1    **Characterization of Primed Adaptation in the *Escherichia coli* type I-E CRISPR-Cas**

2    **System**

3

4

5    Anne M. Stringer[1], Lauren A. Cooper[2], Sujatha Kadaba[1], Shailab Shrestha[1] and Joseph T.

6    Wade[1,2*]

7

8

9    [1]Wadsworth Center, New York State Department of Health, Albany, New York, USA

10   [2]Department of Biomedical Sciences, School of Public Health, University at Albany, Albany,

11   New York, USA.

12

13

14   [*]Corresponding author:

15   joseph.wade@health.ny.gov

16

**ABSTRACT**

CRISPR-Cas systems are bacterial immune systems that target invading nucleic acid. The hallmark of CRISPR-Cas systems is the CRISPR array, a genetic locus that includes short sequences known as "spacers", that are derived from invading nucleic acid. Upon exposure to an invading nucleic acid molecule, bacteria/archaea with functional CRISPR-Cas systems can add new spacers to their CRISPR arrays in a process known as "adaptation". In type I CRISPR-Cas systems, which represent the majority of CRISPR-Cas systems found in nature, adaptation can occur by two mechanisms: naïve and primed. Here, we show that, for the archetypal type I-E CRISPR-Cas system from *Escherichia coli*, primed adaptation occurs at least 1,000 times more efficiently than naïve adaptation. By initiating primed adaptation on the *E. coli* chromosome, we show that spacers can be acquired across distances of >100 kb from the initially targeted site, and we identify multiple factors that influence the efficiency with which sequences are acquired as new spacers. Thus, our data provide insight into the mechanism of primed adaptation.

31 **INTRODUCTION**

32

33 CRISPR-Cas systems are adaptive immune systems that are widespread throughout the bacterial

34 and archaeal kingdoms and protect cells from invading nucleic acid molecules. Although the

35 mechanistic details vary widely across different CRISPR-Cas systems, the basic mechanism is

36 shared across all systems (Wright et al., 2016). A short RNA, known as a CRISPR RNA

37 (crRNA) associates with a surveillance CRISPR-associated (Cas) protein, or complex of Cas

38 proteins. This protein-RNA complex binds to the invading nucleic acid molecule, at least in part

39 due to base-pairing between the crRNA and complementary sequence in the invading nucleic

40 acid (the "protospacer"). This leads to cleavage of the invading nucleic acid either by the

41 surveillance protein-RNA complex, or by an additional protein that is recruited to this complex.

42 The most abundant type of CRISPR-Cas system in nature is the type I system (Koonin et al.,

43 2017); the archetypal type I system is the type I-E system of *Escherichia coli*, which employs a

44 large complex of five Cas proteins (some proteins present in multiple copies), "Cascade", as the

45 surveillance complex, and recruits a separate nuclease protein, Cas3 (Brouns et al., 2008).

46

47 A hallmark of all CRISPR-Cas systems is the CRISPR array, a genetic locus consisting of

48 multiple units of short (~20-40 nt) "repeat" and "spacer" sequences (Wright et al., 2016). Within

49 a given CRISPR array, repeat sequences are identical, and each repeat is separated from the next

50 by a spacer. Spacer sequences are variable, and are derived from invading nucleic acid

51 molecules. The CRISPR array is typically transcribed as a single, long RNA, which is then

52 processed by a Cas protein to yield individual crRNAs. Each crRNA contains a single spacer

53 flanked by partial repeats, and the spacer sequence provides the base-pairing complementarity to

54    a protospacer in the nucleic acid target. For CRISPR-Cas systems targeting DNA, an additional

55    sequence in the DNA target is typically required for CRISPR-Cas immunity. This short sequence

56    (typically 2-4 bp), known as the protospacer-adjacent motif (PAM), is located close to the

57    protospacer, and serves as a binding site for one of the Cas proteins in the surveillance complex

58    (Leenay and Beisel, 2017).

59

60    A key feature of all CRISPR-Cas systems is their ability to acquire new immunity elements, in a

61    process known as adaptation (Jackson et al., 2017; Sternberg et al., 2016). During adaptation,

62    one or more spacer+repeat units are added to the end of a CRISPR array. For type I CRISPR-Cas

63    systems, there are two mechanisms of adaptation: naïve and primed. Naïve adaptation occurs

64    when an organism is invaded by a nucleic acid sequence that is not already matched by an

65    existing spacer. Naïve adaptation typically requires only the Cas1 and Cas2 proteins, and has

66    been proposed to occur by Cas1-2 acquisition of DNA fragments generated by resection of

67    double-strand breaks (Levy et al., 2015). By contrast, primed adaptation requires invasion of the

68    cell by a nucleic acid sequence for which there is already a partially or completely matching

69    spacer. The mechanism of primed adaptation is not well understood, but requires all Cas proteins

70    (Datsenko et al., 2012). Targeting of a protospacer by Cascade with a partially or fully

71    complementary spacer leads to recruitment of Cas3, which is both a nuclease and a helicase.

72    Primed adaptation likely involves translocation of Cas3 away from the protospacer-bound

73    Cascade (Dillard et al., 2018; Redding et al., 2015), with Cas3 generating the substrates for

74    Cas1-2 to integrate new spacers into the CRISPR array (Dillard et al., 2018; Künne et al., 2016;

75    Semenova et al., 2016). However, the mechanism by which Cas3 and/or other proteins generate

76    substrates for Cas1-2 is poorly understood (Musharova et al., 2017).

77

78    For the purposes of this study, we refer to the source of a spacer as the "pre-spacer" location, and

79    we refer to the site of Cascade binding as the protospacer. The sequence context of spacers

80    acquired during primed adaptation has a large impact on the efficiency with which a given

81    sequence is acquired (Datsenko et al., 2012; Fineran et al., 2014; Li et al., 2014; Musharova et

82    al., 2018, 2018; Rao et al., 2017; Richter et al., 2014; Savitskaya et al., 2013; Shmakov et al.,

83    2014; Staals et al., 2016; Strotskaya et al., 2017; Swarts et al., 2012). Previous studies have

84    shown that pre-spacers are located predominantly upstream of the protospacer, consistent with

85    the unidirectional helicase activity of Cas3. For type I-F systems, pre-spacers are equally

86    abundant on both the target and non-target strands (Richter et al., 2014; Staals et al., 2016),

87    whereas for type I-B, I-C, and I-E systems, pre-spacers are located predominantly on the non-

88    target strand (Datsenko et al., 2012; Fineran et al., 2014; Li et al., 2014; Rao et al., 2017;

89    Savitskaya et al., 2013; Shmakov et al., 2014; Strotskaya et al., 2017; Swarts et al., 2012). For

90    type I-F systems, pre-spacers are located over a ~5 kb region upstream of the protospacer (Staals

91    et al., 2016), whereas for type I-E systems, acquisition has been reported over distances of up to

92    ~10 kb (Strotskaya et al., 2017), although most studies involved plasmids of <10 kbp for which

93    the distance of primed adaptation cannot be measured since it likely exceeds the size of the

94    plasmid (Datsenko et al., 2012; Fineran et al., 2014; Savitskaya et al., 2013; Shmakov et al.,

95    2014). For all type I systems tested to date, most pre-spacers are flanked by an optimal PAM

96    sequence (Datsenko et al., 2012; Fineran et al., 2014; Li et al., 2014; Rao et al., 2017; Richter et

97    al., 2014; Savitskaya et al., 2013; Shmakov et al., 2014; Staals et al., 2016; Strotskaya et al.,

98    2017; Swarts et al., 2012); indeed, pre-spacers that are not immediately flanked by an optimal

99    PAM sequence are often the result of (i) "slipping" errors, where an optimal PAM is located one

100   or two bases away from the expected location, suggesting incorrect processing of the DNA by

101   Cas1-2 (Li et al., 2014; Rao et al., 2017; Shmakov et al., 2014; Staals et al., 2016), or (ii)

102   "flipping" errors, where the PAM is on the opposite side and strand of the pre-spacer, suggesting

103   that Cas1-2 has incorrectly inserted the spacer into the CRISPR array in the reverse orientation

104   (Li et al., 2014; Rao et al., 2017; Shmakov et al., 2014; Staals et al., 2016).

105

106   Here, we further investigate the process of primed adaptation in the type I-E system of *E. coli*.

107   Previous studies have shown that primed adaptation is considerably more efficient than naïve

108   adaptation in type I-B, I-E and I-F systems (Datsenko et al., 2012; Li et al., 2014; Staals et al.,

109   2016). We quantify this difference for the *E. coli* system, showing that primed adaptation is at

110   least 1,000 times more efficient than naïve adaptation. We observe primed adaptation from the *E.*

111   *coli* chromosome over distances >100 kb. Using the rich source of pre-spacers inferred from

112   these data, we identify 7 features of pre-spacers that determine the efficiency with which they are

113   acquired, supporting and extending previous observations: (i) the DNA strand of the pre-spacer

114   relative to the protospacer (i.e. target vs non-target strand), (ii) the position of the pre-spacer

115   upstream or downstream of the protospacer, (iii) the presence and position of an AAG PAM

116   immediately adjacent to the pre-spacer, (iv) distance of the pre-spacer from the protospacer over

117   a ~100 kb region, (v) the presence of an AAG within the pre-spacer, (vi) whether or not the pre-

118   spacer is acquired from within ~200 nt of the protospacer, (vii) the sequence of positions 30-33

119   of the pre-spacer. Based on these factors, we propose a model of primed adaptation in which a

120   Cas3-containing complex travels large distances along the DNA from the protospacer, cutting

121   the DNA at or near every AAG, contributing to the generation of substrates for Cas1-2.

122

123    **RESULTS AND DISCUSSION**

124

125    **Primed adaptation is >1,000-fold more efficient than naïve adaptation in *E. coli***

126    Studies of naïve adaptation in *E. coli* have involved high-level overexpression of Cas1 and Cas2

127    (Levy et al., 2015; Yosef et al., 2013). Hence, it is unclear how the efficiency of naïve adaptation

128    compares to that of primed adaptation. We previously used a highly sensitive, fluorescence-

129    based reporter (Amlinger et al., 2017) to precisely quantify primed adaptation in *E. coli* (Cooper

130    et al., 2018). For protospacers with mismatches in the PAM-proximal "seed" region, or extensive

131    mismatches in the PAM-distal region, we failed to detect adaptation, strongly suggesting that

132    naïve adaptation is extremely inefficient in these cells. To directly compare the efficiencies of

133    primed and naïve adaptation in genetically similar strains, we used a Δ*cas3* strain of *E. coli* in

134    which all other *cas* genes are constitutively expressed from their native locus. This strain also

135    contains a minimal CRISPR array (i.e. one repeat, and a leader sequence) associated with an out-

136    of-frame copy of *yfp*. Expansion of this artificial array by a single repeat/spacer unit puts *yfp* in

137    frame, leading to YFP fluorescence (Figure 1A). We then introduced a plasmid expressing *cas3*

138    from an arabinose-inducible promoter, or an equivalent empty vector. We also introduced a

139    plasmid expressing a crRNA that targets a protospacer on the same plasmid ("self-targeting

140    plasmid"; stp), or an equivalent empty vector. Thus, we were able to measure the level of

141    adaptation in four different strains: (i) $cas3^+$ $crRNA^+$, (ii) $cas3^-$ $crRNA^+$, (iii) $cas3^+$ $crRNA^-$, (iv)

142    $cas3^-$ $crRNA^-$. We detected robust adaptation in cells expressing *cas3* and containing the stp

143    ($cas3^+$ $crRNA^+$). By contrast, we detected background levels of fluorescence in cells lacking

144    *cas3* and/or the stp crRNA (Figure 1B; 0-2 $YFP^+$ cells from a total of ~80,000 each). Thus, we

145    detected robust primed adaptation, but no naïve adaptation in a genetically similar strain. We

146   cannot rule out the possibility that naïve adaptation occurs at a frequency below detection using

147   this highly sensitive assay; nonetheless, we can conclude that primed adaptation is at least 1,000-

148   fold more efficient than naïve adaptation in cells expressing *cas* genes at equivalent levels.

