

1 **Supporting Information**

2

3 *Analysing the effect of type IIs restriction cloning scars on promoter sequence*
4 *activity.*

5

6 The use of type IIs restriction cloning to join DNA parts resulted in the
7 insertion of a 4 bp scar sequence (ACCT) between the putative Distal Regulatory
8 Sequence (DRS) and RBS regions of the bioinformatically identified 100 bp
9 *Geobacillus* promoters. Previous studies have shown that novel cryptic functionality
10 can potentially arise through the fusion of previously characterised genetic
11 parts^{45,47,84}. In particular, any alterations to the mRNA secondary structure arising
12 from scar sequences located between the RBS and CDS can potentially negatively
13 impact the efficiency of translation initiation⁸⁵, thereby altering protein production. To
14 assess the impact of the type IIS scar sequence on the activity of *Geobacillus*
15 promoters, a subset of 24 bioinformatically identified sequences were characterised
16 upstream of GFP both with and without scars (Supporting Figure 2).

17

18 For the majority of the analysed *promoter::GFP* constructs, the insertion of
19 scar sequences between DRS and RBS had no statistically significant impact on
20 GFP fluorescence (Supporting Figure 2A). Scar sequences significantly changed the
21 regulatory activity of 4 sequences. Activity was significantly decreased in 3 promoters
22 (P_34, adjusted P value = 0.000; P_18, adjusted P value = 0.022; P_56, adjusted P
23 value = 0.024) and increased in 1 sequence (P_79, adjusted P value = 0.009) when
24 scar sequences were inserted.

25

26 The observed differences in fluorescence between scarred and un-scarred
27 *promoter::GFP* constructs were hypothesised to be the result of alterations to mRNA
28 secondary structure. To test this hypothesis, the free-energy associated with mRNA
29 folding was calculated using the mFold zipfold server⁸⁶. The default settings were
30 used, with RNA 2.3 energy rules. Folding energies were returned in kcal/mol. The
31 temperature at which folding was simulated was set to 60 °C, to reflect the
32 temperature at which *Geobacillus* cultures were incubated. The sequence window in
33 which folding was analysed stretched from -20 to +20, relative to the adenine residue
34 of the GFP start codon.

35

36 3 of the promoters for which fluorescence output was significantly altered by
37 the inclusion of scar sequences also showed the greatest changes in mRNA

38 secondary structure free energy (Supporting Figure 2B). *P_79::GFP*, for which GFP
39 fluorescence was statistically significantly increased by the inclusion of scar
40 sequences, showed the greatest relaxation of mRNA secondary structure out of the
41 24 analysed sequences. *P_18::GFP* also showed an increase in fluorescence output
42 after scar insertion and relaxation of the mRNA secondary structure. Conversely,
43 *P_34::GFP* and *P_56::GFP* showed a statistically significant decrease in
44 fluorescence and an increase in mRNA secondary structure.

45

46 The correlation between relaxed mRNA secondary structure at the RBS-CDS
47 junction was corroborated by the literature^{85,87-89}. The presence of significant
48 secondary structure surrounding the RBS is hypothesised to negatively impact on the
49 ability of an mRNA transcript to sequester ribosomes, thereby reducing the rate at
50 which transcripts are translated, although the strength of this correlation may be
51 sensitive to genetic context and cellular concentrations of amino acids and
52 tRNAs^{90,91}.

53

54 *Artificial Neural Network sequence-function modelling*

55

56 ANNs were fit using the multilayer perceptron algorithm of JMP software. All
57 ANNs used a single hidden layer, as a single hidden layer had previously been
58 shown to be sufficiently complex to describe 224 bp *E. coli* promoters⁵⁷. Furthermore,
59 the increase in model complexity arising from additional hidden layers carried the risk
60 of increasing model variance and overfitting the model to the training data⁹².

61

62 A screening design was used to identify which of the ANN parameters were
63 having the greatest impact on model performance. The ANN parameters that were
64 included in the screening design are summarised in Supporting Table 4. The
65 specified parameters were combined in a full-factorial manner, which resulted in 81
66 ANN architectures being specified. As the result to which an individual ANN
67 converges is dependent on the random seed used to generate the starting network
68 weights⁹², each of the 81 architectures was fit 500 times, with each fit using a unique,
69 specified random seed. All models used the squared penalty method. For each of the
70 81 ANN architectures, the single ANN that returned the highest R² value when
71 applied to the test set was identified.

72

73 The results of the screening experiment were subjected to statistical analysis.
74 A standard least squares model with effect screening emphasis showed that both the

75 number of sequence positions included in the ANN and the personality of the
76 activation function were having a significant effect on model performance at the 0.05
77 significance level (number of sequence positions LogWorth = 7.632, P = 0.000,
78 activation function personality LogWorth = 6.298, P = 0.000). The number of nodes in
79 the hidden layer and the K value used for cross validation did not have a significant
80 effect (P = 0.582 and P = 0.671, respectively).

81

82 The results of the screening design were also analysed using a PLS model,
83 using the NIPALS algorithm and KFold cross validation where K = 7. The R² value
84 returned when the models were applied to the test set was used as the y variable,
85 and the parameters summarised in Supporting Table 4 were used as x variables. The
86 resulting model extracted a single latent variable and was capable of explaining
87 12.5% of the cumulative variation in X and 52.525% of the cumulative variation in Y.

88

89 In PLS models, the Variable Importance in Projection (VIP) statistic can be
90 used to determine the importance of a given variable in determining model output⁵⁴.
91 Variables with a VIP score above the threshold value of 0.8 are commonly accepted
92 to have a significance impact on model output⁹³. In this instance, 4 variables
93 exceeded the threshold VIP value (Supporting Figure 10). The Linear and TanH
94 activation functions were shown to have a significant effect on PLS model output, as
95 did the effect of including 20 or 100 sequence positions in the ANNs. Analysis of the
96 model coefficients showed that the TanH activation function was predicted to
97 positively contribute to ANN predictive power, whereas the linear activation function
98 was predicted to negatively impact predictive power, suggesting that linear modelling
99 provided an inadequately sophisticated mathematical abstraction to adequately
100 describe the promoter design space.

101

102 Taken together, the results of the screening design suggested that the TanH
103 activation function and the number of promoter sequence positions included in the
104 ANN were likely to contribute most to predictive power. A second iteration of ANN
105 design was therefore undertaken using only TanH activation functions. The number
106 of promoter nucleotide sequence positions that were included in the models ranged
107 from 10 to 100, increasing in increments of 10 positions, and the number of nodes in
108 the hidden layer was between 3 and 15, increasing in increments of 2. The
109 parameters were combined in a full factorial manner, and each of the 70 resulting
110 architectures was fit 1,000 times, with each run using a unique, defined random

111 seed. 500 runs of each ANN architecture used K = 4 KFold cross validation, and the
112 remaining 500 runs used K = 5. All ANNs used the squared penalty method.

113

114 The individual ANN that returned the highest R² value when applied to the
115 test set was identified for each of the 70 architectures. The best performing ANN that
116 was obtained returned an R² value of 0.9304 when applied to the test set, and
117 modelled GFP fluorescence as a function of 20 promoter nucleotide sequence
118 positions.

119

120 In the case of ANNs that modelled GFP fluorescence as a function of
121 complete 104 bp promoter sequences, the optimum network obtained used 9 nodes
122 in the hidden layer and returned an R² value of 0.8199 when applied to the test set.
123 To test if ANNs of increased complexity could better model complete promoter
124 sequences, additional single layer ANNs were trained using 17, 19 or 21 nodes.
125 1,000 ANNs were trained for each of the defined architectures, but none performed
126 better than the 9-node model when applied to the test set, either in terms of R² or
127 RASE (Supporting Figure 11).

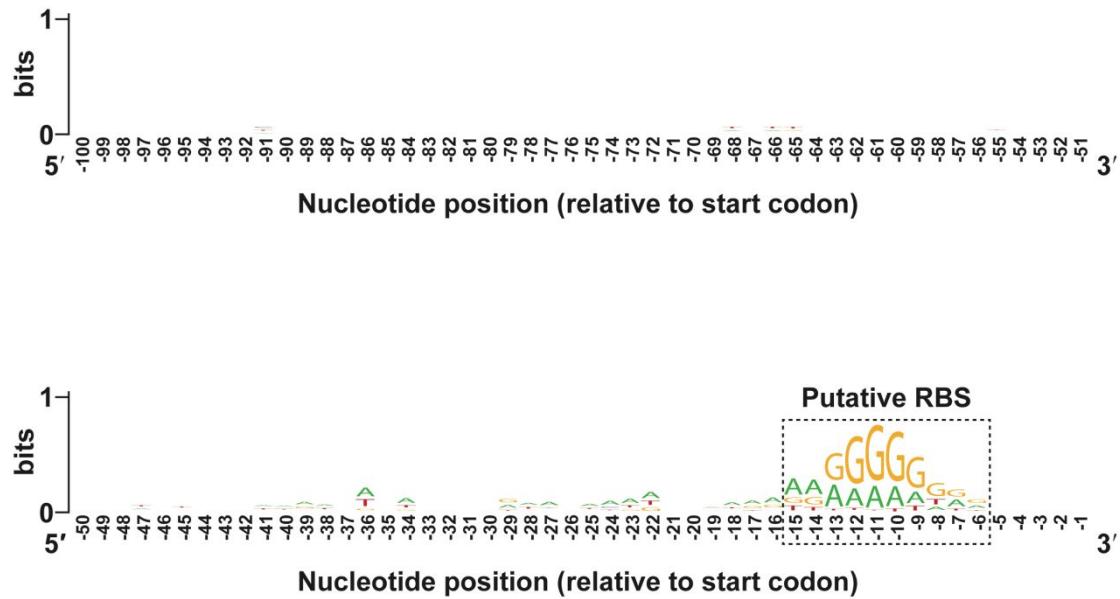
128

129 Although single “best” performing models are most often presented in support
130 of conclusions, this approach falsely assumes that only 1 model explains the
131 response surface of interest⁹⁴. Model ensembling aims to avoid this issue by
132 combining the outputs of multiple individual models trained on the same task⁹⁵.
133 Theoretically, individual constituents of the ensemble can counteract deficiencies in
134 other members, thereby improving the generalisation performance of the ensemble
135 as compared to the individual base models⁹⁵. The simplest method for ensembling
136 multiple quantitative models is to calculate the arithmetic mean of the predicted
137 values for a given set of x variables (e.g. the average predicted activity for a given
138 promoter nucleotide sequence) and subsequently assess the accuracy of the resulting
139 predictions.

140

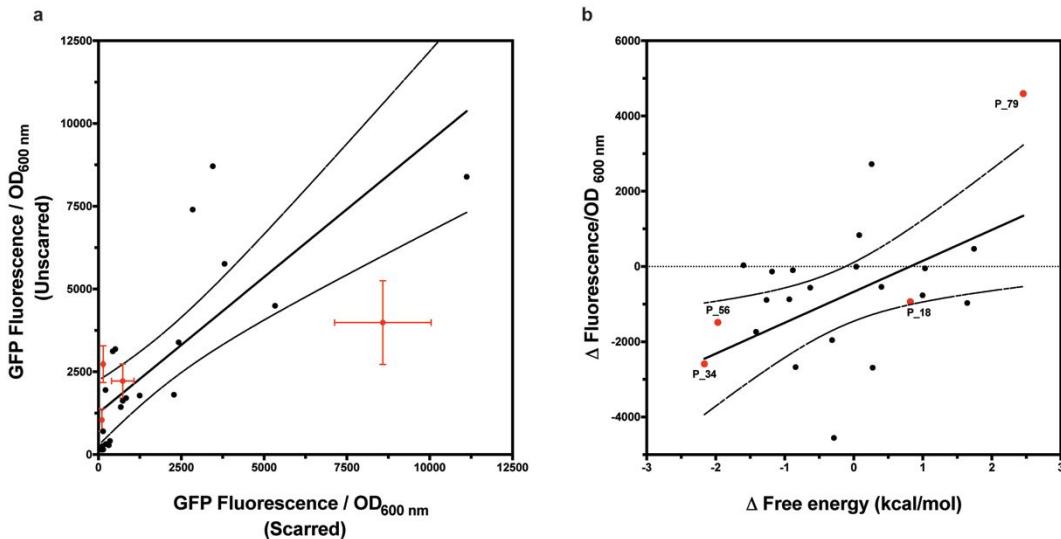
141 Model ensembling was therefore applied in attempt to increase predictive
142 power. For each number of promoter nucleotide sequence positions that were
143 modelled, the network architecture that returned the single ANN with the highest test
144 set R² value was identified. The top 10 highest performing ANNs with the chosen
145 architecture (as judged by the test set R² value) were used to create 9 progressively
146 larger ensembles, with ANNs being included in descending order of their test set R²
147 value. In each ensemble, the mean predicted value for each of the 10 promoter

148 sequences in the test set was calculated and the accuracy of the resulting predictions
149 was assessed. The ensemble that was obtained with the lowest predictive error
150 contained 2 constituent ANNs, each of which modelled promoter activity as a function
151 of 20 sequence positions using 5 nodes in the hidden layer and TanH activation
152 functions. The optimal model returned an R^2 value of 0.9746 when applied to the
153 training and validation sets and an R^2 value of 0.9691 when applied to the test set.



