1    **Promoter activity of ORF-less gene cassettes isolated from the oral metagenome**

2

3    **Supathep Tansirichaiya[1+], Peter Mullany[1], Adam P. Roberts[1,2]\***

4

5    **Keywords**: integron, ORF-less gene cassettes, promoter activity, oral metagenome

6    **Running Title**: Promoter cassettes

7

8

9

10

11

12    [1]Department of Microbial Diseases, University College London, Eastman Dental Institute, 256 Gray's

13    Inn Road, London, WC1X 8LD, UK.

14    [2]Liverpool School of Tropical Medicine, Pembroke Place, Liverpool, L3 5QA, UK

15    *Corresponding author; Liverpool School of Tropical Medicine, Pembroke Place, Liverpool, L3 5QA, UK

16    T: +44(0)151 705 3247. E: Adam.Roberts@lstmed.ac.uk

17    [+]Present affiliation: Department of Clinical Dentistry, Faculty of Health Sciences, UiT the Arctic

18    University of Norway, Tromsø, Norway

**Abstract**

19

20    Integrons are genetic elements consisting of a functional platform for recombination and expression

21    of gene cassettes (GCs). GCs usually carry promoter-less open reading frames (ORFs), encoding

22    proteins with various functions including antibiotic resistance. The transcription of GCs relies mainly

23    on a cassette promoter ($P_C$), located upstream of an array of GCs. Some integron GCs, called ORF-less

24    GCs, contain no identifiable ORF with a small number shown to be involved in antisense mRNA

25    mediated gene regulation.

26    In this study, promoter sequences were identified, using *in silico* analysis, within GCs PCR amplified

27    from the oral metagenome. The promoter activity of ORF-less GCs was verified by cloning them

28    upstream of a *gusA* reporter, proving they can function as a promoter, presumably allowing bacteria

29    to adapt to multiple stresses within the complex physico-chemical environment of the human oral

30    cavity. A bi-directional promoter detection system was also developed allowing direct identification of

31    clones with promoter-containing GCs on agar plates. Novel promoter-containing GCs were identified

32    from the human oral metagenomic DNA using this construct, called pBiDiPD.

33    This is the first demonstration and detection of promoter activity of ORF-less GCs and the development

34    of an agar plate-based detection system will enable similar studies in other environments.

35

**Introduction**

36

37 Integrons are bacterial genetic elements able to integrate and express genes present on gene cassettes

38 (GCs) [1-3]. They consist of two main components; a functional platform and a variable array of GCs.

39 The functional platform, located on the 5' end of an integron, consists of an integrase gene (*intI*), and

40 its promoter (P$_{intI}$), an *attI* recombination site and a constitutive cassette promoter (P$_C$) for the

41 expression of GCs [4]. IntI is a site-specific tyrosine integrase that catalyses the insertion and excision

42 of GCs via recombination mainly at *attI* and the *attC*, the latter located at the joint of excised,

43 circularised GCs. The integrase gene; *intI,* is normally transcribed in the opposite direction to GCs

44 within an integron (Fig 1A). However, some integrons have integrase genes transcribed in the same

45 directions as their GCs. These are called unusual integrons or reverse integrons (Fig 1B), and have been

46 identified in *Treponema denticola*, *Chlorobium phaeobacteroides* and *Blastopirellula marina* [5, 6].

47 The second part of an integron is an array of GCs. Each usually contains a single promoterless open

48 reading frame (ORF) and an *attC* recombination site [7]. The proteins encoded by GCs are diverse and

49 include those associated with antibiotic resistance, virulence, and metabolism [2, 8]. When a GC is

50 excised from integron, it forms a non-replicative mobile genetic element, which can be a substrate for

51 integrase mediated recombination between *attI* (on the integrons) and *attC* (on the circular GC). This

52 directionality of recombination is favoured over *attC*:*attC* recombination, resulting in the usual

53 insertion of a newly integrated GC immediately next to the P$_C$ promoter in the first position of GC

54 array, ensuring maximal expression [9-11].

55 The expression of integron integrases is controlled via the SOS response; there is a LexA-binding site

56 located in the P$_{intI}$ [12]. In the absence of stress, the transcriptional repressor LexA binds to P$_{intL}$ and

57 prevents the transcription of *intI*. The SOS response is activated upon the accumulation of single-

58 strand DNA (ssDNA), generated during DNA damage, DNA repair, transformation, conjugation and

59 certain antibiotic exposure e.g. trimethoprim and fluoroquinolones [13-15]. RecA recognises ssDNA

60 and polymerises into RecA nucleofilaments, which induce autocleavage of LexA, releasing P$_{intI}$ from

61 repression and allowing *intI* transcription [12, 16]. By controlling the expression of IntI, bacteria can

62 reshuffle their GCs at the precise moments of need (stress), thereby avoiding accumulation of random

63 recombination events that could be deleterious to the host cell [17, 18].

64 As most of the GCs do not contain a promoter, their expression usually relies on the $P_C$ promoter. The

65 level of expression of GCs varies depending on the distance from $P_C$, as the strength of expression

66 decreases when GCs are located further from $P_C$ [19]. This ensures that a recently acquired GC will be

67 immediately expressed ensuring rapid adaptation due to stress-induced repositioned gene within the

68 integron GC array. There are also some GCs that contain their own promoters, ensuring constitutive

69 expression of their genes regardless of the $P_C$ promoter and their position within the integron array;

70 examples include *cmlA1* (chloramphenicol resistance), *qnrVC1* (quinolone resistance), *ere(A)*

71 (erythromycin resistance) and many of the GCs encoding toxin-antitoxin (TA) systems [20-23].

72 Integron GCs have been identified from environments such as soils, marine sediments, seawater and

73 more recently from human oral metagenomes [24-28]. In our previous study on the detection of

74 integron GCs in the human oral metagenome, we found 13 ORF-less GCs out of 63 identified GCs (20%)

75 [28]. ORF-less GCs have been shown to encode regulatory RNAs, for example the trans-acting small

76 RNA (sRNA)-Xcc1, encoded by the ORF-less GC of a *Xanthomonas campestris* pv. *campestris* integron,

77 which is involved in regulation of virulence [29]. Whilst promoter activity of ORF-less GCs has been

78 discussed, this has not been experimentally demonstrated [8].

79 In this study, we performed *in silico* analysis to identify promoter sequences in the GCs identified in

80 our previous study on the oral metagenome. Promoter activity was experimentally determined by

81 cloning the selected GCs upstream of the *gusA* reporter gene and measuring β-glucuronidase enzyme

82 activity. Furthermore, we devised a GC-based promoter detection strategy utilising PCR and

83 subsequent cloning between divergently orientated reporter genes. With this system, the successful

84 cloning of amplicons from promoter-containing GCs can be visualised directly on agar plates, allowing

85 the direct isolation of GC PCR amplicons with promoter activity from metagenomic DNA.

86  **Results**

87  ***in silico* analysis of the promoter sequences on the ORF-less GCs.**

88  Among 63 GCs previously identified from human oral metagenomic DNA, 13 were predicted to be ORF-

89  less GCs [28]. Using BPROM promoter prediction software, all ORF-less GCs were predicted to contain

90  promoters on both strands, suggesting that these GCs can transcribe genes in flanking GCs (Table 1).

91  In this study, we have defined the sense strand as the same strand containing the $P_C$ promoter (Fig 1).

92  **Determination of promoter activity of the ORF-less GCs using the β-glucuronidase assay.**

93  Five GCs were chosen for experimental expression analysis. GC TMB4 (amplified with primers targeting

94  *intI* and *attC*) was selected as it is ORF-less and located in the first position of the integron array [28].

95  ORF-less GCs MMU23 and MMB37 were selected as they had the highest overall score predicted by

96  BPROM. Finally, GCs SSU17 and MMB3 were selected as controls, to represent GCs with an ORF.

97  As BPROM predicted putative promoter sequences on both strands, promoter activity of the selected

98  GCs was determined by directionally cloning upstream of a promoterless β-glucuronidase (*gusA*) gene

99  on pCC1BAC-*lacZα-gusA* (Fig 2) in both directions. For the first position GC; TMB4, three different

100  constructs were made: TMB4 $P_C$ promoter, TMB4 GC, and TMB4 $P_C$-GC constructs. As the TMB4 $P_C$

101  promoter was not identical to the $P_C$ of *T. denticola* integron [30], the $P_C$ of another integron GC; TMB1

102  [28], which was identical to it [30], was included. As the selected GCs were likely derived from

103  *Treponema* spp., two experimentally verified *T. denticola* promoters, $P_{TdTro}$ and $P_{Fla}$, were also included

104  as controls showing that *T. denticola* promoters can be recognised in our *E. coli* host [31, 32]. $P_{Fla}$ and

105  $P_{Tdtro}$ were selected as they rely on different sigma factors. $P_{Tdtro}$ is recognised by sigma factor 70 ($σ^{70}$)

106  that is responsible for the transcription of most genes during growth in both *E. coli* and *Treponema*

107  spp. [31, 33], while $P_{Fla}$ is recognised by sigma factor 28 ($σ^{28}$), involved in the expression of flagella-

108  related genes in motile bacteria [32, 34]. This will determine the limitations of our assay in recognising

109  promoters associated with different types of sigma factors. The results are shown in figure 3. MMB37

5

110  and MMB3 had promoter activity on one strand, while MMU23 and SSU17 had no promoter activity,

111  compared to the negative control. The TMB4-$P_C$, TMB4 GC, TMB4 $P_C$-GC and TMB1- $P_C$ constructs, all

112  showed promoter activities on both strands. The $P_{Tdtro}$ from *T. denticola* showed strong promoter

113  activity on both sense and antisense strands, verifying that promoters from *T. denticola* are recognised

114  by *E. coli*. As the $P_C$ promoter sequences on TMB1 and TMB4 samples were different at several

115  nucleotides, it was shown that TMB4-$P_C$ had higher promoter activities than the TMB1-$P_C$ in both

116  directions (Fig 3).