149   Studies of other type I systems suggest that naïve adaptation is similarly inefficient, with primed

150   adaptation being >500-fold more efficient than naïve adaptation in a type I-F system (Staals et

151   al., 2016), and naïve adaptation being undetectable in a type I-B system (Li et al., 2014). This

152   suggests that CRISPR immunity is extremely inefficient for prokaryotes with type I systems

153   when encountering an invading DNA molecule that is not already at least a partial match to an

154   existing spacer. Hence, CRISPR immunity in such a situation would likely only be effective in

155   the context of a large, clonal population of bacteria/archaea.

156

157   **Mapping newly acquired spacers suggests primed adaptation can occur from chromosomal**

158   **sites**

159   Previous studies have shown that the majority of pre-spacers acquired by primed adaptation in *E.*

160   *coli* are located on the non-target strand, and are associated with an AAG PAM (Datsenko et al.,

161   2012; Fineran et al., 2014; Savitskaya et al., 2013; Swarts et al., 2012). To determine whether the

162   same features are true for our experimental set-up, we induced primed adaptation using the stp,

163   and determined the location of pre-spacers by PCR-amplification and sequencing of the

164   expanded CRISPR arrays. As expected, the majority of acquired spacers mapped to the stp

165   (Figure 2A), and the distribution of pre-spacers on the stp was independent of the CRISPR array

166   in which the spacer was acquired (*E.* coli has two CRISPR arrays) (Figure 2B). 75.3% of pre-

167   spacers were located on the same DNA strand as the protospacer, and 94.4% of spacers were

168   associated with an AAG PAM. Thus, our data are consistent with earlier studies of primed

169    adaptation. Intriguingly, many pre-spacers did not map to the stp, suggesting another source for

170    the acquired spacers. We determined whether any of the unmapped pre-spacers map to the *cas3*-

171    expressing plasmid or to the chromosome. Many of the reads map uniquely to the *cas3*-

172    expressing plasmid (Figure 2C), although it is important to note that much of the sequence of this

173    plasmid is identical to that of the stp, and reads that match shared regions were mapped to the

174    stp. Unexpectedly, a small proportion of the reads mapped uniquely to the chromosome. We

175    speculated that pre-spacers mapping uniquely to the *cas3*-expressing plasmid or the chromosome

176    represent a second round of primed adaptation. This is possible because the strain we used

177    contains both the CRISPR-I and CRISPR-II arrays. Hence, even though we sequenced the first

178    acquired spacer in each array, this spacer may have been acquired as a result of primed

179    adaptation using a newly acquired spacer in the other array. Consistent with this hypothesis, pre-

180    spacers mapping uniquely to the chromosome were found adjacent to sequences shared with the

181    stp, i.e. the *araC* gene (Figure 2D), or the transcription terminators of rRNA loci. Thus, our data

182    indicate that primed adaptation from chromosomal pre-spacers can be readily detected using this

183    approach.

184

185    **Primed adaptation occurs over a distance of >100 kb from the protospacer**

186    Previous studies suggest that primed adaptation in *E. coli* occurs over relatively long distances,

187    exceeding the size of a typical plasmid (Strotskaya et al., 2017). Given that we observed primed

188    adaptation from chromosomal sites when targeting the stp with a crRNA, we reasoned that

189    targeting a chromosomal site with a crRNA would lead to extensive primed adaptation from the

190    chromosome, allowing us to infer the distance over which primed adaptation occurs. In two

191    independent experiments, we initiated primed adaptation from protospacers located (i)

192    immediately upstream of the *lacZ* gene, with the target strand of the protospacer on the plus

193    strand of the genome ("*lacZ+*" protospacer), and (ii) ~10 kb from the *lacZ* gene, inside the *mhpT*

194    gene, with the target strand of the protospacer on the minus strand of the genome ("*mhpT-*"

195    protospacer). Both protospacers have an AGG PAM, which can induce both interference and

196    primed adaptation (Cooper et al., 2018; Fineran et al., 2014; Musharova et al., 2019; Xue et al.,

197    2015). Moreover, we previously showed that targeting the *lacZ+* protospacer with a crRNA

198    leads to robust association of Cascade (Cooper et al., 2018). Following induction of primed

199    adaptation for each of the two protospacers, we inferred the location of pre-spacers by PCR-

200    amplification and sequencing of the first acquired spacer in the expanded CRISPR-II array.

201    Thus, we observed robust primed adaptation on the chromosome, centered at each of the

202    protospacers (Figure 3). Consistent with primed adaptation studies using plasmids, the majority

203    of pre-spacers were located on the non-target strand, upstream of the protospacer (Figure 3).

204    Remarkably, we detected primed adaptation over a distance of >100 kb upstream of each of the

205    two protospacers, with the extent of primed adaptation decreasing as a function of distance from

206    the protospacer, albeit with considerable local variation (Figure 3). Of note, a recent study

207    described priming on the *E. coli* chromosome over a similar distance (Shiriaeva et al., 2019). Our

208    data strongly suggest that Cas3 can translocate >100 kb from the protospacer during primed

209    adaptation. Since the strain we used has both CRISPR arrays intact, it is possible that some

210    spacers acquired in CRISPR-II represent two rounds of primed adaptation, with the first round of

211    primed adaptation leading to spacer acquisition in CRISPR-I. However, we believe this is a rare

212    occurrence. There is an insertion element ~15 kb away from the *lacZ+* protospacer that contains

213    many frequently used pre-spacers. This insertion element sequence is duplicated at several other

214    genomic locations, but we observed very few pre-spacers adjacent to the duplicated regions.

215     Hence, we conclude that very few spacers acquired in CRISPR-I from this insertion element led

216     to a second round of primed adaptation in CRISPR-II. Nonetheless, even if a substantial

217     proportion of the pre-spacers we observed result from two rounds of primed adaptation, a single

218     round of primed adaptation must still be able to occur over a distance of >50 kb.

219

220     The large distance over which we observed primed adaptation explains why the efficiency of

221     primed adaptation from a plasmid does not appear to decrease as a function of distance from the

222     protospacer; presumably, Cas3 translocates many times around the circular plasmid. It also

223     suggests that primed adaptation would facilitate acquisition of spacers from any region of even

224     large bacteriophage genomes or plasmid, which may increase the effectiveness of the immune

225     response. The large distance over which primed adaptation occurs on the chromosome also

226     means that there are many more unique pre-spacers than used during primed adaptation from a

227     plasmid. We reasoned that this much larger set of pre-spacers would allow for a more in-depth

228     analysis of the pre-spacer features that influence the efficiency of primed adaptation.

229

230     **Off-target Cascade binding sites are not associated with primed adaptation**

231     We previously showed that Cascade association with a DNA target requires only limited base-

232     pairing between the crRNA spacer and the protospacer. Hence, Cascade binds to many

233     chromosomal sites that typically have between 5 and 10 nt matches to the 5' end of a crRNA

234     spacer, immediately flanked by an optimal (AAG) PAM (Cooper et al., 2018). When targeting

235     the *lacZ*+ protospacer with a crRNA, we observed Cascade binding to ~100 off-target sites

236     (Cooper et al., 2018). We determined the level of primed adaptation from the non-target strand in

237     the 10 kb upstream each of the 76 off-target sites associated with a seed sequence resembling

238    that of the on-target *lacZ+* protospacer. The level of Cascade association and the local pre-spacer

239    usage frequency were not significantly correlated (Figure 4; Spearman's Correlation Test, two-

240    tailed, $p = 0.72$), with pre-spacer usage near off-target sites typically being due to duplicated

241    sequences with identical copies close to the *lacZ+* protospacer. Moreover, the local pre-spacer

242    usage frequency in regions adjacent to off-target Cascade binding sites was not significantly

243    higher than that of 1000 randomly selected genomic locations (Mann-Whitney U Test $p = 0.054$).

244    Thus, our data suggest that off-target Cascade binding events do not lead to primed adaptation,

245    consistent with our previous observation that extensive PAM-distal mismatches between a

246    crRNA spacer and its cognate protospacer prevent primed adaptation (Cooper et al., 2018).

247

248    **Uneven distribution of chromosomal pre-spacers relative to the protospacer**

249    The vast majority of pre-spacers we identified are 33 nt long (98.6% for both the *lacZ+* and

250    *mhpT-* protospacers; note that pre-spacers include the last base of the PAM, which is known to

251    be added to the CRISPR array by Cas1-2; (Goren et al., 2012)). Unless specifically stated, the

252    analyses described below focus exclusively on the 33 nt pre-spacers. We also excluded genomic

253    regions covered by insertion elements, repetitive sequences that confound alignment of DNA

254    sequence reads. Lastly, since the independent replicate datasets for each protospacer were highly

255    similar (Figure 5A + B), all analyses described below use datasets where the two replicates were

256    combined.

257

258    Although the majority of pre-spacers were located on the non-target strand, upstream of the

259    protospacers, we also observed a low level of primed adaptation on the target strand, upstream of

260    the protospacer (Figure 3 + 6A). As for pre-spacers on the non-target strand, the frequency of

261    pre-spacer usage on the target-strand decreased as a function of distance from the protospacer,

262    albeit with considerable local variation. Strikingly, the fraction of pre-spacers on the target

263    strand, upstream of the protospacer, was ~4-fold higher for the *mhpT-* protospacer than for the

264    *lacZ+* protospacer, suggesting that the protospacer sequence context determines the bias in

265    primed adaptation towards the non-target strand (Figure 3 + 6A). Since the two protospacers we

266    used are only ~10 kb apart, and are in opposite DNA orientations, many of the *target* strand pre-

267    spacers for the *mhpT-* protospacer are from the same genomic region and orientation as *non-*

268    *target* strand pre-spacers for the *lacZ+* protospacer (pink box in Figure 3). Indeed, we observed

269    very similar profiles of pre-spacer location and abundance across this region for the two datasets

270    (Figure 6B), indicating that primed adaptation on the target strand likely occurs by the same

271    mechanism as primed adaptation on the non-target strand.

272

273    The strong bias in pre-spacer usage to the region on the non-target strand, upstream of the

274    protospacer, presumably reflects a bias in the direction of Cas3 translocation away. We propose

275    that Cas3 translocation is strongly biased to the non-target strand because of the inherent

276    asymmetry in the Cascade-bound protospacer. In type I-F systems, the only characterized type I

277    systems where Cas3 translocation appears to occur with equal efficiency on both the target and

278    non-target strands, Cas2 is fused to Cas3, such that there are two Cas3 subunits in the Cas1-2-3

279    complex (Fagerlund et al., 2017). We suggest that this allows translocation of Cas2-3 in either

280    strand with roughly equal efficiency, presumably using a different Cas3 subunit for each of the

281    two strands. However, there are likely to be other factors that influence the degree to which Cas3

282    selects strand, since our data show that the ratio of target:non-target strand usage varies between

283    protospacers (Figure 6A).

284

285 **Acquisition of spacers with non-AAG PAMs is frequently due to slipping**

286 We examined all pre-spacers associated with the *lacZ+* and *mhpT-* protospacers on the non-target

287 strand, in the 10 kb region upstream of the protospacer. 95.0% and 96.0% of pre-spacer usage

288 was associated with an AAG PAM for the *lacZ+* and *mhpT-* protospacers, respectively. While

289 the large majority of pre-spacer usage is associated with an AAG PAM, we were interested to

290 determine the basis for selecting pre-spacers with a non-AAG PAM. Previous studies of type I-E,

291 type I-B, type I-C and type I-F CRISPR-Cas systems have shown that pre-spacers with non-

292 canonical PAMs can be acquired due to a phenomenon known as "slipping", whereby the pre-

293 spacer with a non-AAG PAM is positioned one or two nucleotides from a canonical, AAG-

294 associated pre-spacer (Li et al., 2017; Rao et al., 2017; Shmakov et al., 2014; Staals et al., 2016).