155
 156 **Supporting Figure 1: Visualisation of a sequence alignment of 100 putative**
 157 **promoters used to identify the putative location of the Ribosome Binding Site**
 158 **(RBS).**

159
 160 The dashed box highlights the location of the putative RBS. The overall height of
 161 individual stacks indicates the degree of sequence conservation at a given position,
 162 and the height of the nucleotide symbols indicates the conservation of each nucleic
 163 acid at that position. Position numbering is relative to the start codon of the
 164 downstream coding sequence. Sequences were aligned and visualised using
 165 WebLogo version 2.8.2⁹⁶.
 166



167
168

169 **Supporting Figure 2: The effect of cloning scar sequences on promoter
170 activity.**

171

172 Fluorescence and absorbance measurements after 24 h incubation in 96-well plate
173 format. In both instances, points that are coloured red represent promoters for which
174 GFP fluorescence was statistically significantly changed when the cloning scar was
175 inserted in the DNA sequence. Significance was determined using multiple t-tests,
176 using the Holm-Šidák method to correct for multiple comparisons and a significance
177 level of 0.05.

178

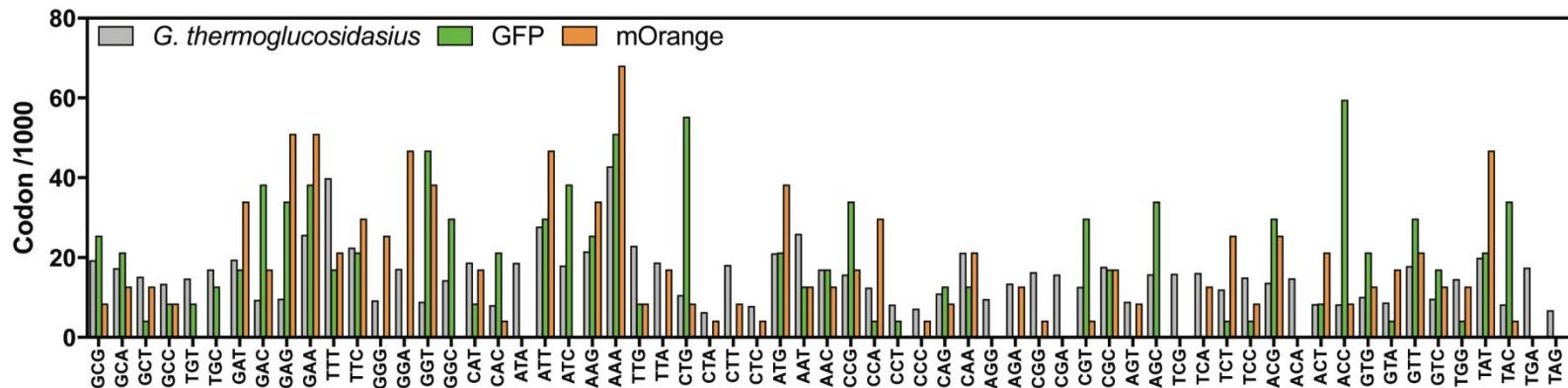
179 A) Fluorescence activity levels of scarred and unscarred promoters. Points represent
180 the mean GFP output of each promoter, from $3 \leq n \leq 9$ starter cultures arising from
181 independent transformation events. Standard deviation error bars are shown only on
182 the statistically significant points. The solid line represents a linear regression of the
183 data, with 95% confidence limits represented by dashed lines. The linear regression
184 had an R^2 value of 0.5216.

185

186 B) Comparing the change in GFP fluorescence and the change in free energy of the
187 mRNA secondary structure of the *promoter::GFP* fusions when cloning scar
188 sequences were inserted. Free energies were calculated using the mFold zipfold
189 server⁸⁶, using default settings and RNA 2.3 energy rules. The sequence window for
190 which secondary structure was calculated stretched from -20 to +20 relative to the
191 adenine residue of the GFP start codon. The temperature at which folding was
192 simulated was set to 60 °C, to reflect the temperature at which *G. thermoglucosidasius*
193 cultures were incubated. The solid line represents a linear

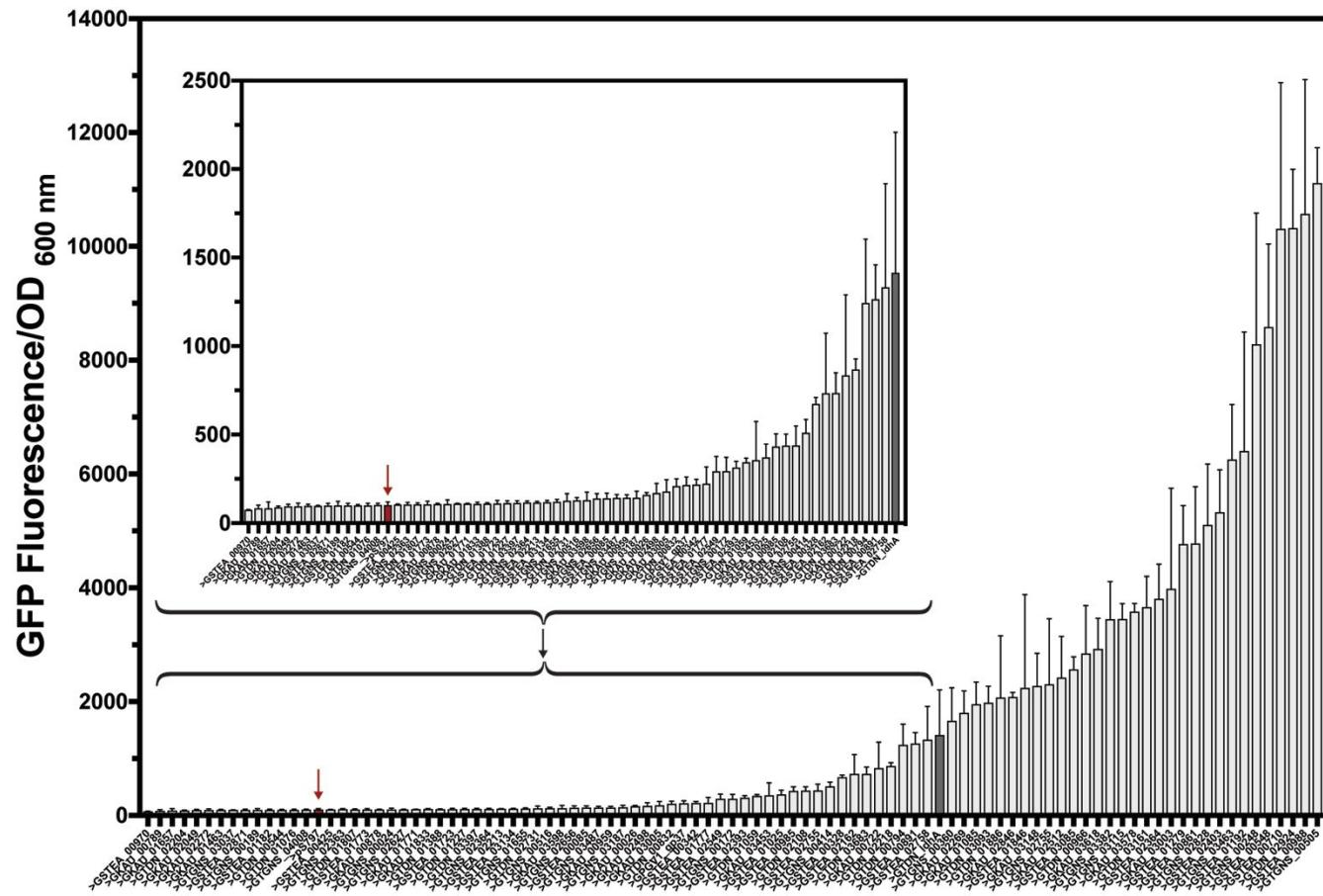
194 regression of the data, with 95% confidence limits represented by dashed lines. The
195 linear regression had an R^2 value of 0.2407.

196



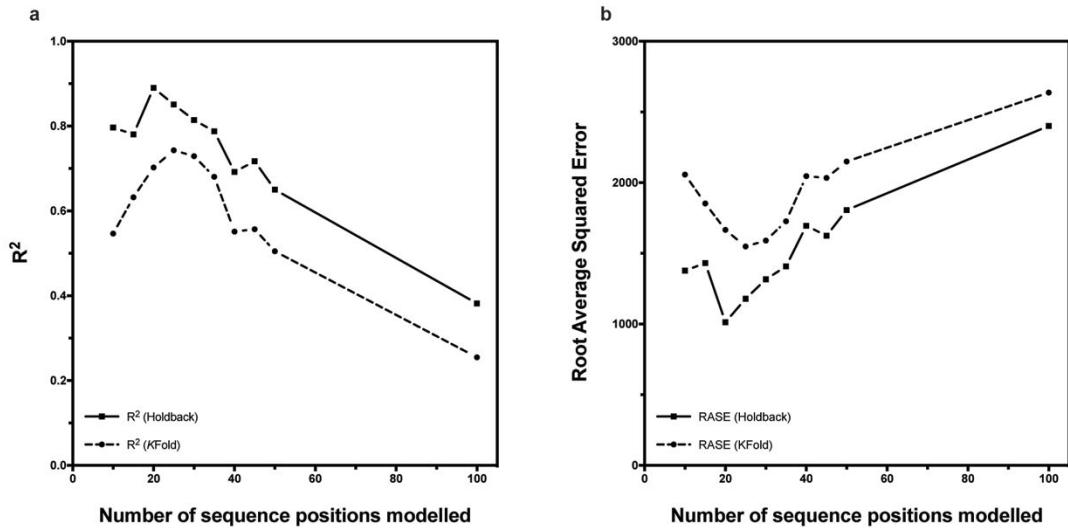
197 Supplementary Figure 3: A comparison of codon usage in the GFP and mOrange reporter sequences.

198



Supporting Figure 4: Activity levels of putative promoter sequences characterised upstream of GFP in *G. thermoglucosidasius*.

Fluorescence and absorbance measurements after 24 h incubation. The positive control, the *G. thermodenitrificans* *IdhA* promoter, is represented by the dark grey bar. The negative control, *G. thermoglucosidasius* transformed with the empty pS797 vector, is shown in red. Standard deviation error bars shown, unless hidden by the bar. Bars represent the mean of n = 3 starter cultures, arising from independent transformation events, with the following exceptions: n = 4 (GSTEA_02871, GTGNS_02828); n = 5 (GSTEA_02755); n = 6 (GTDN_00832, GTDN_01227, GSTEA_02755); n = 11 (pos. control); n = 14 (neg. control).



