST101XII encoded on a tRNA-*pheU*

416   integrated GI is present only in the two draft ST101 genomes MS6201 and MS6203. Other ST101

417   MTases however, show a variable distribution. For example, M.EcoST101I is encoded on Phi2 and

418   shows a distribution consistent with the variability of this element in Clades 1 and 2 as well as the

419   ExPEC complete genomes ED1a and CFT073. In MS6192, MS6193 and MS6194, M.EcoST101II is

420   encoded on Phi6, however a homolog is also found in HVH 98 with several Phi6 gene modules also

18

421      conserved in HVH 98. Interestingly, while M.EcoST101VI and M.EcoST101VII are both encoded on

422      Phi7 in the complete genomes MS6192 and MS6193, their distribution differs throughout the

423      ST101 lineage, which is likely due to differences in Phi7 gene content across the lineage.

424



425
426
427
428      **Figure 4. Distribution of MTases in the _E. coli_ ST101 lineage.** MTases conserved in _E. coli_ MS6192
429      and MS6193 (A: purple), MTases encoded only in MS6192 (B: blue), MTases encoded only in
430      MS6193 (C: pink), MTases encoded on the ST101 draft genomes only (D: green) and accessory
431      ST101 MTases not encoded in _E. coli_ MS6192 or MS6193 (E: orange) are shown along the X-axis.
432      Strain identifiers are listed on the Y-axis, with ST101 strains ordered according to phylogenetic
433      relationship. Black shading indicates a match of >=90% nucleotide identity with a minimum of 95%
434      query coverage. Calculated by comparing the query sequences of _E. coli_ MS6192 and MS6193
435      MTases and MTases defined as the gold standard from the REBASE database [14] to the complete
436      genomes or draft assemblies of _E. coli_ ST101 strains, as implemented in Seqfindr
437      (http://github.com/mscook/seqfindr). Blastn results can also be found in Supplementary Dataset,
438      Table S6.
439

440      Plasmid-borne MTases are also variably conserved throughout Clade 1 strains. The active, plasmid

441      encoded Type I RM system RM.EcoST101IX is conserved across all Clade 1 strains that harbour

442      Tn_6601_, carried on the pIP1206-like F-type plasmid. However, the Clade 1 isolates NA099,

443    MS 79-10 and MS 107-1 also encode a homologous MTase. Further investigations reveal that

444    NA099 shares an identical RM system to RM.EcoST101IX and contains an F-type plasmid, also

445    encoding the $bla_{NDM-1}$ locus. The MTase in MS 79-10 and MS 107-1 however shares only 97%

446    nucleotide identity to the MTase M.EcoST101IX, with these genomes not encoding a $bla_{NDM-1}$-

447    positive F-type plasmid. Similarly, the Type II MTase M.EcoST101X is conserved across all Clade 1

448    strains that encode the pGUE-NDM-like FII plasmid, with an M.EcoST101X homolog (100%

449    nucleotide ID) also present in MS6201. M.EcoST101X homologs are also present in several publicly

450    available *E. coli* and *K. pneumoniae* plasmid sequences (Supplementary Dataset, Table S7),

451    suggesting that despite inactivity under normal laboratory growth conditions, this MTase may

452    have an important biological function. Lastly, homologs of M.EcoST101XI encoded on the IncI1

453    plasmid pMS6192C are present in the majority of ST101 strains, even in genomes that do not carry

454    IncI1 plasmids.

455

456    Using the REBASE Gold Standard database (MTases that have been experimentally validated) and

457    removing redundancy, we identified three additional accessory MTase genes that were not

458    encoded within the complete genomes MS6192 and MS6193 (Figure 4). One of these MTases,

459    M.EcoRII is part of a Type II RM system, present on an FII plasmid in MS6201 and is a predicted

460    Dcm homolog. The Clade 2 isolates KTE184, KTE91 and KTE90 contain a Type I MTase similar (96%

461    nucleotide ID) to M.EcoMIII from the ExPEC complete genome EC958 [32]. Lastly, the Clade 2

462    strain HVH 98 also encodes a Type I MTase, homologous (92% nucleotide ID) to M.Sen11997I from

463    the *Salmonella enterica* subsp. *enterica* serovar Chester strain ATCC 11997 [33].