**Detection of promoter-containing GCs from oral metagenome.**

118  The pCC1BAC-*lacZα-gusA* plasmid, developed for the above enzymatic assay,  had the potential to be

119  used in an agar plate-based detection strategy to detect amplified integron GCs with promoter activity

120  on either strand of DNA. This construct is called the Bi-Directional Promoter Detection plasmid

121  (pBiDiPD). To verify the utility of pBiDiPD, and also to detect novel GCs containing promoter sequences

122  in the human oral metagenome, integron GCs were amplified with SUPA4-*Nsi*I/SUPA3-*Nhe*I and

123  MARS5-*Nsi*I/MARS2-*Nhe*I primers [28], and cloned into pBiDiPD. By spreading transformants on LB

124  plates containing X-gal/IPTG and 4-methylumbelliferyl β-D-glucuronide (MUG), clones  containing

125  inserts with promoter activity in either direction could be identified. The clones with GCs containing a

126  promoter on the sense strand showed blue fluorescence when visualised under UV light, reflecting the

127  activity of β-glucuronidase enzymes catalysing MUG to yield the blue-fluorescent 4-

128  methylumbelliferyl. Clones with promoter activity on the antisense strand resulted in blue colonies as

129  a result of β-galactosidase enzymes catalysing X-Gal into a blue insoluble pigment 5,5'-dibromo-4,4'-

130  dichloro-indigo (Fig 4).

131  After screening clones from both amplicon libraries (amplified with SUPA3-SUPA4 and MARS2-MARS5

132  primers), 23 different GCs with promoter activities were identified (Table 2). Fourteen of these were

133  similar to the GCs identified in the previous study with >86% nucleotide identity [28]. Among the

6

134    recovered promoter-containing GCs, 9 out of 23 were novel including sample SSU-Pro-20, SSU-Pro-27,

135    SSU-Pro-32, SSU-Pro-46, SSU-Pro-65, MMU-Pro-5, MMU-Pro-24, and MMU-Pro-53. Artefactual PCRs

136    were discounted by detecting the consensus R' (1R) core sites [GTTRR(Y)R(Y)Y(R)] and the

137    complementary R'' (1L) core sites [R(Y)Y(R)Y(R)YAAC] of *attC* located downstream from the *attC*

138    forward primers and upstream from the *attC* reverse primers, respectively (Supplementary Table 1)

139    [35].

140    The GCs can be categorised into two groups, one predicted to encode toxin-antitoxin systems in 12

141    out of 23 GCs, including plasmid stabilization protein (toxin)-prevent-host-death protein (antitoxin),

142    BrnT (toxin)-BrnA (antitoxin), VapC (toxin)-AbrB/MazE/SpoVT family protein (antitoxin), RelE/ParE

143    family (toxin)-XRE transcriptional regulator (antitoxin). The second group contained ORF-less GCs,

144    which could be found in 7 samples, all reported in the previous study, except sample MMU-Pro-53.

145    Most of the samples (14 out of 23 GCs) showed the promoter activity only on the sense strand.

146    Samples with promoter activity only on the antisense strand were MMU-Pro-6, MMU-Pro-63, and

147    MMU-Pro-65, while 6 out of 23 GCs exhibited promoter activity on both strands.

148    **Discussion**

149    Integrons are important disseminators of antimicrobial resistance genes and therefore, it is important

150    to understand the diversity of GCs and how their expression is controlled. Even though most of the

151    GCs contained a single ORF, ORF-less GCs have also been found [24, 27, 28, 36, 37].

152    In this study, we determined promoter activity from GCs isolated by PCR from metagenomic DNA by

153    measuring promoter activity from multiple GC containing constructs. As the ORF-less GCs were

154    recovered from the oral metagenome, there is little information regarding the original host. Therefore,

155    we chose to test the promoter activities by using an *E. coli* surrogate. Nucleotide sequence analysis

156    suggested that these GCs were likely to be derived from *Treponema* spp., therefore, the ability of *E.*

157    *coli* to utilise *T. denticola* promoter sequences was determined by including the experimentally verified

7

158    *T. denticola* promoter, $P_{TdTro}$ [31] which showed high activity on both sense and antisense strands,

159    providing confidence that *E. coli* could be used. However, as no promoter activity was detected from

160    $P_{Fla}$, it suggested that our enzymatic assay cannot detect promoters associated with $\sigma^{28}$ from

161    *Treponema* spp., which could be due to an inability for the *E. coli* host to recognise the *Treponema* $\sigma^{28}$

162    promoter sequence.

163    Promoter activities of the ORF-less GCs were confirmed and quantified by using a β-glucuronidase

164    enzymatic assay. This is the first time that the promoter activity of ORF-less GCs has been

165    demonstrated *in vitro*, as shown by the activity on the sense strand of the MMB37 and both strands

166    of the TMB4. A study on the *Vibrio* integron, containing a 116-cassette array, showed that most of the

167    GCs are transcribed [38]. Therefore, ORF-less GCs could be responsible for transcription of the other

168    GCs not transcribed by $P_C$.

169    For the TMB4 GC (ORF-less GC in the first position), it was initially hypothesised that the promoter

170    could help increase the expression of the downstream GCs. This was shown when $P_C$ promoter was

171    coupled with a second promoter ($P_2$) (found in 10% of class 1 integron and located 119 bp downstream

172    from $P_C$), could result in a significantly higher expression of GCs [39, 40]. The constructs of TMB4 $P_C$

173    and TMB4 $P_C$+GC were therefore included in the assay to determine whether having a promoter GC at

174    the first position could help promote the expression of downstream GCs. The results show that

175    coupling promoter GC in the first position slightly increased the expression of reporter genes (Fig. 3).

176    However, this was not significant (*p*-value >0.99 by using ordinary one-way ANOVA followed by

177    Bonferroni's post-hoc).

178    The lack of additive promoter activity can be explained by more competition for enzymes involved in

179    transcription such as RNA polymerases (RNAP) or sigma factors to be available for transcription from

180    each promoter, resulting in lower transcriptional level [41]. Another, not mutually exclusive possibility

181    is transcriptional interference (TI) between the four promoters on the TMB4 Pc+GC construct. We have

182 experimentally shown promoter activity of TMB4 $P_C$ and TMB4 GC constructs on both strands,

183 indicating convergent TI is a possibility.

184 In usual integrons, $P_C$ is in *intI*, which is convergent to the integron integrase promoter $P_{IntI}$ downstream

185 (Fig 1), resulting in TI. The TI between $P_C$ and $P_{IntI}$ has been shown to control the expression of integrase

186 and the subsequent recombination of GCs. The weaker strength of $P_C$ could result in higher expression

187 of integrase, which increases recombination of GCs [42, 43]. This relationship of $P_C$ and *intI* might also

188 apply to the reverse integrons found in *T. denticola,* even though their position and direction of $P_{intI}$,

189 $P_C$ and *intI* gene are in a different orientation compared to the usual integrons (Fig 1).

190 Due to the lack of additive promoter activity when Pc and an ORF-less GC promoter were tested in

191 tandem we hypothesised that there is an alternative selective advantage for having an ORF-less,

192 promoter-containing GC in the first position on an integron GC array.

193 The expression level of cassette genes located further down in the array normally decreases due to

194 the formation of a stem-loop structure on mRNA at *attC* sites, which impede the progression of the

195 ribosome [44]. It was previously shown that the level of streptomycin resistance was reduced four

196 times, when the *aadA2*-containing GC was located in the second position [45]. However, our data

197 shows that the insertion of an ORF-less, promoter-containing GC in the first position did not decrease

198 the *gusA* expression significantly (considered as a proxy for the expression of gene(s) in the second

199 GC), i.e. comparing the data for TMB4 $P_C$ and TMB4 $P_C$+GC. Therefore, we hypothesised that promoter-

200 containing GCs could act as a genetic clutch, where the expression of the original first GC is partially

201 disengaged from the $P_C$ promoter and replaced by the one on the ORF-less promoter containing GC

202 (Fig 5A). This can prevent a significant change in expression of the first GC while a new, first GC is

203 sampled from the pool of GCs in order to adapt to an additional stress concurrent with the selective

204 pressure requiring expression of the first GC. This system would work as a genetic clutch with the

205 insertion of any GC containing a promoter in the same direction as $P_C$, so it could be the insertion of

206 either ORF-less GCs such as TMB4 GC, or other promoter-containing GCs such as the multiple TA-

207      containing GCs we have identified; providing another selective advantage to retaining them and

208      explaining their varied position within the GC array.

209      A genetic clutch within an integron can be of benefit to bacteria when they are exposed to multiple

210      environmental stresses such as two different antibiotics simultaneously. The first resistance gene

211      (green ORF in Fig 5Biii) can be expressed by the $P_C$ promoter, while the second resistance gene (blue

212      GC), located in the third position, is expressed by $P_C$ and the promoter GC. Therefore, allowing bacteria

213      to survive in both the presence of both drugs.

214      As the other ORF-less GC MMU23 showed no promoter activity it may have other functions or carry a

215      promoter that can be recognised in its native host but not in *E. coli*, or require other sigma factors. For

216      the ORF-containing GC MMB3 sample, the promoter activity was found on the sense strand. This GC

217      was predicted to carry toxin-antitoxin (TA) ORFs, including the PIN toxin and ribbon-helix-helix

218      antitoxin domains, which were shown to contain their own promoter. Sample SSU17 and MMU23,

219      which showed no promoter activity, can be considered as a control; illustrating that not all of GCs

220      amplified from the oral metagenome exhibited promoter activity within our assay.

221      The pCC1BAC-*lacZα-gusA* plasmid, developed for the enzymatic assay, also had potential to be used

222      for the detection of promoter activity in either direction from GCs. The clones with promoters on the

223      sense strand can be detected under UV light and showed blue fluorescence because β-glucuronidase

224      can cleave the substrate, MUG, on the plate, which produces a fluorescence compound called

225      methylumbelliferone. For the clones carrying promoters on the antisense strand, they can be detected

226      by blue-white screening as β-galactosidase can cleave X-gal, producing an intensively blue product

227      called 5,5'-dibromo-4,4'-dichloro-indigo, which can be viewed by eye under normal light.

228      To verify the application of pCC1BAC-*lacZα-gusA* plasmids as promoter detection system, integron GCs

229      were amplified from the human saliva metagenome by using SUPA3-SUPA4 and MARS2-MARS5

230      primers, which amplify integron GCs from the oral metagenome [28]. After cloning the amplified GCs

231    between both reporter genes, two groups of GCs were identified with promoter activities: ORF-less

232    GCs and TA-containing GCs. By detecting 7 clones containing ORF-less GCs with promoter activity it

233    further supported that one of the functions of ORF-less GCs in integrons is to provide promoter

234    activities.