223

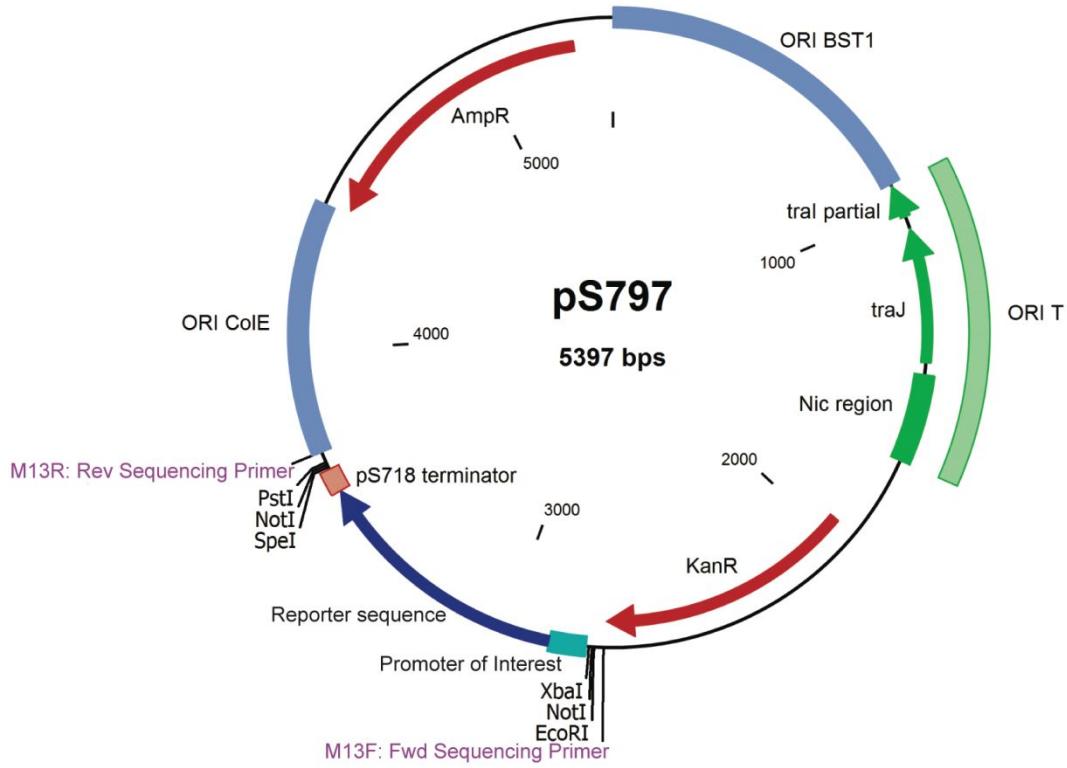
224 **Supporting Figure 5: R^2 (a) and Root Average Squared Error (RASE, b) values**
 225 **returned by PLS sequence-function models when applied to a test data set.**

226

227 PLS models were trained using the NIPALS algorithm, with a maximum of 10 latent
 228 variables permitted per model. For each number of sequence positions modelled, 4
 229 PLS regressions were fit using KFold cross validation with K values of 4, 5, 7 and 10.
 230 2 models were also fit at each position using holdback cross validation, with 20% or
 231 33% of the training data set withheld to act as a validation set. Points represent the
 232 individual model that returned the highest R^2 value for the given number of promoter
 233 sequence positions. The circular points and dashed lines represent models trained
 234 using KFold cross validation. The square points and solid lines represent models
 235 trained using holdback cross validation.

236

237



238

239 **Supporting Figure 6: Plasmid map of the pS797 expression vector used for *in***
 240 **vivo characterisation of putative promoter sequences.**

241

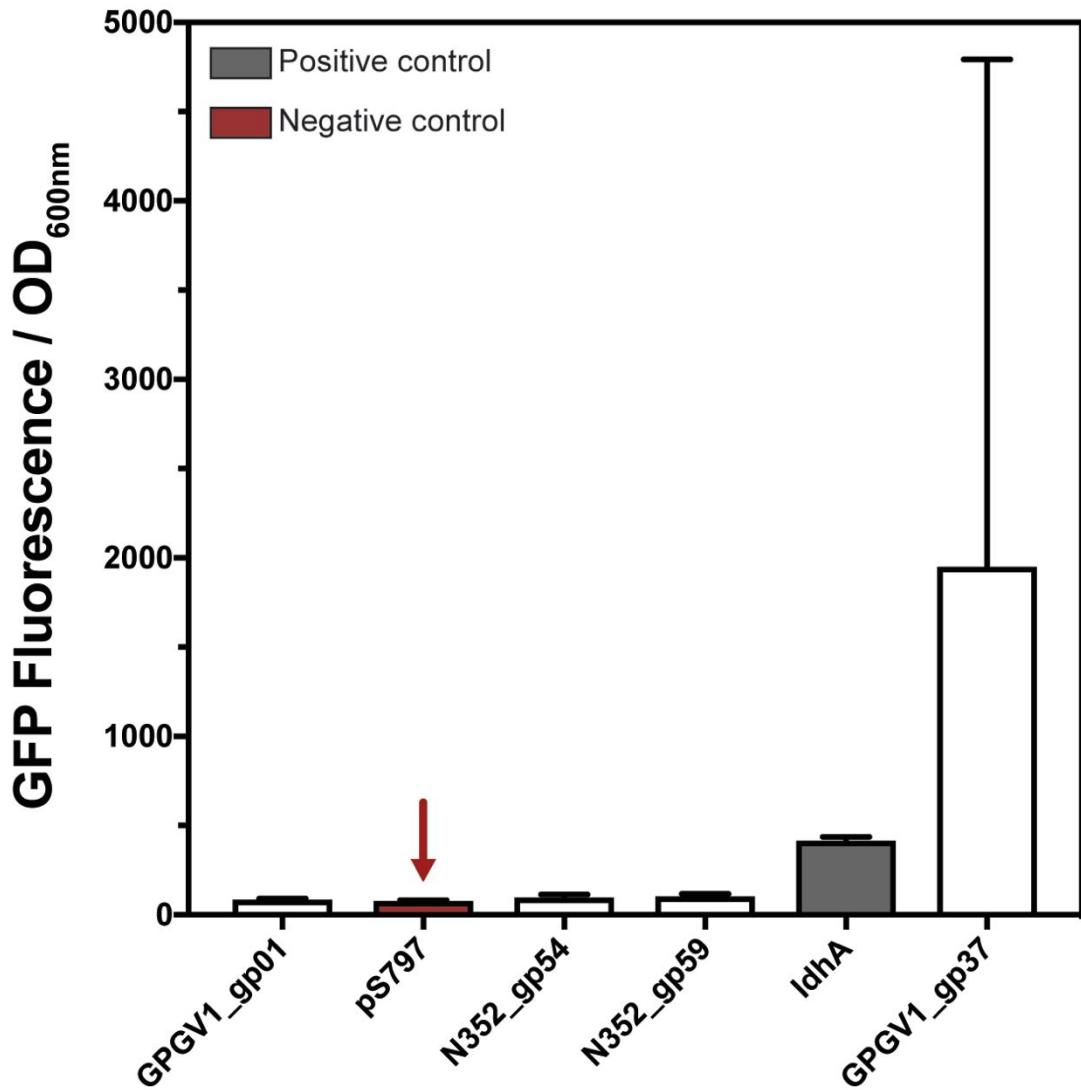
242 The origin of transfer, ORI T, which contains the machinery necessary for
 243 mobilisation of the vector during conjugal transformation, is shown in green. Antibiotic
 244 resistance genes, conferring resistance to ampicillin (AmpR) and kanamycin (KanR)
 245 are shown in red. 2 origins of replication are present: ORI CoIE for replication in *E.*
 246 *coli*, and ORI BST1 for replication in *Geobacillus*. Both origins are shown in blue. The
 247 binding sites of primers used for sequence verification of the promoter and reporter
 248 sequences are shown in purple.

249

250 The promoter of interest, shown in cyan, is located between multiple cloning sites
 251 (MCS) containing the listed restriction sites. The reporter protein used for promoter
 252 characterisation is also located between the MCS and is followed by the S718
 253 terminator sequence. The reporter and terminator sequences are shown in dark blue
 254 and red, respectively.

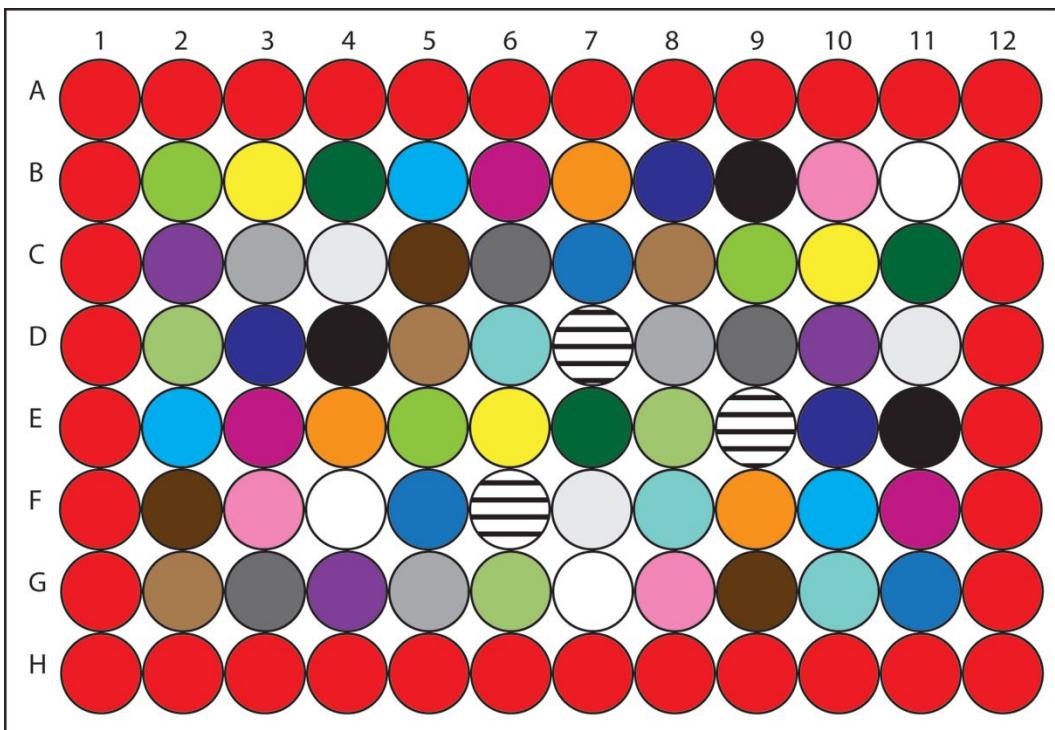
255

256 The conjugal plasmid pS797 was developed by ZuvaSyntha Ltd. (Hertfordshire, UK)
 257 from the electroporation plasmid pNC2, which in turn was developed from plasmid
 258 pUCG18 (GenBank EU547236)⁹².