464

465    **Discussion**

466    We have previously characterised the role of MGEs in the carriage of $bla_{NDM-1}$, conferring

467    carbapenem resistance in the two *E. coli* ST101 complete genomes (MS6192 and MS6193) and five

468     draft genomes (MS6194, MS6201, MS6203, MS6204 and MS6207) [10]. In the present study, we

469     used these genomes and the kinetics of PacBio SMRT sequencing to bioinformatically characterise

470     DNA MTases, assign recognition sites with their cognate MTase and to define the genomic context

471     of MTases within our collection, facilitating the first comprehensive methylome analysis of the

472     ST101 lineage.  We identified 14 DNA MTases and eight $^{6m}$A recognition sites, including one novel

473     site that could not be assigned to its cognate MTase. We also showed that eight MTases shared by

474     MS6192 and MS6193 were either inactive under the growth conditions tested or responsible for

475     $^{5m}$C methylation, which was not characterised in this study. Transcriptionally silent MTases may be

476     active under specific circumstances such as stress induction or changes in environment. It is

477     possible that cloning and expression of these genes via a plasmid system in a MTase-free strain of

478     *E. coli*, as has been performed previously for other MTases [34], could reveal their target

479     specificity. Overall, our capacity to resolve complex MGEs and define the genomic context of

480     MTases within the ST101 lineage has revealed strain, clade and lineage-wide methylome

481     heterogeneity.

482

483     There is an almost identical methylation profile between the two complete ST101 genomes

484     MS6192 and MS6193, however we show that the acquisition of a single, active RM system

485     (RM.EcoST101IV) encoded on the Tn7-like Tn (present only in MS6193) resulted in 788

486     differentially methylated sites. While more than 96% of sites were within intragenic regions of the

487     genome, 27 sites were within intergenic regions, with all but one located in putative promoter

488     regions (which we defined as ≤300 bp from a start codon). Thus, it is possible that methylation of

489     the RM.EcoST101IV site 5'-$^{6m}$**A**CGN$_5$G**T**TG-3' could result in an indirect role in gene regulation.

490     While the gene regulatory role of orphan MTases such as Dam has previously been demonstrated

491     [11], there are also examples of acquisitions of RM systems that have caused differential

492     methylation patterns and thus differential gene regulation. For example, comparisons of the

21

493     knockout mutant of the Type II RM system RM.EcoGIII encoded on the Shiga toxin phage, to the

494     wild-type *E. coli* C227-11 strain led to more than a third of all genes differentially expressed [34],

495     indicating that acquired MTases encoded on MGEs can result in significant changes to gene

496     expression. Future work will involve analysing the intersection of the methylome and

497     transcriptome via RNA sequencing methods.

498

499     Currently, it is unknown whether the additional RM system RM.EcoST101IV could generate

500     barriers of DNA exchange and influence the gene pool available to MS6193 however, RM systems

501     do have a role in maintaining species identity and restricting horizontal gene transfer in some

502     species. For example, in *Neisseria meningitidis*, the distribution of RM systems is consistent with

503     its phylogenetic clade structure [35]. Intraclade HGT was significantly more likely than interclade

504     HGT, highlighting that RM systems generate barriers to DNA exchange and are involved in the

505     evolution of distinct lineages [35]. In *Staphylococcus aureus*, a mutation in the restriction subunit

506     (*hsdR*) of the Type I RM system Sau1 is vital for plasmid transformation of the laboratory strain

507     *S. aureus* RN4220, allowing uptake of foreign DNA [36]. Additionally, distinct variants of two

508     specificity units (*hsdS*) encoded on GIs were identified across the different lineages of *S. aureus*,

509     indicating lineage-specific sequence specificity [36]. In *Burkholderia pseudomallei*, each lineage

510     contained a distinct complement of RM systems, which caused clade-specific methylation patterns.