235    TA-containing GCs are abundant in chromosomal integrons (CIs), which were suggested to have a role

236    in preventing random deletion of GCs and stabilising the large arrays CIs [22, 46, 47]. TA systems

237    normally encode a stable toxin and a labile antitoxin [48], therefore TA cassettes have to carry their

238    own promoters to ensure their expression. These were found in CIs of *Treponema* spp., such as the

239    HicA-HicB TA-containing GC in the fourth position within the GC array (Accession number NC_002967)

240    in the CI from *T. denticola* [30]. As most of the GCs amplified with our primers were homologous with

241    *Treponema* spp., these TA-containing GCs should be present in our oral metagenome and were

242    detected by our pBiDiPD based on their promoter activities.

243    Two of the GCs, SSU-Pro-9 and MMU-Pro-18, were similar to the MMB3 and MMB37 GCs, respectively,

244    which were shown by the β-glucuronidase enzyme assay to have promoter activity on the sense

245    strand. The phenotypes of SSU-Pro-9 and MMU-Pro-18 colonies also showed only a blue fluorescence

246    phenotype, reflecting the promoter activity on the sense strand, which corresponded with the

247    enzymatic assay results of MMB3 and MMB37.

248    To summarise, the promoter activities of the ORF-less integron GCs were experimentally

249    demonstrated by using a robust β-glucuronidase enzyme assay, confirming that one of the functions

250    of ORF-less GCs is to provide promoters for the expression of ORF containing GCs, in addition to

251    expression from $P_C$. The dual reporter plasmid; pBiDiPD, was developed for the direct visualisation of

252    clones containing gene cassettes with promoter activity on agar plates. This can be applied as a

253    detection system for promoter activity for any other DNA fragments.

254

255 **Materials and methods**

256 *in silico* **analysis of the human oral cavity gene cassettes and the construction of pCC1BAC-*lacZα*-GC-**

257 ***gusA* constructs.**

258 All of the ORF-less GCs and some of the GCs containing ORFs, identified in the previous study [28],

259 were predicted for putative promoter sequences by using the web-based software BPROM in the

260 Softberry package [49].

261 **Construction of pUC19-GC-*gusA* and pCC1BAC-*lacZα*-GC-*gusA* constructs.**

262 To determine the promoter activity of the selected GCs, the constructs were initially cloned into the

263 EcoRI and KpnI restriction sites on pUC19-P*tet*(M)-*gusA* plasmid [50]. The selected GCs were amplified

264 from the pGEM-T easy vector containing the GC amplicon from a previous study [28], as shown in

265 Supplementary Fig 1, by using primer listed in Supplementary Table 2.

266 Due to a significant difference in the plasmid copy number in some constructs of the pUC19-GC-*gusA*,

267 new constructs were prepared based on a low copy number CopyControl™ pCC1BAC™ vector

268 (Epicenter, UK) as it will be maintained in *E. coli* cell as one plasmid per cell and enable us to control

269 the plasmid copy number to be similar between each construct. The construct was designed to contain

270 two reporter genes, β-galactosidase *lacZα* and β-glucuronidase *gusA* genes (Fig 2 and Supplementary

271 Fig 2). As *lacZα* on pCC1BAC contained T7 promoter sequence, it was first deleted by using Q5® Site-

272 Directed Mutagenesis Kit (New England Biolabs, UK). The backbone of pCC1BAC was amplified with

273 pCC1BAC-del*LacZ*-F1 and pCC1BAC-del*LacZ*-R1, and the amplified products were treated with a

274 Kinase-Ligase-DpnI (KLD) enzyme mix, following the instructions from the manufacturer. The KLD-

275 treated product was then transformed into *E. coli* α-Select Silver Efficiency competent cells (Bioline,

276 UK) following the instructions from the manufacturer. The pCC1BAC-del*LacZ* plasmid was then

277 extracted from *E. coli* by using QIAprep Spin Miniprep Kit (Qiagen, UK), following the manufacturer's

278 instructions.

12

279    The *lacZα* reporter gene was amplified from the pUC19 vector (New England Biolabs, UK) with *LacZ-*

280    F1 and *LacZ*-R1 primers. For *gusA* reporter gene, it was amplified from pUC19-P*tet*(M)-*gusA* with *gusA-*

281    *F1* and *gusA-R1* primers. A bidirectional terminator, modified from *lux* operon, was added to *LacZ*-F1

282    and *gusA*-R1 primers, resulting in two bi-directional terminators flanking the *lacZα-gusA* reporter

283    genes [51]. This was done to prevent transcriptional read-through from the promoter in the plasmid

284    backbone and to also prevent promoters from the inserts interfering with the expression of genes on

285    the plasmid backbone. The *lacZα* and *gusA* amplicons were digested with NsiI restriction enzymes

286    (New England Biolabs, UK) and ligated together by using T4 DNA ligase (New England Biolabs, UK). The

287    *lacZα-gusA* ligated product was directionally cloned into the pCC1BAC-del*LacZ* plasmid by digesting

288    them with AatII and AvrII restriction enzymes and ligated together, resulting in pCC1BAC-*lacZα-gusA*

289    plasmid.

290    The selected GCs were amplified from each pUC19-GC-*gusA* constructs by using primer listed in

291    Supplementary Table 1. The amplicons were double digested with *Nsi*I and *Nhe*I and directionally

292    cloned into a pre-digested pCC1BAC-*lacZα-gusA* plasmid, then transformed into *E. coli* α-Select Silver

293    Efficiency competent cells.

294    **Determination of β-glucuronidase enzymatic activity.**

295    The β-glucuronidase enzymatic assay was performed to measure the promoter activity based on the

296    expression of *gusA*, following the protocol described previously with some modifications [52]. The

297    overnight cultures of *E. coli* containing the reporter constructs were prepared in LB broth

298    supplemented with 12.5 μg/mL chloramphenicol. The $OD_{600}$ of each overnight culture was measured.

299    An aliquot of 1 mL of the overnight culture was centrifuged at 3000 x *g* for 10 min and discarded the

300    supernatant. The cell pellets were incubated at -70°C for 1 hr and resuspended in 800 μl of pH 7 Z

301    buffer (50 mM 2-mercaptoethanol, 40 mM $NaH_2PO_4 \cdot H_2O$, 60 mM $Na_2HPO_4 \cdot 7H_2O$, 10 mM KCl, and 1mM

302    $MgSO_4 \cdot 7H_2O$) and 8 μl of toluene. The mixture was transferred to a 2 ml cryotube containing glass

303    beads (150–212 μm in diameter) (Sigma, UK) and vortexed twice for 5 min each with an incubation on

304    ice for 1 min in between. The glass beads were then removed by centrifugation at 3000 x $g$ for 3 min.

305    One-hundred microliters of cell lysate were mixed with 700 µl of Z-buffer, then incubated at 37°C for

306    5min. One-hundred sixty microliters of 6 mM ρ-nitrophenyl-β-D-glucuronide (PNPG) was then added

307    to the reaction and incubated at 37°C for 5 min. The reactions were stopped by adding 400 µl of 1 M

308    $Na_2CO_3$ and centrifuged at 3000 x $g$ for 10 min to remove cell debris and glass beads. The absorbance

309    of the supernatant was measured with a spectrophotometer at the wavelength of 405 nm. Three

310    biological replicates of the β-glucuronidase enzymatic assay were performed. The β-glucuronidase

311    Miller units were calculated from $\frac{A_{405}\times1000}{OD_{600}\times time\text{ (min)}\times1.25\times volume\text{(mL)}}$ [53].

**Statistical analysis.**

313    The average and standard deviation of β-glucuronidase concentration were calculated from three

314    biological replicates, which were used for the columns and error bars in figure 3, respectively. The

315    statistical comparisons between the negative control (pCC1BAC-*lacZ-gusA*) to the other constructs

316    were performed by using ordinary one-way ANOVA with either Dunnett's post-hoc test (to compare

317    each construct with a negative control) or Bonferroni's post-hoc test (to compare constructs between

318    themselves). The groups with statistically significantly difference from the control had the *p*-value of

319    less than 0.05.

**Recovery of promoter-containing GCs from the human oral metagenome**

321    The integron GCs were amplified from the human oral metagenome by using as described previously

322    with SUPA4-NsiI-SUPA4-NheI and MARS5-NsiI-MARS2-NheI primers [28]. The human oral

323    metagenomic DNA was previously extracted from the saliva samples collected from 11 volunteers in

324    the Department of Microbial Diseases, UCL Eastman Dental Institute [28]. The Ethical approval for the

325    collection and uses of saliva samples was obtained from University College London (UCL) Ethics

326    Committee (project number 5017/001).

327 The amplified products were purified and digested with NsiI and NheI and ligated with the pre-digested

328 pCC1BAC-*lacZα-gusA* plasmid. The ligated products were transformed into *E. coli* α-Select Silver

329 Efficiency competent cells by heat shock. Cells were spread on LB agar supplemented with 12.5 μg/mL

330 chloramphenicol, 80 μg/mL X-gal, 50μM IPTG, and 70 μg/mL 4-methylumbelliferyl-β-D-glucuronide

331 (MUG). After incubation at 37°C for 18 hr, the colonies with β-galactosidase activity from *lacZ* was

332 detected by blue-white screening on the agar plate, and the β-glucuronidase activity from *gusA* was

333 visualisation under UV light. Colonies exhibiting either activity were selected and subcultured on fresh

334 agar plates. The inserts were amplified by colony PCR using *lacZ*-F2 and *gusA-F2* primers and

335 sequenced by sequencing service from Genewiz (Genewiz, UK).