259
 260 **Supporting Figure 7: Initial characterisation of putative promoter sequences**
 261 **isolated from bacteriophage genomes.**

262
 263 Promoter activity in *G. thermoglucosidasius* was characterised after 24 h growth. The
 264 positive control, the *G. thermodenitrificans* *IdhA* promoter is shown in dark grey and
 265 the negative control, *G. thermoglucosidasius* transformed with an empty pS797
 266 vector, is shown in red. Bars represent the mean of n = 3 independent starter
 267 cultures arising from independent transformation events, except in the case of
 268 GPGV1_gp37, where n = 9. Standard deviation error bars shown, unless hidden by
 269 the bar.
 270



271

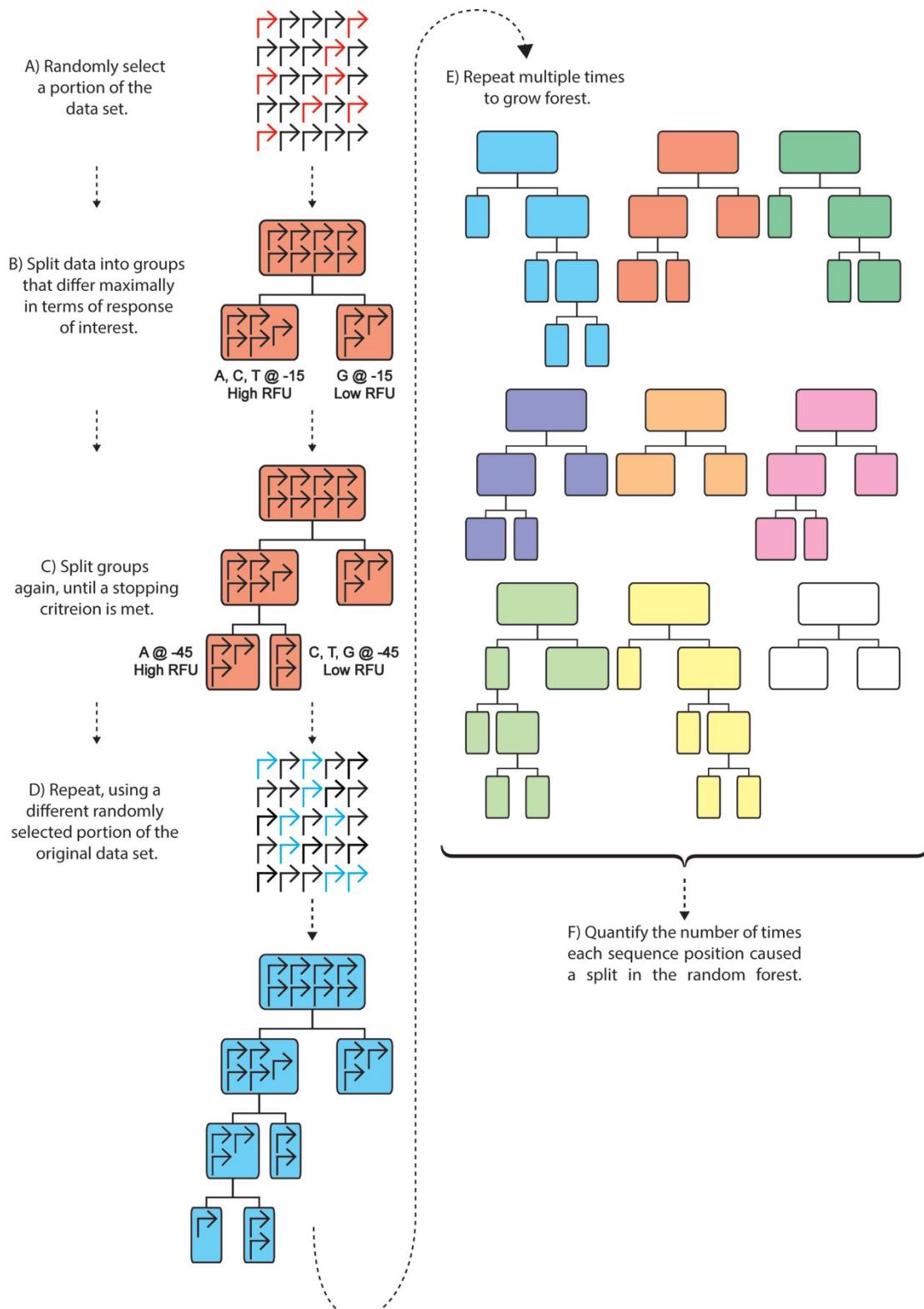
272 **Supporting Figure 8: Schematic representation of Latin rectangle 96-well plate**
 273 **layout used in promoter characterisation.**

274

275 96-well plates contained 3x 200 µl aliquots from each of 20 starter cultures.
 276 Disparate colours and patterns represent aliquots taken from the same starter
 277 culture. The outermost wells, highlighted in red, contained sterile growth media to
 278 reduce edge effects.

279

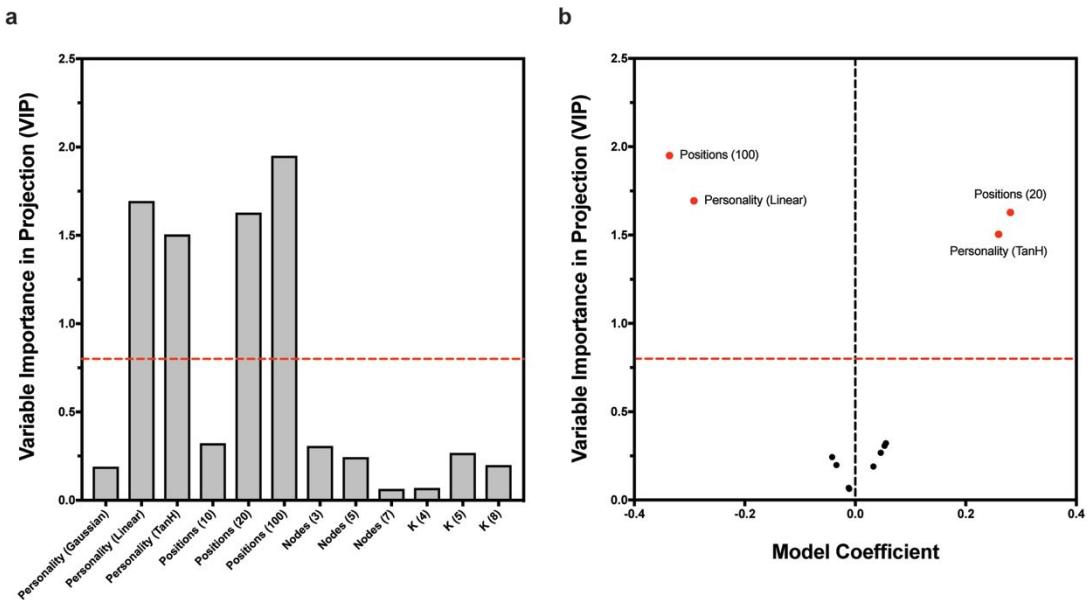
280



281

282 **Supporting Figure 9: Schematic representation of a random forest partition**
283 **model, as applied to promoter sequences.**

284



285

286 **Supporting Figure 10: Assessing the contribution of ANN model parameters to**
 287 **determining predictive power using a PLS model.**

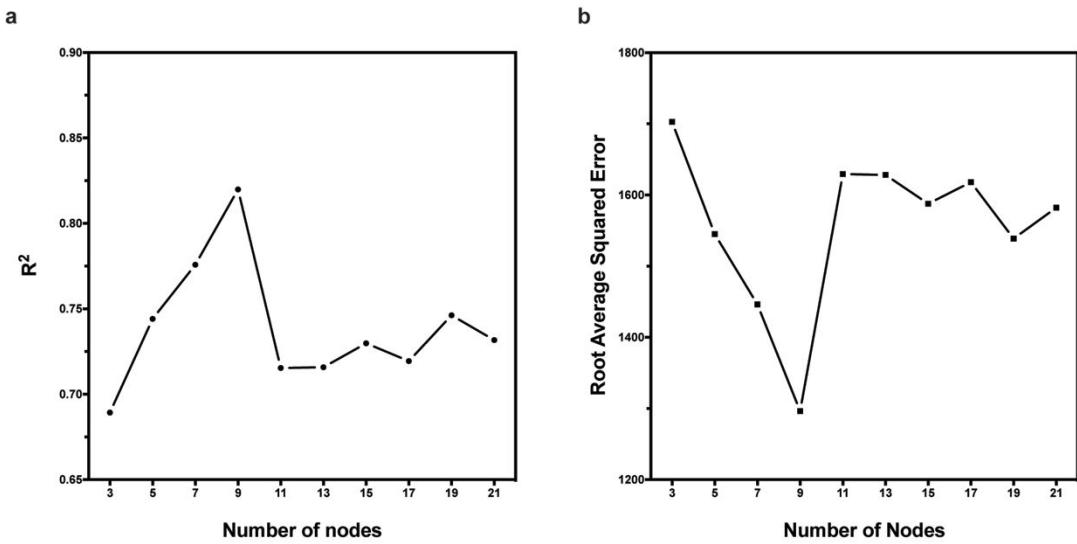
288

289 a) Variable Importance in Projection (VIP) plot and b) VIP v. model coefficient plot.

290

291 In both panels, the dashed red line represents the VIP threshold value of 0.8, above
 292 which variables are predicted to have a significant impact on model output. VIP and
 293 coefficient values were returned by a PLS regression that modelled the R^2 value
 294 returned by ANNs when applied to a test data set as a function of activation function
 295 personality (Personality), the number of promoter nucleotide sequence positions
 296 included in the ANN (Positions), the number of nodes included in the hidden layer
 297 (Nodes) and the value of K used in ANN KFold cross validation (K).

298



299

300 **Supporting Figure 11: Model performance statistics for ANNs modelling GFP**
 301 **fluorescence as a function of complete 104 bp promoters.**

302

303 R^2 (A) and Root Average Squared Error (RASE, B) values were returned when ANNs
 304 were applied to an independent test set of 10 promoter sequences. For each number
 305 of nodes specified, 1,000 ANNs were fit, with each fit having a unique, specified
 306 random seed. For each number of nodes, 500 models were fit using KFold cross
 307 validation with $K = 4$, and 500 models were fit with $K = 5$. Points represent the R^2 and
 308 RASE values returned by the single highest performing ANN for each number of
 309 nodes.

310

	Sequence		Sequence		Sequence		Sequence
ada	ATAAGAAT	cpxR	AGTGTAAA	cysB	TCTTGCAT	fruR	TGAATCGA
araC	AAGATTAG	cpxR	GGAGTAAC	cysB	TGTATATA	fur	ATAATTGT
araC	TAACAAAA	crp	TGTCGTTAA	cytR	TTCTGTAA	fur	ATTGCAT
araC	TCCGACCT	crp	TTTTGTTA	cytR	ACTTGTAA	fur	TAATGCTT
araC	GACATAAG	crp	AAGTGTGA	cytR	ATAGAACT	fur	ATGATAAT
araC	TCCTTGTT	crp	TAGACACT	deoR	AATTCTAA	fur	ATAATGAT
araC	ACAGAGGG	crp	AGATCACA	deoR	TTAGAATA	fur	ATCATTAA
arcA	ATTGTAA	crp	TAATGTGA	deoR	AATTTTAT	fur	CATTATAG
arcA	TTAATTAA	crp	TTAAATTG	dnaA	ATTCACAA	galR	TGTAAGCG
arcA	TAATTAAA	crp	ATGCGAGG	dnaA	CCACAAAGT	galR	ATGTAACC
arcA	TTAACTAA	crp	TCACACTT	dnaA	AGTTATCC	gcvA	AACTAATT
arcA	TAACCTAA	crp	AGATGTGA	dnaA	TTGTTATC	gcvA	TTATATTT
arcA	TAATTATA	crp	TTAGATTA	dnaA	TATACACA	gcvA	ATAAGCTA
arcA	TCATGTAA	crp	GATGTGTA	dnaA	TTTGGATA	gcvA	ACTAATAG
arcA	AATAAAAA	crp	AAATGTAA	fadR	CCGACCTA	glpR	TTCAAAAT
arcA	AAAAGGGA	crp	AATAAATT	fadR	GGACTTGT	glpR	TCGAATTA
arcA	TAACAATT	crp	TCACATT	fadR	AACTCATC	glpR	CACACATT
arcA	AATTGTAA	crp	TATACATA	farR	TGTATTAT	glpR	ACGATAAG
arcA	TCACAAAA	crp	ACACACAT	fhlA	TCATTTTC	glpR	AATTTGAG
argR	ATAAAAAT	crp	CTCGTTT	fis	TCTTAACT	hipB	CCCTTAAG
argR	ATAATCAT	crp	AAATAACA	fis	TCAGAGGA	hipB	AAAGGATA
argR	TTTTTTAT	crp	AACCCCTC	fis	CTCATT	hipB	TAAGATGT
argR	TTATAATT	crp	GTTAGCTC	fis	TGTAAATT	hipB	CCTTTTA
argR	ATATTCA	crp	ATAACAAT	fis	TTATCTAA	hns	TAGGCTGA
argR	AATAATT	crp	TGTGATCT	fis	ACAGTTGT	hns	ACAGAGTA
argR	AATTCAAT	crp	CGCACATA	fis	TATACTTA	hns	AAAGGAAT
argR	CATAAAAA	crp	TTTGCAA	fis	TATTCTAT	hns	TAATTAA
argR	CATCCATA	crp	TTAGTTCA	fis	AATTATT	hns	ATTAATT
argR	TAATTCA	crp	CACAGTGC	fis	ACAATTAT	ihf	TTTTATT
argR	AATTCA	crp	GTGCTCCC	fis	AAATGTGA	ihf	GCTTAGAG
argR	CACAATAA	crp	ATCACAA	fis	CCCTCCGT	ihf	TTGGAGTC
argR	AATTAATA	crp	TTTGTGAG	fis	AAAAATAA	ihf	ATCATACA
argR2	ATACACTA	crp	TCACAATT	fis	TCTTTAAT	ihf	TAGGATAA
argR2	CATATT	crp	AGGTAACA	flhCD	TGAGGGGT	ihf	AGATATAT
argR2	TTTTTATT	crp	ATGCCGTA	flhCD	GGGCTTT	ihf	CCGGCCTA
argR2	ACATATAA	crp	TCATATCA	fnr	ATCAATT	ihf	CCATGTAG
argR2	TTTATT	crp	TCAAACCT	fnr	TTTTTGAA	ihf	AATAAAAT
argR2	ATATAAAAT	crp	GTTCAATT	fnr	ATCAATAA	ihf	TGTAGGCC
argR2	GCAAATAA	crp	ATTGTGA	fnr	TTATACAA	ihf	AAAGTATA
carP	CACTTTT	crp	GTGTGTAA	fnr	ACAATT	ihf	TAGATTAA
carP	CTGTAAAA	crp	GAAACATA	fnr	TAAATTGT	ihf	TATACTTT
cpxR	GTAAATT	crp	TAACCAAT	fnr	ATTGT	ihf	ACAGACAA
cpxR	CTTGTAAG	cspA	GCAAGACA	fnr	ATAAATAT	ihf	TGATATGA
cpxR	TAAAATT	cynR	TAAGGTAA	fnr	TGTAAATC	ihf	AGTATACA
cpxR	TAAAACAA	cynR	ATAAGTAA	fnr	TCAAGAGT	ihf	CTCCTGA
cpxR	TAAAAAGA	cysB	TATAGTTA	fnr	CCTACCTC	ihf	TAATCTAT