295 To determine the frequency of slipping in our data, we calculated all pairwise distances for non-

296 AAG pre-spacers to AAG pre-spacers on the non-target strand, in the 10 kb upstream of each of

297 the protospacers. We observed a strong enrichment of pre-spacers that could be assigned to "-1",

298 "+1", or "+2" slips (Figure 7A + B), with +1 slips being the most frequent, accounting for 37%

299 of all non-AAG pre-spacer usage. Our data indicate that together, +1 slips and -1 slips represent

300 the majority (60%) of primed spacer acquisition on the non-target strand for pre-spacers not

301 associated with an AAG PAM. The remaining pre-spacers that lack an AAG PAM and are not

302 associated with slipping events typically have suboptimal PAMs, with the third base of the PAM

303 strongly enriched for G, and the second base moderately enriched for A or C (Figure 7C). Thus,

304 our data suggest that while the specificity for an AAG PAM is high during primed adaptation,

305 there is a hierarchy of selectivity within the PAM, with the third position being the most

306 important, and the first position being the least important.

307

308    Analysis of slipping in a type I-F CRISPR-Cas system demonstrated that slipped pre-spacers are

309    associated with PAMs that lead to efficient primed adaptation (Jackson et al., 2019). Our data

310    indicate that the most common form of slipping is a +1 slip, which would lead to an AGN PAM

311    (Figure 7B). The next most common form of slipping is a -1 slip, which would lead to a NAA

312    PAM (Figure 7B). Previous studies of PAM sequences that lead to interference and/or primed

313    adaptation indicate that of all the possible AGN and NAA PAMs, only AGT, and possibly CAA,

314    fail to lead to interference and/or primed adaptation (Musharova et al., 2019; Xue et al., 2015).

315    Thus, slipping in the *E. coli* CRISPR-Cas system likely leads to acquisition of a functional

316    spacer in most cases.

317

318    Previous studies of naïve adaptation showed only ~25% of pre-spacer usage is associated with

319    AAG PAMs (Levy et al., 2015; Yosef et al., 2013), in contrast to the >95% observed for primed

320    adaptation in our study and other studies (Savitskaya et al., 2013). One possible explanation is

321    that the high level of Cas1-2 overexpression required to detect naïve adaptation causes a

322    reduction in PAM specificity. An alternative hypothesis is that Cas1-2 association with the

323    translocating Cas3-containing complex during priming leads to an increase in specificity for

324    AAG; Cas3 cuts DNA with some sequence specificity (Künne et al., 2016), although this low

325    level of specificity alone is unlikely to account for the large difference in PAM specificity

326    between naïve and primed adaptation. Cas8e, a subunit of Cascade, binds specifically to AAG

327    PAMs. It is possible that Cas8e translocates with Cas3 and facilitates nicking at AAG sequences.

328    In support of this idea, Cas8 is fused to Cas3 in some bacterial species (Westra et al., 2012).

329    Moreover, Cas8e associates only weakly with the rest of the Cascade complex, and frequently

330      dissociates from the complex (Jore et al., 2011), in particular when there are suboptimal bases in

331      the PAM or in the PAM-proximal region of the protospacer (Jung et al., 2017). However, *in vitro*

332      reconstitution of Cas3 translocation for the *Thermobifida fusca* Cas3 suggests that translocating

333      Cas3 does not associate with Cas8e (Dillard et al., 2018).

334

335      **Acquisition of spacers from the target strand is frequently due to flipping, often following a**

336      **slipping event**

337      We examined all pre-spacers associated with the *lacZ+* and *mhpT-* protospacers on the target

338      strand, in the 10 kb downstream of the protospacer, i.e. on the opposite strand to where the most

339      frequently used pre-spacers are located. 2.4% and 2.1% of pre-spacer usage within 10 kb of the

340      protospacer, on this side (i.e. upstream of the protospacer on the non-target strand, or

341      downstream of the protospacer on the target strand) is on the target strand for the *lacZ+* and

342      *mhpT-* protospacers, respectively. While the large majority of pre-spacer usage is on the non-

343      target strand, we were interested to determine the basis for selecting pre-spacers from the target

344      strand. Previous studies of type I-E, type I-B, type I-C and type I-F CRISPR-Cas systems have

345      shown that pre-spacers with non-canonical PAMs can be acquired due to a phenomenon known

346      as "flipping", whereby a pre-spacer is selected from the non-target strand but then inserted into

347      the CRISPR array in the opposite orientation (Figure 8A) (Li et al., 2017; Rao et al., 2017;

348      Shmakov et al., 2014; Staals et al., 2016). To determine the frequency of flipping in our data, we

349      calculated all pairwise distances for target-strand pre-spacers to AAG pre-spacers on the non-

350      target strand, in the 10 kb upstream of each of the protospacers. We reasoned that a distance of

351      32 nt would represent a flipping event involving a canonical pre-spacer with an AAG PAM. We

352      observed an enrichment of distances that support flipping events involving spacers with a -1, +1,

353   or +2 nt slip, in addition to flipping events involving canonical (i.e. no slip) pre-spacers (Figure

354   8B). Thus, relative to the frequency of pre-spacer usage for the canonical pre-spacers, flipped

355   pre-spacers with an additional slip occur more often. We presume that pre-spacers from the non-

356   target strand with an AAG PAM are unlikely to be inserted into the CRISPR array in the

357   incorrect orientation, whereas pre-spacers generated by slipping events are more likely to be

358   flipped because they have a suboptimal PAM.

359

360   **Pre-spacers with non-canonical lengths are often associated with slipping**

361   While the large majority of pre-spacers are 33 nt long, we were interested to determine the basis

362   for selecting pre-spacers with non-canonical lengths. 0.3% of both *lacZ+* and *mhpT-* pre-spacer

363   usage was associated with pre-spacers that are 32 nt long, and 1.1% of pre-spacer usage for both

364   *lacZ+* and *mhpT-* pre-spacers was associated with pre-spacers that are 34 nt long. We chose to

365   focus on the region <10 kb upstream of the protospacer, on the non-target strand, where the

366   usage of pre-spacers of non-canonical lengths is highest (89.8% and 82.9% of all non-canonical

367   length pre-spacers for *lacZ+* and *mhpT-* protospacers, respectively). The majority of 32 nt pre-

368   spacer usage for both *lacZ+* and *mhpT-* was associated with AAG PAMs; hence, these pre-

369   spacers are equivalent to efficiently used 33 nt pre-spacers, but lack the most PAM-distal base.

370   These AAG-associated pre-spacers account for 87.7% (lacZ+) and 95.7% (*mhpT-*) of all 32 nt

371   pre-spacer usage. Similarly, the majority of 34 nt pre-spacer usage was associated with AAG

372   PAMs; hence, these pre-spacers are equivalent to efficiently used 33 nt pre-spacers, but with an

373   additional base at the PAM-distal end. These AAG-associated pre-spacers account for 79.7%

374   (*lacZ+*) and 86.7% (*mhpT-*) of all 34 nt pre-spacer usage. There is also a significant enrichment

375   for 32 nt pre-spacer usage associated with +1 slips, and for 34 nt pre-spacer usage associated

376  with -1 slips (Figure 9), indicating that slipping sometimes involves maintaining the PAM-distal

377  boundary while either shortening or lengthening the spacer. The equivalent observation has been

378  made for type I-B, type I-C and type I-F CRISPR-Cas systems (Li et al., 2017; Rao et al., 2017;

379  Staals et al., 2016).

380

381  **The efficiency of primed adaptation is reduced for pre-spacers within 200 bp of the**

382  **protospacer**

383  To determine if the pattern of pre-spacer usage is dependent upon the location of the protospacer,

384  we initiated primed adaptation by targeting a location ~10 kb downstream of the *lacZ+*

385  protospacer, within the *codA* gene; we refer to this as the *codA+* protospacer. Replicate datasets

386  for *codA+* were highly correlated (Figure 5C); hence we combined the two datasets for further

387  analysis. As expected, we detected pre-spacer usage over a ~100 kb distance from the *codA+*

388  protospacer, with the majority of pre-spacer positions overlapping those used when targeting the

389  *lacZ+* protospacer (Figure 10A). We compared the frequency of pre-spacer usage for pre-spacers

390  <10 kb upstream of the *lacZ+* protospacer on the non-target strand when targeting either the

391  *lacZ+* or *codA+* protospacer. In almost all cases, the relative usage of pre-spacers was similar for

392  each of the protospacers targeted (Figure 10B), strongly suggesting that pre-spacer usage

393  frequency is an inherent property of the pre-spacer, and is not related to the location of the

394  protospacer. However, the three pre-spacers located within 200 nt of the *lacZ+* protospacer were

395  used proportionally less frequently when targeting the *lacZ+* protospacer than when targeting the

396  *codA+* protospacer (Figure 10B). Thus, our data suggest that pre-spacers very close to the

397  protospacer are used inefficiently. Biochemical studies of Cas3 suggested that Cas3 generates a

398  single-stranded region adjacent to the Cascade-bound protospacer of <300 nt (Redding et al.,

399    2015) or <500 nt (Dillard et al., 2018), for Cas3 from *E. coli* or *T. fusca*, respectively. We

400    propose that Cas3 generates a single-stranded region of a similar length *in vivo*, and that no

401    spacers can be acquired from the single-stranded DNA. Our data suggest that in the region <200

402    bp from the protospacer, the relative efficiency of pre-spacer usage frequency increases as a

403    function of increasing distance from the protospacer (Figure 10B). Hence, we propose that Cas3

404    transitions stochastically between a state where it nicks DNA frequently, generating an extended

405    region of single-stranded DNA, and a state where it nicks DNA infrequently, with the average

406    distance traveled until the transition being 100-200 bp.

407

408    **The efficiency of primed adaptation is reduced for pre-spacers with an internal AAG**

409    It is well established that the frequency of pre-spacer usage varies considerably between pre-

410    spacer sequences (Musharova et al., 2017, 2018; Savitskaya et al., 2013). However, the basis for

411    this variability is poorly understood. As described above, the frequency of pre-spacer usage

412    decreases globally as a function of increasing distance from the protospacer (Figure 3 + 11A+B).

413    We fitted these data to an exponential decay model, relating pre-spacer usage to distance from

414    the protospacer. We then selected "enriched" pre-spacers that have a >4-fold higher usage

415    frequency than predicted by the model, and "de-enriched" pre-spacers that have a >4-fold lower

416    usage frequency than predicted by the model. We determined the frequency of all tetranucleotide

417    and trinucleotide sequences within the sets of enriched and de-enriched pre-spacers (Figure 11C-

418    D). Almost all tetranucleotide and trinucleotide sequences were found at similar frequencies

419    within the enriched and de-enriched pre-spacers. However, AAG-containing pre-spacers were

420    found far less frequently in enriched pre-spacers than in de-enriched pre-spacers. We compared

421    the frequency of pre-spacer usage for AAG-containing and non-AAG-containing pre-spacers

422     (Figure 11A-B). The frequency was, on average, ~3.5-fold higher for non-AAG-containing pre-

423     spacers than for AAG-containing pre-spacers. We conclude that the presence of an AAG within

424     a pre-spacer substantially reduces the frequency with which that pre-spacer is selected during

425     primed adaptation.

426

427     To experimentally test the effect of an AAG within a pre-spacer, we generated a second self-

428     targeting plasmid, "stp2", that includes a site for easy introduction of potential pre-spacers

429     flanked by an AAG PAM (see Materials and Methods for a detailed description of the plasmid

430     and assessment of its ability to cause primed adaptation). The introduced pre-spacer in

431     unmodified stp2 has AAA at positions 19-21. We also constructed a derivative of stp2 in which

432     the introduced pre-spacer has been modified to contain AAG rather than AAA at positions 19-

433     21. We refer to this mutant stp2 plasmid as stp2-mut1 (Figure 12A). We initiated primed

434     adaptation in cells containing either stp2 or stp2-mut1. We determined the location of pre-

435     spacers for each of the two plasmids by PCR-amplification and sequencing of the expanded

436     CRISPR arrays. We compared the relative usage frequency for each pre-spacer on the non-target

437     strand that is associated with an AAG PAM (these represent the most frequently used pre-

438     spacers; Figure 2A) for each of the two plasmids. As expected, pre-spacers with the same

439     sequence were used at the same frequencies for each plasmid (Figure 12B). By contrast, the

440     AAG-containing pre-spacer in stp2-mut1 was used at a substantially lower frequency than the

441     equivalent AAA-containing pre-spacer in wild-type stp2 (Figure 12B). Thus, our experimental

442     data are consistent with the bioinformatic analysis, and show that pre-spacers with an internal

443     AAG are used at a reduced frequency during primed adaptation. We note also that a recent,

444    independent study reached a similar conclusion about the impact of AAG sequences within pre-

445    spacers (Musharova et al., 2018).