511     Transformation with reporter plasmids that carried specific restriction sites impeded the ability of

512     the *E. coli* strains encoding distinct *B. pseudomallei* RM systems to be transformed [37]. It is

513     therefore predicted that these lineage-specific RM systems partition the species by restricting HGT

514     and inhibiting uptake of non-self-DNA [37]. Whether RM systems within ST101 present a

515     significant barrier to HGT between lineages has yet to be elucidated and represents an area of

516     future research interest.

517

518    In the seven PacBio sequenced genomes, we characterised only a single ST101 MTase capable of

519    $^{5m}$C methylation (*dcm* encoded by M.EcoST101Dcm, which methylates 5'-C$^{5m}$C<u>WG</u>G-3' sites) that

520    has previously been characterised in *E. coli* [38]. However, our ability to detect $^{m5}$C methylation

521    was limited. The kinetics of $^{5m}$C methylation are subtle and spread over several bases as the

522    modification is hidden in the major groove of the DNA, limiting the effectiveness of the detection

523    algorithm [20]. This could have been overcome by increasing the number of SMRT cells used and

524    thus throughput, increasing the sequencing coverage to 250X [22]. Alternatively, enzymatic

525    conversion via Ten-eleven translocation (Tet) treatment to convert $^{5m}$C to 5-carboxylcytosine

526    increases the size of the modification, enhancing the kinetic signal [39]. However, this conversion

527    is sometimes incomplete and even with complete conversion, $^{5m}$C isn't always detected at

528    complete levels [40]. Thus, we focused our study on the dominant $^{6m}$A modifications in *E. coli*.

529

530    To date, eleven *E. coli* methylomes have been published, including the *E. coli* strains DH5α,

531    BL21(DE3) and Bal225 [41], C277-11 [34], RM13514 and RM13516 [42], EC958 [32], CFT073 and

532    K-12 substr. MG1655 [23] and 95NR1 and 95JB1 [43]. These studies have highlighted the diversity

533    of MTases across *E. coli*, with the MTase complement and site specificities varying significantly

534    even between members of the same phylogroup and ST. However, in general, each study was

535    restricted to a very small number of genomes, limiting our knowledge of MTase conservation

536    across whole lineages. Currently, this is only the second study of the distribution of MTases within

537    strains of an *E. coli* lineage, where we first noted the importance of MGEs in the distribution of

538    MTases and showed lineage-specific methyltransferase patterns in the UPEC ST131 clone [32]. By

539    characterising the genomic context of all MTases in our two ST101 complete genomes and active

540    MTases in our five draft genomes, we showed that the majority are encoded on MGEs. Including

541    an additional 13 published and publicly available draft genomes confirmed that variation in

542    MTases within the ST101 lineage was mostly due to MGE differences between genomes.

543    Furthermore, there were limited numbers of accessory MTases identified that were not encoded

544    within either MS6192 or MS6193. While our identification of accessory MTases was restricted as

545    we only compared against MTases that have been experimentally shown to possess methylation

546    activity (REBASE Gold Standard) [14], these limited numbers of accessory MTases may indicate

547    that MTases act as a barrier to HGT within the lineage.

548

549    Our analysis of the ST101 methyome shows that even within a single clade, substantial differences

550    in MTase content can occur, highlighting the need for multiple PacBio genomes across all clades to

551    reveal the full extent of epigenomic diversity within a lineage. Additionally, our findings

552    demonstrate the significant role of MGEs in enabling very similar bacterial strains to rapidly

553    acquire genome-wide differences in their methylome, highlighting the expanding role of MGEs in

554    *E. coli* evolution. Further work studying the intersection between the methylome and

555    transcriptome will expand our understanding of the functional roles of DNA methylation in

556    bacteria and provide new insights into how strain and lineage-specific methylome changes drive

557    host adaptation.