	Sequence		Sequence		Sequence		Sequence
ihf	TTTCAAAA	melR	TCCCATAA	oxyR	ATTAGTGT	rpoD15	TAATGTTT
ihf	AAATAAAA	metJ	ATAAGCGT	oxyR	CGGAGTAA	rpoD15	CTAAATAA
ihf	TTAATTAA	metJ	ATACATCT	oxyR	GATTAATT	rpoD15	TCCTCTGT
ihf	GCTAATAG	metJ	ACATCTAA	oxyR	TAGAATAG	rpoD15	AGAGTCAA
ihf	CTTCGGGA	metJ	AGACATCC	pdhR	TTGTTAAA	rpoD15	TTTGTGTA
ihf	CTAAGGGC	metJ	TGTAAACAA	phoB	TTTAATTAA	rpoD15	TTAACACAA
ihf	TGTAAGAA	metJ	GCTAAAAT	phoB	CTGTCATA	rpoD15	CCTCTCCC
ihf	ACAAAAAA	metJ	CTTTACAT	phoB	TTTTCATA	rpoD15	CCCTTAAA
ilvY	ATATATCA	metJ	TAAACGGA	phoB	ACGCATAA	rpoD15	TATAAGTC
lacI	AGTAACAA	metJ	TAGATGTG	phoB	TAATATAT	rpoD15	GGTATACT
lexA	TATATAAA	metR	ATTTTTCC	phoB	TAATCTGT	rpoD15	AGGATTAA
lexA	TATACAGT	metR	TTTTTTCA	phoB	AATAAAAG	rpoD15	TGTCTCAC
lexA	TCCATACA	metR	TTCTTTTC	phoB	TTTATTAA	rpoD15	TTTTAACAA
lexA	ATAAATAA	metR	CAAATTTT	phoB	TGTCATCA	rpoD15	GTGAATAG
lexA	TATATACT	modE	TATACAAG	phoB	TCATAAAA	rpoD15	AAGGTAGA
lexA	TATATACA	modE	CTATATAC	phoB	CTTACATA	rpoD15	TTTCGCA
lexA	AAACCACAA	modE	CCTACATA	phoB	CATCTTTC	rpoD16	CACATTTC
lexA	TACTGTAT	nagC	CATTTCAC	phoB	ATGTAACA	rpoD16	GCTTTTTA
lexA	ATATAAAA	nagC	GAAATAAG	phoB	TAATTCGA	rpoD16	TTGCAAAT
lexA	TGTACACA	nagC	ATATTTTA	phoB3	TCCTTACA	rpoD16	TCATAACAA
lexA	ATATATAC	nagC	ATTTTAGA	purR	AATAAAAG	rpoD16	CGTATAAT
lexA	ATAGATAA	nagC	TTTAATTAA	purR	ATTTCAAG	rpoD16	ACCTAAAT
lexA	TAAAAACAA	nagC	CTTATTAA	purR	TAATCGTT	rpoD16	ACGGTATA
lexA	TATATTCA	narL	TGCTCCTT	purR	CGCGAGGT	rpoD16	TCGCCCT
lexA	TTTTTTTA	narL	TAACTCTT	purR	TGAGGAAA	rpoD16	TTTCAAAT
lexA	TATATATA	narL	GTAAGGGT	purR	TTTTTAAG	rpoD16	TAGCGTAA
lrp	ATCTTTTA	narL	TCCCCATG	purR	CAAGGAAA	rpoD16	CAGATATA
lrp	TATTTTTT	narL	TCCTAAAG	purR	CTTTTTCT	rpoD16	ATTCATA
lrp	TTTAGTGT	narL	AAGTAGTA	purR	TTTTCGTT	rpoD16	TTTTTATA
lrp	ATTTTTTT	narL	TCGGGGTA	purR	CTTTCCCT	rpoD16	TGTCAAGA
lrp	TTGTCTTA	narL	ACTCCTTA	purR	GAAACGAG	rpoD16	CCCCATAAA
lrp	TGCATTTT	narL	TAACTCTA	purR	TTTCGTTT	rpoD16	TTCAATCT
lrp	TGTAGAAT	narP	TTGAGGTA	purR	CGTTTTTT	rpoD16	TAGTCGG
lrp	ATTTATTA	narP	AAGGAGTA	rhaS	TCCTGTCA	rpoD16	TCTTTAGG
lrp	TTTCTTTT	narP	TTTAGAGT	rpoD15	AGCCTAAA	rpoD16	TAGAATGC
lrp	TATTCTTA	ntrC	TGCACTAA	rpoD15	AACGCATA	rpoD16	CTCCTGAT
lrp	TTGACAAT	ompR	AAATCACA	rpoD15	ACATACTA	rpoD16	CCTTATAA
lrp	TATTTATT	ompR	TCATATTAA	rpoD15	CATGATAG	rpoD16	AAATAATT
lrp	TAAGATAA	ompR	GTAACATA	rpoD15	TAGCTTAT	rpoD16	TGTAGACT
lrp	AATACATA	ompR	TGTTTACA	rpoD15	TTTTGTAA	rpoD16	GTTAGGGT
lrp	ATAACTAA	ompR	TTTACATT	rpoD15	TTTTGTAA	rpoD16	GAGTATAA
malT	AGGGAGGA	ompR	GTTACATA	rpoD15	TAAGGTAA	rpoD16	ACAGTATA
malT	GGGGAGGA	ompR	TTCTTTTT	rpoD15	TTTTCACG	rpoD16	TTATAAAA
marR	TATACTTG	ompR	TTGTAGCA	rpoD15	GCTCAAGA	rpoD16	TAATTAGA
marR	CAACTAAT	ompR	ATATTACA	rpoD15	TGTGTAGG	rpoD16	TAGAGCAC

Sequence	Sequence	Sequence	Sequence
rpoD16 TGTATAAT	rpoD17 CCCGTAGG	rpoD17 GTGTCATA	rpoH2 TTTTTTTT
rpoD16 AGTCTCAA	rpoD17 AATTTCT	rpoD17 ATTAGAGT	rpoH2 TCTCCCTT
rpoD16 AGTGTAAAT	rpoD17 CTCCCTTT	rpoD17 TGCAATAA	rpoH2 CTCTCCCA
rpoD16 CCTACAAT	rpoD17 CCAAATAG	rpoD17 TTGTCTGA	rpoH2 AAATAATG
rpoD16 AACTAGTT	rpoD17 TTTTACTT	rpoD18 GATAGAAT	rpoH2 ACAAAAGA
rpoD16 CATAGTGT	rpoD17 CTAAAACCA	rpoD18 ACTTGT TA	rpoH2 ATTCTACC
rpoD16 TGATATAA	rpoD17 TTTTATAG	rpoD18 ATCATTGT	rpoH2 CCCTTTAA
rpoD17 AACTAACAC	rpoD17 CTCTGTAG	rpoD18 AATTGAGG	rpoH3 CTCCCCCT
rpoD17 CGTTGTAA	rpoD17 CAAGAGGG	rpoD18 AGTGTAGT	rpoH3 CAAAAAAA
rpoD17 ACTTTTGT	rpoD17 AGTAGAAT	rpoD18 TTGCTTTT	rpoH3 CTTCCCTT
rpoD17 TTTGTTTT	rpoD17 ATACTTAA	rpoD18 AGGCCCTC	rpoH3 CTTGAATA
rpoD17 TATAGATT	rpoD17 AATCTTTA	rpoD18 ACGACCCC	rpoH3 CTGATAAG
rpoD17 ATTTTGTA	rpoD17 GTTATACT	rpoD18 AAATATAT	rpoH3 CCCCCCAT
rpoD17 CCAAATTG	rpoD17 GAGTAAAA	rpoD18 TTGAGATA	rpoH3 AATAACTC
rpoD17 CTTAAAAT	rpoD17 CTAACCCT	rpoD18 GTGAGGGA	rpoN AAATTGTA
rpoD17 ACTTTTAG	rpoD17 TAGCAACA	rpoD18 TTACACTT	rpoS17 TTATGTTT
rpoD17 TGATAATT	rpoD17 TAGCCTTT	rpoD18 TTGATAGG	rpoS17 CCCCCCTC
rpoD17 CCCATAAC	rpoD17 GCATAATG	rpoD18 TGTATGAT	rpoS17 CTATTATA
rpoD17 GAGACACA	rpoD17 CTTACTTT	rpoD18 CTTATACT	rpoS17 GGAGGGTG
rpoD17 TAATGTAA	rpoD17 TTTTCCTT	rpoD18 AGTACGGG	rpoS17 CTCACAAA
rpoD17 CCTATAGT	rpoD17 TTCCCTGT	rpoD18 GGCATAAT	rpoS17 CAGACATA
rpoD17 ATTAGTTA	rpoD17 CCCCTTTT	rpoD18 TTAATT TT	rpoS17 CGGCTAGT
rpoD17 TTAGGCTA	rpoD17 TTGTGAAT	rpoD18 GTACAATC	rpoS17 AGAGGGAG
rpoD17 ATACTATA	rpoD17 ATAATAAT	rpoD18 TTTAAAAT	rpoS17 TTATATT A
rpoD17 ACAAGTGT	rpoD17 CAGCATAAC	rpoD18 AGACTACT	rpoS18 GTAGTCTA
rpoD17 GACCCCAC	rpoD17 TCCCTTTT	rpoD18 GGTAGACT	rpoS18 TGTCA TGA
rpoD17 AATAGTTA	rpoD17 TAGGAATA	rpoD19 ACGTGCTA	rpoS18 TAGTTTAG
rpoD17 CAATTCTA	rpoD17 GCTACAAT	rpoD19 ATTGTTTT	rpoS18 CAGCTAGT
rpoD17 AGTTAATA	rpoD17 CATGCTAA	rpoD19 TGTGTTAG	soxS GATAAGCG
rpoD17 CTTATAAG	rpoD17 ACAATCTT	rpoD19 CGGAGTAG	soxS TATCATT T
rpoD17 TACCTAAA	rpoD17 TTGGAATA	rpoD19 ACAGATTT	soxS TTGCTCCT
rpoD17 CAATTCTA T	rpoD17 GGAGTGTA	rpoD19 TTTCATAT	soxS AACTGTAA
rpoD17 AATAAGTA	rpoD17 CGTACT	rpoD19 GCTATAAT	soxS ATTGTTA
rpoD17 AATAAAATA	rpoD17 ATAACAGG	rpoD19 TTACCTCA	soxS AACCCCCG
rpoD17 TGCTTGTA	rpoD17 CCCCTGGA	rpoD19 GAATGTAG	soxS ATAATTCA
rpoD17 CATATAGT	rpoD17 CCTATACT	rpoD19 ACTAGCGA	soxS CTTTCCC
rpoD17 AGATCGGA	rpoD17 ACAGGAAT	rpoD19 CACCTAAA	soxS CAAGTATA
rpoD17 TTCCTCCT	rpoD17 TGTTATAA	rpoD19 GTGATCTA	soxS AAACCATA
rpoD17 AAAAATAG	rpoD17 GTTCGAGG	rpoD19 TTAACATG	soxS AATTCTTT
rpoD17 TTGATATA	rpoD17 CCCTAAAA	rpoD19 TACTTAAA	torR GTTCATAT
rpoD17 TAACGGAT	rpoD17 TTATGATA	rpoD19 CAAGCTTG	trpR TTAACTAG
rpoD17 AAATCAGA	rpoD17 TTTATAAT	rpoD19 TGTTGTGT	trpR TACTCTTT
rpoD17 TTTGTAT	rpoD17 CTCCCATAG	rpoD19 TCCTGCTA	tus TAGTATGT
rpoD17 CGCCTTTT	rpoD17 CTTAACCT	rpoE ACTTTACA	tus ATAAGTAT
rpoD17 TGTAAAAT	rpoD17 GTCCTACA	rpoE GTCTGATA	tus TAACTAAT

