1  **TReSR: A PCR-compatible DNA sequence design method for engineering proteins containing**

2  **tandem repeats**

3  James A Davey* & Natalie K Goto*

4  Department of Chemistry and Biomolecular Sciences, University of Ottawa, 10 Marie-Curie, Ottawa

5  Ontario, K1N 6N5, Canada

6  * Correspondence to JAD: jamesa_davey@dfci.harvard.edu & NKG: natalie.goto@uottawa.ca

7

8  **ABSTRACT.** Protein tandem repeats (TRs) are motifs comprised of near-identical contiguous sequence

9  duplications. They are found in approximately 14% of all proteins and are implicated in diverse biological

10  functions facilitating both structured and disordered protein-protein and protein-DNA interactions. These

11  functionalities make protein TR domains an attractive component for the modular design of protein

12  constructs. However, the repetitive nature of DNA sequences encoding TR motifs complicates their

13  synthesis and mutagenesis by traditional molecular biology workflows commonly employed by protein

14  engineers and synthetic biologists. To address this challenge, we developed a computational protocol to

15  significantly reduce the complementarity of DNA sequences encoding TRs called TReSR (for **T**andem

16  **Re**peat DNA **S**equence **R**edesign). The utility of TReSR was demonstrated by constructing a novel

17  constitutive repressor synthesized by duplicating the LacI DNA binding domain into a single-chain TR

18  construct by assembly PCR. Repressor function was evaluated by expression of a fluorescent reporter

19  delivered on a single plasmid encoding a three-component genetic circuit. The successful application of

20  TReSR to construct a novel TR-containing repressor with a DNA sequence that is amenable to PCR-based

21  construction and manipulation will enable the incorporation of a wide range of TR-containing proteins for

22  protein engineering and synthetic biology applications.

23

24  **INTRODUCTION**

25  The ability to rapidly construct, evaluate, and sequence libraries of protein variants is essential to the

26  workflow employed by protein engineers and synthetic biologists who strive to create proteins and genetic

27  circuits with new and improved properties and functions [1, 2]. A useful category of biomacromolecular

28  components can be derived from tandem repeat (TR) amino acid sequence motifs found in approximately

29     14% of all proteins [3] facilitating an array of both structured and disordered protein-protein and protein-

30     nucleic acid interactions [4, 5]. TR sequence motifs encompass a number of modular protein components,

31     including smaller intra-domain motifs forming fibrous structures (e.g., collagen and α-helical coiled-coils) [6,

32     7], intermediate sized motifs (e.g., WD40, leucine-rich, armadillo, ankyrin, Kelch, and HEAT repeat

33     domains) forming elongated solenoid and closed toroid structures [8−13], in addition to larger bead-on-a-

34     string multi-domain motifs [14]. Despite their utility and abundance, TR sequence motifs remain unexploited

35     as a class of modular components for the purposes of protein and genetic circuit engineering precisely

36     because they are encoded by repetitive DNA sequences that prohibit the routine application of PCR-based

37     molecular biology techniques [15−18].

38

39     Although it is possible to synthesize TR sequences by full-length gene synthesis, the downstream PCR-

40     based manipulations that are routinely employed in protein engineering and synthetic biology workflows will

41     be complicated by the presence of repetitive DNA sequences in these constructs. To circumvent this short-

42     coming, the degenerate nature of DNA codons encoding the canonical amino acids (with the exception of

43     Trp and Met) can be exploited to construct a TR protein encoded by a gene designed to have reduced DNA

44     sequence complementarity, thereby rendering it compatible with downstream PCR-based manipulations.

45     Furthermore, this DNA design strategy should also make it possible to employ the assembly polymerase

46     chain reaction (aPCR) to construct TR-encoding genes more cost-effectively than full-length gene

47     synthesis, since aPCR utilizes oligonucleotide primers as the only template-donating reagent in the reaction

48     [16]. However, the task of redesigning a TR DNA sequence to make it suitable for aPCR is not trivial, as

49     the probability of generating misassembled products increases with the number of primers used in the

50     reaction. Consequently, aPCR approaches in gene synthesis are limited to DNA sequences that have

51     relatively low complementary between non-overlapping segments of primers.

52

53     To create a TR-encoding DNA sequence that would be amenable to both aPCR synthesis and PCR-based

54     mutagenesis, we have devised a DNA sequence redesign strategy called TReSR (for **T**andem **Re**peat DNA

55     **Se**quence **R**edesign) to introduce silent mutations that would allow for gene construction by aPCR while

56     preserving the amino acid identity of the translated TR construct (Fig 1).  To test this methodology, we

57      designed a novel 178 amino acid residue bead-on-a-string TR protein containing a duplication of the N-

58      terminal DNA binding domain (DBD) of the bacterial repressor LacI [19, 20], a construct that has the

59      potential to expand the toolbox of DNA-binding proteins for synthetic biology applications. Application of

60      TReSR allowed for the design of a TR-encoding DNA template having reduced sequence identity (66%)

61      compared to an initial 100% sequence identity between targeted regions of the TRs. This reduction in

62      sequence complementarity enabled synthesis of the full DNA sequence by aPCR and splicing by overlap

63      extension (SOE) [18]. This template was also compatible with PCR-based site-directed mutagenesis, which

64      was used to introduce domain selective triple-mutations designed to specifically bind a variant of the *lac*

65      operator or abolish its DNA-binding activity [21].

66

67      **Fig 1. Overview of the TR DNA sequence redesign strategy implemented in TReSR.**

68      The design strategy presented in this study is schematically outlined for the construction of a TR containing

69      two identical 20 amino acid segments from the N-terminus of the LacI repressor.  (A) The TReSR protocol

70      is initiated by dissection of the 20-amino acid target sequence into contiguous 5-residue segments (labelled

71      with upper-case roman numerals) for DNA sequence redesign. (B) This is followed by the generation of a

72      sequence list (with individual sequence entries labelled with lower-case roman numerals) constructed from

73      combinations of synonymous codons that encode the amino acid sequence for each segment. An example

74      codon combination encoding the amino acid sequence for segment I is shown that uses the codons

75      highlighted in red. The label for this codon combination is given by a number for each amino acid that

76      corresponds to the list position for the codon used (e.g., for the sequence shown, the first codon is used for

77      all amino acids except for the last one which used the 6th codon in the list). Melting temperatures (Tm) of

78      all codon combinations are then calculated using the UNAfold web server24 to provide a measure of

79      homodimerization affinities for the forward ($T_{FF}$) and reverse complement ($T_{RR}$) sequences, along with the

80      Tm of heterodimerization for the forward sequence with its reverse complement ($T_{FR}$) and with the reverse

81      complement of the wild-type sequence ($T_{WT}$). Sequences are then filtered and discarded based on

82      computed hybridization metrics (described in detail in the Materials and Methods), favouring codon

83      combinations that maximize the Tm of heterodimization ($T_{FR}$) while minimizing the Tm of homodimerization

84      (TFF and TRR) and hybridization with the wild-type sequence ($T_{WT}$). (C) The third step of the TReSR

85    protocol assigns codon combinations to groups according to sequence similarity. All pair-wise percent

86    sequence identities are calculated (shown as a heat map for codon combinations (i) to (iv)) and used to

87    identify pairs of codon combinations having high sequence complementarity (e.g., codon combination (i) is

88    similar to (ii) and (iii), and dissimilar to (iv) – (vi)). These are plotted in an interaction graph of codon

89    combination space which is shown for the six codon combinations partitioned into four unique clusters

90    (shaded portions) according to their sequence similarity. (Red arrows indicate codon combinations that

91    share a high degree of percent identity and would therefore be assigned to the same group, while green

92    lines indicate codon combinations that are distinct, and consequently assigned to different groups.) (D)

93    After group assignment, the fourth step involves the assembly of sequences from two adjacent codon

94    combinations from different groups (shown for codon combinations from orange, purple and blue groups

95    from interaction graph). Hybridization metrics are calculated for the joined adjacent segments and then the

96    list of paired codon combinations is filtered (as was done in the second step, B) to eliminate paired

97    segments which are predicted to have problematic homodimerization behaviours (i.e., high $T_{FF}$ and $T_{RR}$).

98    (E) The TReSR algorithm is concluded following a depth-first-search of remaining adjacent codon

99    combinations to identify sequence paths joining contiguous segments. An example TR sequence path is

100    shown with adjacent segment codon combination pairs connected by green arrows for the first domain and

101    continued with red arrows for the path encoding the second domain. A randomly selected sequence

102    resulting from an assembled path is then evaluated as described in the Materials and Methods to confirm

103    that the DNA sequence would be suitable for aPCR construction of the target gene.