446

447    A previous study mapped pre-spacers for naïve adaptation in a strain of *E. coli* that is not able to

448    undergo primed adaptation (Yosef et al., 2013). The authors sequenced 934,202 newly acquired

449    spacers that came from the chromosome, covering 84,951 pre-spacer sites. Although the

450    frequency of pre-spacers with AAG PAMs is considerably lower for naïve adaptation than

451    primed adaptation, many pre-spacers are associated with an AAG PAM (Yosef et al., 2013).

452    Using these data, we determined the frequency of usage for all potential chromosomal pre-

453    spacers associated with an AAG PAM. 18,872 sequences were identified as pre-spacers at least

454    once, whereas 108,207 sequences were not. We then compared the number of AAG-containing

455    (i.e. an internal AAG) pre-spacers within each of these groups to determine whether the impact

456    of an internal AAG is limited to primed adaptation. Although the proportion of AAG-containing

457    pre-spacers was similar for each group, a significantly larger fraction (38.8%) of the unused pre-

458    spacers had an internal AAG than the fraction of used pre-spacers (34.4%; Chi Squared Test with

459    Yates' Continuity Correction $p < 1e^{-7}$). As a control, we analyzed the proportion of AGA-

460    containing pre-spacers for the same groups and found no significant difference (Chi Squared

461    Test with Yates' Continuity Correction $p = 0.08$). We then divided the 18,872 used pre-spacer

462    sequences into AAG-containing and non-containing groups and compared the usage frequencies

463    for each group. While the distribution of usage frequencies was similar between the two groups

464    (Figure 12C), it was slightly and significantly higher for the AAG-lacking group (Mann Whitney

465    U Test $p = 2.5e^{-4}$). We saw no significant difference when comparing AGA-containing and

466    AGA-lacking pre-spacers (Mann Whitney U Test $p = 0.63$). We conclude that AAG-containing

467     pre-spacers are acquired at a lower efficiency during naïve adaptation. However, the magnitude

468     of the effect appears very small, in contrast to the large effect we observed for primed adaptation

469     (Figures 11 + 12B). Consistent with this, a previous study did not detect an effect of AAG within

470     pre-spacers on the efficiency of naïve adaptation (Musharova et al., 2018).

471

472     Since AAG is the PAM associated with >95% of pre-spacer usage during primed adaptation, we

473     presume that the impact of AAGs within pre-spacers is connected to the selection of pre-spacers

474     with AAG PAMs. One possible explanation for the reduced frequency of usage for AAG-

475     containing pre-spacers is that multiple Cas3-containing complexes could translocate from a

476     single Cascade-bound protospacer, such that stable association of a Cas3-containing complex

477     with an AAG PAM sequence prevents association of a second Cas3-containing complex with a

478     nearby AAG PAM. However, this would not explain the size of the decrease in pre-spacer usage

479     frequency associated with an internal AAG (3.5-fold), or its impact over a >100 kb region; for

480     the effect of an internal AAG to be explained by competition within Cas3-continaing complexes,

481     most of the >1,000 AAG sequences within 100 kb of the protospacer would need to be occupied

482     at any given time. A more likely explanation is that AAG sequences are frequently nicked, which

483     would represent a semi-stable change in the DNA that can be caused by a transiently positioned

484     Cas3-containing complex. Recent studies suggest that pre-spacers are nicked within the PAM

485     during primed adaptation (Musharova et al., 2017; Shiriaeva et al., 2019). We propose that the

486     Cas3-containing complex nicks all AAG sequences that it encounters as it translocates away

487     from the Cascade-bound protospacer (after transitioning from the highly active state within 200

488     nt of the protospacer), and that pre-spacers with an internal nick cannot be used as substrates for

489     adaptation by Cas1-2.

490

491    Intriguingly, the impact of an AAG within a pre-spacer appears to diminish as the AAG is

492    positioned further from the PAM (Figure 13A) (Musharova et al., 2018). Based on our nicking

493    model, we hypothesized that closely positioned pairs of AAG sequences would lead to closely

494    spaced nicks, and hence single-stranded DNA gaps that may then prevent the ability of these

495    sequences to be integrated into a CRISPR array. To test this hypothesis, we repeated the stp2

496    primed adaptation assay described above (Figure 12B) using a derivative of the stp2-mut1

497    plasmid where the sequence between the PAM and the internal AAG of the modified pre-spacer

498    was changed to a more G/C-rich sequence with a higher melting temperature. We reasoned that

499    this would reduce the propensity for the DNA between the PAM and the internal AAG to

500    become single-stranded in the event that the PAM AAG and the internal AAG were both nicked.

501    We refer to this plasmid as stp2-mut2 (Figure 12A). We also constructed an equivalent plasmid,

502    stp2-mut3, where the AAG within the pre-spacer was changed to an AAA (Figure 12A).

503    Although the relative usage frequency of the modified pre-spacer was higher for stp2-mut2 than

504    for stp2-mut1 (Figure 13B), it was also higher for stp2-mut3 than for stp2 (Figure 13C),

505    suggesting that this effect is independent of the internal AAG, and arguing against our

506    hypothesis. Moreover, these data suggest that the sequence of the PAM-proximal region of a pre-

507    spacer impacts the frequency with which it is used during primed adaptation.

508

509    **The efficiency of primed adaptation is modulated by a PAM-distal motif in the pre-spacer**

510    Although the presence of an AAG within the pre-spacer is clearly an important factor in

511    determining the frequency of pre-spacer usage, there is considerable local variability in pre-

512    spacer usage for pre-spacers that lack an internal AAG (Figure 11A + B). We used the frequency

513    of pre-spacer usage for AAG-flanked pre-spacers in the 100 kb region upstream of the *lacZ+* and

514    *mhpT-* protospacers to fit exponential decay models, estimating pre-spacer usage as a function of

515    distance from the protospacer. Based on the modeled exponential decay, we estimate that the

516    frequency of pre-spacer usage drops by 50% every ~11,500 – 13,900 kb. Strikingly, an *in vitro*

517    study of *E. coli* Cas3 translocation indicated that 50% of Cas3 molecules dissociate from DNA

518    after ~12 kb (Redding et al., 2015), strongly suggesting that the decrease in pre-spacer usage we

519    observe as a function of distance from the protospacer reflects the extent of Cas3 translocation.

520

521    We compared the frequency of pre-spacer usage to the predicted frequency based on the

522    exponential decay model. We then selected all pre-spacers for which the frequency of usage was

523    >4-fold greater than that predicted by the model. Alignment of these pre-spacer sequences

524    revealed sequence bias at positions 28, 30, 32 and 33 (Figure 14A). We next selected all pre-

525    spacers for which the frequency of usage was >4-fold lower than that predicted by the model.

526    Alignment of these pre-spacer sequences revealed sequence bias at positions 30, 31, 32 and 33

527    (Figure 14B). Strikingly, the hierarchy of preferred bases at positions 30, 32 and 33 was

528    precisely inverted for the enriched and de-enriched pre-spacers. To determine whether the

529    identified sequence biases are meaningful, we determined the information content of the position

530    weight matrices (Schneider et al., 1986) for the enriched and de-enriched pre-spacers. We then

531    determined the information content of the same number of randomly selected pre-spacer

532    sequences, and repeated this analysis 100,000 times. Based on the distribution of information

533    content scores for the randomly selected sequences, the information content scores of the

534    position weight matrices from the enriched and de-enriched pre-spacers were significantly higher

535    than expected by chance ($Z = 10.4$ and $p < 1e^{-14}$ for enriched pre-spacers; $Z = 4.9$ and $p = 5e^{-7}$ for

536 de-enriched pre-spacers). We conclude that the sequence of base positions at the 3' end of the

537 pre-spacer influences usage frequency, with positions 30, 32 and 33 being the most important.

538

539 To experimentally test the impact of pre-spacer positions 30-33 on primed adaptation, we

540 constructed derivatives of stp2 that contain a pre-spacer where positions 30-31 were modified

541 from TAAA to CCAA (stp2-mut4), or positions 30-33 were modified to CCCG (stp2-mut5;

542 Figure 14C). Based on the enrichment/de-enrichment of pre-spacers in the chromosomal primed

543 adaptation experiment (Figure 14A-B), we reasoned that changing positions 30-33 to CCAA

544 would not have a large effect on the frequency of pre-spacer usage, since our data suggest that

545 only one of the two modified bases (position 30) would generate a sub-optimal pre-spacer. By

546 contrast, we reasoned that changing positions 30-33 to CCCG would have a substantial negative

547 impact on the frequency of pre-spacer usage, since it has the most disfavored base at three PAM-

548 distal positions. Consistent with our expectation, changing positions 30-31 to CC had no effect

549 on pre-spacer usage frequency (Figure 14D), whereas changing positions 30-33 to CCCG

550 substantially reduced pre-spacer usage frequency (Figure 14E).

551

552 A previous study suggested that the sequence of positions 32 and 33 of the pre-spacer influences

553 the efficiency of naïve adaptation. Specifically, a C at position 32 or T at position 33 resulted in

554 inefficient adaptation whereas an A at position 32 and or an A at position 33 resulted in efficient

555 adaptation (Yosef et al., 2013). Our data are consistent with the idea that the sequence of

556 positions 28-33 are important, and moreover are consistent with the specific sequence effects

557 observed for naïve adaptation. We conclude that the sequence of the PAM-distal portion of the

558 pre-spacer influences both naïve and primed adaptation, likely due to effects on Cas1-2

559   association and/or integration of the pre-spacer into the CRISPR array. Our data suggest that the

560   number of important positions at the PAM-distal end of the pre-spacer is larger than previously

561   suggested (Yosef et al., 2013), encompassing perhaps as many as six positions.

562

563   **A model for primed adaptation**

564   We propose that once recruited to a protospacer-bound Cascade, Cas3 translocates along the

565   DNA, upstream from the PAM, consistent with *in vitro* studies (Dillard et al., 2018; Redding et

566   al., 2015). The initial translocating state of Cas3 is likely to be different to its final translocating

567   state; we propose that this initial state is associated with frequent, sequence-independent nicking

568   of one strand of the DNA. Our data suggest that the transition from the frequent nicking state to

569   the infrequent nicking state is stochastic, and occurs in all cases within ~200 nt of the

570   protospacer. Biochemical data suggest that translocating Cas3 initially remains associated with

571   Cascade (Dillard et al., 2018; Redding et al., 2015). We propose that the interaction between

572   Cas3 and Cascade causes Cas3 to be in the frequent nicking state, and that it is the dissociation

573   of Cas3 from Cascade as tension builds in the DNA that causes Cas3 to transition to the

574   infrequent nicking state. After the transition, we propose that Cas3, or possibly Cas1 that is

575   associated with Cas3, cuts DNA at AAG sequences. Given the magnitude of the defect in primed

576   adaptation associated with the presence of an AAG within a spacer, it is likely that this nicking

577   occurs at almost every AAG encountered by Cas3. It is unclear why the degree to which an AAG

578   within a pre-spacer affects primed adaptation is dependent on the position of the AAG relative to

579   the PAM. Given that an AAG within a pre-spacer has a small but significant negative impact on

580   naïve adaptation, we suggest that Cas1 is responsible for the DNA cut at AAG sequences.