558

559    **Acknowledgements**

566

575

576    **Conflicts of interest**

577    The authors declare that there are no conflicts of interest.

**578**  **Table 1. Summary of DNA methyltransferases and Restriction-Modification systems identified in the *E. coli* ST101 complete genomes MS6192 and**
**579**  **MS6193.**

| Protein or gene | Coordinates MS6192/MS6193 | Genomic Location | Gene Cluster[1] | GC gene cluster (%) | Type | R-M/ Orphan | Mod Type | Motif Site | Comments |
|---|---|---|---|---|---|---|---|---|---|
| M.EcoST101I | 983527-984384/ 983525-984382 | Phi2 | M | 47.55 | II alpha | Orphan | $^{6m}$A | Unknown | Dam-like MTase; phage origin. Unknown activity |
| M.EcoST101Dcm | 2135755-2137173/ 2160691-2162109 | Core-Chr | M | 51.86 | II | Orphan | $^{5m}$C | 5'-CCWGG-3' | Dcm. Active MTase, sites not quantified in this study |
| M.EcoST101II | 2641213-2641740/ 2666149-2666676 | Phi6 | M | 51.89 | II alpha | Orphan | $^{6m}$A | Unknown | Dam-like MTase; phage origin. Unknown activity |
| M.EcoST101III | 3601556-3602440/ 3626492-3627376 | Core-Chr | M | 48.02 | II beta | Orphan | $^{6m}$A | 5'-ATGCAT-3' | Not active in this study. 99% aa identity to YhdJ |
| M.EcoST101Dam | 3684643-3685479/ 3709579-3710415 | Core-Chr | M | 49.58 | II alpha | Orphan | $^{6m}$A | 5'-GATC-3' | Dam. Active MTase. |
| RM.EcoST101IV | 3835989-3846759[2] | Tn7-like GI | M-x-S-x-x-R | 54.79 | I gamma | R-M | $^{6m}$A | 5'-ACGN$_5$GTTG-3' | Active R-M system. 100% ID, 100% query cover to M.KpnAATIII in *K. pneumoniae* AATZP |
| *mcrCB* | 4718391-4721806/ 4762204-4765615 | GI-*leuX* | R-S | 46.01 | IV | R-M | $^{5m}$C | 5'-RmC(N40–3000)RmC-3' | Activity undetermined. Cleaves DNA containing $^{m5}$C on one or both strands |
| RM.EcoST101V | 4761209-4721802/ 4805022-4810741 | Chr – RD12[3] | S-M-R | 48.06 | I | R-M | $^{6m}$A | 5'-CCAN$_9$TTGG-3' | Active RM system. *hsdS, hsdM, hsdR*. 100% ID, 100% query cover to RM.EcoG089ORF25920P in *E. coli* GB089 |
| M.EcoST101VI | 4811797-4812450/ 4855610-4856263 | Phi7 | M | 54.28 | II | Orphan | $^{6m}$A | Unknown | Unknown activity. 99.694% ID, 100% query cover to M.EcoACNORF4826P in *E. coli* ACN001 |

| | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|
| M.EcoST101VII | 4817088-4818140/ 4860901-4861953 | Phi7 | M | 47.86 | II beta | Orphan | $^{6m}A/^{4m}C$ | Unknown | Unknown activity. 100% ID, 100% query cover to M.EcoACNORF4834P in *E. coli* ACN001 |
| M.EcoST101VIII | 8143-8826/ 8148-8831 | pMS6192A-NDM/ pMS6193A-NDM | M | 56.58 | II | Orphan | $^{6m}A/^{4m}C$ | Unknown | Unknown activity. 100% ID, 100% query cover to M.KpnKF3ORF164P in *K. pneumoniae* strain KF3 |
| RM.EcoST101IX | 124879-130953/ 124883-130957 | pMS6192A-NDM/ pMS6193A-NDM | M-S-R | 43.52 | I gamma | R-M | $^{6m}A$ | 5'-CC**A**GN₆R**T**TG-3' | Active R-M system. *hsdS, hsdM, hsdR*. 99.936% ID, 100% query cover to M.Sen33676ORF4987P in *Salmonella enterica* subsp. enterica servoar Typhimurium 33676 |
| M.EcoST101X | 10410-11093/ 10410-11093 | pMS6192B/ pMS6193B | M | 55.4 | II beta | Orphan | $^{6m}A/^{4m}C$ | Unknown | Unknown activity. 93% ID, 100% query cover to M.Eco29787ORF26870P in *E. coli* plasmid pCFSAN029787_02 |
| M.EcoMS6192XI | 33422-34105 [3] | pMS6192C | M | 61.73 | II beta | Orphan | $^{6m}A/^{4m}C$ | Unknown | Unknown activity. 97% ID, 100% query cover to M.Eco6409ORF23710P in *E. coli* plasmid p6409-151.583kb |