104

105    The function of the designed repressor was evaluated using a three-component genetic circuit where

106    expression of enhanced green fluorescent protein (eGFP) could be inhibited by expression of our TR

107    repressor construct binding to a unique operator element incorporated in the reporter protein promoter

108    sequence. Measurement of density-normalized culture fluorescence in the absence or presence of an

109    expression-inducing agent for repressor expression demonstrated that only those repressor constructs

110    containing a functional DNA binding sequence in both DBDs could repress expression of eGFP. This

111    genetic circuit was also used to demonstrate that a 19-residue C-terminal truncation of the duplicated DBD,

112    corresponding to the linker helix hinge of LacI (residues 61−89) also served as a functional repressor [22,

113    23]. These results demonstrate the utility of TReSR to manipulate DNA sequence components encoding

114    TRs and create a new DNA binding module that can be used as a repressor in a genetic circuit.

115

116    **MATERIALS AND METHODS**

117    **Calculation and design of the tandem repeat DNA sequences.** The target DNA sequence for the

118    development of our TreSR protocol was a new scDBD containing a TR of two consecutive LacI DBDs

119    (residues 1 through 89). The goal of the DNA sequence redesign procedure was to introduce silent

120    mutations that would allow the duplicated DNA sequence to be constructed via site-selective

121    oligonucleotide assembly by reducing the sequence similarity between the TR-encoding regions while

122    preserving the amino acid identity of the construct. A summary of the TReSR workflow is shown in Figure

123    1. Segment lengths between 5 and 7 amino acid residues were chosen to reduce the total combinatorial

124    space of silent mutations to be evaluated (Fig 1A). Thermodynamic parameter evaluation was conducted

125    using the DINAMelt two-state melting hybridization application made available through the UNAfold web

126    server [24] to predict melting temperatures (Tm) for homodimeric pairs of forward ($T_{FF}$) and reverse

127    complement ($T_{RR}$) sequences, as well as heterodimerization between the forward and reverse complement

128    ($T_{FR}$) sequences (Fig 1B). Codon segment combinations were then filtered, rejecting segments that form

129    undesired stable homodimers ($T_{FF}$ and $T_{RR}$) or which hybridize with the wild-type LacI sequence ($T_{WT}$), while

130    preferentially selecting for codon combinations that have strong heterodimerization ($T_{FR}$) potentials using a

131    percentile-based threshold calculated for each segment. Specifically, a more stringent 50th percentile was

132    used to set parameter thresholds which minimized potential off-target segment assemblies ($T_{FF}$, $T_{RR}$, and

133    $T_{WT}$) while a less stringent 10th percentile was employed to establish parameter thresholds favouring

134    hybridization with the target segment ($T_{FR}$). The values for these percentile-based thermodynamic

135    parameter thresholds are reported in Table S1. A comparison of similarity between codon combinations

136    encoding the same protein segment was performed by computing pair-wise percent sequence identities

137    (Fig 1C). Codon combinations were grouped according to whether they shared high sequence

138    complementarity (percent sequence identity ≥ 80.0%), and whether the codon pair shared similar profiles

139    for percent sequence identity values with respect to the other codon combinations for the segment, as

140    evaluated by a cosine similarity comparison (cos ≥ 0.9975). Due to the large number of remaining codon

141 combinations, the proceeding steps in the TReSR protocol were limited to codon combinations from four

142 randomly selected groups for each segment, excluding all other codon combinations from further

143 consideration (Table S1). All combinations of adjacent pairs of codon combinations were joined, and $T_{FF}$,

144 $T_{RR}$, and $T_{FR}$ values evaluated (Fig 1D) and filtered (Table S2) to discard segment pairs predicted to have

145 high $T_{FF}$, $T_{RR}$ and low $T_{FR}$ values which would prove potentially problematic during aPCR. Again, a

146 percentile-based threshold was employed to discard adjacent codon combination pairs with high

147 homodimerization affinities (20th percentile) while a fixed value for heterodimerization ($T_{FR}$ = 80.0 °C) was

148 employed to select for an appropriate set of adjacent codon combination pairs to carry forward in the

149 protocol. The TReSR protocol was concluded using a depth-first search to design the single-chain construct

150 template (Fig 1E), selecting 100 paths which visit distinct codon combinations from different groupings

151 thereby ensuring that the duplicated DNA sequences would be dissimilar and thus amenable to aPCR

152 synthesis. A single path of segments was selected to serve as the TR DNA sequence template reported in

153 Table S3. This TR DNA sequence template was then partitioned into oligonucleotide primers for aPCR

154 synthesis [25]. Refer to the TReSR Computer Code and Documentation section in the Supplementary

155 Information for the Python3.8 computer code and documentation for the TReSR protocol program.

156

157 **Construction of the genetic circuit.** The pET-11a plasmid (Novagen) was used as the genetic vector to

158 host all three components constituting our genetic circuit. These three components include Cloning Site I

159 whose genetic insert is expressed by the pDBD promoter regulated by the LacI mutant W220F (LacI$_{W220F}$),

160 Cloning Site II serving as the reporter protein expression cassette under the control of the pGFP promoter,

161 and Cloning Site III providing LacI$_{W220F}$ under the control of its native promoter pLacI. This plasmid was

162 transformed and propagated in electrocompetent *Escherichia coli* DH10B [26] via the ColE1 origin with a

163 copy number estimated at 25 to 30 plasmids per cell [27] and AmpR selection marker conferring ampicillin

164 resistance with working concentration of 100 µg·mL$^{-1}$ [28]. Combinations of the pLacI and pLacI$^Q$ promoters

165 [29] were paired with the LacI repressor and its variant W220F [30], constructed by successive quick-

166 change PCR reactions. Cloning Sites I and II were incorporated into the pET-11a vector via circular

167 polymerase extension cloning (CPEC) [31] of linear insert cassettes synthesized by aPCR [16] of the pDBD

168 and pGFP promoters and their fusion to the eGFP gene [32] by SOE PCR [18]. Unique restriction enzyme

169    sequences for BamHI and NdeI were introduced to the flanking regions of Cloning Site I, and XhoI and

170    NheI flanking Cloning Site II to enable insertion of gene cassettes at these sites, and both Cloning Sites I

171    and II are flanked by an identical T7 terminator sequence [33]. Insertion of gene cassettes into Cloning Site

172    III were performed using the NdeI restriction sequence belonging to cloning Site I and an EcoRI site 42 BP

173    downstream of the T7 terminator sequence of Cloning Site II.

174

175    **Construction of the eGFP and dGFP genes.** A copy of *Aequorea victoria* $GFP_{S65T}$ [34] was provided to

176    us by the Chica laboratory, incorporated into the cloning site of the pET-11a. This gene was used as the

177    template to produce eGFP by introducing F64L and H231L mutations by site-directed mutagenesis, yielding

178    eGFP ($avGFP_{F64L/S65T/H231L}$). Notably, our eGFP gene lacks the M1_S2insV insertion mutation

179    corresponding to the NcoI cloning scar found in the originally reported construct [32]. The decoy fluorescent

180    protein (dGFP) used to mimic eGFP expression burden while masking fluorescence output was produced

181    by introducing the R96A mutation [35] into eGFP by quick-change mutagenesis producing

182    $avGFP_{F64L/S65T/R96A/H231L}$.

183

184    **Molecular Biology Reagents and Sequencing Service.** All aPCR, SOE, CPEC, and quick change PCR

185    reactions were performed using Vent DNA polymerase (purchased from New England Biolabs, NEB) and

186    oligonucleotide primers purchased from Eurofins Genomics. Quick change PCR reactions were adapted

187    by replacing the DpnI digestion step with a gel extraction protocol (QIAquick Gel Extraction Kit, Qiagen).

188    Restriction digestion reactions for preparation of vector and insert DNA was processed using BamHI-HF,

189    EcoRI-HF, NdeI, NheI-HF, and XhoI enzymes (NEB). Vector DNA was dephosphorylated using quick cow

190    intestinal phosphatase (QCIP) and ligation reactions conducted using T7 DNA ligase (NEB). Purification of

191    DNA products was done by PCR cleanup (E.Z.N.A Cycle Pure Kit, Omega Bio-Tek) or gel extraction

192    (QIAquick Gel Extraction Kit, Qiagen). Assembled plasmids were transformed into *E. coli* DH10B by

193    electroporation and harvested by miniprep (E.Z.N.A. Plasmid DNA Mini Kit II, Omega Bio-Tek). All culturing

194    was performed in LB Lenox media (BioShop) spiked with 100 µg·mL$^{-1}$ of ampicillin (BioShop). Solid media

195    support was produced by dissolving agar (at a concentration of 15 g·L$^{-1}$, BioShop) in LB (Lenox) liquid

196    media preparation. Induction of Cloning Site I was controlled through isopropyl-β-D-thiogalactopyranoside

197  (IPTG, BioShop) added to achieve a working concentration of 10 mM. Evaluation of *in vivo* reporter protein

198  expression was performed using a SpectraMax M2 plate reader (Molecular Devices) to record absorbance

199  ($\lambda_{ab}$ = 600 nm) and fluorescence ($\lambda_{ex}$ = 485 nm, $\lambda_{em}$ = 510 nm, fixed gain = medium, 30 flashes per read)

200  from 100 µL aliquots of cell culture in black-walled and clear-bottomed 96-well format microplates (Greiner).