336 **Sequence analysis and nomenclature of promoter-containing GC amplicons.**

337 DNA sequences were visualised and analysed by using BioEdit version 7.2.0

338 (http://www.mbio.ncsu.edu/bioedit/bioedit.html). The contigs from sequencing reactions were

339 combined by using CAP contig function in the software. The sequences were then matched to the

340 nucleotide and protein database by using BlastN and BlastX from the National Centre for

341 Biotechnology Information (NCBI), respectively. The criteria for the sequence analysis of integron GC

342 were the same as described in the previous study [28]. Two additional criteria for the verification of

343 GCs detected with pCC1BAC-*lacZα-gusA* were included. Any clones containing incomplete GCs, caused

344 by digestion at internal NsiI and NheI restriction sites on the GCs, were excluded from the dataset.

345 Also chimeric inserts, which were the ligation products between digested amplicons, were also

346 excluded.

347 The promoter-containing GCs were named as described in the previous study [28]. The first and second

348 letters represented the forward primer and reverse primer used in the amplification. The third letter

349 represents the source of the human oral metagenomic DNA which is U for the United Kingdom. This

350 was followed by term "Pro", indicating the presence of promoter activity, and the number of the clone.

351 For instance, SSU-Pro-1 stands for the first clone amplified from the UK oral metagenome by using

352    SUPA3 and SUPA4 primers. The sequences of these GCs were deposited in the DNA database with the

353    accession number from MH536747 to MH536769.

354    **Conflict of interest**

355    Nothing to Declare

356    **References**

357    1.      Hall RM, Collis CM. Antibiotic resistance in gram-negative bacteria: the role of gene cassettes

358    and integrons. Drug resistance updates : reviews and commentaries in antimicrobial and anticancer

359    chemotherapy. 1998;1(2):109-19. Epub 2006/08/15. PubMed PMID: 16904397.

360    2.      Cambray G, Guerout AM, Mazel D. Integrons. Annual review of genetics. 2010;44:141-66. Epub

361    2010/08/17. doi: 10.1146/annurev-genet-102209-163504. PubMed PMID: 20707672.

362    3.      Michael CA, Gillings MR, Holmes AJ, Hughes L, Andrew NR, Holley MP, et al. Mobile gene

363    cassettes: a fundamental resource for bacterial evolution. The American naturalist. 2004;164(1):1-12.

364    Epub 2004/07/22. doi: 10.1086/421733. PubMed PMID: 15266366.

365    4.      Escudero JA, Loot* C, Nivina A, Mazel D. The Integron: Adaptation On Demand. Microbiology

366    spectrum. 2015;3(2). doi: doi:10.1128/microbiolspec.MDNA3-0019-2014.

367    5.      Boucher Y, Labbate M, Koenig JE, Stokes HW. Integrons: mobilizable platforms that promote

368    genetic diversity in bacteria. Trends in microbiology. 2007;15. doi: 10.1016/j.tim.2007.05.004.

369    6.      Wu Y-W, Doak TG, Ye Y. The gain and loss of chromosomal integron systems in the *Treponema*

370    species. BMC evolutionary biology. 2013;13(1):1-9. doi: 10.1186/1471-2148-13-16.

371   7.      Mazel D. Integrons: agents of bacterial evolution. Nature reviews. 2006;4.

372   8.      Gillings MR. Integrons: past, present, and future. Microbiology and molecular biology reviews

373   : MMBR. 2014;78(2):257-77. Epub 2014/05/23. doi: 10.1128/mmbr.00056-13. PubMed PMID:

374   24847022; PubMed Central PMCID: PMCPMC4054258.

375   9.      Collis CM, Hall RM. Gene cassettes from the insert region of integrons are excised as covalently

376   closed circles. Molecular microbiology. 1992;6(19):2875-85. Epub 1992/10/01. PubMed PMID:

377   1331702.

378   10.     Partridge SR, Recchia GD, Scaramuzzi C, Collis CM, Stokes HW, Hall RM. Definition of the *attI1*

379   site of class 1 integrons. Microbiology (Reading, England). 2000;146 ( Pt 11):2855-64. Epub

380   2000/11/07. PubMed PMID: 11065364.

381   11.     Collis CM, Recchia GD, Kim M-J, Stokes HW, Hall RM. Efficiency of Recombination Reactions

382   Catalyzed by Class 1 Integron Integrase IntI1. Journal of bacteriology. 2001;183(8):2535-42. doi:

383   10.1128/JB.183.8.2535-2542.2001. PubMed PMID: PMC95170.

384   12.     Guerin E, Cambray G, Sanchez-Alberola N, Campoy S, Erill I, Da Re S, et al. The SOS response

385   controls integron recombination. Science (New York, NY). 2009;324(5930):1034. Epub 2009/05/23.

386   doi: 10.1126/science.1172914. PubMed PMID: 19460999.

387   13.     Baharoglu Z, Bikard D, Mazel D. Conjugative DNA transfer induces the bacterial SOS response

388   and promotes antibiotic resistance development through integron activation. PLOS Genetics.

389   2010;6(10):e1001165. doi: 10.1371/journal.pgen.1001165.

390   14.    Baharoglu Z, Krin E, Mazel D. Connecting environment and genome plasticity in the

391   characterization of transformation-induced SOS regulation and carbon catabolite control of the *Vibrio*

392   *cholerae* integron integrase. Journal of bacteriology. 2012;194(7):1659-67. Epub 2012/01/31. doi:

393   10.1128/jb.05982-11. PubMed PMID: 22287520; PubMed Central PMCID: PMCPMC3302476.

394   15.    Baharoglu Z, Mazel D. *Vibrio cholerae* triggers SOS and mutagenesis in response to a wide

395   range of antibiotics: a route towards multiresistance. Antimicrobial agents and chemotherapy.

396   2011;55(5):2438-41. Epub 2011/02/09. doi: 10.1128/aac.01549-10. PubMed PMID: 21300836;

397   PubMed Central PMCID: PMCPMC3088271.

398   16.    Little JW. Mechanism of specific LexA cleavage: autodigestion and the role of RecA coprotease.

399   Biochimie. 1991;73(4):411-21. Epub 1991/04/01. PubMed PMID: 1911941.

400   17.    Harms K, Starikova I, Johnsen PJ. Costly Class-1 integrons and the domestication of the the

401   functional integrase. Mobile genetic elements. 2013;3(2):e24774. Epub 2013/08/06. doi:

402   10.4161/mge.24774. PubMed PMID: 23914313; PubMed Central PMCID: PMCPMC3681742.

403   18.    Starikova I, Harms K, Haugen P, Lunde TT, Primicerio R, Samuelsen O, et al. A trade-off between

404   the fitness cost of functional integrases and long-term stability of integrons. PLoS Pathog.

405   2012;8(11):e1003043. Epub 2012/12/05. doi: 10.1371/journal.ppat.1003043. PubMed PMID:

406   23209414; PubMed Central PMCID: PMCPMC3510236.

407   19.    Coyne S, Guigon G, Courvalin P, Perichon B. Screening and quantification of the expression of

408   antibiotic resistance genes in *Acinetobacter baumannii* with a microarray. Antimicrobial agents and

409   chemotherapy. 2010;54(1):333-40. Epub 2009/11/04. doi: 10.1128/aac.01037-09. PubMed PMID:

410   19884373; PubMed Central PMCID: PMCPMC2798560.

411    20.     Stokes HW, Hall RM. Sequence analysis of the inducible chloramphenicol resistance

412    determinant in the Tn*1696* integron suggests regulation by translational attenuation. Plasmid.

413    1991;26(1):10-9. Epub 1991/07/01. PubMed PMID: 1658833.

414    21.     da Fonseca EL, Vicente AC. Functional characterization of a Cassette-specific promoter in the

415    class 1 integron-associated *qnrVC1* gene. Antimicrobial agents and chemotherapy. 2012;56(6):3392-4.

416    Epub 2012/03/07. doi: 10.1128/aac.00113-12. PubMed PMID: 22391535; PubMed Central PMCID:

417    PMCPMC3370728.

418    22.     Szekeres S, Dauti M, Wilde C, Mazel D, Rowe-Magnus DA. Chromosomal toxin-antitoxin loci

419    can diminish large-scale genome reductions in the absence of selection. Molecular microbiology.

420    2007;63(6):1588-605. Epub 2007/03/21. doi: 10.1111/j.1365-2958.2007.05613.x. PubMed PMID:

421    17367382.

422    23.     Biskri L, Mazel D. Erythromycin esterase gene *ere(A)* is located in a functional gene cassette in

423    an unusual class 2 integron. Antimicrobial agents and chemotherapy. 2003;47(10):3326-31. Epub

424    2003/09/25. PubMed PMID: 14506050; PubMed Central PMCID: PMCPMC201170.

425    24.     Stokes HW, Holmes AJ, Nield BS, Holley MP, Nevalainen KM, Mabbutt BC, et al. Gene cassette

426    PCR: sequence-independent recovery of entire genes from environmental DNA. Applied and

427    environmental microbiology. 2001;67(11):5240-6. Epub 2001/10/27. doi: 10.1128/aem.67.11.5240-

428    5246.2001. PubMed PMID: 11679351; PubMed Central PMCID: PMCPMC93296.

429    25.     Elsaied H, Stokes HW, Nakamura T, Kitamura K, Fuse H, Maruyama A. Novel and diverse

430    integron integrase genes and integron-like gene cassettes are prevalent in deep-sea hydrothermal

431    vents. Environmental microbiology. 2007;9(9):2298-312. Epub 2007/08/10. doi: 10.1111/j.1462-

432    2920.2007.01344.x. PubMed PMID: 17686026.

433    26.    Koenig JE, Sharp C, Dlutek M, Curtis B, Joss M, Boucher Y, et al. Integron Gene Cassettes and

434    Degradation of Compounds Associated with Industrial Waste: The Case of the Sydney Tar Ponds. PloS

435    one. 2009;4(4):e5276. doi: 10.1371/journal.pone.0005276.


436    27.    Elsaied H, Stokes HW, Kitamura K, Kurusu Y, Kamagata Y, Maruyama A. Marine integrons

437    containing novel integrase genes, attachment sites, *attI*, and associated gene cassettes in polluted

438    sediments from Suez and Tokyo Bays. The ISME journal. 2011;5(7):1162-77. Epub 2011/01/21. doi:

439    10.1038/ismej.2010.208. PubMed PMID: 21248857; PubMed Central PMCID: PMCPMC3146285.


440    28.    Tansirichaiya S, Rahman MA, Antepowicz A, Mullany P, Roberts AP. Detection of Novel

441    Integrons   in   the   Metagenome   of   Human   Saliva.   PloS   one.   2016;11(6):e0157605.   doi:

442    10.1371/journal.pone.0157605.