Sequence	
tus	CATTAGTA
tus	TAACTAAG
tus	TTGTAACG
tyrR	TAAATAAA
tyrR	TATACAGA
tyrR	GTGTAAAA
tyrR	TTTACAAT
tyrR	TGAGATT
tyrR	TGTAATT
tyrR	CTTACACA
tyrR	TTAATACA
tyrR	AATATATA
tyrR	TATGTAAC
tyrR	ATTGTACA

311

312 **Supporting Table 1: List of Transcription Factor Binding Sites (TFBS) used in**
 313 **promoter identification.**

314

315

316

317

318

Part	Prefix sequence	Suffix sequence
5'-Homologous region	CAGGAAACAGCTATGACCATGGA	TACTTGAGACCACGAAGTTACTA
	ATTCGCGGCCGCTTCTAGAGACT	GTAGCGGCCGCTGCAGGACTGGC
	CTGTGGTCTCAGGAG	CGTCGTTTACA
Distal Regulatory Sequence (DRS)	CAGGAAACAGCTATGACCATGGA	ACCTTGAGACCACGAAGTTACTA
	ATTCGCGGCCGCTTCTAGAGACT	GTAGCGGCCGCTGCAGGACTGGC
	CTGTGGTCTCATACT	CGTCGTTTACA
RBS	CAGGAAACAGCTATGACCATGGA	AATGTGAGACCACGAAGTTACTA
	ATTCGCGGCCGCTTCTAGAGACT	GTAGCGGCCGCTGCAGGACTGGC
	CTGTCTGGAGTCACTGGTCTCAA	CGTCGTTTACA
	CCT	
CDS	CAGGAAACAGCTATGACCATGGA	CTGCTGAGACCAGTGGCTCCAGA
	ATTCGCGGCCGCTTCTAGAGACT	CGAAGTTACTAGTAGCGGCCGCT
	CTGTGGTCTCAAATG	GCAGGACTGGCCGTCGTTTACA
Terminator	CAGGAAACAGCTATGACCATGGA	TTCGTGAGACCACGAAGTTACTA
	ATTCGCGGCCGCTTCTAGAGACT	GTAGCGGCCGCTGCAGGACTGGC
	CTGTGGTCTCACTGC	CGTCGTTTACA
3'-Homologous region	CAGGAAACAGCTATGACCATGGA	GCTTGAGACCACGAAGTTACTA
	ATTCGCGGCCGCTTCTAGAGACT	GTAGCGGCCGCTGCAGGACTGGC
	CTGTGGTCTCATTCG	CGTCGTTTACA
pS797 vector	GGTCTCTTCGTAAGTAGCGG	ATGGAATTCGCGGCCGCTTCTAG
	CCGCTGCAGG	AGTACTAGAGACC

319

320 **Supporting Table 2: DNA sequences of cloning affixes used in type IIs**
 321 **restriction cloning.**

322

323 Prefixes were joined to the 5' end of the relevant parts and suffixes were joined to the
 324 3' end.

325

Promoter	Source	GFP RFU	mOrange RFU	Putative Identity	Gene Ontology Biological process / Molecular function
ATPase/GTPase					
42	GSTEA	432.51	9.59	nirQ Denitrification regulatory protein	ATPase activity
44	GSTEA	674.39	43.08	typA GTP-binding protein	GTPase activity
47	GTDN	2071.08	13.69	ychF Ribosome and GTP binding protein	ATPase activity
51	GSTEA	222.90	179.52	cmpC Bicarbonate transport ATP-binding protein	ATPase activity
55	GTGNS	510.41	105.22	hflX GTPase	GTPase activity
60	GSTEA	1265.28	864.58	sugC Sugar transport ATP-binding protein	ATPase activity
64	GKAU	2083.20	543.27	ravA ATPase	ATPase activity
Biosynthetic processes-related					
4	GTDN	88.81	14.38	clsB Cardiolipin synthase	Cardiolipin biosynthetic process
9	GSTEA	99.76	12.53	ilvC Ketol-acid reductoisomerase	Isoleucine biosynthetic process
18	GKAU	107.35	9.76	hemN Oxygen-independent coproporphyrinogen III oxidase	Porphyrin-containing compound biosynthetic process
39	GSTEA	292.71	90.63	hemL Glutamate-1-semialdehyde 2,1-aminomutase	Protoporphyrinogen IX biosynthetic process
41	GSTEA	371.83	92.95	plsY Glycerol-3-phosphate acyltransferase	Phospholipid biosynthetic process
53	GTDN	344.88	126.45	modB Molybdenum transport system permease	Mo-molybdopterin cofactor biosynthetic process
65	GSTEA	2243.22	459.15	moaB Molybdenum cofactor biosynthesis protein B	Mo-molybdopterin cofactor biosynthetic process
69	GKAU	2928.95	124.90	dapH 2,3,4,5-tetrahydropyridine-2,6-dicarboxylate N-acetyltransferase	Diaminopimelate biosynthetic process
Catabolic processes-related					
6	GKAU	96.15	9.17	lon2 Lon protease 2	Protein catabolic process
8	GTGNS	98.59	8.76	deoB Phosphopentomutase	Deoxyribonucleotide catabolic process
43	GSTEA	439.98	82.08	lytC N-acetylmuramoyl-L-alanine amidase	Peptidoglycan catabolic process

Promoter	Source	GFP RFU	mOrange RFU	Putative Identity	Gene Ontology Biological process / Molecular function
Cell membrane-related					
15	GTGNS	106.07	17.17	Integral membrane protein TerC family protein	Integral membrane component
58	GKAU	835.97	158.54	mcpA Methyl-accepting chemotaxis protein	Integral membrane component
72	GTDN	3662.86	354.38	ydjM Inner membrane protein	Integral membrane component
77	GSTEA	5331.79	223.34	5-bromo-4-chloroindolyl phosphate hydrolysis protein	Integral membrane component
Miscellaneous					
7	GKAU	98.28	11.75	pepA Cytosol aminopeptidase	Aminopeptidase activity
12	GTDN	101.12	9.72	hmuU Hemin transport system permease	Transport
14	GSTEA	105.80	10.65	mtnK Methylthioribose kinase	L-methionine salvage
16	GTDN	106.67	25.62	Ferredoxin--NADP reductase	Cellular iron ion homeostasis
38	GPGV1	215.12	12.08	YqaJ viral recombinase	Recombinase activity
52	GSTEA	314.30	164.02	cobB NAD-dependent protein deacylase	Glutamine metabolic process
61	GSTEA	1332.51	288.34	amaA N-acyl-L-amino acid amidohydrolase	Aminoacylase activity
66	GKAU	2422.96	274.73	yhdN Aldo-keto reductase	Stress response
70	GTGNS	3451.39	175.18	zapA Cell division protein	Cell Division
73	GSTEA	3808.74	142.70	hcnC D-amino-acid-oxidase	Oxireductase activity
74	GKAU	3983.82	389.54	Putative heme peroxidase	Oxidation-reduction process
79	GSTEA	8585.09	135.90	Enterobactin/ferric enterobactin esterase	Iron ion transport
80	GSTEA	10571.15	379.11	hmp Flavohemoprotein	Response to nitrosative stress
Nucleic acid-related					
2	GKAU	83.99	14.59	dnaG DNA primase	DNA replication
17	GSTEA	106.99	13.00	dps DNA protection during starvation protein	Cellular iron ion homeostasis
21	GKAU	109.52	16.29	dps DNA protection during starvation protein	Cellular iron ion homeostasis
46	GTGNS	11115.08	61.66	glcR HTH-type transcriptional repressor GlcR	Transcription Regulation
54	GTDN	438.67	103.75	yfhQ Adenine DNA glycosylase	Base-excision repair
63	GTDN	1982.95	153.12	sigW RNA polymerase sigma factor	Transcription
76	GSTEA	4778.75	169.79	gltX Glutamate-tRNA ligase	Glutamyl-tRNA aminoacylation

Promoter	Source	GFP RFU	mOrange RFU	Putative Identity	Gene Ontology Biological process / Molecular function
Protein folding and secretion-related					
24	GSTEA	112.45	16.50	groL Chaperonin	Protein refolding
48	GKAU	2279.05	19.69	SecG preprotein translocase subunit	Protein secretion
50	GSTEA	218.37	301.62	Putative bifunctional phosphatase/peptidyl-prolyl cis-trans isomerase	Protein folding
Sporulation-related					
3	GKAU	84.71	11.38	yabP Sporulation protein	Sporulation resulting in formation of a cellular spore
10	GTGNS	100.32	9.87	divIB Cell division protein	Sporulation resulting in formation of a cellular spore
11	GSTEA	100.94	9.26	Gamma-D-glutamyl-L-diamino acid endopeptidase	Sporulation resulting in formation of a cellular spore
19	GTGNS	108.24	12.04	sspD Small, acid-soluble spore protein D	Sporulation resulting in formation of a cellular spore
20	GTGNS	108.83	9.52	cwlJ Spore coat protein, cell wall hydrolase	Sporulation resulting in formation of a cellular spore
25	GTDN	114.89	8.04	YabP family Sporulation protein	Sporulation resulting in formation of a cellular spore
26	GTGNS	116.61	10.85	mecA Adapter protein	Negative regulation of sporulation
30	GTDN	126.93	11.70	ctpB Carboxy-terminal processing protease	Sporulation resulting in formation of a cellular spore
32	GTGNS	139.46	83.92	spo0F Sporulation initiation phosphotransferase F	Sporulation resulting in formation of a cellular spore
40	GKAU	356.39	23.72	cotJA Spore coat associated protein	Sporulation resulting in formation of a cellular spore
59	GTDN	1244.59	182.38	ecsA ABC-type transporter ATP-binding protein	Sporulation resulting in formation of a cellular spore
Transferases					
36	GKAU	160.56	93.60	Aminoalkylphosphonic acid N-acetyltransferase	Transferase activity
37	GKAU	179.38	76.61	lipL Octanoyl-[GcvH]:protein N-octanoyltransferase	Transferase activity
75	GSTEA	4768.29	355.45	yrrT Methyltransferase	Transferase activity