201  The ColE1 origin and all Cloning Site I, II and LacI genotypes were confirmed by Sanger Sequencing

202  services contracted through Génome Québec (centre d'expertise et de services Génome Québec).

203

204  **Construction of the single-chain tandem repeat DNA binding domain repressors.** The first single-

205  chain tandem repeat DNA binding domain (scDBD) repressor architecture constructed involved the full-

206  length duplication of the N-terminal LacI DBD sequence (residues 1 through 89), incorporating the triple-

207  mutation DFT (Y17**D**/Q18**F**/R22**T**) producing the non-functional scDBD$_{DFT/DFT}$ construct [21]. This construct

208  was produced in a three-step synthesis where two N-terminal (DBD.N) and two C-terminal (DBD.C)

209  fragments (DBD: residues 1 through 29, and LNK: resides 60 through 89) were first produced by aPCR and

210  gel purified (described in detail in Supplementary Information Section 2). The DBD.N and DBD.C fragment

211  pairs were independently fused to an unchanged intermediate PCR fragment (residues 30 through 59) by

212  SOE PCR and subjected to PCR cleanup, prior to a third SOE reaction producing the full-length construct.

213  This construct was digested (BamHI and NdeI) and gel extracted for insertion into the genetic circuit

214  (prepared by gel extraction of restriction digestion reaction with QCIP, BamHI, and NdeI). The plasmid was

215  harvested by miniprep, followed by confirmation of scDBD$_{DFT/DFT}$ construct identity by sequencing. This

216  template was then subjected to site-directed mutagenesis introducing the functional triple-mutation IAN

217  (D17**I**/F18**A**/T22**N**) at N-terminal and C-terminal duplicated DBDs producing three additional constructs by

218  SOE: scDBD$_{IAN/DFT}$, scDBD$_{DFT/IAN}$, and scDBD$_{IAN/IAN}$. These constructs were similarly inserted into the

219  genetic circuit by digestion, gel purification, and ligation, and harvested by miniprep prior to sequencing.

220  Lastly, C-terminal truncations omitting the duplicated 60 through 89 residue segment of each of the four

221  constructs were produced by PCR: scDBD$_{DFT/DFT/\Delta CT}$, scDBD$_{IAN/DFT/\Delta CT}$, scDBD$_{DFT/IAN/\Delta CT}$, and

222  scDBD$_{IAN/IAN/\Delta CT}$. Likewise, these constructs were inserted into the genetic circuit, harvested by miniprep,

223  and sequenced.

224

225     *In vivo* **evaluation of genetic circuit output.** The protocol employed to assay and evaluate genetic circuit

226     outputs was adapted from a previous publication assessing the burden imposed upon endogenous

227     expression factors by exogenous genetic circuits (36). 1 mL LB Lenox pre-cultures were seeded with

228     transformants plated onto a solid LB agar medium, under ampicillin selection. These pre-cultures were

229     grown to stationary phase (37 °C, 300 rpm, 16 hours) and their densities recorded and normalized to 2.6

230     units ($\lambda_{ab}$ = 600 nm). Density-normalized pre-cultures were passaged into a 2× concentration of LB Lenox

231     (2.6 mL volume) spiked with a 2× concentration of ampicillin (5.2 µL volume). Passaged cultures were

232     distributed in 0.3 mL aliquots into deep 96-well format culture plates across 8 wells (−IPTG: 0.3 mL $H_2O$,

233     +IPTG: 0.3 mL 20 mM IPTG). The resulting culture plate setup allows for twelve separate transformants to

234     be grown in quadruplicate at a 0.6 mL culture volume in the presence and absence of 10 mM IPTG with a

235     starting density of 0.005 absorbance units ($\lambda_{ab}$ = 600 nm). Cultures were grown with shaking (37 °C, 300

236     rpm) and sampled in 100 µL aliquots at the 6, 7, 8, and 9-hour time-points, recording their density and

237     fluorescent output. Linear regression analysis of culture density (Y: $\lambda_{ab}$ = 600 nm) as a function of time (X:

238     hours), fluorescence (Y: $\lambda_{ex}$ = 485 nm, $\lambda_{em}$ = 510 nm, gain = medium, 30 flashes per read) as a function of

239     time (X: hours), and fluorescence (Y: $\lambda_{ex}$ = 485 nm, $\lambda_{em}$ = 510 nm, gain = medium, 30 flashes per read) as

240     a function of absorbance (X: $\lambda_{ab}$ = 600 nm) demonstrate that all measurements, regardless of absence or

241     presence of 10 mM IPTG, were linear and recorded at steady-state (Figures S13−S20). Genetic circuit

242     expression output (F) is reported by taking the quotient of fluorescence (GFP) and density (ABS) from each

243     culture measurement (Equation 1).

244    
$$F = \frac{\text{GFP}}{\text{ABS}}$$
Eq. 1

245   All plotted data are reported as the arithmetic average across four measurements, with error bars indicating

246     the standard deviation of the sample. To determine whether changes to genetic circuit outputs are

247     statistically significant, two-tailed homoscedastic t-tests were performed with p-values reported for

248     populations exhibiting statistically significant difference (*i.e.,* p-value ≤ 0.001).

249

250   **RESULTS**

251   **Overview of DNA sequence redesign protocol.** To enable the synthesis of a construct containing TR

252     elements we developed a DNA sequence design protocol called TReSR (for **T**andem **R**epeat **Se**quence

253 **R**edesign) which was implemented as a script run in Python3.8 (Supplementary Information). Using the

254 strategy outlined in Figure 1, TReSR introduces silent mutations into TR sequences to reduce the potential

255 for off-target primer hybridization in PCR reactions, such as those required in aPCR and site directed

256 mutagenesis protocols.  The design protocol is conducted in five steps, beginning with the dissection of TR-

257 encoding gene cassette regions into contiguous segments encoding 5 to 7 amino acid residues each (Fig

258 1A). In the *second* step (Fig 1B), all codon combinations of silent mutations are generated for each segment

259 and melting temperatures (Tm) calculated for the forward ($T_{FF}$) and reverse complement ($T_{RR}$) homodimers,

260 along with that of the forward sequence with its reverse complement ($T_{FR}$) and for the forward sequence

261 with the reverse complement belonging to the wild-type gene ($T_{WT}$). These Tm values are used to exclude

262 sequences prone to homodimerization or hybridization with the wild-type sequence, and include sequences

263 predicted to have strong self-hybridization values. Sequences that do not meet percentile-based thresholds

264 for these thermodynamic parameters are then discarded before moving to the next step. In the third step

265 (Fig 1C) Tm values are calculated for heterodimers formed between pairs of remaining segment sequences

266 to build an interaction graph and identify a set of compatible sequences for each segment (i.e., segments

267 with unique codon combinations having minimal heterodimerization Tm's). In the fourth step (Fig 1D),

268 unique sequences are paired with adjacent segment sequences and again filtered based on calculated Tm

269 values following the same protocols performed in the second step. The fifth and final step (Fig 1E) involves

270 a depth-first search joining randomly selected paths from contiguous segment pairs. One assembled path

271 is then chosen at random and primer hybridization parameters evaluated to ensure that the chosen

272 sequence does not have significant off-target hybridization propensity that would complicate its construction

273 by aPCR.

274

275 To test the ability of TReSR to design an aPCR-compatible DNA sequence for an engineered TR-protein,

276 we chose to construct a single-chain (sc) repressor containing two identical copies of the DBD of the lactose

277 repressor (LacI), called scDBD. Native LacI interacts with DNA as a dimer, with each subunit donating an

278 N-terminal DBD that binds one half of the nearly symmetrical *lacO* operator sequence, followed by a linker

279 region and lactose-binding regulatory domain that inhibits DNA binding when bound to 1,6-allolactose or

280 its analogue isopropyl β-D-1-thiogalactopyranoside (IPTG) [20]. Our experimental construct contains two

281 copies of LacI amino acid residues 1 – 89 organized as a bead-on-a-string tandem repeat. This region of

282 LacI was selected because previously published studies have demonstrated that the duplication of helix-

283 turn-helix domains was sufficient to confer DNA binding capabilities to single-chain repressors [37, 38].

284 Thus, the scDBD repressor construct is designed to bind the operon constitutively to block transcription of

285 the downstream gene. Given the well-characterized suite of DBD-operator sequence combinations that

286 have been identified for this family, the designed scDBD repressor has the potential to expand the range

287 of transcriptional regulators that can be used in synthetic biology applications.