580 [1]M=MTase, R=REase, S=Specificity subunit, x=any other gene. [2]Only found in MS6193. [3]ST101 Region of difference (RD) 12 - in complete genomes
581 MS6192 and MS6193. [4]Only found in MS6192.

582

583

584

585

586

587

588

589 **Table 2. Summary of DNA methyltransferase recognition sites identified in the PacBio sequenced *E. coli* ST101 strains in this study.**

| Motif[1] | Name | Position of MTase | Mod Type | Mod Pos | Strains | No. sites in genome | No. sites detected | Methylated (%) | Mean QV | Mean IPD Ratio | Mean Site coverage |
|---|---|---|---|---|---|---|---|---|---|---|---|
| GA*T*C | M.EcoST101Dam | Core-genome | 6mA | 2 | MS6207 | 43154 | 38809 | 89.93 | 35.07 | 5.97 | 16.25 |
| | | | | | MS6201 | 42732 | 39639 | 92.76 | 36.62 | 5.94 | 17.32 |
| | | | | | MS6203 | 42144 | 39517 | 93.77 | 37.26 | 5.89 | 17.91 |
| | | | | | MS6204 | 43182 | 37737 | 87.39 | 34.34 | 5.99 | 15.67 |
| | | | | | MS6192 | 41748 | 41664 | 99.8 | 95.05 | 6.09 | 53.82 |
| | | | | | MS6194 | 41956 | 39948 | 95.21 | 38.88 | 5.99 | 18.87 |
| | | | | | MS6193 | 41646 | 41458 | 99.55 | 73.68 | 5.89 | 41.53 |
| **A**CGN$_5$GTTG/ CAA**C**N$_5$CGT | RM.EcoST101IV | Tn7-like GI | 6mA | 1/3 | MS6193 | 788 | 783/777 | 99.37/98.6 | 69.53/68.31 | 6.71/6.36 | 41.22/41.26 |
| CC**A**N$_9$TTGG/ CCA**A**N$_9$TGG | RM.EcoST101V | RD12[2] | 6mA | 3/4 | MS6203 | 518 | 487/468 | 94.02/90.35 | 38.02/35.67 | 7.39/6.43 | 18.35/18.08 |
| | | | | | MS6192 | 502 | 502/501 | 100/99.8 | 92.37/86 | 7.73/6.73 | 53.47/53.04 |
| | | | | | MS6194 | 499 | 476/466 | 95.39/93.39 | 39.26/36.96 | 6.42/3.22 | 18.95/20.51 |
| | | | | | MS6193 | 492 | 491/485 | 99.8/99.55 | 71.84/68.77 | 7.67/6.87 | 40.89/41.53 |
| CC**A**GN$_6$RTTG/ CA**A**YN$_6$CTGG | RM.EcoST101IX | pIP1206-like FII plasmid | 6mA | 3/3 | MS6204 | 853 | 756/698 | 88.63/81.83 | 35.05/32.33 | 8.79/6.1 | 15.67/15.43 |
| | | | | | MS6192 | 821 | 821/815 | 100/99.27 | 95.12/86.69 | 8.99/5.97 | 53.72/52.39 |
| | | | | | MS6194 | 833 | 807/765 | 96.88/91.84 | 39.3/36.16 | 8.78/5.91 | 18.62/18.57 |
| | | | | | MS6193 | 820 | 819/808 | 99.88/95.54 | 73.75/69.92 | 9.06/5.95 | 40.99/40.55 |
| GC**A**N$_5$GTTC/ GA**A**CN$_5$TGC | RM.EcoST101XII | GI-*pheU* | 6mA | 3/3 | MS6201 | 593 | 557/529 | 93.93/89.21 | 36.42/34.56 | 7.39/6.35 | 17.37/16.83 |
| | | | | | MS6203 | 577 | 544/519 | 94.28/89.95 | 36.44/34.82 | 7.34/6.08 | 17.57/17.35 |
| AGG**A**NTT | N/A | Unknown | 6mA | 4 | MS6203 | 1986 | 1803 | 90.79 | 35.53 | 5.82 | 17.88 |