581   Presumably, during naïve adaptation a second cut by Cas1 is less likely than during primed

582    adaptation, because Cas1 would not be associated with a translocating Cas3, and hence is less

583    likely to be suitably positioned on the DNA.

584

585    It is unclear how substrates are generated for Cas1-2 during primed adaptation, although nicking

586    of at least one DNA strand by Cas3 is required (Datsenko et al., 2012). There is strong evidence

587    for a nick between the second and third positions of the PAM (Musharova et al., 2017; Shiriaeva

588    et al., 2019), consistent with our proposal that every AAG is cut by the translocating Cas3-

589    containing complex. However, this would only account for one of the four cuts required for

590    substrate generation. Future work is required to determine whether Cas3, Cas1, and/or non-Cas

591    proteins are required for this process. Once substrates have been generated, the efficiency of

592    integration into the array is dependent on the sequence. Our data indicate that the six most PAM-

593    distal positions of the pre-spacer contribute to the efficiency of primed adaptation (Figure 14).

594    Given that the importance of the two most PAM-distal positions has been recognized to be

595    important for naïve adaptation (Yosef et al., 2013), we propose that the impact of these positions

596    is solely on Cas1-2. These positions may be important for stable association of Cas1-2, or for the

597    process of integration into the CRISPR array. Our data also suggest that sequences within pre-

598    spacers, more proximal to the PAM, play an important role in primed adaptation (Figure 13), but

599    the mechanism for this is unknown.

600

601

602 **MATERIALS AND METHODS**

603

604 **Strains and Plasmids**

605 Strains, plasmids, and oligonucleotides are listed in Tables 1, 2, and 3, respectively. All strains

606 are derivatives of MG1655 (Blattner et al., 1997). CB386, AMD536 and AMD688 have been

607 previously described (Cooper et al., 2018; Luo et al., 2014). 1XDNAi, as described previously

608 (Caliando and Voigt, 2015), contains a chromosomally integrated actuator, pACT-01. AMD671

609 is a derivative of 1XDNAi that we constructed using FRUIT recombineering (Stringer et al.,

610 2012). Specifically, *thyA* was deleted in 1XDNAi, using oligonucleotides JW472 and JW473

611 and strain AMD052 (MG1655 Δ*thyA*) as a template. *thyA* was inserted downstream of the

612 pACT-01 actuator using FRUIT (Stringer et al., 2012) with oligonucleotides JW9016 and

613 JW9017. Oligonucleotides JW9009 and JW9010 were used to amplify *cas1-cas2* from MG1655

614 to replace *thyA* downstream of pACT-01. Lastly, the *lacZ* gene was replaced by *thyA* using

615 FRUIT (Stringer et al., 2012) with oligonucleotides JW9066 and JW9067.

616

617 The self-targeting plasmid (stp; pAMD189), the Cas3-expressing plasmid (pAMD191), and the

618 plasmid that expresses the crRNA targeting the *lacZ+* protospacer (pCB380), have been

619 described previously (Cooper et al., 2018; Luo et al., 2014). All other crRNA-expressing

620 plasmids are derivatives of pAMD179 (Cooper et al., 2018). To clone individual spacers, pairs of

621 oligonucleotides were annealed, extended, and cloned using In-Fusion (Clontech) into the <u>*Xho*</u>I

622 and <u>*Sac*</u>II sites of pAMD179 to generate pAMD211 (targets *codA+* protospacer; with

623 oligonucleotides JW8010 and JW6518), and pAMD212 (targets *mhpT-* protospacer; with

624 oligonucleotides JW8011 and JW6518).

625

626      The second self-targeting plasmid (stp2) is a derivative of pSDS009. pSDS009 was generated by

627      cloning a synthesized dsDNA fragment (gBlock JW9076; Integrated DNA Technologies, Inc.;

628      Table 3) using In-Fusion (Clontech) into the *Cla*I and *Hin*dIII sites of pBAD24 amp (Guzman et

629      al., 1995). stp2 and its derivatives were generated by annealing and extending pairs of

630      oligonucleotides, and cloning using In-Fusion (Clontech) into the <u>*Hin*</u>dIII site of pSDS009 to

631      generate pSK014 (stp2; cloned with oligonucleotides JW10011 and JW10012), pSK013 (stp2-

632      mut1; cloned with oligonucleotides JW10009 and JW10010), pSK015 (stp2-mut2; cloned with

633      oligonucleotides JW10013 and JW10014), pSK035 (stp2-mut3; cloned with oligonucleotides

634      JW10015 and JW10016), pSK017 (stp2-mut4; cloned with oligonucleotides JW10021 and

635      JW10022), and pSK016 (stp2-mut5; cloned with oligonucleotides JW10017 and JW10018).

636

637      **Quantification of adaptation using a YFP reporter**

638      AMD688 was transformed with two plasmids: the first plasmid was either pAMD189 or empty

639      pBAD24, while the second was either pAMD191 or empty pBAD33. Cells were grown

640      overnight in LB supplemented with 100 ug/mL ampicillin and 30 ug/mL chloramphenicol and

641      were sub-cultured the next day 1:100 for 6 hours in LB supplemented with 0.2% arabinose and

642      30 ug/mL chloramphenicol. Cells were pelleted by centrifugation and resuspended in M9

643      minimal medium in twice the original volume ($OD_{600} \approx 1.0$). Samples were transferred to 5 mL

644      polystyrene round-bottom tubes and analyzed by flow cytometry for single-cell detection of YFP

645      expression using the BD FACSAria IIU Cell Sorter. We recorded 100,000 events for each

646      sample.

647

648 **PCR and Sequencing to Assess Primed Adaptation**

649 To determine the set of pre-spacers acquired during priming, pAMD191 was transformed into

650 AMD536 along with either pAMD189, pCB380, pAMD211 or pAMD212. Cells were grown

651 overnight in LB supplemented with 100 μg/mL ampicillin, 30 μg/mL chloramphenicol, and

652 0.2% glucose at 37 °C with aeration, and sub-cultured the next day in LB supplemented with

653 0.2% arabinose at 37 °C with aeration for one hour.

654

655 To experimentally assess the impact of specific sequences on pre-spacer usage, pSK013,

656 pSK014, pSK015, pSK035, pSK016, or pSK017 were transformed into AMD671. Cells were

657 grown overnight in LB supplemented with 100 μg/mL ampicillin, and 0.2% glucose at 37°C with

658 aeration, and sub-cultured the next day in LB supplemented with 0.2% arabinose at 37°C with

659 aeration for six hours.

660

661 For both sets of cultures described above, cells were pelleted from 1 mL of culture by

662 centrifugation, and cell pellets were frozen at -20 °C. PCRs were then performed on the cell

663 pellets, amplifying the CRISPR arrays using oligonucleotides JW7816 and JW7817 for CRISPR-

664 I and JW7818 and JW7819 for CRISPR-II. PCR products were visualized on acrylamide gels

665 and the first expansion band was extracted and purified. A second round of amplification was

666 performed using a universal forward oligonucleotide (JW7820 or JW8053), and one of a set of

667 reverse oligonucleotides containing Nextera (Illumina) indices (JW8054, JW8057, JW8062,

668 JW8476, JW8477, JW8478, JW8479, JW8480, JW8481 or JW8485). Samples were purified

669 using a MinElute kit (Qiagen) and pooled. Samples were sequenced using an Illumina MiSeq

670 Instrument (Wadsworth Center Applied Genomic Technologies Core).

671

672 Plasmids pSK013, pSK014, pSK015, pSK035, pSK016, or pSK017 serendipitously function as

673 self-targeting plasmids despite the lack of a cloned spacer. These plasmids are derivatives of

674 pSDS009, which is designed to express crRNAs from a cassette that includes a spacer with a full

675 repeat downstream, but only 8 nt of repeat sequence upstream. pSDS009 does not itself contain a

676 canonical spacer between the partial and complete repeat sequences. Similarly, pSK013,

677 pSK014, pSK015, pSK035, pSK016, and pSK017 do not include a canonical spacer sequence.

678 Nonetheless, multiple lines of evidence indicate that these plasmids drive efficient primed

679 adaptation in cells expressing *cas3*: (i) the majority of pre-spacer usage is from the plasmid; (ii)

680 96% of pre-spacer usage on the plasmid is associated with an AAG PAM; (iii) for AAG-flanked

681 pre-spacers that are shared with the stp, are found on the stp non-target strand, and are not shared

682 with pAMD191, pre-spacer usage from pSK014 correlates well with that from the stp

683 (correlation coefficient of 0.93). We suspect that the downstream repeat sequence in pSDS009

684 and its derivatives also functions as an upstream repeat sequence for a cryptic crRNA. In this

685 scenario, the sequence downstream of this repeat serves as a spacer, with the end of the crRNA

686 presumably being determined by the position of the downstream transcription terminator. We

687 previously showed that the exact same sequence results in expression of a cryptic crRNA from a

688 derivative of pAMD179 (Cooper et al., 2018). The presumed spacer sequence is duplicated in the

689 same orientation, a short distance downstream, which would facilitate interference and primed

690 adaptation.

691

692 **Spacer sequence extraction and alignment to reference genomes**

693   Newly acquired spacers found directly adjacent to the leader proximal endogenous spacer (the

694   first primed adaptation acquired spacer) were extracted from .fastq files using custom Python

695   scripts for each CRISPR array being analyzed. The scripts first selected sequences that included

696   the first spacer of the unexpanded CRISPR array (spacer #1) and at least one repeat sequence.

697   The script then removed sequences for which the gap between spacer #1 and the next upstream

698   repeat differed from the expected value. The script then extracted the sequence between the first

699   repeat and the repeat immediately upstream of spacer #1. This was presumed to consist of

700   complete spacer sequences, or combinations of spacers and repeats (i.e. the result of >1 spacer

701   being added to the CRISPR array). Sequences were discarded if they were shorter than 30 nt or

702   longer than 36 nt, ensuring that only arrays with a single expansion were analyzed. Spacers were

703   mapped to their respective genomes (stp, pAMD191, or *E. coli* K-12 MG1566 U00096.3) using

704   CLC Genomics Workbench v10.1.1, requiring perfect matches. For experiments involving

705   mapping reads to multiple genomes, we first mapped to stp, then mapped all unmapped reads to

706   pAMD191, and then mapped all unmapped reads to *E. coli* K-12 MG1655.

707

**Analysis of primed adaptation in regions adjacent to off-target Cascade binding sites**

709   We previously described a set of 76 off-target Cascade binding sites whose sequence is

710   consistent with Cascade binding in association with the crRNA targeting the *lacZ+* protospacer

711   (Cooper et al., 2018). We counted the frequency of pre-spacer usage for each of these off-target

712   protospacers, searching on the non-target strand in the 10 kb upstream of the protospacer. We

713   repeated the analysis for 1,000 randomly selected genomic positions equally distributed across

714   the forward and reverse strands.

715

**Analysis of slipping for 33 nt pre-spacers**

716

717 We selected all 33 nt pre-spacers on the non-target strand, in the 10 kb regions upstream of the

718 *lacZ*+ and *mhpT*- protospacers, for which the PAM was not AAG. We calculated all pairwise

719 distances for these pre-spacers to AAG-associated pre-spacers on the non-target strand, in the 10

720 kb upstream of each of the protospacer. For the analysis represented in Figure 7A, only unique

721 pre-spacer positions were considered, so information on the frequency with which different pre-

722 spacers are used was lot. Values plotted in Figure 7A represent the number of instances of the

723 indicated slipping distance divided by the total number of comparisons between pre-spacer

724 positions and AAGs, multiplied by 1000. Values listed in Figure 7B take into account the fact

725 that many pre-spacer positions were detected more than once, and hence these numbers control

726 for the relative use of the different sequences. Pre-spacers with a non-AAG PAM that did not

727 derive from a -1, +1 or +2 slip were selected; pre-spacers that were sequenced multiple times

728 were represented at the frequency with which they were detected. The associated PAM

729 sequences were then converted into the logos shown in Figure 7C using Weblogo (Crooks et al.,

730 2004).