288

289 **Implementation of the TReSR protocol.** To make the problem of DNA sequence redesign more tractable,

290 it was necessary to first divide the targeted sequence into smaller segments encoding between 5 and 7

291 amino acid residues. The choice of a maximum of 7 residues per segment reduced the combinatorial

292 sequence space to a manageable size and also made evaluation of thermodynamic parameters more

293 efficient. This choice of segment length was also convenient for design of a sequence that would be

294 compatible with aPCR, since the DNA encoding these segments would be half the length of a typical

295 oligonucleotide primer needed for this method.

296

297 For the redesign of the scDBD TR-encoding sequences, LacI DBD residues 1 through 29 was divided into

298 5 contiguous segments (labeled A through E), and residues 60 through 89, divided into six contiguous

299 segments (F through K). For each segment (Table S1), all possible codon combinations containing silent

300 mutations were generated, yielding between 128 (segment B) and 3,072 (segment E) different DNA

301 sequences per segment. The UNAFold web server was employed to calculate Tm hybridization values ($T_{FF}$,

302 $T_{RR}$, $T_{FR}$, and $T_{WT}$) for all DNA segments, and the list of segments pruned using percentile-based thresholds

303 ($T_{FF} = 0.5$, $T_{RR} = 0.5$, $T_{FR} = 0.2$, and $T_{WT} = 0.5$). The values for these thresholds, tabulated in Table S1,

304 pruned approximately 70 to 86% of the total DNA codon combinations belonging to each segment (Fig 1B).

305 For each segment, pairwise percent sequence identity values were calculated for the filtered set of codon

306 combinations. These values were used to construct an interaction graph with vertices, representing

307 individual codon combinations, connected by edges, indicting the percent sequence identity between pairs

308 of codon combinations belonging to the graph (schematically illustrated in Fig 1C). This graph was used to

309    group codon combinations based on their shared sequence identity. A pair of codon combinations were

310    assigned the same group designation if they shared 80% sequence identity and if they had similar percent

311    identity profiles with the remaining codon combinations in the graph, as determined by a cosine similarity

312    comparison (threshold ≥ 0.9975). This grouping procedure produced between 6 (segment B) and 129

313    (segment H) distinct groupings for the set of eleven protein segments (Table S1). To reduce the total

314    number codon combinations carried forward for the remainder of the TReSR protocol, the space of codon

315    combinations was constrained to those belonging to four randomly selected groups from each segment.

316

317    To determine which codon combinations originating from adjacent protein segments were compatible for

318    synthesis by PCR, the thermodynamic parameters ($T_{FF}$, $T_{RR}$, and $T_{FR}$) for DNA sequences constructed from

319    pairs of adjacent codon combinations were compiled, specifically: codon combinations for segment A were

320    paired with those for segment B (A+B), as well as B+C, C+D, D+E, F+G, G+H, H+I, I+J, J+K and K+A. This

321    list of adjacent codon combinations was then filtered using the thresholds reported in Table S2 (20th

322    percentile for $T_{FF}$ and $T_{RR}$ with a fixed value of 80 °C for $T_{FR}$), reducing the number of paired DNA sequences

323    by approximately 13 to 40%. A depth-first search was then performed to assemble scDBD DNA template

324    sequences from the set of filtered adjacent segments using contiguous segments belonging to distinct

325    groupings. The resulting template was designed to encode LacI residues 1 through 29 (segments A to E)

326    joined with LacI residues 60 through 89 (segments F to K) to make the N-terminal DBD (DBD.N) joined to

327    a template encoding the same DBD sequence at the C-terminus (DBD.C).

328

329    We analyzed the first DNA template produced by TReSR (out of 100 templates generated) and found that

330    42 and 46 unique silent mutations were introduced into the DBD.N and DBD.C templates across a design

331    space of 59 total codons (Fig 2). The resulting DBD.N and DBD.C templates have reduced sequence

332    identity with each another, measured at 65% between the pair, and the wild-type LacI sequence, recorded

333    at 66% and 63%, respectively. Thermodynamic parameters calculated for segment sequences (Table S3)

334    showed hybridization values that are compatible with aPCR synthesis, with reduced homodimerization Tm

335    (°C) across all segments ($-58.9 \leq T_{FF} \leq 18.7$ and $-43.7 \leq T_{RR} \leq 19.6$ °C). All segment hybridization affinities

336    ($62.5 \leq T_{FR} \leq 74.4$ °C) were within the range typically required for annealing and extension steps employed

337      during PCR. An additional advantage of the TReSR-generated sequences are the low hybridization Tm

338      values with the wild-type LacI sequence ($-15.4 \leq T_{WT} \leq 43.7$ °C) suggesting that downstream PCR

339      manipulation of the scDBD sequence should be possible even with the presence of the LacI gene on the

340      same plasmid.

341

342      **Synthesis of the Tandem Repeat Repressor.** The template sequence shown in Figure 2 was partitioned

343      into twelve primers for each domain (N.1 – N.12 and C.1 – C.12 for DBD.N and DBD.C, respectively) such

344      that the 3′-termini of all primers were comprised of at least one G/C base-pair and sequence overlap with

345      adjacent primers was designed to ensure efficient assembly of primers (Tm > 64 °C) while limiting length

346      to no more than 44 bases (Table 1).  Residues 17, 18, and 22 directing the operator specificity for each

347      DBD are delivered on primers 4 and 5, named according to their triple-mutation identity (DFT: **D**17/**F**18/**T**22

348      and IAN: **I**17/**A**18/**N**22). According to predictions of thermodynamic parameters shown in Table 1, all

349      primers are expected to adopt linear secondary structures in solution ($\Delta G_F^{72°C}$ and $\Delta G_R^{72°C} > 0.0$ kcal·mol$^{-1}$)

350      preferentially favouring hybridization with their reverse complement sequences ($T_{FR} \geq 75.2$ °C) over

351      formation of undesired homodimers ($T_{FF} \leq 40.7$ °C and $T_{RR} \leq 46.3$ °C). Lastly, a comparison of predicted

352      hybridization affinities between pairs of primers was conducted to ensure successful assembly of the target

353      template DNA sequences for scDBD$_{DFT/DFT}$ and scDBD$_{IAN/IAN}$ (Fig S1). Analysis of predicted hybridization

354      Tm values suggests that all primers will hybridize with sufficient affinity to their adjacent counterparts under

355      reaction conditions employed during PCR ($T_{HYB} \geq 70$ °C) without forming side-products that result from

356      hybridization between pairs of non-adjacent primers.

357 **Table 1.** Assembly PCR oligonucleotide primers for the TReSR designed tandem repeat repressor