443    29.    Chen X-L, Tang D-J, Jiang R-P, He Y-Q, Jiang B-L, Lu G-T, et al. sRNA-Xcc1, an integron-encoded

444    transposon- and plasmid-transferred trans-acting sRNA, is under the positive control of the key

445    virulence regulators HrpG and HrpX of *Xanthomonas campestris* pathovar *campestris*. RNA Biology.

446    2011;8(6):947-53. doi: 10.4161/rna.8.6.16690.


447    30.    Coleman N, Tetu S, Wilson N, Holmes A. An unusual integron in *Treponema denticola*.

448    Microbiology   (Reading,   England).   2004;150(Pt   11):3524-6.   Epub   2004/11/06.   doi:

449    10.1099/mic.0.27569-0. PubMed PMID: 15528643.


450    31.    Brett PJ, Burtnick MN, Fenno JC, Gherardini FC. *Treponema denticola* TroR is a manganese-

451    and iron-dependent transcriptional repressor. Molecular microbiology. 2008;70(2):396-409. Epub

452    2008/09/03. doi: 10.1111/j.1365-2958.2008.06418.x. PubMed PMID: 18761626; PubMed Central

453    PMCID: PMCPMC2628430.

454    32.    Limberger RJ, Slivienski LL, Izard J, Samsonoff WA. Insertional inactivation of *Treponema*

455    *denticola tap1* results in a nonmotile mutant with elongated flagellar hooks. Journal of bacteriology.

456    1999;181(12):3743-50. Epub 1999/06/15. PubMed PMID: 10368149; PubMed Central PMCID:

457    PMCPMC93852.


458    33.    Paget MSB, Helmann JD. The $\sigma^{70}$ family of sigma factors. Genome Biology. 2003;4(1):203-.

459    PubMed PMID: PMC151288.


460    34.    Koo B-M, Rhodius VA, Campbell EA, Gross CA. Mutational analysis of *Escherichia coli* $\sigma^{28}$ and

461    its target promoters reveal recognition of a composite −10 region, comprised of an "extended −10

462    motif" and a core-10 element. Molecular microbiology. 2009;72(4):830-43. doi: 10.1111/j.1365-

463    2958.2009.06691.x. PubMed PMID: PMC2756079.


464    35.    Stokes HW, O'Gorman DB, Recchia GD, Parsekhian M, Hall RM. Structure and function of 59-

465    base element recombination sites associated with mobile gene cassettes. Molecular microbiology.

466    1997;26(4):731-45. doi: 10.1046/j.1365-2958.1997.6091980.x.


467    36.    Boucher Y, Nesbø CL, Joss MJ, Robinson A, Mabbutt BC, Gillings MR, et al. Recovery and

468    evolutionary analysis of complete integron gene cassette arrays from Vibrio. BMC evolutionary

469    biology. 2006;6(1):3. doi: 10.1186/1471-2148-6-3.


470    37.    Li X, Shi L, Yang W, Li L, Yamasaki S. New array of *aacA4-catB3-dfrA1* gene cassettes and a

471    noncoding cassette from a class-1-integron-positive clinical strain of *Pseudomonas aeruginosa*.

472    Antimicrobial    agents    and    chemotherapy.    2006;50(6):2278-9.    Epub    2006/05/26.    doi:

473    10.1128/aac.01378-05. PubMed PMID: 16723608; PubMed Central PMCID: PMCPMC1479156.

21

474   38.    Michael CA, Labbate M. Gene cassette transcription in a large integron-associated array. BMC

475   genetics. 2010;11:82. Epub 2010/09/17. doi: 10.1186/1471-2156-11-82. PubMed PMID: 20843359;

476   PubMed Central PMCID: PMCPMC2945992.


477   39.    Papagiannitsis CC, Tzouvelekis LS, Miriagou V. Relative strengths of the class 1 integron

478   promoter hybrid 2 and the combinations of strong and hybrid 1 with an active p2 promoter.

479   Antimicrobial    agents    and    chemotherapy.    2009;53(1):277-80.    Epub    2008/11/13.    doi:

480   10.1128/aac.00912-08. PubMed PMID: 19001114; PubMed Central PMCID: PMCPMC2612174.


481   40.    Lévesque C, Brassard S, Lapointe J, Roy PH. Diversity and relative strength of tandem

482   promoters for the antibiotic-resistance genes of several integron. Gene. 1994;142(1):49-54. doi:

483   http://dx.doi.org/10.1016/0378-1119(94)90353-0.


484   41.    Singh SS, Singh N, Bonocora RP, Fitzgerald DM, Wade JT, Grainger DC. Widespread suppression

485   of intragenic transcription initiation by H-NS. Genes & development. 2014;28(3):214-9. Epub

486   2014/01/23. doi: 10.1101/gad.234336.113. PubMed PMID: 24449106; PubMed Central PMCID:

487   PMCPMC3923964.


488   42.    Jové T, Da Re S, Denis F, Mazel D, Ploy M-C. Inverse correlation between promoter strength

489   and    excision    activity    in    class    1    integrons.    PLOS    Genetics.    2010;6(1):e1000793.    doi:

490   10.1371/journal.pgen.1000793.


491   43.    Guerin E, Jove T, Tabesse A, Mazel D, Ploy MC. High-level gene cassette transcription prevents

492   integrase expression in class 1 integrons. Journal of bacteriology. 2011;193(20):5675-82. Epub

493   2011/08/23. doi: 10.1128/jb.05246-11. PubMed PMID: 21856858; PubMed Central PMCID:

494   PMCPMC3187215.

495    44.    Jacquier H, Zaoui C, Sanson-le Pors MJ, Mazel D, Bercot B. Translation regulation of integrons

496    gene cassette expression by the *attC* sites. Molecular microbiology. 2009;72(6):1475-86. Epub

497    2009/06/03. doi: 10.1111/j.1365-2958.2009.06736.x. PubMed PMID: 19486293.


498    45.    Collis CM, Hall RM. Expression of antibiotic resistance genes in the integrated cassettes of

499    integrons. Antimicrobial agents and chemotherapy. 1995;39(1):155-62. Epub 1995/01/01. PubMed

500    PMID: 7695299; PubMed Central PMCID: PMCPMC162502.


501    46.    Guerout AM, Iqbal N, Mine N, Ducos-Galand M, Van Melderen L, Mazel D. Characterization of

502    the *phd-doc* and *ccd* toxin-antitoxin cassettes from *Vibrio* superintegrons. Journal of bacteriology.

503    2013;195(10):2270-83. Epub 2013/03/12. doi: 10.1128/jb.01389-12. PubMed PMID: 23475970;

504    PubMed Central PMCID: PMCPMC3650543.


505    47.    Rowe-Magnus DA, Guerout AM, Biskri L, Bouige P, Mazel D. Comparative analysis of

506    superintegrons: engineering extensive genetic diversity in the Vibrionaceae. Genome research.

507    2003;13(3):428-42. Epub 2003/03/06. doi: 10.1101/gr.617103. PubMed PMID: 12618374; PubMed

508    Central PMCID: PMCPMC430272.


509    48.    Van Melderen L, Saavedra De Bast M. Bacterial toxin-antitoxin systems: more than selfish

510    entities? PLOS genetics. 2009;5(3):e1000437. Epub 2009/03/28. doi: 10.1371/journal.pgen.1000437.

511    PubMed PMID: 19325885; PubMed Central PMCID: PMCPMC2654758.


512    49.    Solovyev V, Salamov A. Automatic Annotation of Microbial Genomes and Metagenomic

513    Sequences. In: Li RW, editor. In Metagenomics and its Applications in Agriculture, Biomedicine and

514    Environmental Studies: Nova Science Publishers; 2011. p. 61-78.

23

515    50.    Seier-Petersen MA, Jasni A, Aarestrup FM, Vigre H, Mullany P, Roberts AP, et al. Effect of

516    subinhibitory concentrations of four commonly used biocides on the conjugative transfer of Tn*916* in

517    *Bacillus subtilis*. The Journal of antimicrobial chemotherapy. 2014;69(2):343-8. Epub 2013/10/05. doi:

518    10.1093/jac/dkt370. PubMed PMID: 24092655; PubMed Central PMCID: PMCPMC3886932.


519    51.    Swartzman A, Kapoor S, Graham AF, Meighen EA. A new *Vibrio fischeri lux* gene precedes a

520    bidirectional termination site for the *lux* operon. Journal of bacteriology. 1990;172(12):6797-802.

521    Epub 1990/12/01. PubMed PMID: 2254256; PubMed Central PMCID: PMCPMC210795.


522    52.    Dupuy B, Sonenshein AL. Regulated transcription of *Clostridium difficile* toxin genes. Molecular

523    microbiology. 1998;27(1):107-20. doi: 10.1046/j.1365-2958.1998.00663.x.