731

**Analysis of flipping**

732

733 We selected all 33 nt pre-spacers on the target strand, in the 10 kb regions downstream of the

734 *lacZ*+ and *mhpT*- protospacers. We calculated all pairwise distances for these pre-spacers to

735 AAG-associated pre-spacers on the non-target strand, in the 10 kb upstream of each of the

736 protospacer. For the analysis represented in Figure 8B, only unique pre-spacer positions were

737 considered, so information on the frequency with which different pre-spacers are used was lot.

738 Values plotted in Figure 8B represent the number of instances of the indicated slipping distance

739    divided by the total number of comparisons between pre-spacer positions and AAGs, multiplied

740    by 1000.

741

742    **Analysis of slipping for spacers of non-canonical length**

743    We selected all 32 and 34 nt pre-spacers on the non-target strand, in the 10 kb regions upstream

744    of the *lacZ+* and *mhpT-* protospacers. We calculated all pairwise distances for these pre-spacers

745    to AAG-associated pre-spacers on the non-target strand, in the 10 kb upstream of each of the

746    protospacer. For the analysis represented in Figure 9C-D, only unique pre-spacer positions were

747    considered, so information on the frequency with which different pre-spacers are used was lot.

748    Values plotted in Figure 9C-D represent the number of instances of the indicated slipping

749    distance divided by the total number of comparisons between pre-spacer positions and AAGs,

750    multiplied by 1000.

751

752    **Identification of pre-spacers with higher-than-expected or lower-than-expected usage**

753    **frequency, independent of internal AAG sequences**

754    We took the frequency of pre-spacer usage for all AAG-flanked 33 nt pre-spacers in the 100 kb

755    region upstream of the *lacZ+* and *mhpT-* protospacers. We then fit an exponential decay model to

756    each dataset using Microsoft Excel ("best fit line" function). Using the parameters associated

757    with these models, we determined the expected usage frequency for pre-spacers at every position

758    in the 100 kb windows, and we then took the ratio of the actual usage frequency to the expected

759    usage frequency. Pre-spacers with ratios >4 or <0.25 were selected for analysis of tetranucleotide

760    and trinucleotide content.

761

762 **Identification of pre-spacers with higher-than-expected or lower-than-expected usage**

763 **frequency, independent of internal AAG sequences**

764 We took the frequency of pre-spacer usage for all AAG-flanked 33 nt pre-spacers in the 100 kb

765 region upstream of the *lacZ+* and *mhpT-* protospacers, excluding pre-spacers with an internal

766 AAG. We then fit an exponential decay model to each dataset using Microsoft Excel ("best fit

767 line" function). Using the parameters associated with these models, we determined the expected

768 usage frequency for pre-spacers at every position in the 100 kb windows, and we then took the

769 ratio of the actual usage frequency to the expected usage frequency. Pre-spacers with ratios >4 or

770 <0.25 were used as input for Weblogo (Crooks et al., 2004) to create the sequence logos shown

771 in Figure 14A-B.

772

# REFERENCES

787

788

789      Amlinger, L., Hoekzema, M., Wagner, E.G.H., Koskiniemi, S., and Lundgren, M. (2017).

790      Fluorescent CRISPR Adaptation Reporter for rapid quantification of spacer acquisition. Sci. Rep.

791      *7*, 10392.

792      Blattner, F.R., Plunkett, G., Bloch, C.A., Perna, N.T., Burland, V., Riley, M., Collado-Vides, J.,

793      Glasner, J.D., Rode, C.K., Mayhew, G.F., et al. (1997). The complete genome sequence of

794      Escherichia coli K-12. Science *277*, 1453–1462.

795      Brouns, S.J.J., Jore, M.M., Lundgren, M., Westra, E.R., Slijkhuis, R.J.H., Snijders, A.P.L.,

796      Dickman, M.J., Makarova, K.S., Koonin, E.V., and van der Oost, J. (2008). Small CRISPR

797      RNAs guide antiviral defense in prokaryotes. Science *321*, 960–964.

798      Caliando, B.J., and Voigt, C.A. (2015). Targeted DNA degradation using a CRISPR device

799      stably carried in the host genome. Nat. Commun. *6*, 6989.

800      Cooper, L.A., Stringer, A.M., and Wade, J.T. (2018). Determining the Specificity of Cascade

801      Binding, Interference, and Primed Adaptation In Vivo in the Escherichia coli Type I-E CRISPR-

802      Cas System. MBio *9*.

803      Crooks, G.E., Hon, G., Chandonia, J.-M., and Brenner, S.E. (2004). WebLogo: a sequence logo

804      generator. Genome Res. *14*, 1188–1190.

805   Datsenko, K.A., Pougach, K., Tikhonov, A., Wanner, B.L., Severinov, K., and Semenova, E.

806   (2012). Molecular memory of prior infections activates the CRISPR/Cas adaptive bacterial

807   immunity system. Nat. Commun. *3*, 945.

808   Dillard, K.E., Brown, M.W., Johnson, N.V., Xiao, Y., Dolan, A., Hernandez, E., Dahlhauser,

809   S.D., Kim, Y., Myler, L.R., Anslyn, E.V., et al. (2018). Assembly and Translocation of a

810   CRISPR-Cas Primed Acquisition Complex. Cell *175*, 934-946.e15.

811   Fagerlund, R.D., Wilkinson, M.E., Klykov, O., Barendregt, A., Pearce, F.G., Kieper, S.N.,

812   Maxwell, H.W.R., Capolupo, A., Heck, A.J.R., Krause, K.L., et al. (2017). Spacer capture and

813   integration by a type I-F Cas1–Cas2-3 CRISPR adaptation complex. Proc. Natl. Acad. Sci. *114*,

814   E5122–E5128.

815   Fineran, P.C., Gerritzen, M.J.H., Suárez-Diez, M., Künne, T., Boekhorst, J., Hijum, S.A.F.T.

816   van, Staals, R.H.J., and Brouns, S.J.J. (2014). Degenerate target sites mediate rapid primed

817   CRISPR adaptation. Proc. Natl. Acad. Sci. *111*, E1629–E1638.

818   Goren, M.G., Yosef, I., Auster, O., and Qimron, U. (2012). Experimental definition of a

819   clustered regularly interspaced short palindromic duplicon in Escherichia coli. J. Mol. Biol. *423*,

820   14–16.

821   Guzman, L.M., Belin, D., Carson, M.J., and Beckwith, J. (1995). Tight regulation, modulation,

822   and high-level expression by vectors containing the arabinose PBAD promoter. J. Bacteriol. *177*,

823   4121–4130.

824   Jackson, S.A., McKenzie, R.E., Fagerlund, R.D., Kieper, S.N., Fineran, P.C., and Brouns, S.J.J.

825   (2017). CRISPR-Cas: Adapting to change. Science *356*.

826     Jackson, S.A., Birkholz, N., Malone, L.M., and Fineran, P.C. (2019). Imprecise Spacer

827     Acquisition Generates CRISPR-Cas Immune Diversity through Primed Adaptation. Cell Host

828     Microbe *25*, 250-260.e4.

829     Jore, M.M., Lundgren, M., van Duijn, E., Bultema, J.B., Westra, E.R., Waghmare, S.P.,

830     Wiedenheft, B., Pul, Ü., Wurm, R., Wagner, R., et al. (2011). Structural basis for CRISPR RNA-

831     guided DNA recognition by Cascade. Nat. Struct. Mol. Biol. *18*, 529–536.

832     Jung, C., Hawkins, J.A., Jones, S.K., Xiao, Y., Rybarski, J.R., Dillard, K.E., Hussmann, J.,

833     Saifuddin, F.A., Savran, C.A., Ellington, A.D., et al. (2017). Massively Parallel Biophysical

834     Analysis of CRISPR-Cas Complexes on Next Generation Sequencing Chips. Cell *170*, 35-

835     47.e13.

836     Koonin, E.V., Makarova, K.S., and Zhang, F. (2017). Diversity, classification and evolution of

837     CRISPR-Cas systems. Curr. Opin. Microbiol. *37*, 67–78.

838     Künne, T., Kieper, S.N., Bannenberg, J.W., Vogel, A.I.M., Miellet, W.R., Klein, M., Depken,

839     M., Suarez-Diez, M., and Brouns, S.J.J. (2016). Cas3-Derived Target DNA Degradation

840     Fragments Fuel Primed CRISPR Adaptation. Mol. Cell *63*, 852–864.

841     Leenay, R.T., and Beisel, C.L. (2017). Deciphering, Communicating, and Engineering the

842     CRISPR PAM. J. Mol. Biol. *429*, 177–191.

843     Levy, A., Goren, M.G., Yosef, I., Auster, O., Manor, M., Amitai, G., Edgar, R., Qimron, U., and

844     Sorek, R. (2015). CRISPR adaptation biases explain preference for acquisition of foreign DNA.

845     Nature *520*, 505–510.

846    Li, M., Wang, R., Zhao, D., and Xiang, H. (2014). Adaptation of the Haloarcula hispanica

847    CRISPR-Cas system to a purified virus strictly requires a priming process. Nucleic Acids Res.

848    *42*, 2483–2492.

849    Li, M., Gong, L., Zhao, D., Zhou, J., and Xiang, H. (2017). The spacer size of I-B CRISPR is

850    modulated by the terminal sequence of the protospacer. Nucleic Acids Res. *45*, 4642–4654.

851    Luo, M.L., Mullis, A.S., Leenay, R.T., and Beisel, C.L. (2014). Repurposing endogenous type I

852    CRISPR-Cas systems for programmable gene repression. Nucleic Acids Res. gku971.

853    Musharova, O., Klimuk, E., Datsenko, K.A., Metlitskaya, A., Logacheva, M., Semenova, E.,

854    Severinov, K., and Savitskaya, E. (2017). Spacer-length DNA intermediates are associated with

855    Cas1 in cells undergoing primed CRISPR adaptation. Nucleic Acids Res. *45*, 3297–3307.

856    Musharova, O., Vyhovskyi, D., Medvedeva, S., Guzina, J., Zhitnyuk, Y., Djordjevic, M.,

857    Severinov, K., and Savitskaya, E. (2018). Avoidance of Trinucleotide Corresponding to

858    Consensus Protospacer Adjacent Motif Controls the Efficiency of Prespacer Selection during

859    Primed Adaptation. MBio *9*.

860    Musharova, O., Sitnik, V., Vlot, M., Savitskaya, E., Datsenko, K.A., Krivoy, A., Fedorov, I.,

861    Semenova, E., Brouns, S.J.J., and Severinov, K. (2019). Systematic analysis of Type I-E

862    Escherichia coli CRISPR-Cas PAM sequences ability to promote interference and primed

863    adaptation. Mol. Microbiol. *111*, 1558–1570.

864    Rao, C., Chin, D., and Ensminger, A.W. (2017). Priming in a permissive type I-C CRISPR-Cas

865    system reveals distinct dynamics of spacer acquisition and loss. RNA N. Y. N *23*, 1525–1538.

866    Redding, S., Sternberg, S.H., Marshall, M., Gibb, B., Bhat, P., Guegler, C.K., Wiedenheft, B.,

867    Doudna, J.A., and Greene, E.C. (2015). Surveillance and Processing of Foreign DNA by the

868    Escherichia coli CRISPR-Cas System. Cell *163*, 854–865.

869    Richter, C., Dy, R.L., McKenzie, R.E., Watson, B.N.J., Taylor, C., Chang, J.T., McNeil, M.B.,

870    Staals, R.H.J., and Fineran, P.C. (2014). Priming in the Type I-F CRISPR-Cas system triggers

871    strand-independent spacer acquisition, bi-directionally from the primed protospacer. Nucleic

872    Acids Res. *42*, 8516–8526.

873    Savitskaya, E., Semenova, E., Dedkov, V., Metlitskaya, A., and Severinov, K. (2013). High-

874    throughput analysis of type I-E CRISPR/Cas spacer acquisition in E. coli. RNA Biol. *10*, 716–

875    725.