Promoter	Source	GFP RFU	mOrange RFU	Putative Identity	Gene Ontology Biological process / Molecular function
Unknown Function					
1	GSTEA	75.01	8.74	Hypothetical protein	n/a
5	GKAU	95.81	12.39	Hypothetical protein	n/a
13	GTGNS	101.81	9.00	Hypothetical protein	n/a
22	GKAU	109.63	12.29	Hypothetical protein	n/a
23	GTDN	110.48	23.59	Hypothetical protein	n/a
27	GSTEA	117.47	11.07	Hypothetical protein	n/a
28	GSTEA	118.11	14.48	YugN-like family protein (uncharacterised)	n/a
29	GTGNS	121.19	13.72	Hypothetical protein	n/a
31	GKAU	130.44	75.44	Hypothetical protein	n/a
33	GSTEA	140.90	33.40	Hypothetical protein	n/a
34	GKAU	143.87	44.62	Hypothetical protein	n/a
35	GTGNS	144.61	35.58	Hypothetical protein	n/a
45	GTDN	868.61	86.09	Hypothetical protein	n/a
49	GTDN	113.02	154.27	Hypothetical protein	n/a
56	GSTEA	734.10	233.05	Hypothetical protein	n/a
57	GTDN	735.03	164.92	Hypothetical protein	n/a
62	GKAU	1950.40	155.49	Hypothetical protein	n/a
67	GSTEA	2571.00	452.36	Hypothetical protein	n/a
68	GTDN	2846.23	102.78	Hypothetical protein	n/a
71	GKAU	3582.83	193.07	yugN-like family protein, unknown function	n/a
78	GSTEA	6405.40	557.41	Hypothetical protein	n/a

326

327

Supporting Table 3: Analysis of the native genes with which the characterised promoters were originally associated.

328

329

330

Promoters highlighted in yellow are those which were identified as functioning consistently independent of the downstream CDS. Source organisms: GKAU = *G. kaustophilus*; GTDN = *G. thermodenitrificans*; GTGNS = *G. thermoglucosidasius*; GSTEA = *G. stearothermophilus*. Gene Ontology information was obtained from the UniProt database.

331

Parameter	Levels specified
Activation function personality	Gaussian, Linear or TanH
Number of nodes in the hidden layer	3, 5 or 7
Cross validation methodology	KFold with K = 4, K = 5 or K = 8
Number of sequence positions modelled	10, 20 or 100

332

333 **Supporting Table 4: Artificial Neural Network parameters included in the**
334 **architecture optimisation screening design, and the values specified for each**
335 **parameter.**

336

Promoter	DNA Sequence
1	GTCATGCAGCCATCTTTCTCCTTTCTCCATTATAATGCCGCTCGATCGAT TCCATAATGTAAATGGATAGAACGAACAGGACCTGGAAGGATGCCGCCA
2	TACGGGCAAGGGGGCGAAACGAACGTTGAAAAATTGTCGAAATGAAAGGA AAATCGGCTTGCTGCCGAGAATAATACGAACCTACGGGTTGTTCGT
3	AACCATCCCCGCTGACATACATGGTACAAACAACGGATGAGAAAGGCGTT GGCCTTGCCGGCTTCTATACTAAATGAAACCTATGCGGGGATGAAC
4	TTCGGTAAGATGGGAGAGAACGCGGCTTCAGCTCTGCTATTTGTGCGCTGT AGGATACATTAATAAGTACTGAAGCGACGCACCTACAGGAGGGGAAGGC
5	CATCAAGAAATAACTGAAATCCGTGAGAGATGCCTATACTTCGCTTCGCCTCCTT GCATAGTCTAACAGTGCAACTAGGCGAAAGACCTGAAGGGGGAGGGC
6	ATGCCCTGGCGCCGGTCGGCGCTTTTTGTTATTCCGCCAGTCGCGGAGA TACTATCATACTACATACCATCTAAAGGCAACCTGGAGGGATATCGGGGC
7	CCGTTTCAGCTGATTTTGATGACAGCAAGGAGAGAAAAGTTAGATAGTCGA AAGAATCGCTTGCTTCGCTATAATGAAGGACCTGGGAATGAATCAAC
8	GTCAGACTTCTGACGAGTGCAATGCTCAGGAGTTCTATCAGTATTCGGAATG ATAAAACAGAAAGATCACGTTATGAAGACAACCTGGAGGAATTATGATT
9	CGATCGTGTCAAAACATTGTCGATACGAATAACCGCGCTGTGCGCAATGAGGG CGCATGGCGGCCATTACACTACACGACAACCTAAGGAGAGGGTAATC
10	AAAAACATGTTGCTCGCTTAGGAACGTAATATAACGGCATCAGATGTGCCG TTTGTGTTATAAACATGCTTCGAGTTAACCTTCGAAAGCGGTGA
11	CCGCAACGGCTGCTCCCTGCTGGAGCAAGCCGTTGTTCACTTTGGCGCC TGCGCATATGTTGAAGGTAGAAGAAAGGGACCTAAGGGAGATGGCGC
12	GTTAGGTTGCGCTCGGTTGTCACTTGATTGAGTTTATGAGGCGGTTGAACGAG TGGGATGAGCGCAGATAAAAGAACACACCTGGTAGGATGAGGACA
13	TTCCCTCCATTGTCATTGTCATTGTTAGTTGACGGTTCGGCGGATTCAATTG GTGGAAAATTTCGCAATTGGATGTTAAAGGACCTGAAGGGAAACGGTG
14	TCTGAAAGATAAGAGGCGGAAGACGATGCATCTCAAGCCTTTCCGTATGCG AAAGAGGCTTTCAATTGCCAAACCTAACCTAGGGGAACGTGAAA
15	AAAGAAGGGAAATAAAAAACCACGTTTTAAAATTGGTATGAAATGAAACAT GCCCTACTATAATGGTGGTAAAACAAAGACCTAAATGAGGGGCACGC
16	AACCGTGGACGTTGAAAATCGCAACATTCTATTATCATAGTTGGCAACGGAA CGAAAAACATTGTCCTGTTAGACAACGAACCTAAGGATGGAATCCAT
17	TCATATTCCGGTGAATCGGGCAAGCTATCCATCGACACGTTCCATGTTCA TCAACAAGCCAGCTGGGTACAAGAACCTAACCTAGGAGGAGTGAGAGG
18	CCCCCCTCGTTAAAGCAAATGAAAGCTTCCAATTGCATCTATAAATTATAGAC AAAGGAGTTGCTGCTATATTAGACAAACCTAGGGGTTGTTGCG
19	GCAGAAAACAGCATAATTTCATCAACTTCCGTTATAGAGGATTCTGATTGACA CATACTAACATTGCCAGACAACACTATGACCTTAGGAGGAAATACA
20	ATCAGAAATTCACTTCAAAATGGATTGTCCAATTGTCACCATTTCCTTATA

	TCAATAGCCTAATATTGTTGAGAAAAAGACACCTAGTTGAGGTGACAAT
21	CTTCTGGGAGTTGCAAGCCCCACTTCAGCGTCGGAAGTGGGGTAGTTGA CCAACAAGCCAGCTGGGTACAAGAACATAACCTAGGAGGAGTGAGAGG
22	AAGCACCAAAGAAGTTCATGCCAACCTGGACGGCGCATAGGACAATGGCG TTTCAATAGGCTATAGTAACAGGCCATCACCTGAAAGGGAAATCGGG
23	ATAATAGGAAAATCGCGCAGGATATCCGGACGAAATCGAACGCATATTGAGGATG TCTGGCATATGCTCGAACAAACAGAGGAAGACCTAAAGGGAGAGAGGC
24	TAAACATTCCGAAAAGTACATATCCTTAACCTTATGACACGAATACTTACCTT TTTCACGGCATGCAGTGAACAAATTACATACCTAGGAGGTAACGGGGT
25	GCAATGGGCCTGCCGACAGGCAGGCCGTTGTTTGAAGCGGTGCTAGAAC CCCCTTTGTCATACATATGAGATGAAAACCTGGGGTTCTTTTAA
26	AATAAAAGTACAAGATGTTATCTGTATTAAATAGAATGTTGGGGTAAAAA TCCCTTCAACTTTATTCTTATTAGAACCTGGAGAGTCTGCGCT
27	CCCCTTTTGCAAAGTTCACTATATGATGGAGCCATGCGTTCATGC CTGTTGCCATGACGAATAATGAGGATCGGACCTGCTAAGGAGGAACC
28	CTATATGCCCTGCCGTTTCAACCAGCGTCTGAAGTCGCCGCGTCAAATA ATCCCCATCAAACGGAAACGATAAAAAGACCTAAGGAGGGTTGTCT
29	ACTTCAACCACATGGAGGCGGAAAACCATTCTGTTCTCGCATACTTATG AAAGGGATATGTTCTATTCTCAGCCAAGAACCTAAGGAGGCAAAACA
30	GAAGCATACTCAGTACAAGCCTGACATATGATGGAGATAAGGGCATGGCACGGCA TGATGCGGTGCTTTCTTATGTATAGGACCTTGAGGAGGAACAAAC
31	GGCAAAC TGCTGTCTATTATATCACAAACCGATAGAAAATCGGCATGTTCCGG TACACTAGAAAAGACGCCACTTGCCAGACCTAAAGGAGTCACCCA
32	ATTATCCAATGCAAAAAAAATGGCAGGAAATAACTGAAGGATGCGAAACTTA TAACATGGAACAAAATAATGAAAGAATTGACCTAGGTGTGTCTATAAG
33	CACCGTCGGTCTGCTAGTTGTTCAAAAAATGATTGATTTCGCCGCTTTG CCCCTACAATGATAAGAACACTGCACGGACCTAAAGGAACAGCACGA
34	AAAAGAGGAATTCTTCTTAAACGTCGATTTGTTGATTCTTCTCGCGAT GGGACCCGTTTATGGTATAGTGGCTACACCTGTAGGAGGGACTGTG
35	CTTCTTGCGAAATGGTCTAGTTGTCGTGCCGAAGGGAGTCAGGTTT ACCGCGAATTGATTATCATCATTGAAAGACCTATGGGGAGATGACA
36	CTTCGAAACGATGATCAACGTTGCCAGGGAGTTCCATTGGACAAA CATGGATATCATGAATGTAGAAATAGAAGAACCTAGGGAGATGGAGTC
37	GAGGGAGGCATGCCCTGTCCTTTGCCATTGCTAATGGCGATTATGT TATAATGAATAGAGTTGGAATGATCTATGCACCTAAGGGTTGCTCGCG
38	ATAACTCATCTTATAAATACCTCATCAAGTCATCAATTCCCTCCTTCCAATGT AAATATATGTATAATTCTGTATCAAAACCTAGGGGGGAATTGT
39	TACATAATCGTGCCTTTCACCTTTGTTCCAGACAAGTGTGCGCG TATGGTAAAATAAGAGGAGTCATGGAAACCTCAGGAGGATGAAGCG
40	TTCCTCCTTATTGCCATGATTGCCATGACAAACAGGCATGAACGGTACTGC TCTATCATAATAGACTGAACCTTGCCAAACCTAAAAGGAGGGTTTC