524    53.    Miller JH. Experiments in molecular genetics. NY: Cold Spring Harbor Laboratory Press: Cold

525    Spring Harbor; 1972.

526 **Tables**

527 **Table 1;** The putative promoters for ORF-less GCs and GCs with an ORF (SSU17 and MMB3) predicted

528 using BPROM.

| Clones | Strand | -10 box | -35 box | Score | | |
|--------|--------|---------|---------|-------|------|------------------------------|
| | | | | -10 box | -35 box | Linear discriminant function (LDF)* (Overall score) |
| TMB4 | + | AGGTATAAT | ATAAGA | 89 | -10 | 9.78 |
| | - | CATTATTTT | TTGACA | 41 | 66 | 7.60 |
| SSU9 | + | AATTATAAT | TAAAAA | 74 | 0 | 7.04 |
| | - | TAGTATAAT | TTTATT | 80 | 34 | 7.11 |
| MMU2 | + | AATTATAAT | TTAAAA | 74 | 37 | 8.36 |
| | - | TAGTATAAT | TTTATT | 80 | 34 | 8.90 |
| MMU11 | + | ATGTAAAAT | TTGCTG | 75 | 47 | 11.34 |
| | + | AACTATACT | AGGAAA | 59 | -7 | 5.99 |
| | - | AAATAAAAT | TTTTCA | 56 | 34 | 6.96 |
| | - | CTATAAATT | TTTCAA | 44 | 36 | 3.24 |
| MMU19 | + | AGGTATAAT | TAGAAA | 89 | 23 | 9.07 |
| | + | TTGAAAAAT | TTGCGG | 44 | 32 | 3.43 |
| | - | TATTATAAT | TTTCCT | 79 | 37 | 9.10 |
| MMU23 | + | AATTATAAT | TAAAAG | 74 | -6 | 9.84 |
| | + | TTTTATTAT | TTGATG | 72 | 52 | 6.05 |
| | - | TATTATAAT | TTTCCT | 79 | 37 | 8.66 |
| | - | TAGTATAAT | TTTATT | 80 | 34 | 8.05 |
| MMB2 | + | AATTATAAT | TATAAG | 74 | -2 | 8.71 |
| | + | TATTATAAT | TTGATG | 79 | 52 | 7.88 |
| | - | TATTATAAT | TTTCCT | 79 | 37 | 9.10 |

25

| Clones | Strand | -10 box | -35 box | Score | | |
|--------|--------|---------|---------|-------|-------|------------------------------------------------------|
| | | | | -10 box | -35 box | Linear discriminant function (LDF)* (Overall score) |
| | - | TATTATAAT | TTTATT | 79 | 34 | 8.84 |
| MMB5 | + | AATTATAAT | TTAAAA | 74 | 37 | 8.36 |
| | - | TAGTATAAT | TTTATT | 80 | 34 | 7.95 |
| MMB20 | + | AATTATAAT | TAAAAG | 74 | -6 | 9.09 |
| | - | TATTATAAT | TTTCCT | 79 | 37 | 9.10 |
| MMB32 | + | TATTATAAT | TTGATG | 79 | 52 | 6.28 |
| | + | AGATATAAA | GTGTAA | 39 | 14 | 4.84 |
| | - | TATTATAAT | TTGATT | 79 | 53 | 6.61 |
| | - | TTTTATTTT | TTAAAA | 52 | 37 | 5.11 |
| MMB36 | + | AATTATAAT | TTAAAA | 74 | 37 | 6.94 |
| | + | TATTATAAT | TTGATG | 79 | 52 | 6.45 |
| | - | TATTATAAT | TTTATT | 79 | 34 | 7.44 |
| | - | TTTTAAAAT | TTGACT | 79 | 61 | 6.13 |
| MMB37 | + | AATTATAAT | TAAAAG | 74 | -6 | 9.11 |
| | + | TTATATAAT | TTGATG | 75 | 52 | 8.55 |
| | - | TAGTATTAT | TTTATT | 66 | 34 | 10.48 |
| | - | TATTATAAT | TTTCCT | 79 | 37 | 9.10 |
| SSU17 | + | CTTTATAAT | ATGAAT | 82 | 25 | 7.80 |
| | + | TGATAAAAT | GTGAAA | 75 | 27 | 4.62 |
| | - | TGATATAAT | TTTATT | 82 | 34 | 9.34 |
| | - | TGATTAGAT | TTTATG | 21 | 33 | 5.10 |
| MMB3 | + | CTGTATATT | TTGATA | 63 | 58 | 6.74 |
| | + | ATTTATGAT | ATGAAA | 65 | 30 | 5.18 |
| | - | ATGTATTGT | TTGATG | 44 | 52 | 6.64 |

| Clones | Strand | -10 box | -35 box | Score | | |
|---|---|---|---|---|---|---|
| | | | | -10 box | -35 box | Linear discriminant function (LDF)* (Overall score) |
| | - | GCATATAAT | TTCTCT | 65 | 28 | 4.75 |

* The LDF takes into account motifs found in promoters: -10 and -35 boxes, a distance between -10 and -35 boxes, and frequencies of certain nucleotides represented in transcription start sites. It can be approximated as log(<likelihood of a site being promoter>/<likelihood of a site not being promoter>) [49].
** The selected samples for the enzymatic assay are highlighted in yellow.

529

530 **Table 2.** Characterisation of the human oral integron GCs containing promoter sequences detected by pBiDiPD.

| Gene cassettes | Primer pair | Cassette Size (bp) | Orientation* | BlastN | | BlastX | | | | Promoter activity | | Accesion number |
| | | | | Closest homologue | Percentage identity (%)/ Coverage (%) | Closest homologue | ORF size (bp) | Percentage identity (%)/Coverage (%) | Accession number of the homologous proteins (BlastX) | Sense Strand (*gusA*) | Antisense strand (*lacZ*) | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| SSU-Pro-7 | SUPA3-SUPA4 | 1001 | | SSU22 | 98/95 | Prevent-host-death protein (Phd_YefM antitoxin superfamily) [*Treponema vincentii*] | 264 | 97/100 | WP_006188308.1 | Y | N | MH536747 |
| | | | | | | XRE family transcriptional regulator [*Treponema vincentii*] | 441 | 98/100 | WP_006188306.1 | | | |
| SSU-Pro-9 | SUPA3-SUPA4 | 834 | | MMB3 | 98/99 | Hypothetical protein (antitoxin, ribbon-helix-helix domain protein) [*Treponema putidum*] | 246 | 67/100 | WP_044978234.1 | Y | N | MH536748 |
| | | | | | | Twitching motility protein PilT (PIN toxin domain) [*Treponema putidum*] | 414 | 71/100 | AIN93467.1 | | | |
| SSU-Pro-13 | SUPA3-SUPA4 | 855 | | MMB39 | 98/99 | Toxin RelE [*Treponema medium*] | 357 | 95/100 | WP_016522532.1 | Y | N | MH536749 |
| | | | | | | Transcriptional regulator (Antitoxin, XRE family) [*Treponema medium*] | 330 | 95/100 | WP_016522533.1 | | | |
| SSU-Pro-16 | SUPA3-SUPA4 | 925 | | SSU28 | 98/100 | AbrB/MazE/SpoVT family DNA-binding domain-containing protein (Antitoxin) [*Treponema putidum*] | 231 | 96/100 | WP_044979179.1 | Y | N | MH536750 |

28

| Gene cassettes | Primer pair | Cassette Size (bp) | Orientation* | BlastN | | BlastX | | | | Promoter activity | | Accesion number |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | | | Closest homologue | Percentage identity (%)/ Coverage (%) | Closest homologue | ORF size (bp) | Percentage identity (%)/Coverage (%) | Accession number of the homologous proteins (BlastX) | Sense Strand (*gusA*) | Antisense strand (*lacZ*) | |
| | | | | | | Endoribonuclease MazF (Toxin) [*Treponema denticola*] | 336 | 99/100 | WP_010694033.1 | | | |
| SSU-Pro-20 | SUPA3-SUPA4 | 1263 | | MMU28 | 77/42 | Prevent-host-death protein (Phd_YefM antitoxin superfamily) [*Treponema sp.* JC4] | 249 | 75/88.3 | WP_009103386.1 | Y | N | MH536751 |
| | | | | | | Plasmid stabilization protein (ParE toxin superfamily) [*Treponema sp. JC4]* | 147 | 57/43.8 | WP_009104800.1 | | | |
| SSU-Pro-24 | SUPA3-SUPA4 | 425 | | SSU9 | 99/100 | - | - | - | - | Y | Y | MH536752 |
| SSU-Pro-27 | SUPA3-SUPA4 | 753 | | *Treponema putidum* strain OMZ 758 | 93/100 | BrnT family toxin [*Treponema sp.*] | 273 | 99/100 | WP_002666393.1 | Y | N | MH536753 |
| | | | | | | CopG family transcriptional regulator (BrnA antitoxin) [*Treponema denticola*] | 288 | 99/100 | WP_044909778.1 | | | |
| SSU-Pro-32 | SUPA3-SUPA4 | 972 | | No significant similarity found. | - | RelE/ParE family toxin [*Treponema denticola*] | 354 | 98/100 | WP_002683264.1 | Y | N | MH536754 |
| | | | | | | XRE family transcriptional regulator [*Treponema denticola*] | 273 | 100/100 | WP_002683262.1 | | | |

| Gene cassettes | Primer pair | Cassette Size (bp) | Orientation* | BlastN | | BlastX | | | | Promoter activity | | Accesion number |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | | | Closest homologue | Percentage identity (%)/ Coverage (%) | Closest homologue | ORF size (bp) | Percentage identity (%)/Coverage (%) | Accession number of the homologous proteins (BlastX) | Sense Strand (*gusA*) | Antisense strand (*lacZ*) | |
| SSU-Pro-34 | SUPA3-SUPA4 | 832 |  | SSU5 | 99/100 | Hypothetical protein (antitoxin, ribbon-helix-helix domain protein) [*Treponema putidum*] | 246 | 67/100 | WP_044978234.1 | Y | N | MH536755 |
| | | | | | | PIN domain-containing protein [*Treponema putidum*] | 414 | 71/100 | WP_044978236.1 | | | |
| SSU-Pro-39 | SUPA3-SUPA4 | 1137 |  | MMU25 | 99/99 | Hypothetical protein [uncultured bacterium] | 462 | 99/100 | ANC55535.1 | Y | N | MH536756 |
| | | | | | | Hypothetical protein [*Treponema maltophilum*] | 213 | 88/100 | WP_016526060.1 | | | |
| | | | | | | PemK/MazF family toxin [*Fibrobacter sp. UWCM*] | 342 | 80/100 | WP_022932935.1 | | | |
| SSU-Pro-46 | SUPA3-SUPA4 | 971 |  | No significant similarity found | - | Hypothetical protein [*Treponema socranskii*] | 267 | 80/100 | WP_021329686.1 | Y | N | MH536757 |
| | | | | | | Hypothetical protein [*Treponema socranskii*] | 228 | 84/100 | WP_021329641.1 | | | |
| | | | | | | DUF4160 domain-containing protein [*Treponema sp. C6A8*] | 276 | 67/100 | WP_027729334.1 | | | |
| SSU-Pro-65 | SUPA3-SUPA4 | 811 |  | *Treponema sp.* OMZ 838 | 91/21 | AbrB/MazE/SpoVT family DNA-binding domain-containing protein (Antitoxin) [*Treponema denticola*] | 228 | 93/100 | WP_010693782.1 | Y | N | MH536758 |