876    Schneider, T.D., Stormo, G.D., Gold, L., and Ehrenfeucht, A. (1986). Information content of

877    binding sites on nucleotide sequences. J. Mol. Biol. *188*, 415–431.

878    Semenova, E., Savitskaya, E., Musharova, O., Strotskaya, A., Vorontsova, D., Datsenko, K.A.,

879    Logacheva, M.D., and Severinov, K. (2016). Highly efficient primed spacer acquisition from

880    targets destroyed by the Escherichia coli type I-E CRISPR-Cas interfering complex. Proc. Natl.

881    Acad. Sci. 201602639.

882    Shiriaeva, A.A., Savitskaya, E., Datsenko, K.A., Vvedenskaya, I.O., Fedorova, I., Morozova, N.,

883    Metlitskaya, A., Sabantsev, A., Nickels, B.E., Severinov, K., et al. (2019). Detection of spacer

884    precursors formed in vivo during primed CRISPR adaptation. Nat. Commun. *10*, 4603.

885    Shmakov, S., Savitskaya, E., Semenova, E., Logacheva, M.D., Datsenko, K.A., and Severinov,

886    K. (2014). Pervasive generation of oppositely oriented spacers during CRISPR adaptation.

887    Nucleic Acids Res. *42*, 5907–5916.

888    Staals, R.H.J., Jackson, S.A., Biswas, A., Brouns, S.J.J., Brown, C.M., and Fineran, P.C. (2016).

889    Interference-driven spacer acquisition is dominant over naive and primed adaptation in a native

890    CRISPR–Cas system. Nat. Commun. *7*, 12853.

891    Sternberg, S.H., Richter, H., Charpentier, E., and Qimron, U. (2016). Adaptation in CRISPR-Cas

892    Systems. Mol. Cell *61*, 797–808.

893    Stringer, A.M., Singh, N., Yermakova, A., Petrone, B.L., Amarasinghe, J.J., Reyes-Diaz, L.,

894    Mantis, N.J., and Wade, J.T. (2012). FRUIT, a scar-free system for targeted chromosomal

895    mutagenesis, epitope tagging, and promoter replacement in Escherichia coli and Salmonella

896    enterica. PloS One *7*, e44841–e44841.

897    Strotskaya, A., Savitskaya, E., Metlitskaya, A., Morozova, N., Datsenko, K.A., Semenova, E.,

898    and Severinov, K. (2017). The action of Escherichia coli CRISPR-Cas system on lytic

899    bacteriophages with different lifestyles and development strategies. Nucleic Acids Res *45*, 1946–

900    1957.

901    Swarts, D.C., Mosterd, C., van Passel, M.W.J., and Brouns, S.J.J. (2012). CRISPR Interference

902    Directs Strand Specific Spacer Acquisition. PLoS ONE *7*, e35888.

903    Westra, E.R., van Erp, P.B.G., Künne, T., Wong, S.P., Staals, R.H.J., Seegers, C.L.C., Bollen, S.,

904    Jore, M.M., Semenova, E., Severinov, K., et al. (2012). CRISPR immunity relies on the

905    consecutive binding and degradation of negatively supercoiled invader DNA by Cascade and

906    Cas3. Mol. Cell *46*, 595–605.

907    Wright, A.V., Nuñez, J.K., and Doudna, J.A. (2016). Biology and Applications of CRISPR

908    Systems: Harnessing Nature's Toolbox for Genome Engineering. Cell *164*, 29–44.

909    Xue, C., Seetharam, A.S., Musharova, O., Severinov, K., Brouns, S.J., Severin, A.J., and

910    Sashital, D.G. (2015). CRISPR interference and priming varies with individual spacer sequences.

911    Nucleic Acids Res *43*, 10831–10847.

912    Yosef, I., Shitrit, D., Goren, M.G., Burstein, D., Pupko, T., and Qimron, U. (2013). DNA motifs

913    determining the efficiency of adaptation into the Escherichia coli CRISPR array. Proc. Natl.

914    Acad. Sci. U. S. A. *110*, 14396–14401.

915

916

917    **Table 1. Strains used in this study.**

| Name | Description | Source |
|------|-------------|--------|
| MG1655 | Wild-type *E. coli* | (Blattner et al., 1997) |
| CB386 | MG1655 [Δ*cas3* P*cse1*]::[cat PJ23119] | (Luo et al., 2014) |
| AMD536 | MG1655 [Δ*cas3* P*cse1*]::[ PJ23119] | (Cooper et al., 2018) |
| AMD688 | MLS1003 [Δ*cas3* P*cse1*]::[ PJ23119] | (Cooper et al., 2018) |
| 1XDNAi | MG1655 Δ*araC-araBAD*, ΔP*lacI:lacI*, Δ*cas-*CRISPR::pACT-01 (Kan-) | (Caliando and Voigt, 2015) |
| AMD671 | MG1655 Δ*araC-araBAD*, ΔP*lacI:lacI*, Δcas-CRISPR::pACT-01, *cas1*$^{+}$ *cas2*$^{+}$, Δ*lacZ::thyA* | This Study |
| AMD052 | MG1655 Δ*thyA* | (Stringer et al., 2012) |

918

919

920    **Table 2. Plasmids used in this study.**

| Name | Description | Source |
|------|-------------|--------|
| pBAD24 amp | Empty pBAD24 amp | (Guzman et al., 1995) |
| pBAD33 cam | Empty pBAD33 cam | (Guzman et al., 1995) |
| pAMD179 | Parent vector for cloning crRNAs | (Cooper et al., 2018) |
| pAMD189 | stp; pAMD179 expressing a self-targeting crRNA | (Cooper et al., 2018) |
| pAMD191 | pBAD33-*cas3* | (Cooper et al., 2018) |
| pCB380 | pcrRNA.con-*lacZ* | (Luo et al., 2014) |
| pAMD211 | pAMD179 expressing a crRNA targeting *mhpT* | This Study |
| pAMD212 | pAMD179 expressing a crRNA targeting *codA* | This Study |
| pSDS009 | Parent vector for cloning pre-spacers | This Study |
| pSK013 | pSDS009 stp2-mut1 pre-spacer | This Study |
| pSK014 | stp2; pAMD179 with an added pre-spacer | This Study |
| pSK015 | pSDS009 stp2-mut2 pre-spacer | This Study |
| pSK035 | pSDS009 stp2-mut3 pre-spacer | This Study |
| pSK016 | pSDS009 stp2-mut5 pre-spacer | This Study |
| psK017 | pSDS009 stp2-mut4 pre-spacer | This Study |

921

922

923 **Table 3. Oligonucleotides and synthesized dsDNA used in this study.**

| Name | Sequence |
|------|----------|
| JW472 | CCGACGCGCAGTTTA |
| JW473 | CACGTTGTGTTTTCATGC |
| JW6518 | CAGCGGGGATAAACC |
| JW7816 | TCGTCGGCAGCGTCAGATGTGTATAAGAGACAGAAGGTTGGTGGGTTGTTTTTATGGG |
| JW7817 | GTCTCGTGGGCTCGGAGATGTGTATAAGAGACAGTTATCAATTACAACCGACAGGGAGCC |
| JW7818 | TCGTCGGCAGCGTCAGATGTGTATAAGAGACAGAAAGTTGGTAGATTGTGACTGGC |
| JW7819 | GTCTCGTGGGCTCGGAGATGTGTATAAGAGACAGCAACAGCAGCACCCATGAC |
| JW7820 | AATGATACGGCGACCACCGAGATCTACACGCGTAAGATCGTCGGCAGCGTC |
| JW8010 | GCGCGGGGAACTCGAGCAGTCGCGCTTTGTCGAAACCGTTGCTGCCGGTTTATCCCCGC |
| JW8011 | GCGCGGGGAACTCGAGGTGCCAGGGCATACAAAACGCTTTGCCCACGGTTTATCCCCGC |
| JW8053 | AATGATACGGCGACCACCGAGATCTACACTCCAGGTATCGTCGGCAGCGTCAGATGTG |
| JW8054 | CAAGCAGAAGACGGCATACGAGATGGATTCACGTCTCGTGGGCTCGGAGATGTG |
| JW8057 | CAAGCAGAAGACGGCATACGAGATCGCATTAGGTCTCGTGGGCTCGGAGATGTG |
| JW8062 | CAAGCAGAAGACGGCATACGAGATATGGACTCGTCTCGTGGGCTCGGAGATGTG |
| JW8476 | CAAGCAGAAGACGGCATACGAGATTAAGGCGAGTCTCGTGGGCTCGGAGATGTG |
| JW8477 | CAAGCAGAAGACGGCATACGAGATCGTACTAGGTCTCGTGGGCTCGGAGATGTG |
| JW8478 | CAAGCAGAAGACGGCATACGAGATAGGCAGAAGTCTCGTGGGCTCGGAGATGTG |
| JW8479 | CAAGCAGAAGACGGCATACGAGATTCCTGAGCGTCTCGTGGGCTCGGAGATGTG |
| JW8480 | CAAGCAGAAGACGGCATACGAGATGGACTCCTGTCTCGTGGGCTCGGAGATGTG |
| JW8481 | CAAGCAGAAGACGGCATACGAGATTAGGCATGGTCTCGTGGGCTCGGAGATGTG |
| JW8485 | CAAGCAGAAGACGGCATACGAGATCGAGGCTGGTCTCGTGGGCTCGGAGATGTG |
| JW9009 | ATTGGGCCAGCTAAATCG |
| JW9010 | GGCTCATTATACCAGTCAGGACGTTGGGAAGAGGCCGCTCAAACAGGTAAAAAAGACACC |

| JW9016 | GCTAAATCGATGGGATGTGGCTTGCTATCTTTGGCTCCACTGTGATAGACAGCTGCATGCAT |
|--------|-------------------------------------------------------------------|
| JW9017 | AGAACTGGCTCATTATACCAGTCAGGACGTTGGGAAGAGGCCGCGTGTAGGCTGGAGCTG |
| JW9066 | TGTGGAATTGTGAGCGGATAACAATTTCACACAGGAAACAGCTTAGACAGCTGCATGCAT |
| JW9067 | TTCCTTACGCGAAATACGGGCAGACATGGCCTGCCCGGTTATTAGTGTAGGCTGGAGCTG |
| JW9076 (dsDNA) | TGTTTGACAGCTTATCATCGATTTGACAGCTAGCTCAGTCCTAGGTATAATGCTAGCATAAACCGCGGTACTTAGCTCCTCAGATTAGGATTGCGGAGAATAACAACCGCCGTTCTCATCGAGTAATCTCCGGATATCGACCCATAACGGGCAATGATAAAAGGAGTAACCTGTGAAAAAGATGCAATCTATCGTACTCGCACTTTCCCTGGTTCTGGTCGCTCCCATGGCAGCACAGGCTGCGGAAATTACGTTAGTCCCGTCAGTAAAATTACAGATAGGGAGGAGCGATCCTGCAGTGTTCCCCGCGCCAGCGGGGATAAACCGAGGGAACTGCCAGGCATCAAATAAAACGAAAGGCTCAGTCGAAAGACTGGGCCTTTCGTTTTATCTGTTGTTTGTCGGTGAACGCTCTCCTGAGTAGGACAAGCTTGGCTGTTTTGGCGGA |
| JW10009 | TGAGTAGGACAAGCTTACTAATTATAATAGAAGCCAGATACTAAAAAGCTTGGCTGTTTT |
| JW10010 | AAAACAGCCAAGCTTTTTAGTATCTGGCTTCTATTATAATTAGTAAGCTTGTCCTACTCA |
| JW10011 | TGAGTAGGACAAGCTTACTAATTATAATAGAAACCAGATACTAAAAAGCTTGGCTGTTTT |
| JW10012 | AAAACAGCCAAGCTTTTTAGTATCTGGTTTCTATTATAATTAGTAAGCTTGTCCTACTCA |
| JW10013 | TGAGTAGGACAAGCTTGCCGGCTGCAGCAGAAGCCAGATACTAAAAAGCTTGGCTGTTTT |
| JW10014 | AAAACAGCCAAGCTTTTTAGTATCTGGCTTCTATTATATATAGTAAGCTTGTCCTACTCA |
| JW10015 | TGAGTAGGACAAGCTTGCCGGCTGCAGCAGAAACCAGATACTAAAAAGCTTGGCTGTTTT |
| JW10016 | AAAACAGCCAAGCTTTTTAGTATCTGGCTTCTATTATATATAGTAAGCTTGTCCTACTCA |
| JW10017 | TGAGTAGGACAAGCTTACTAATTATAATAGAAACCAGATACCCCGAAGCTTGGCTGTTTT |
| JW10018 | AAAACAGCCAAGCTTCGGGGTATCTGGTTTCTATTATAATTAGTAAGCTTGTCCTACTCA |
| JW10021 | TGAGTAGGACAAGCTTACTAATTATAATAGAAACCAGATACCCAAAAGCTTGGCTGTTTT |
| JW10022 | AAAACAGCCAAGCTTTTGGGTATCTGGTTTCTATTATAATTAGTAAGCTTGTCCTACTCA |