| Primer Name | Sequence (5′ → 3′) | Hybridization (°C) | | | Folding (kcal·mol⁻¹) | |
|---|---|---|---|---|---|---|
| | | $T_{FF}$ | $T_{RR}$ | $T_{FR}$ | $\Delta G_F^{345K}$ | $\Delta G_R^{345K}$ |
| pDBD.F | CCAGTAGTAGGTTGAGGC | -29.0 | -20.9 | 71.2 | 2.62 | 2.89 |
| pDBD.R | CAGGCTTCATTTTTTTCCTCCTTCTAGTTTAAACAAAATTATTTG | 40.3 | 38.0 | 79.1 | 1.48 | 1.63 |
| N.1 | CTAGAAGGAGGAAAAAAATGAAGCCTGTTACCCTG | -5.0 | -15.2 | 81.2 | 1.68 | 1.48 |
| N.2 | CTCTGCCACGTCATACAGGGTAACAGGCTTCATTTTTTTC | 26.7 | 30.1 | 84.7 | 2.03 | 1.93 |
| N.3 | CCTGTATGACGTGGCAGAGTATGCAGGAGTGAGC | 40.7 | 11.2 | 86.5 | 0.84 | 1.43 |
| DFT.N.4 | CTACCGTGCTAACCGTAAAATCGCTCACTCCTGCATACTCTG | -1.2 | -0.1 | 86.7 | 2.54 | 1.02 |
| IAN.N.4 | CTACATTAGAGACCGTGGCAATGCTCACTCCTGCATACTCTG | 36.6 | 37.3 | 86.7 | 0.94 | 0.93 |
| DFT.N.5 | GATTTTACGGTTAGCACGGTAGTAAACCAAGCCTCCCATG | 28.6 | 36.6 | 85.5 | 1.27 | 1.26 |
| IAN.N.5 | CATTGCCACGGTCTCTAATGTAGTAAACCAAGCCTCCCATG | -20.5 | -3.1 | 86.5 | 1.97 | 1.24 |
| N.6 | CATGGGAGGCTTGGTTTACTAC | -3.1 | -1.6 | 75.2 | 2.29 | 2.61 |
| LacI.N.F | GTAGTAAACCAAGCCTCCCATGTTTCTGCGAAAACGC | 26.2 | 27.6 | 85.7 | 1.82 | 0.22 |
| LacI.N.R | GACTCCAATGAGCAGGGATTGTTTGCCCGCCAGTTG | 25.4 | 29.5 | 88.8 | 1.64 | 1.24 |
| N.7 | CAATCCCTGCTCATTGGAGTCGCTACATCGTC | 10.6 | 14.0 | 84.4 | 1.52 | 1.57 |
| N.8 | GCGTGTAAGGCAAGGGACGATGTAGCGACTCCAATG | 26.3 | 10.6 | 87.8 | 1.57 | 1.72 |
| N.9 | GTCCCTTGCCTTACACGCCCCCTCTCAAATC | -17.0 | -4.3 | 86.7 | 1.72 | 2.26 |
| N.10 | CCTTGACTTTATGGCAGCTACGATTTGAGAGGGGGCGTG | -1.6 | 18.8 | 88.3 | 1.24 | 1.65 |
| N.11 | CGTAGCTGCCATAAAGTCAAGGGCTGACCAAATG | 18.8 | 23.4 | 84.8 | 1.34 | 1.09 |
| N.12 | CATACAAAGTGACGGGTTTCATTTGGTCAGCCCTTGAC | 5.7 | -0.6 | 85.3 | 1.44 | 1.39 |
| C.1 | CAAATGAAACCCGTCACTTTGTATGATGTAG | -4.3 | -2.6 | 77.6 | 2.04 | 1.66 |
| C.2 | GCATATTCGGCTACATCATACAAAGTGACGGGTTTC | -2.6 | -12.4 | 82.7 | 1.66 | 2.04 |
| C.3 | CTTTGTATGATGTAGCCGAATATGCAGGCGTAAG | 29.8 | 32.3 | 81.5 | 1.49 | 1.94 |
| DFT.C.4 | CTACAGTAGAGACGGTGAAGTCACTTACGCCTGCATATTCGG | 33.4 | 35.1 | 86.3 | 1.62 | 0.72 |
| IAN.C.4 | CTACATTGCTCACTGTAGCGATACTTACGCCTGCATATTCGG | 34.4 | 46.3 | 85.7 | 1.38 | 1.10 |
| DFT.C.5 | CTTCACCGTCTCTACTGTAGTCAATCAGGCGAGTCATG | 40.2 | 37.1 | 84.8 | 1.89 | 1.58 |
| IAN.C.5 | CGCTACAGTGAGCAATGTAGTCAATCAGGCGAGTCATG | 47.6 | 39.6 | 85.7 | 1.98 | 1.84 |
| C.6 | CATGACTCGCCTGATTGACTAC | -0.4 | -6.4 | 74.9 | 2.12 | 2.10 |
| LacI.C.F | GTAGTCAATCAGGCGAGTCATGTTTCTGCGAAAACGCG | 30.1 | 31.4 | 86.4 | 1.60 | 0.22 |
| LacI.C.R | GCTCTGTTTGCCCGCCAGTTG | -3.2 | 19.0 | 81.5 | 2.35 | 1.79 |
| C.7 | CTGGCGGGCAAACAGAGCCTTTTGATAGGGGTAGCAACG | 25.1 | 26.8 | 90.0 | 0.93 | 1.57 |
| C.8 | GAACTCGTTGCTACCCCTATCAAAAGG | 8.7 | 1.2 | 78.9 | 1.96 | 1.79 |
| C.9 | GATAGGGGTAGCAACGAGTTCATTGGCACTC | 1.2 | 29.8 | 83.2 | 1.79 | 1.68 |
| C.10 | CTATCTGGGAAGGTGCATGGAGTGCCAATGAACTCGTTG | 30.1 | 22.1 | 87.1 | 1.03 | 0.78 |
| C.11 | CCATGCACCTTCCCAGATAGTGGCAGCAATCAAATCGAG | 15.0 | 17.8 | 87.4 | 1.27 | 1.03 |
| C.12 | CCTATCATTACTGGTCCGCTCTCGATTTGATTGCTGCCAC | 17.8 | 15.0 | 86.7 | 1.39 | 0.41 |
| T7T.F | GAGCGGACCAGTAATGATAGGGATCC | 29.3 | 25.6 | 80.0 | 0.41 | 2.00 |
| T7T.R | GCAGCCGGATCCCTATCATTACTGGTCCG | 34.3 | 30.3 | 85.1 | 1.62 | 0.41 |
| ΔCT.R | GCAGCCGGATCCCTATCATTAACTCGTTGCTACCCCTATCAAAAGG | 34.3 | 30.3 | 88.3 | 2.00 | 1.96 |

358

359   To create the DNA cassettes encoding scDBD tandem repeat repressor containing either the DFT or IAN

360   set of mutations, PCR reactions were performed to synthesize 4 fragments for each domain, with DBD.N

361   produced by SOE of fragments 1 through 4 and DBD.C produced by SOE of fragments 4 through 7, using

362   the reaction schematic illustrated in Figure S2. Specifically, aPCR was conducted using primers: N.1

363   through N.6 to produce fragment 2 (encoding LacI residues 1 to 29 for DBD.N),  N.7 through N.12 to

364   produce fragment 4 (encoding LacI residues 60 to 89 for DBD.N), C.1 through C.6 to produce fragment 5

365   (encoding LacI residues 1 to 29 for DBD.C), and DBD.C.7 through DBD.C.12 with T7T.F and T7T.R to

366   produce fragment 7 (encoding LacI residues 60 to 89 for DBD.C). Regions of the scDBD-encoding

367   sequence that were not part of the TR targeted by TReSR were amplified using conventional PCR to make

368   fragments 1, 3 and 6. Successful production of all fragments was supported by agarose gel analysis which

369   all showed a single band at the expected molecular length (Fig S3). SOE was then used to construct the

370   larger fragments encoding the DBD.N domain (composed of fragments 1 to 4) and DBD.C domain

371   (composed of fragments 4 to 7).  To construct the full-length tandem repeat constructs scDBD$_{DFT/DFT}$ and

372   scDBD$_{IAN/IAN}$ (sequences provided in Fig S4A, B), fragments encoding DBD.N and DBD.C containing the

373   appropriate triple mutant were joined using a final SOE reaction. Agarose gel electrophoresis demonstrated

374   that all SOE reactions were successful in producing the target fragments (Fig S3).

375

376   As intended by the TReSR design protocol, the full-length template containing this TR-encoding sequence

377   could also be used to site-selectively mutate a single DBD without interference from the other DBD in the

378   TR.  Moreover, those PCR mutagenesis reactions were performed using DNA templates in a plasmid also

379   carrying the gene for the native LacI repressor.  No cross-reactivity with the primers targeting one the TR

380   domains was detected, with only the targeted PCR product being observed by agarose gel electrophoresis

381   (Figure S3) and confirmed by DNA sequencing (Fig S4). This was expected since the TReSR algorithm

382   was designed to exclude sequences that might hybridize with the WT gene (*i.e.,* $T_{WT} \leq 43.7$ °C for all

383   segment sequences in Table S3). Together, these results demonstrate that application of the TReSR

384   protocol enabled the design of TR DNA sequence templates suitable for assembly and manipulation by

385   PCR.

386

387 **Design of a three-component genetic circuit to evaluate function of scDBD constructs *in vivo*.** To

388 evaluate the function of our scDBD constructs, a three-component genetic circuit was designed (Fig 3A)

389 placing expression of the experimental scDBD repressor under inducible control to evaluate its function

390 reported by a cell-based fluorescence assay. We chose to construct our genetic circuit on a single plasmid

391 (sequence provided in Fig S5) since this was expected to reduce its burden on the fitness of its biological

392 hosts by reducing the number of replication origins and selection markers required to propagate and select

393 for the genetic circuit [36]. The genetic circuit contains three components, identified as Cloning Sites I, II,

394 and III, each responsible for delivery of the experimental scDBD repressor, eGFP reporting protein, and

395 LacI repressor protein responsible for regulation of scDBD expression, respectively. Cloning Site III

396 incorporates the gene encoding the W220F variant of the lac repressor ($LacI_{W220F}$) under the control of the

397 constitutive pLacI promoter (pLacI), labelled $pLacI(LacI_{W220F})$. This variant of the LacI repressor was

398 selected after testing a set of plasmids with combinations of repressors (LacI or $LacI_{W220F}$) paired with

399 promoters (pLacI or $pLacI^Q$), confirming the superior ability of $pLacI(LacI_{W220F})$ to repress transcription from

400 the pDBD promoter bearing $lacO^{sym}$ operator sequences in the absence of inducer while simultaneously

401 maximing output expression upon induction (Fig S6−S9 and Table S4−S8) [30]. The ability of the

402 $pLacI(LacI_{W220F})$ regulatory component to minimize the occurrence of inducer-free (*i.e.,* 'leaky') expression