590 [1]By convention, bold bases and underlined bases indicate methylation on the forward and reverse strand, respectively. [2]ST101 Region of Difference
591 (RD) 12 – defined in our complete genomes MS6192 and MS6193. QV = Quality Value

592     **References**

593     1.      **Peirano G, Mulvey GL, Armstrong GD, Pitout JD**. Virulence potential and adherence
594     properties of Escherichia coli that produce CTX-M and NDM beta-lactamases. *Journal of medical*
595     *microbiology* 2013;62(Pt 4):525-530.
596     2.      **Mora A, Blanco M, Lopez C, Mamani R, Blanco JE et al.** Emergence of clonal groups
597     O1:HNM-D-ST59, O15:H1-D-ST393, O20:H34/HNM-D-ST354, O25b:H4-B2-ST131 and ONT:H21,42-
598     B1-ST101 among CTX-M-14-producing Escherichia coli clinical isolates in Galicia, northwest Spain.
599     *International journal of antimicrobial agents* 2011;37(1):16-21.
600     3.      **Hertz FB, Nielsen JB, Schonning K, Littauer P, Knudsen JD et al.** Erratum to: Population
601     structure of Drug-Susceptible,-Resistant and ESBL-producing Escherichia coli from community-
602     acquired urinary tract infections. *BMC microbiology* 2016;16(1):114.
603     4.      **Wagner S, Gally DL, Argyle SA**. Multidrug-resistant Escherichia coli from canine urinary
604     tract infections tend to have commensal phylotypes, lower prevalence of virulence determinants
605     and ampC-replicons. *Veterinary microbiology* 2014;169(3-4):171-178.
606     5.      **Mushtaq S, Irfan S, Sarma JB, Doumith M, Pike R et al.** Phylogenetic diversity of
607     Escherichia coli strains producing NDM-type carbapenemases. *The Journal of antimicrobial*
608     *chemotherapy* 2011;66(9):2002-2005.
609     6.      **Yoo JS, Kim HM, Koo HS, Yang JW, Yoo JI et al.** Nosocomial transmission of NDM-1-
610     producing Escherichia coli ST101 in a Korean hospital. *The Journal of antimicrobial chemotherapy*
611     2013;68(9):2170-2172.
612     7.      **Poirel L, Savov E, Nazli A, Trifonova A, Todorova I et al.** Outbreak caused by NDM-1- and
613     RmtB-producing Escherichia coli in Bulgaria. *Antimicrobial agents and chemotherapy*
614     2014;58(4):2472-2474.
615     8.      **Toleman MA, Bugert JJ, Nizam SA**. Extensively drug-resistant New Delhi metallo-beta-
616     lactamase-encoding bacteria in the environment, Dhaka, Bangladesh, 2012. *Emerging infectious*
617     *diseases* 2015;21(6):1027-1030.
618     9.      **Ranjan A, Shaik S, Mondal A, Nandanwar N, Hussain A et al.** Molecular Epidemiology and
619     Genome Dynamics of New Delhi Metallo-beta-Lactamase-Producing Extraintestinal Pathogenic
620     Escherichia coli Strains from India. *Antimicrobial agents and chemotherapy* 2016;60(11):6795-
621     6805.
622     10.     **Ashcroft MM, Forde BM, Phan M-D, Peters KM, Henderson A et al.** Genomic
623     characterisation and context of the $bla_{NDM-1}$ carbapenemase in *Escherichia coli* ST101. *bioRxiv*
624     2020:860726.
625     11.     **Casadesus J, Low D**. Epigenetic gene regulation in the bacterial world. *Microbiology and*
626     *molecular biology reviews : MMBR* 2006;70(3):830-856.
627     12.     **Korlach J, Turner SW**. Going beyond five bases in DNA sequencing. *Current opinion in*
628     *structural biology* 2012;22(3):251-261.
629     13.     **Vasu K, Nagaraja V**. Diverse functions of restriction-modification systems in addition to
630     cellular defense. *Microbiology and molecular biology reviews : MMBR* 2013;77(1):53-72.
631     14.     **Roberts RJ, Vincze T, Posfai J, Macelis D**. REBASE--a database for DNA restriction and
632     modification: enzymes, genes and genomes. *Nucleic acids research* 2015;43(Database issue):D298-
633     299.
634     15.     **Murray NE**. Type I restriction systems: sophisticated molecular machines (a legacy of
635     Bertani and Weigle). *Microbiology and molecular biology reviews : MMBR* 2000;64(2):412-434.
636     16.     **Tock MR, Dryden DT**. The biology of restriction and anti-restriction. *Current opinion in*
637     *microbiology* 2005;8(4):466-472.