41	CAAAACAACCGCTTTCTTCATTCAGTCTACCGTTAACGGCGATGATAA AATAGAAAACATGCACAAACTGTTGAGTAAACCTAGGGTGTGATTG
42	TTTACATAACATTATGGTAATTCTGTTGATCCACCGAACGATT TAGTAGAATGGGAGAAGGACACGACAACAGACCTAGAGGAAGACGTAGG
43	CGCCTGTTCTGAGCGGTTTCTGCCGACAGACAGGAAATGCCATTGTG CGCGAACCATATAGTAGAATGAGAGGAAACACCTAGGAGGGAGCGCG
44	ACGGCTCTCACCAAACGTTATGTTGCATCTATATGCTTGATTCTATGTT TGTGATATAATAGACGAGCAGCTGTCTAACACCTCAGGAGTGACCAACA
45	AGGATTGGCAAATTGCAAGAATGAAAGCCGCCGCTGCAATTTCCTGCGAAAG CGGTAGAATAACACAAAGCAAACGCTGTGGAACCTAAGGGATGACAAGTG
46	TTTTTTGAGTGTGTTGATTGAAAATGATTGATTAAATGTTTCATTAT GATAAGGGTACAAAGTAATCATTCCATCCACCTAAGGAGGTGGGATCG
47	AAAAACGCTTTCCAGAACAGACTGGCGGTACGCTGTCATTCTGCCAGC CCATATAATGAGGGAAGATGCGAACGGAAACCTAGAGGAGTGGGACGA
48	TGTTCCCTTCGTGTTGAGACAGGCCAGTTGCCAACGGCGATCTAT GGTATAGTGAATTATTGCAACGTCGGTCGACCTTGGAGGTGTCAGGC
49	CCAGCTGCCACCCTGCCAGTCGTTAAAAATGATTGATTCCGCCGTTTG CCCCTACAATGATAACAGGAAACAGCATGGACCTAAAGGAACAGCGAGA
50	TCATGGGAAAAATGGAACGAACCGTGATTGGTTGACATTCAATTGTTT CAGTACGGTAAACAGTGAATAACATTTCGGACCTCAAGGAGAAGATAAC
51	ATGTCAGTCCACTCACTGCCAATAGTGACAAAATCATTCTCGTTGGACGAT TTCATTATAATAACCATACTAATGATTACACCTTACGGAGGGAAACCA
52	CTAGAAACGAGTGAAGAGGAAAGTAAAATGATTGGTGGAAAGAGATGAAGCATC AGAAATGAGAGTCGCTATGGTGCGACAGTACCTAAGGAGGGAAACGGA
53	TTTCCCACTCCTATGATAACAATCATTGTGAGTAGTGTACAGCGTTGTCTGCT TGTGTTCCGTATACTAGATGGGATGTGGCACCTAACAGAGGGTTGATT
54	TGTAGTGGATCACCCAGCCGTTGTCTCCTGCATGCTGCAAGTCATGTTGGGG TTTTTGGTATAATGGGATGGACAATAGACACCTACCATGGGGGACATA
55	ATATAATAACAAACTGTATGCCGCTATATAGCGATTGCAAGCAAGATTGCTAT GATGGAGGCAGATTGAGTCATGCAGATGGAACCTAACAGAGGTGCTGCAT
56	GCCTTGCCAAATCATCGCGCTCCCGGCCATCATTGCGCTGGCGGACAAGTTC CTCTATAATGAAGACGCCGACTATTTCACCTGGACGGAGGGACGAT
57	GCCAGTCGTCAATTCCCTGCTTCGACCATCGTGCCTGCGCAGGCAAGTCC TCTATAATTAGGACGCCGACTATTGTACCTGGACGGAGGGACGAT
58	ATTATGCCAGGAAGGAAATGAATGATCATTACGAATCTATTAAAGCGGTGAA GATAAAGAAAAGACGATTGCATAATGGAAAACCTAGGAGGAGAACGACT
59	CATTATACAGGACGACGAAAGAAACTTTGGTAACGATAGACGTTCTTG AAAATGGGACAAACGTGATTCAAGATAACGACCTGAAAGGAAGCGAACG
60	ACGGGCGACTTTTGTCACCTGCTATTCAAGATAACGACCTGAAAGGAAGCGAACG ATAATATGAGTGTGAGAAAACGATTCAATACCTAAAGGGGGATTG
61	AAACGTGCAGCGGGCTGCCTTCTTCATGAATATAGAGAAATGAAAATTCAAAAA

	TTCCGTTATAATATTAGGAAAGCGAAAACCTAGAGGAGGAAACGTG
62	AGCGTAAGCGAACGCTCCTTTGCGCCCTGGCGTACAACAGCGCGAAA ATCATGGTATGATGGAAGAAAAGAATGACACCTAGGGAGACGTGGCA
63	CCTCCCCTTTTCAGTTCTATGATACAATAAACATGTTGATGAAACTTTC TTCATCAAGCGCCGTAATATAAGATGAGCACCTGTAAGGGCGGGGTA
64	TTCCAAAAGTTGTTATTCGATATTATAATAGTATAAACTAGATATTGTC TACTGCTAAAAGTCCCATTGATCACCATAAACCTAAGGAGTTGATGTA
65	GTCTTTCCTTCGTTGCCAGCGAGAAGCAGCCGGCAAGAGTTGGTAATT TGCTATACTATAGCTAACACACAATGACAACCTAAAGGAGATGAGGGA
66	ATTGATGGACTCTTCATTGTAACATGTTGCCATTTCGCTCCGTT GTACAATGAACATATCATCATCGATCGCAACCTAAGGGCGCAAAG
67	ATACATGGCTCGTCCTCCTTTCTCCCCGCTCCCTTTCGCGTTC GTCTGCTATGATGGACATAGGCAGAAAATAACCTGAGGGGGATGGAG
68	TTTTACATAATCTTTTGCCATCCTCCCTCTGACATCTCTCGCTTCCAT TTTCTTCATTGATGGTATGATGAATAATACCTGGAGGTGCCAGATT
69	AGATATTACAGCATTGCTCGATGAATTGTGCTTGTACAATGATCTGTT TTGCTAAAGTAAAAAGAATAACAGGTGTAACCTAGGAGGATATGCC
70	CATCGTTTTCTAGGCATATCTAGTAGTCGTATCTTCATTATATCAGCATT CATGTTATGATAAACAAATAGGATTTGAAGACCTAATGGGGAAATCAT
71	AACAATTGTTGAAAAGCAAACCGTTCTGTACTATTATAATGAGGGATGATT TCCTGTTGGCGCTTTACATACTGTGTTGACCTTAGGAGGAAGAAAG
72	GTCACATTTGGTGCATTGTCACAAATTGTCACAAAGCGAACAAATATTCTT GTTATGGTATAATGGGGAAACAAGGAAGTACCTTGAGGGATGTTCG
73	CGCCTTTGACGCCCTGCCGGCTTCGGGCCGGCCACACCGTTCCGATTGGC GCCGTCTATGGTATACTAGTAGAAAAAAGGACCTAAAGGTGTTCATGA
74	ACAATTCACTACTCCACTTTGATACGCTGATAGCCGTTTTGGCAAAC GTGATATAGTGTACAAATTGTTGACCTTAAGGAGTGACGATG
75	ATCATTGAAACGGCGGAAGAAGTTATGGTACAGTGGGAGCAGGTTATCATCTAAAG AATACTATGTTATAGGAATTATAAACCTACCTATGGAGTGAGGCCAT
76	ACCGTTAGGCGCCGAATGGGGCGAGCCGTCCGTTGAACTTGCGCGATAAAC GGTACAATAAGGTACAAGTTGGATGAAAACCTGTGGAGGGTTGGACG
77	TGCGCGCCGTCTAGTTGGTGTGAAACTCTTGGCGCTTGTACGTATCATGGC GTGGACATGTTATAATGGAAGAAAAGATGAACCTAGAAAGGGATGGACA
78	GATTCCACCTCTAAATTCTTTCCCTCCATTATCTTCCACAAGACCGG CCACAAATAGAATGCGATATGATATAGTGAACCTAAAGAGGGAGGGATT
79	TGTTCGGCAAGTTCACGAAAAATTGAGCCCTCGTTAACATCTTGCCTAAA TATGTTACAATAACAAGTGAACCTGCGCTTACCTCATGGAGGGATGATC
80	CAATTATTAAAGATGTATTGAAATGTTGCTTTCACTATGCGTATGATTGAA TAAAACATGTATTAAAATACATGTTTACCTGAGAGGAGATGAAAC
<i>IdhA</i>	CTGCCTCGTCCATTGGCTTAATGGAGGTTGTCATGAAAATGACAAACACG TCCAAACAATTGCCATAATCGTTACGCATAGTTGATTCATCGCGTAAATA

ATTTGTGAATGTATTCACACCTATAAGAAGGGAGAATAGT

337

338 **Supporting Table 5: DNA sequences of the characterised promoters.**

339

340 Supplementary References

- 341 84. Yao, A. I., Fenton, T. A., Owsley, K., Seitzer, P., Larsen, D. J., Sit, H., Lau, J., Nair, A.,
342 Tantiongloc, J., Tagkopoulos, I., and Facciotti, M. T. (2013) Promoter Element Arising
343 from the Fusion of Standard BioBrick Parts. *ACS Synth. Biol.* 2, 111–120.
- 344 85. Mirzadeh, K., Martínez, V., Todd, S., Guntur, S., Herrgård, M. J., Elofsson, A., Nørholm,
345 M. H. H., and Daley, D. O. (2015) Enhanced Protein Production in *Escherichia coli* by
346 Optimization of Cloning Scars at the Vector–Coding Sequence Junction. *ACS Synth.
347 Biol.* 4, 959–965.
- 348 86. Zuker, M. (2003) Mfold web server for nucleic acid folding and hybridization prediction.
349 *Nucleic Acids Res.* 31, 3406–3415.
- 350 87. Kudla, G., Murray, A. W., Tollervey, D., and Plotkin, J. B. (2009) Coding-sequence
351 Determinants of Gene Expression in *Escherichia coli*. *Science* 324, 255–258.
- 352 88. Bentele, K., Saffert, P., Rauscher, R., Ignatova, Z., and Blüthgen, N. (2013) Efficient
353 translation initiation dictates codon usage at gene start. *Mol. Syst. Biol.* 9, 1–10.
- 354 89. Mortimer, S. A., Kidwell, M. A., and Doudna, J. A. (2014) Insights into RNA structure and
355 function from genome-wide studies. *Nat. Rev. Genet.* 15, 469–479.
- 356 90. Welch, M., Govindarajan, S., Ness, J. E., Villalobos, A., Gurney, A., Minshull, J., and
357 Gustafsson, C. (2009) Design Parameters to Control Synthetic Gene Expression in
358 *Escherichia coli*. *PLoS ONE* 4, e7002.
- 359 91. Tuller, T., and Zur, H. (2014) Multiple roles of the coding sequence 5' end in gene
360 expression regulation. *Nucleic Acids Res.* 43, 13–28.
- 361 92. Hastie, T., Tibshirani, R., and Friedman, J. (2009) Neural Networks. In *The Elements of
362 Statistical Learning: Data Mining, Inference and Prediction*, 2nd ed., pp 389–416,
363 Springer, New York.
- 364 93. Eriksson, L., Andersson, P. L., Johansson, E., and Tysklind, M. (2006) Megavariate
365 analysis of environmental QSAR data. Part I – A basic framework founded on principal
366 component analysis (PCA), partial least squares (PLS), and statistical molecular design
367 (SMD). *Mol. Divers.* 10, 169–186.
- 368 94. Clyde, M. (2002) Model Averaging. In *Subjective and Objective Bayesian Statistics*
369 (Press, S. J., Ed.), John Wiley & Sons, Inc., New York.
- 370 95. Yang, J., Zeng, X., Zhong, S., and Wu, S. (2013) Effective Neural Network Ensemble
371 Approach for Improving Generalization Performance. *IEEE Trans. Neural Netw.
372 Learning Syst.* 24, 878–887.
- 373 96. Crooks, G. E., Hon, G., Chandonia, J. M., and Brenner, S. E. (2004) WebLogo: A
374 Sequence Logo Generator. *Genome Res.* 14, 1188–1190.
- 375