403 events was required to evaluate scDBD function since the output of the genetic circuit must be reported in

404 the absence and presence of scDBD expression. Details concerning this engineering effort are included in

405 the Supplementary Information section: Engineering and Optimization of the three-Component Genetic

406 Circuit.

407

408 Cloning Site I delivers the experimental scDBD repressor constructs that were synthesized by aPCR and

409 SOE (Fig S3) using primers designed by TReSR (Table 1). The scDBD constructs were inserted into

410 Cloning Site I under the control of the promoter pDBD (Fig 3B), outfitted with a pair of $lacO^{sym}$ operator

411 sequences ($lacO^{sym}$: 5′−AATT**GTG**AGCGCT**CAC**AATT−3′), placed at core [36] and proximal [39] positions

412 relative to the RNA polymerase recruitment sequence [40]. Repression of pDBD by $LacI_{W220F}$ is mediated

413 by a specific DBD-operator interacting pair (*i.e.*, the LacI DBD containing the wild-type triple-residue

414 sequence **Y**17/**Q**18/**R**22 is selectively recruited to the $lacO^{sym}$ operator sequence) [21]. This promoter

415　architecture ensures that scDBD expression can be selectively controlled by the addition of IPTG in a dose-

416　dependent manner, minimizing the basal level of expression in the absence of the inducer (Fig S14).

417

418　The activity of the genetic circuit is reported using a third component, which delivers the genetically encoded

419　reporter, enhanced green fluorescent protein (eGFP) [32], to Cloning Site II whose promoter, pGFP (Fig

420　3B), is regulated by the expression of functional scDBD repressor constructs. Specific recruitment of

421　functional scDBD constructs to pGFP is accomplished by employing the DBD triple-mutation

422　Y17**I**/Q18**A**/R22**N** which selectively binds a variant of the symmetric lac-type operator called lacO$^{TTA}$

423　(sequence 5′−AATT**TTA**AGCGCT**TAA**AATT−3′, with bolded residues indicating site of mutations) [21]. This

424　three-component genetic circuit architecture ensures that repression of pGFP is specifically mediated by

425　functional scDBD repressor without interference from LacI$_{W220F}$ which is incorporated to regulate expression

426　of scDBD constructs. This genetic circuit setup is therefore designed to allow expression of eGFP in the

427　absence of IPTG since expression of scDBD by its promoter (pDBD) is inhibited by LacI$_{W220F}$. Conversely,

428　in the presence of IPTG, LacI$_{W220F}$ dissociates from pDBD, enabling expression of the scDBD repressor

429　candidate which, if functional, would bind to pGFP to repress expression of the reporter protein. Thus, the

430　genetic circuit functions to report on scDBD repressor activity by inverting input induction and output

431　fluorescent signal. To reduce the potential influence of junction interference on expression levels of eGFP,

432　identical T7 terminator sequences [33] were introduced at the 3′-termini of both Cloning Site I and II coding

433　regions, while upstream promoter elements were outfitted with identical riboJ genetic insulator, hairpin, and

434　ribosome binding site sequences [41, 42]. Relative expression levels were measured for pDBD and pGFP

435　promoters in a series of experiments to demarcate the minimum and maximum signal output that can be

436　produced by the genetic circuit in our chosen host expression system, with results described in detail in

437　Supplementary Information section: Expression Controls for the three-Component Genetic Circuit (Fig

438　S10−S13 and Table S9). With this data it was possible to use this single-plasmid genetic circuit to evaluate

439　the function of our designed scDBD repressors in a quantitative manner (Fig S14B).

440

441　**Evaluation of scDBD repressor function.** To evaluate the function of our scDBD repressor, four variants

442　of the genetic circuit were made from a combination of DBDs with the functional IAN (Y17**I**/Q18**A**/R22**N**)

443 and non-functional DFT (Y17**D**/Q18**F**/R22**T**) triple-mutations, incorporated into DBD.N and/or DBD.C

444 domains of our scDBD repressor construct (Fig S15-S22 and Table S10). As native lac repressor binds

445 DNA in a dimeric state, we anticipated that only scDBD repressor constructs incorporating the functional

446 IAN mutation in both DBDs would be able to bind pGFP to repress transcription of the eGFP gene. As

447 shown in Figure 4, a representative sample of the density-normalized fluorescence taken at the 8-hour time

448 point in the presence of 10 mM IPTG resulted in a 5-fold reduction in genetic circuit output signal relative

449 to that obtained for the circuit grown in the absence of IPTG. This repression of eGFP expression by the

450 scDBD repressor was only obtained when both N- and C-terminal DBDs contained the IAN mutation

451 required for recognition of the lacO$^{sym}$ variant operator (lacO$^{TTA}$: G6T/T5/G4A) incorporated in the pGFP

452 promoter. Moreover, the same result was obtained when this combination of scDBDs was truncated

453 (scDBD$_{IAN/IAN/\Delta CT}$) to eliminate the C-terminal copy of the DBD linker region (residues 61 to 89), as only the

454 variant that contained the IAN mutation in both DBDs showed a reduction in fluorescence upon addition of

455 IPTG (3.8 ± 0.2-fold decrease in density normalized fluorescence in the presence of 10 mM IPTG at the 8-

456 hour time point). These results suggest that both scDBD$_{IAN/IAN}$ and its truncated counterpart,

457 scDBD$_{IAN/IAN/\Delta CT}$, act to selectively repress expression from the pGFP promoter without the need for

458 dimerization that is characteristic of the native lac repressor. This demonstration of scDBD repressor

459 function illustrates the ability of TReSR to create new functional proteins containing TR motifs without the

460 need to resort to total gene synthesis.

461

462 **DISCUSSION**

463 We chose to demonstrate the utility of TReSR by duplicating a domain-length sequence, in this case the

464 LacI DBD, since this type of TR construct tends to be one of the most difficult to construct by aPCR methods.

465 The repression of eGFP expression via the action of the scDBD$_{IAN/IAN}$ and scDBD$_{IAN/IAN/-LNK}$ repressors

466 reported by our genetic circuit demonstrates that even for this challenging system, TReSR was able to

467 create a DNA sequence encoding the TR that allowed its cost-effective assembly by aPCR and SOE.

468 Moreover, the DNA sequence produced by TReSR was also compatible with downstream introduction of

469 mutations by PCR-based site-directed mutagenesis. TReSR design of DNA sequences therefore makes it

470 possible to avoid the well-documented difficulties that are normally associated with manipulating repeating

471    DNA sequences [15]. This presents a significant advantage over other approaches that first independently

472    engineer the function of modular domains and then assemble a final TR construct by gene synthesis or

473    DNA ligation. For example, a phage display approach has been employed to identify pairs of zinc finger

474    motifs that could be expressed together as a bead-on-a-string TR-containing protein capable of recognizing

475    DNA sequences in the HIV-1 promoter [43, 44]. However, this method cannot be readily applied to TR

476    constructs comprised of modules that do not have function when expressed as individual domains, like the

477    DBDs that were used to engineer the scDBD in this study. As demonstrated using the scDBD triple-mutant

478    variants containing a single functional DBD *(i.e.,* inactivating DFT triple mutations introduced to one of the

479    DBDs), scDBD repressors required two functional DBDs to achieve repression. The use of TReSR to create

480    TR-containing proteins with DNA sequences that allow aPCR assembly and PCR-based manipulation

481    opens the door to simultaneous screening of more than one module and increases the range of TR-

482    containing proteins that can be designed.

483

484    It is expected that our DNA sequence redesign strategy will be able to successfully accommodate the

485    duplication of at least three domain-length sequences into a single construct, since it was possible to

486    perform site-selective PCR-based manipulation of a single DBD in a plasmid that contained DBD

487    sequences from both native LacI and the scDBD construct. Moreover, TReSR brought the sequence

488    identity between segments encoding each DBD down to 66%, which can be considered a benchmark for

489    predicting the success of future sequence redesign projects (i.e. aPCR and mutagenesis should be possible

490    if a similar level of sequence identity is obtained from TReSR-designed sequences for other engineered TR

491    proteins). Using this benchmark as a guideline, we anticipate that this will likely be possible for the design

492    of a protein containing four TR domains, since most amino acids are encoded by four degenerate codons.

493    Moreover, for TRs where the repeated sequence is shorter than the domain-length sequences targeted

494    here (e.g. heptad repeat of a leucine zipper), it should be possible to use TReSR to create proteins

495    containing a larger number of TRs.

496

497    While the task of computational redesign of DNA sequences to allow PCR-based mutagenesis and

498    manipulation is not unique to this study, this is the first to identify dissimilar DNA sequence fragments prior

499    to assembly of the full-length construct. Previous strategies have been proposed where the targeted

500    sequence is fragmented into oligonucleotides without introducing codon substitutions [25], or where

501    sequence selection is rooted in thermodynamic prediction of oligonucleotide hybridization behaviours [24].