| Gene cassettes | Primer pair | Cassette Size (bp) | Orientation* | BlastN | | BlastX | | | | Promoter activity | | Accesion number |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | | | Closest homologue | Percentage identity (%)/ Coverage (%) | Closest homologue | ORF size (bp) | Percentage identity (%)/Coverage (%) | Accession number of the homologous proteins (BlastX) | Sense Strand (*gusA*) | Antisense strand (*lacZ*) | |
| | | | | | | VapC family toxin [*Treponema denticola*] | 402 | 93/100 | WP_010693784.1 | | | |
| MMU-Pro-4 | MARS5-MARS2 | 520 | | MMU2 | 99/100 | - | - | - | - | Y | Y | MH536759 |
| MMU-Pro-5 | MARS5-MARS2 | 983 | | *Treponema putidum* strain OMZ 758 | 94/78 | Prevent-host-death protein (Phd_YefM antitoxin superfamily) [*Treponema denticola*] | 240 | 98/98.8 | WP_002669519.1 | Y | Y | MH536760 |
| | | | | | | RelE/StbE family addiction module toxin [*Treponema denticola*] | 318 | 94/100 | WP_002688980.1 | | | |
| MMU-Pro-6 | MARS5-MARS2 | 737 | | MMB36 | 86/100 | - | - | - | - | N | Y | MH536761 |
| MMU-Pro-18 | MARS5-MARS2 | 634 | | MMB37 | 95/100 | - | - | - | - | Y | N | MH536762 |
| MMU-Pro-22 | MARS5-MARS2 | 431 | | MMU19 | 91/100 | - | - | - | - | Y | Y | MH536763 |
| MMU-Pro-24 | MARS5-MARS2 | 904 | | No significant similarity found | - | Universal stress protein [*Marinobacter sp.*] | 348 | 30/54.1 | WP_008177208.1 | Y | Y | MH536764 |
| | | | | | | Hypothetical protein [*Methylobacter tundripaludum*] | 213 | 79/100 | WP_031438379.1 | | | |
| | | | | | | Prevent-host-death protein [*Treponema pedis*] | 84 | 76/27.8 | WP_024469914.1 | | | |

| Gene cassettes | Primer pair | Cassette Size (bp) | Orientation* | BlastN | | BlastX | | | | Promoter activity | | Accesion number |
| | | | | Closest homologue | Percentage identity (%)/ Coverage (%) | Closest homologue | ORF size (bp) | Percentage identity (%)/Coverage (%) | Accession number of the homologous proteins (BlastX) | Sense Strand (*gusA*) | Antisense strand (*lacZ*) | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| MMU-Pro-31 | MARS5-MARS2 | 574 | | MMB5 | 88/70 | - | - | - | - | Y | N | MH536765 |
| MMU-Pro-48 | MARS5-MARS2 | 817 | | *Treponema sp.* OMZ 838 | 91/25 | AbrB/MazE/SpoVT family DNA-binding domain-containing protein [*Treponema denticola*] | 228 | 93/100 | WP_010693782.1 | Y | N | MH536766 |
| | | | | | | VapC family toxin [*Treponema denticola*] | 402 | 93/100 | WP_010693784.1 | | | |
| MMU-Pro-53 | MARS5-MARS2 | 430 | | No significant similarity found | - | - | - | - | - | Y | Y | MH536767 |
| MMU-Pro-63 | MARS5-MARS2 | 927 | | SSU8 | 99/99 | Hypothetical protein [*Treponema denticola*] | 531 | 98/93.7 | WP_002692239.1 | N | Y | MH536768 |
| MMU-Pro-65 | MARS5-MARS2 | 896 | | MMU27 | 99/100 | Hypothetical protein [uncultured bacterium] | 399 | 99/84.2 | ANC55539.1 | N | Y | MH536769 |
| | | | | | | Hypothetical protein [uncultured bacterium] | 357 | 99/100 | ANC55540.1 | | | |

* The orange half circles and green arrow boxes are representing *attC* sites and ORFs, respectively.

** The GCs found in this study are highlighted in yellow. Those not highlighted were also detected in Tansirichaiya et al. (2016) [28].

531

**Figure Legends**

**Figure 1:** A generalised structure of (A) usual integrons and (B) unusual, or reverse integrons. The green arrows indicate the primer binding sites on the unusual integron structure of *T. denticola*. The grey and blue open arrowed boxes represent integrase gene (*intI*) and the open reading frames (ORFs), respectively, pointing in the direction of transcription. The promoters, $P_{intI}$ and $P_C$, were represented by black arrows. The recombination sites, *attI* and *attC*, were represented by yellow and orange circles, respectively.

**Figure 2:** The structure of pCC1BAC-*lacZα-gusA* plasmid. The green, blue and orange open arrowed boxes represent *lacZα, gusA* and chloramphenicol resistance gene, respectively, pointing in the direction of transcription. The black lines indicate the position of restriction sites on the plasmid. The red circles indicate bidirectional transcriptional terminators.

**Figure 3:** The promoter activity from pCC1BAC-*lacZα*-GC-*gusA* constructs estimated by β-glucuronidase enzyme assays. Error bars indicate the standard errors of the means from three replicates. The asterisks (*) indicate the constructs were statistically significantly different from the negative control group (pCC1BAC-*lacZα-gusA*) with the *p*-value <0.05 by using ordinary one-way ANOVA followed by Dunnett's multiple comparison tests.

**Figure 4:** The detection of the integron GCs by using pBiDiPD. A.) Blue-white screening to detect for the clones with promoter activity on the antisense strand, B.) Exposing the colonies under the UV light to detect clones with promoter activity on the sense strand. The positive (+) and negative (-) colonies were the *E. coli* containing pCC1BAC-*lacZα*-TMB4-Pc-*gusA* (with experimentally proven promoter activities on either strand of DNA and pCC1BAC-*lacZα-gusA* (no promoter activity), respectively

**Figure 5:** The proposed genetic clutch. (A) When a promoter-containing GC inserts into the first position, it can act as a genetic clutch by disengaging the original first GC (blue arrow) from PC promoter and replaced with the one on promoter GC. When a new GC (green arrow) inserts, it can be

556    expressed by PC promoter, while the blue GC is expressed by promoter-containing GC and PC

557    promoter. (B) The expression level of gene cassettes with and without a genetic clutch.  The estimated

558    levels of expression of the blue ORF in i.) the first, ii.) the second and iii.) the third position were shown

559    in the bar chart. The solid bars represent the situation when promoter-containing GC was inserted

560    upstream of the blue GC, while the gridded bars represent the situation when no promoter-containing

561    GC was inserted. The asterisks indicate the experimentally verified expression level, suggested by the

562    results in Figure 3 (TMB4 PC and TMB4 PC+GC).  The expression of the blue ORF was hypothesised to

563    be decreased when more GCs are inserted without the presence of a promoter-containing GC as a

564    genetic clutch (gridded bars), based on the data from previous study [45].

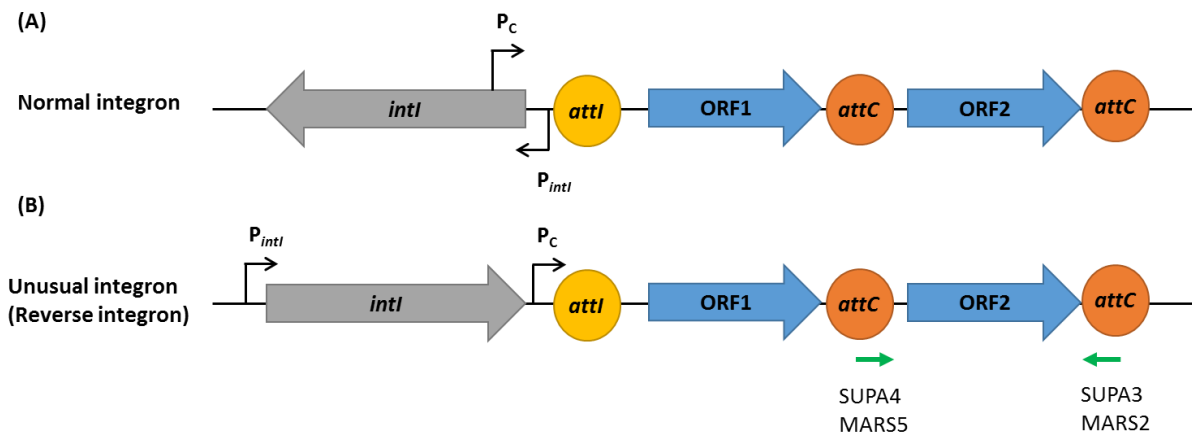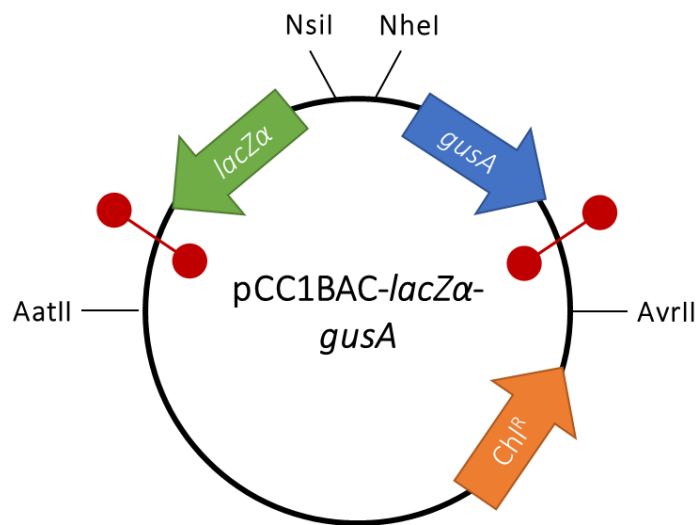565 **Figure 1**



566

567     **Figure 2.**



568

**Figure 3.**



The concentration of β-glucuronidase enzyme in Miller unit (U) from pCC1BAC-*lacZ-gusA* constructs