638     17.     **Roberts RJ, Belfort M, Bestor T, Bhagwat AS, Bickle TA et al.** A nomenclature for
639     restriction enzymes, DNA methyltransferases, homing endonucleases and their genes. *Nucleic*
640     *acids research* 2003;31(7):1805-1812.
641     18.     **Pingoud A, Wilson GG, Wende W**. Type II restriction endonucleases--a historical
642     perspective and more. *Nucleic acids research* 2014;42(12):7489-7527.
643     19.     **Oliveira PH, Touchon M, Rocha EP**. The interplay of restriction-modification systems with
644     mobile genetic elements and their prokaryotic hosts. *Nucleic acids research* 2014;42(16):10618-
645     10631.
646     20.     **Flusberg BA, Webster DR, Lee JH, Travers KJ, Olivares EC et al.** Direct detection of DNA
647     methylation during single-molecule, real-time sequencing. *Nature methods* 2010;7(6):461-465.
648     21.     **Chin CS, Alexander DH, Marks P, Klammer AA, Drake J et al.** Nonhybrid, finished microbial
649     genome assemblies from long-read SMRT sequencing data. *Nature methods* 2013;10(6):563-569.
650     22.     **Biosciences P**. 2015. White Paper - Detecting DNA Base Modifications Using Single
651     Molecule, Real-Time Sequencing. http://www.pacb.com/wp-
652     content/uploads/2015/09/WP_Detecting_DNA_Base_Modifications_Using_SMRT_Sequencing.pdf
653     [accessed 06/03/2017].
654     23.     **Blow MJ, Clark TA, Daum CG, Deutschbauer AM, Fomenkov A et al.** The Epigenomic
655     Landscape of Prokaryotes. *PLoS genetics* 2016;12(2):e1005854.
656     24.     **Davis BM, Chao MC, Waldor MK**. Entering the era of bacterial epigenomics with single
657     molecule real time DNA sequencing. *Current opinion in microbiology* 2013;16(2):192-198.
658     25.     **Quinlan AR, Hall IM**. BEDTools: a flexible suite of utilities for comparing genomic features.
659     *Bioinformatics (Oxford, England)* 2010;26(6):841-842.
660     26.     **Huang da W, Sherman BT, Lempicki RA**. Systematic and integrative analysis of large gene
661     lists using DAVID bioinformatics resources. *Nature protocols* 2009;4(1):44-57.
662     27.     **Huang da W, Sherman BT, Lempicki RA**. Bioinformatics enrichment tools: paths toward the
663     comprehensive functional analysis of large gene lists. *Nucleic acids research* 2009;37(1):1-13.
664     28.     **Broadbent SE, Balbontin R, Casadesus J, Marinus MG, van der Woude M**. YhdJ, a
665     nonessential CcrM-like DNA methyltransferase of Escherichia coli and Salmonella enterica. *Journal*
666     *of bacteriology* 2007;189(11):4325-4327.
667     29.     **Conlan S, Lau AF, Palmore TN, Frank KM, Segre JA**. Complete Genome Sequence of a
668     Klebsiella pneumoniae Strain Carrying blaNDM-1 on a Multidrug Resistance Plasmid. *Genome*
669     *announcements* 2016;4(4).
670     30.     **Poirel L, Bonnin RA, Nordmann P**. Analysis of the resistome of a multidrug-resistant NDM-
671     1-producing Escherichia coli strain by high-throughput genome sequencing. *Antimicrobial agents*
672     *and chemotherapy* 2011;55(9):4224-4229.
673     31.     **Khayi S, Blin P, Chong TM, Chan KG, Faure D**. Complete Chromosome and Plasmid
674     Sequences of Two Plant Pathogens, Dickeya solani Strains D s0432-1 and PPO 9019. *Genome*
675     *announcements* 2018;6(17).
676     32.     **Forde BM, Phan MD, Gawthorne JA, Ashcroft MM, Stanton-Cook M et al.** Lineage-Specific
677     Methyltransferases Define the Methylome of the Globally Disseminated Escherichia coli ST131
678     Clone. *mBio* 2015;6(6).
679     33.     **Timme RE, Pettengill JB, Allard MW, Strain E, Barrangou R et al.** Phylogenetic diversity of
680     the enteric pathogen Salmonella enterica subsp. enterica inferred from genome-wide reference-
681     free SNP characters. *Genome biology and evolution* 2013;5(11):2109-2123.
682     34.     **Fang G, Munera D, Friedman DI, Mandlik A, Chao MC et al.** Genome-wide mapping of
683     methylated adenine residues in pathogenic Escherichia coli using single-molecule real-time
684     sequencing. *Nature biotechnology* 2012;30(12):1232-1239.
685     35.     **Budroni S, Siena E, Dunning Hotopp JC, Seib KL, Serruto D et al.** Neisseria meningitidis is
686     structured in clades associated with restriction modification systems that modulate homologous