502    Similar to our strategy, both DNAWorks [45] and Gene2Oligo [46] redesign DNA sequences by

503    computationally evaluating codon substitutions conferring silent mutations which improve the

504    thermodynamic parameters of select oligonucleotides for PCR synthesis. However, neither of these

505    protocols compare DNA sequences of the fragments to reduce the similarity between them, which is critical

506    for the generation of sequences encoding protein TRs.  Although our strategy does not include codon usage

507    frequency data when redesigning DNA sequences [47], this parameter could be included in the criteria used

508    to prune codon combination lists (Fig 1B). Alternatively, low frequency codons could be removed from the

509    codon table used in the TReSR calculations, or the effect of low-frequency codons could be mitigated by

510    employing a tRNA-overexpression strategy with cell strains developed for this purpose [48]. This was not

511    required for the scDBDs designed in this study, however, since expression levels of the repressor were

512    sufficient for functional repression.

513

514    One of the results arising from our demonstration of TReSR utility is the creation of a new scDBD from the

515    DBD of the lac repressor which was capable of repressing expression from a modified *lacO* promoter. This

516    repressor design is similar to a previously engineered scDBD repressor constructed by duplicating the N-

517    terminal DBD from the bacteriophage 434 cI repressor which recognizes the symmetric 434 operator

518    sequence [37]. DNA sequence recognition of this construct could be predictably altered to produce scDBDs

519    that recognize asymmetric operators to investigate the influence of direct and indirect protein-DNA contacts

520    on repressor-operator binding [49] or to identify cognate and specific protein-DNA interacting pairs [50].

521    This same strategy for constructing scDBD repressors has also been employed with the lambda Cro

522    repressor sequence [38]. Our results with the DBD of the lac repressor show that the same strategy can be

523    extended to another well-characterized family of repressors. While the LacI DBD shares a similar helix-

524    turn-helix motif to these bacteriophage repressors, the DNA recognition helix making direct contacts with

525    the operator sequence is oriented in the opposite direction with respect to those bacteriophage repressors

526    [51]. Despite this distinction, our results demonstrate that the same strategy for constructing scDBD

527     architectures from bacteriophage repressors is readily applicable to the LacI DBD and should follow the

528     functional rules defining DBD recognition of operator DNA that had been defined with full-length LacI [21].

529     In addition, this approach has the potential to be extended to include DBDs belonging to other members of

530     the lac repressor superfamily [19], further increasing the range of promoter sequences that could be

531     recognized.

532

533     While TReSR was created for the purpose of redesigning the DNA sequence of a novel TR protein (*i.e.*,

534     the scDBD repressor constructs), this computational strategy also has the potential to be adapted to

535     applications that do not involve TRs. This would involve modification of the TReSR methodology to compare

536     non-identical protein segments, a task that could be facilitated by replacing the percent sequence identity

537     metric (employed in the *third* step of the TReSR protocol, Fig 1C) with calculation of the Tm for hybridization

538     between pairs of DNA sequences. This strategy has the potential to allow the redesign of DNA templates

539     containing problematic regions to make them amenable to PCR-based manipulations by breaking the

540     sequence down into fragments and generating codon combinations with more favorable hybridization

541     parameters. This type of sequence redesign protocol would be particularly useful for mutagenesis of DNA

542     templates in high-throughput procedures (*e.g.*, deep sequencing mutagenesis). The TReSR computational

543     strategy could also be applied to the design and selection of reliable primers for assembly of DNA barcodes

544     used in genotyping large populations of genetic samples. For this application, TReSR could be adapted to

545     compare and select primer combinations that assemble in a defined order to generate unique DNA

546     sequences (*i.e.*, barcodes) appended to amplicons in a single PCR reaction from isolated samples. These

547     samples can then be pooled for next-generation sequencing thus enabling simultaneous sample

548     identification and genotyping, provided that the amplicon length is amenable to the sequencing

549     methodology employed. Similarly, the TReSR protocol could be applied to ribozyme design strategies to

550     design interacting and non-interacting RNA sequences, which would open the door to the engineering of

551     increasingly complex genetic programs and circuitry directing control over gene expression. Overall, the

552     ability to redesign sequences using smaller segments with defined hybridization parameters lies at the core

553     of the TReSR protocol, and offers opportunities for a wide range of potential applications.

554

555 **CONCLUSION.** The PCR synthesis and manipulation of TR DNA sequences presents a prohibitive

556 challenge interfering with the routine incorporation of TR sequences in engineered proteins. To overcome

557 this barrier, we devised and implemented a DNA sequence redesign protocol (TReSR) to construct TR

558 DNA templates that are amenable to assembly and mutagenesis by PCR. TReSR predictions were

559 validated by the construction of a single-chain tandem repeat repressor, created by duplicating the DNA

560 binding domain of LacI. Experimental characterization of repressor construct function using a three-

561 component genetic circuit confirms that this new repressor is functional. The use of TReSR to create TR-

562 containing proteins with DNA sequences that allow aPCR and PCR-based manipulation opens the door to

563 simultaneous screening of both modules and increases the range of TR-containing proteins that can be

564 designed.

565

570

571 **REFERENCES**

572 1. Sinha R, Shukla P. Current trends in protein engineering: updates and progress. Curr Protein Pept Sci.

573 2019;20: 398−407.

574 2. Nandagopal N, Elowitz MB Synthetic biology: integrated gene circuits. Science. 2011;333: 1244−1248.

575 3. Pellegrini M, Marcotte EM, Yeates TO. A fast algorithm for genome-wide analysis of proteins with repeated

576 sequences. Proteins. 1999;35: 440−446.

577 4. Kajava AV. Tandem repeats in proteins: from sequence to structure. J Struct Biol. 2011;197: 279−288.

578 5. Delucchi M, Schaper E, Sachenkova O, Elofsson A, Anisimova M. A new concensus of protein tandem

579 repeats and their relationship with intrinsic disorder. Genes. 2020;11: 407−425.

580 6. Berisio R, Vitagliano L, Mazzarella L, Zagari A. Crystal structure of the collagen triple helix model [(Pro-

581 Pro-Gly)(10)](3). Protein Sci. 2002;11: 262−270.

7. Liu J, Zheng Q, Deng Y, Cheng C-S, Kallenbach NR, Lu M. A seven-helix coiled coil. Proc Natl Acad Sci USA. 2006;103: 15457−15462.

8. Neer EJ, Schmidt CJ, Nambudripad R, Smith TF. The ancient regulatory-protein family of WD-repeat proteins. Nature. 1994;371: 297−300.

9. Kobe B, Deisenhofer J. The leucine-rich repeat: a versatile binding motif. Trends Biochem Sci. 1994;19: 415−421.

10. Hatzfeld M. The armadillo family of structural proteins. Int Rev Cytol. 1999;186: 179−224.

11. Mosavi LK., Cammett TJ, Desrosiers DC, Peng Z-Y. The ankyrin repeat as molecular architecture for protein recognition. Protein Sci. 2004;13: 1435-1448.

12. Adams J, Kelso R, Cooley L. The kelch repeat superfamily of proteins: propellers of cell function. Trends Cell Biol. 2000;10: 17-24.

13. Yoshimura SH, Hirano T. HEAT repeats - versatile arrays of amphiphilic helices working in crowded environments? J Cell Sci. 2016;129: 3963−3970.

14. Paladin L, Hirsh L, Piovesan D, Andrade-Navarro MA, Kajava AV, Tosatto SCE. RepeatsDB 2.0: improved annotation, classification, search and visualization of repeat protein structures. Nucleic Acids Res. 2017;45: D308−D312.

15. Hommelsheim CM, Frantzeskakis L, Huang M, Ülker B. PCR amplification of repetitive DNA: a limitation to genome editing technologies and many other applications. Sci Rep. 2014;23: 5052−5064.

16. Stemmer WP, Crameri A, Ha KD, Brennan TM, Heyneker HL. Single-step assembly of a gene and entire plasmid from large numbers of oligodeoxyribonucleotides. Gene. 1995;164: 49−53.

17. Bryksin AV, Matsumura I. Overlap extension PCR cloning: a simple and reliable way to create recombinant plasmids. Biotechniques. 2010;48: 463−465.

18. Ho SN, Hunt HD, Horton RM, Pullen JK, Pease LR. Site-directed mutagenesis by overlap extension using the polymerase chain reaction. Gene. 1989;77: 51−59.

19. Swint-Kruse L, Matthews KS. Allostery in the LacI/GalR family: variations on a theme. Curr Opin Microbiol. 2009;12: 129−137.

20. Lewis M. The lac repressor. C R Biol. 2005;328: 521−548.

609  21. Milk L, Daber R, Lewis M. Functional rules for lac repressor-operator associations and implications for

610       protein-DNA interactions. Protein Sci. 2010;19: 1162−1172.