571     **Figure 4.**
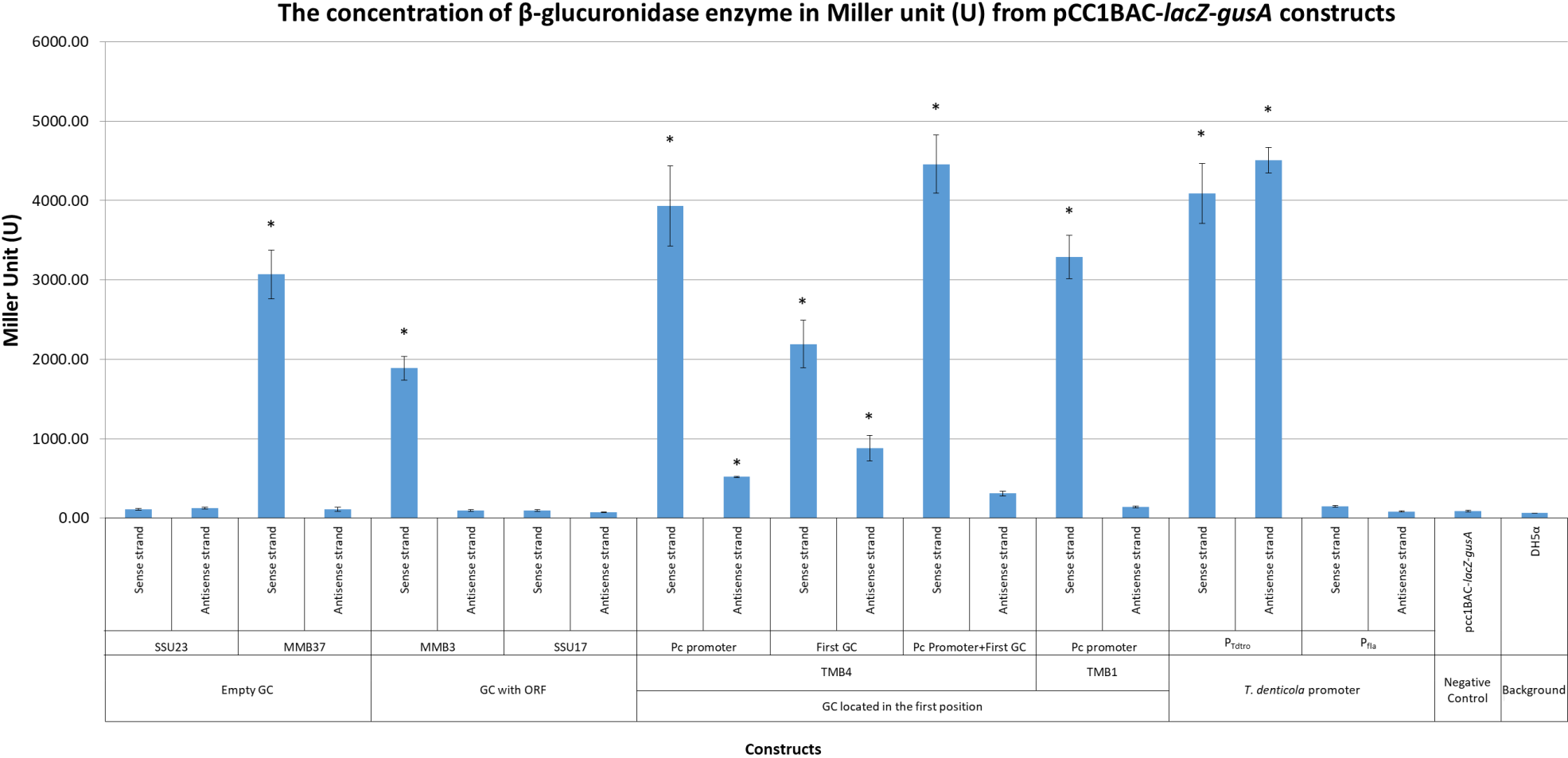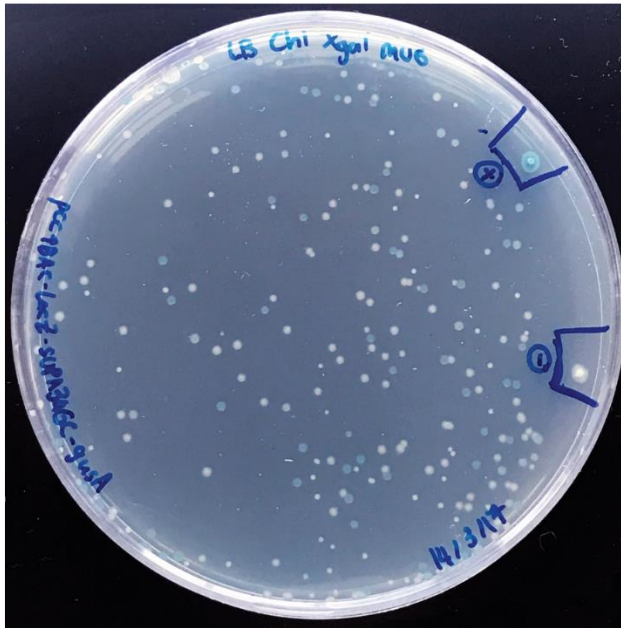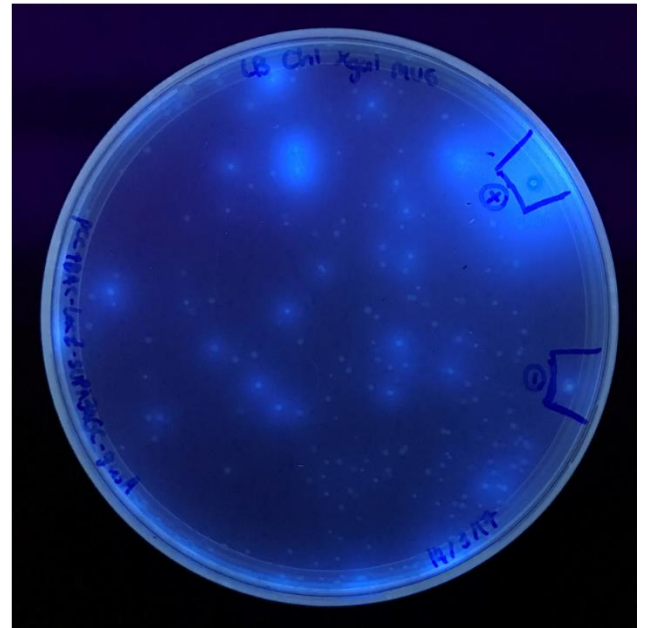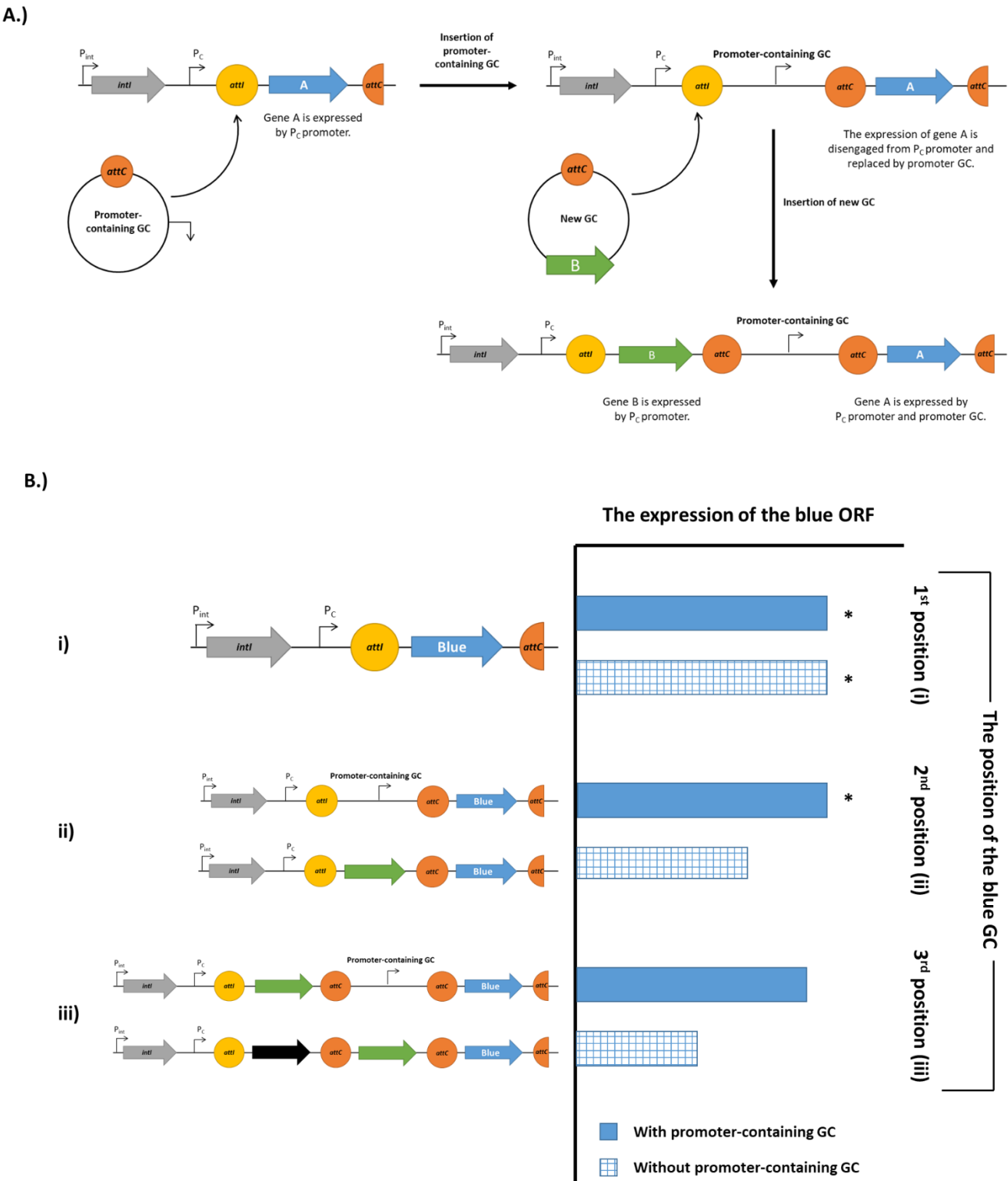
(A)                                        (B)



572

573    **Figure 5.**



574