687    recombination. *Proceedings of the National Academy of Sciences of the United States of America*
688    2011;108(11):4494-4499.
689    36.    **Waldron DE, Lindsay JA**. Sau1: a novel lineage-specific type I restriction-modification
690    system that blocks horizontal gene transfer into Staphylococcus aureus and between S. aureus
691    isolates of different lineages. *Journal of bacteriology* 2006;188(15):5578-5585.
692    37.    **Nandi T, Holden MT, Didelot X, Mehershahi K, Boddey JA et al.** Burkholderia pseudomallei
693    sequencing identifies genomic clades with distinct recombination, accessory, and epigenetic
694    profiles. *Genome research* 2015;25(1):129-141.
695    38.    **Marinus MG, Morris NR**. Isolation of deoxyribonucleic acid methylase mutants of
696    Escherichia coli K-12. *Journal of bacteriology* 1973;114(3):1143-1150.
697    39.    **Clark TA, Lu X, Luong K, Dai Q, Boitano M et al.** Enhanced 5-methylcytosine detection in
698    single-molecule, real-time sequencing via Tet1 oxidation. *BMC biology* 2013;11:4.
699    40.    **Kozdon JB, Melfi MD, Luong K, Clark TA, Boitano M et al.** Global methylation state at
700    base-pair resolution of the Caulobacter genome throughout the cell cycle. *Proceedings of the*
701    *National Academy of Sciences of the United States of America* 2013;110(48):E4658-4667.
702    41.    **Powers JG, Weigman VJ, Shu J, Pufky JM, Cox D et al.** Efficient and accurate whole
703    genome assembly and methylome profiling of E. coli. *BMC genomics* 2013;14:675.
704    42.    **Cooper KK, Mandrell RE, Louie JW, Korlach J, Clark TA et al.** Comparative genomics of
705    enterohemorrhagic Escherichia coli O145:H28 demonstrates a common evolutionary lineage with
706    Escherichia coli O157:H7. *BMC genomics* 2014;15:17.
707    43.    **Forde BM, McAllister LJ, Paton JC, Paton AW, Beatson SA**. SMRT sequencing reveals
708    differential patterns of methylation in two O111:H- STEC isolates from a hemolytic uremic
709    syndrome outbreak in Australia. *Scientific reports* 2019;9(1):9436.

710