611  22. Swint-Kruse L, Larson C, Pettitt BM, Matthews KS. Fine-tuning function: correlation of hinge domain

612       interactions with function distinctions between LacI and PurR. Protein Sci. 2002;11: 778−794.

613  23. Tungtur S, Egan SM., Swint-Kruse L. Functional consequences of exchanging domains between LacI and

614       PurR are mediated by the intervening linker sequence. Proteins. 2007;68: 375−388.

615  24. Markham NR, Zuker M. UNAFold: software for nucleic acid folding and hybridization. Methods Mol Biol.

616       2008;453: 3−31.

617  25. Tian S, Yesselman JD, Cordero P, Das R. Primerize: automated primer assembly for transcribing non-

618       coding RNA domains. Nucleic Acids Res. 2015;43: W522−W526.

619  26. Durfee T, Nelson R, Baldwin S, Plunkette G III, Burland V, Mau B, et al. The complete genome sequence

620       of Escherichia coli DH10B: insights into the biology of a laboratory workhorse. J. Bacteriol. 2008;190:

621       2597−2606.

622  27. Hershfield V, Boyer HW, Yanofsky C, Lovett MA, Helinski DR. Plasmid ColE1 as a molecular vehicle for

623       cloning and amplification of DNA. Proc Natl Acad Sci USA. 1974;71: 3455−3459.

624  28. Steele C, Zhang S, Shillitoe EJ. Effect of different antibiotics on efficiency of transformation of bacteria by

625       electroporation. Biotechniques. 1994;17: 360−365.

626  29. Glascock CB, Weickert MJ. Using chromosomal lacI$^{Q1}$ to control expression of genes on high-copy-number

627       plasmids in Escherichia coli. Gene. 1998;223: 221−231.

628  30. Gatti-Lafranconi P, Dijkman WP, Devenish SRA, Hollfelder F. A single mutation in the core domain of the

629       lac repressor reduces leakiness. Microb Cell Fact. 2013;12: 67−76.

630  31. Quan J, Tian J. Circular polymerase extension cloning of complex gene libraries and pathways. PLoS One.

631       2009;4: e6441−e6446.

632  32. Cormack BP, Valdivia RH, Falkow S. FACS-optimized mutants of the green fluorescent protein (GFP).

633       Gene. 1996;173: 33−38.

634  33. Neff NF, Chamberlin MJ. Termination of transcription by Escherichia coli ribonucleic acid polymerase in

635       vitro. Effect of altered reaction conditions and mutations in the enzyme protein on termination with T7 and

636       T3 deoxyribonucleic acids. Biochemistry. 1980;19: 3005−3015.

34. Prasher DC, Eckenrode VK, Ward WW, Prendergast FG, Cormier MJ. Primary structure of the Aequorea victoria green-fluorescent protein. Gene. 1992;111: 229−233.

35. Barondeau DP, Putnam CD, Kassmann CJ, Tainer JA, Getzoff ED. Mechanism and energetics of green fluorescent protein chromophore synthesis revealed by trapped intermediate structures. Proc Natl Acad Sci USA. 2003;100: 12111−12116.

36. Davey JA, Wilson CJ. Engineered signal-coupled inducible promoters: measuring the apparent RNA-polymerase resource budget. Nucleic Acids Res. 2020;48: 9995−10012.

37. Simoncsits A, Chen J, Percipalle P, Want S, Törö I, Pongor S. Single-chain repressors containing engineered DNA-binding domains of the phage 434 repressor recognize symmetric or asymmetric DNA operators. J Mol Biol. 1997;267: 118−131.

38. Jana R, Hazbun TR, Fields JD, Mossing MC. Single-chain lambda Cro repressors confirm high intrinsic dimer-DNA affinity. Biochemistry. 1998;37: 6446−6455.

39. Garcia HG, Sanchez A, Boedicker JQ, Osborne M, Gelles J, Kondov J, et al. Operator sequence alters gene expression independently of transcription factor occupancy in bacteria. Cell Rep. 2012;2: 150−161.

40. Cox RS III, Surette MG, Elowitz MB. Programming gene expression with combinatorial promoters. Mol Syst Biol. 2007;3: 145−155.

41. Clifton KP, Jones EM, Paudel S, Marken JP, Monette CE, Halleran AD, et al. The genetic insulator RiboJ increases expression of insulated genes. J Biol Eng. 2018;12: 23−28.

42. Lou C, Stanton B, Chen Y-J, Munsky B, Voight CA. Ribozyme-based insulator parts buffer synthetic circuits from genetic context. Nat Biotech. 2012;30: 1137−1142.

43. Isalan M, Klug A, Choo Y. A rapid, generally applicable method to engineer zinc fingers illustrated by targeting the HIV-1 promoter. Nat Biotech. 2001;19: 656−660.

44. Pavletich NP, Pabo CO. Zinc finger-DNA recognition: crystal structure of a Zif268-DNA complex at 2.1 A. Science. 1991;252: 809−817.

45. Hoover DM, Lubkowski J. DNAWorks: an automated method for designing oligonucleotides for PCR-based gene synthesis. Nucleic Acids Res. 2002;30: e43−e39.

46. Rouillard J-M, Lee W, Truan G, Gao X, Zhou X, Gulari E. Gene2Oligo: oligonucleotide design for in vitro gene synthesis. Nucleic Acids Res. 2004;32: W176−W180.

665 47. Nowak RM, Wojtowicz-Krawiec A, Plucienniczak A. DNASynth: A computer program for assembly of

666   artificial gene parts in decreasing temperature. Biomed Res Int. 2015;2015: 413262−413270.

667 48. Lipinszki Z, Vernyik V, Farago N, Sari T, Puskas LG, Blattner FR, et al. Enhancing the translational capacity

668   of E. coli by resolving the codon bias. ACS Synth Biol. 2018;7: 2656−2664.

669 49. Chen J, Pongor S, Simoncsits A. Recognition of DNA by single-chain derivatives of the phage 434

670   repressor: high affinity binding depends on both the contacted and non-contacted base pairs. Nucleic Acids

671   Res. 1997;25: 2047−2054.

672 50. Simoncsits A, Tjörnhammar ML, Wang S, Pongor S. Single-chain 434 repressors with altered DNA-binding

673   specificities. Isolation of mutant single-chain repressors by phenotypic screening of combinatorial mutant

674   libraries. Genetica. 1999;106: 85−92.

675 51. Lehming N, Sartorius J, Oehler S, von Wilcken-Bergmann B, Müller-Hill B. Recognition helices of lac and

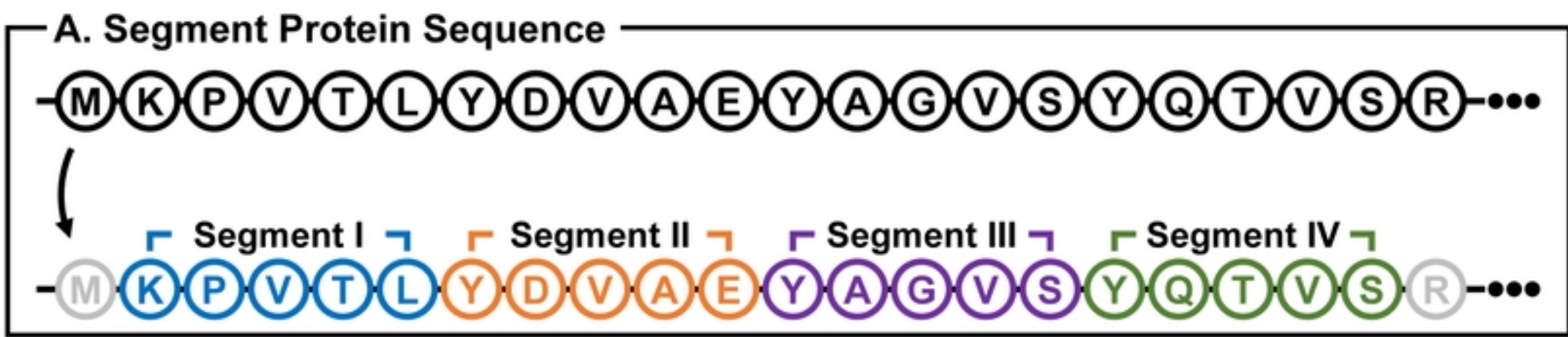676   lambda repressor are oriented in opposite directions and recognize similar DNA sequences. Proc Natl Acad

677   Sci USA. 1988;85: 7947−7951.

Figure 1

Figure 2

Figure 3

**A.**

IAN: I17/A18/N22 (functional DBD triple-mutation)

DFT: D17/F18/T22 (non-functional DBD triple-mutation)

**B.**

$4.5 \times 10^{-10}$

**C.**

$1.3 \times 10^{-9}$

Figure 4