924

925     **FIGURE LEGENDS**

926

927     **Figure 1. Primed adaptation is >1000-fold more efficient than naïve adaptation. (A)**

928     Schematic showing the *yfp* reporter construct of strain AMD688 used to quantify adaptation. The

929     genome of the parent strain has a single CRISPR repeat embedded within a *yfp* open reading

930     frame such that the repeat causes a frame-shift, preventing translation of YFP. Expansion of the

931     CRISPR array by a single repeat and spacer restores the frame, permitting translation of YFP,

932     leading to detectable fluorescence. **(B)** Percentage of YFP⁺ cells for an *E. coli* strain containing

933     the *yfp* reporter (AMD688), either the stp (pAMD189; "crRNA +") or an equivalent empty

934     vector (pBAD24; "crRNA -"), and either a *cas3*-expressing plasmid (pAMD191; "cas3 +") or an

935     equivalent empty vector (pBAD33; "cas3 -"), following addition of arabinose to the cultures.

936

937     **Figure 2. Primed adaptation from a plasmid results in primed adaptation from shared**

938     **chromosomal locations. (A)** Frequency of pre-spacer usage across the stp (pAMD189) during

939     primed adaptation in AMD536 cells containing pAMD189 and pAMD191. Values represent the

940     number of sequenced instances of each pre-spacer, plotted as a function of position on the

941     plasmid. Positive values indicate pre-spacers on the forward strand; negative values indicate pre-

942     spacers on the reverse strand. **(B)** Frequency of pre-spacer usage on the stp compared for spacers

943     acquired in the CRISPR-I and CRISPR-II arrays. Data are only plotted for 33 nt pre-spacers that

944     were sequenced at least once in both datasets. **(C)** Frequency of pre-spacer usage across

945     pAMD191 during primed adaptation in AMD536 cells containing pAMD189 and pAMD191.

946     Only pre-spacers unique to pAMD191 are shown; grey shaded regions indicate shared sequence

947     with pAMD189. **(D)** Frequency of pre-spacer usage across a chromosomal region encompassing

948    the *araC* gene during primed adaptation in AMD536 cells containing pAMD189 and pAMD191.

949    Values are plotted as a function of genome position. Only pre-spacers unique to the chromosome

950    are shown; the grey shaded region indicates shared sequence with pAMD189.

951

952    **Figure 3. Primed adaptation on the chromosome occurs over regions of >100 kb.** Frequency

953    of pre-spacer usage across a chromosomal region encompassing the *lacZ+* and *mhpT-*

954    protospacers for AMD536 cells containing pAMD191 (encodes Cas3), and either pCB380

955    (encodes crRNA targeting *lacZ+*; data shown in black) or pAMD212 (encodes crRNA targeting

956    *mhpT-*; data shown in blue). Values represent the relative number of sequenced instances of each

957    pre-spacer, plotted as a function of genome position. Positive values indicate pre-spacers on the

958    forward strand; negative values indicate pre-spacers on the reverse strand. Green boxes indicate

959    positions of repetitive insertion element sequences. The pink box indicates the region analyzed in

960    Figure 6B.

961

962    **Figure 4. Off-target Cascade binding sites are not associated with detectable primed**

963    **adaptation.** Frequency of pre-spacer usage in the 10 kb upstream of the on-target *lacZ+*

964    protospacer (orange datapoint), or off-target protospacers identified by (Cooper et al., 2018), on

965    the non-target strand, only counting 33 nt pre-spacers. Relative Cascade association with each

966    target sequence, as determined by ChIP-seq (Cooper et al., 2018), is shown on the x-axis.

967

968    **Figure 5. Replicate primed adaptation datasets for chromosomal protospacers are highly**

969    **reproducible.** Frequency of pre-spacer usage across the chromosome for each of two replicates

970    for AMD536 cells containing pAMD191 (encodes Cas3), and **(A)** pCB380 (encodes crRNA

971    targeting *lacZ+*), **(B)** pAMD212 (encodes crRNA targeting *mhpT-*), or **(C)** pAMD211 (encodes

972    crRNA targeting *codA+*). Values represent the number of sequenced instances of each pre-

973    spacer, for every 33 nt pre-spacer that mapped to the chromosome.

974

975    **Figure 6. The majority of primed adaptation occurs on the non-target strand, upstream of**

976    **the protospacer. (A)** Distribution of pre-spacer usage for *lacZ+* (black) and *mhpT-* (blue)

977    between the non-target (positive values) and target (negative values) strands, and between the

978    upstream (dark colors) and downstream (light colors) directions. **(B)** Frequency of pre-spacer

979    usage in the 10 kb upstream of the *mhpT-* protospacer on the target strand for AMD536 cells

980    containing pAMD191 (encodes Cas3), and pCB380 (encodes crRNA targeting *lacZ+*; x-axis) or

981    pAMD212 (encodes crRNA targeting *mhpT-*; y-axis). Values represent the number of sequenced

982    instances of each pre-spacer, for every 33 nt pre-spacer associated with an AAG PAM.

983

984    **Figure 7. Slipping accounts for the majority of pre-spacers on the non-target strand that**

985    **are not associated with an AAG PAM. (A)** Normalized frequency of distances ("slipping

986    distance"; x-axis) between uniquely positioned protospacers with an AAG PAM and uniquely

987    positioned protospacers with a non-AAG PAM, for the 10 kb upstream of the *lacZ+* (black

988    datapoints) and *mhpT-* (blue datapoints) protospacers, on the non-target strand. Values were

989    normalized to the total number of pairwise comparisons. **(B)** Distribution of pre-spacer usage for

990    all pre-spacers with a non-AAG PAM, for the 10 kb upstream of the protospacer on the non-

991    target strand. Values represent averages from the *lacZ+* and *mhpT-* experiments, with error

992    values representing one standard deviation from the mean. **(C)** DNA sequence logos for PAMs

993    from pre-spacers with a non-AAG PAM, for the 10 kb upstream of each of the *lacZ+* and *mhpT-*

994    protospacers on the non-target strand, excluding pre-spacers associated with -2, -1, +1 or +2

995    slipping events. Logos were generated from all pre-spacer usage (i.e. not just unique pre-

996    spacerse).

997

998    **Figure 8. Flipping often occurs following a slipping event. (A)** Schematic showing different

999    possible flipping events. The yellow highlighted sequences indicate the AAG PAM associated

1000   with the non-flipped pre-spacer. **(B)** Normalized frequency of distances ("slipping distance"; x-

1001   axis) between uniquely positioned protospacers with an AAG PAM on the non-target strand, and

1002   uniquely positioned protospacers on the target strand, for the 10 kb upstream of the *lacZ+* (black

1003   datapoints) and *mhpT-* (blue datapoints) protospacers on the non-target strand (i.e. 10 kb

1004   downstream of the protospacers on the target strand). Values were normalized to the total

1005   number of pairwise comparisons.

1006

1007   **Figure 9. Pre-spacers with non-canonical lengths are often associated with slipping. (A)**

1008   Schematic showing pre-spacers of non-canonical lengths associated with an AAG PAM in either

1009   a slipped or non-slipped configuration. **(B)** Normalized frequency of distances ("slipping

1010   distance"; x-axis) between uniquely positioned, 33 nt protospacers with an AAG PAM and

1011   uniquely positioned 32 nt protospacers, or **(C)** 33 nt protospacers, for the 10 kb upstream of the

1012   *lacZ+* (black datapoints) and *mhpT-* (blue datapoints) protospacers, on the non-target strand.

1013   Values were normalized to the total number of pairwise comparisons.

1014

1015   **Figure 10. Pre-spacers with 200 nt of the protospacer are used less frequently. (A)**

1016   Frequency of pre-spacer usage across a chromosomal region encompassing the *lacZ+* and *codA+*

1017    protospacers for AMD536 cells containing pAMD191 (encodes Cas3), and either pCB380

1018    (encodes crRNA targeting *lacZ*+; data shown in black) or pAMD211 (encodes crRNA targeting

1019    *codA*+; data shown in purple). Values represent the relative number of sequenced instances of

1020    each pre-spacer, plotted as a function of genome position. Data are only plotted for pre-spacers

1021    on the reverse strand. The blue box indicates the 10 kb region analyzed in panel B. **(B)**

1022    Frequency of pre-spacer usage in the 10 kb upstream of the *lacZ*+ protospacer on the non-target

1023    strand for AMD536 cells containing pAMD191 (encodes Cas3), and pAMD211 (encodes crRNA

1024    targeting *codA*+; x-axis) or pCB380 (encodes crRNA targeting *lacZ*+; y-axis). Values represent

1025    the number of sequenced instances of each pre-spacer, for every 33 nt pre-spacer associated with

1026    an AAG PAM. Data for pre-spacers within 200 nt of the *lacZ*+ protospacer are circled, and the

1027    distance from the *lacZ*+ protospacer is indicated.

1028

1029    **Figure 11. Pre-spacers that contain an AAG are used less frequently.** Frequency of pre-

1030    spacer usage across the 100 kb upstream of the **(A)** *lacZ*+ and **(B)** *mhpT*- protospacers for

1031    AMD536 cells containing pAMD191 (encodes Cas3), and either pCB380 (encodes crRNA

1032    targeting *lacZ*+) or pAMD212 (encodes crRNA targeting *mhpT*-). One was added to all

1033    frequency values to allow visualization of zero-scoring datapoints. Data are only shown for 33 nt

1034    protospacers on the non-targets strand, associated with an AAG PAM. Orange datapoints

1035    represent pre-spacers with an internal AAG. **(C)** Frequency of tetranucleotide sequence usage

1036    within pre-spacers that were de-enriched (x-axis) or de-enriched (y-axis) relative to an

1037    exponential decay model. Values represent the percentage of all tetranucleotides, averaged for

1038    data from the *lacZ*+ and *mhpT*- protospacer datasets. Error bars represent one standard deviation

1039    from the mean. Tetranucleotide sequences containing AAG are plotted in orange. **(D)** Same as

1040    (C), but for trinucleotide sequences.

1041

1042    **Figure 12. Experimental evidence that pre-spacers containing an AAG are used less**

1043    **frequently in primed adaptation, and to a lesser degree in naïve adaptation.** **(A)** Schematic

1044    showing the unmodified and modified pre-spacer sequences for the stp2, stp2-mut1, stp2-mut2

1045    and stp2-mut3 plasmids. The PAM is indicated by a dashed box. Grey highlighting indicates

1046    changes relative to the unmodified pre-spacer. **(B)** Frequency of pre-spacer usage for all pre-

1047    spacers on the unmodified and mut1 stp2 plasmids, for pre-spacers on the non-target strand with

1048    an AAG PAM. The single modified pre-spacer that differs between the stp2 and stp2-mut1 is

1049    shown in orange. The frequency of pre-spacer usage from the stp2 and stp2-mut1 plasmids are

1050    shown on the *x*-axis and *y*-axis, respectively. **(C)** Tukey boxplot showing the distribution of pre-

1051    spacers usage or all chromosomal sequences with an AAG PAM that were detected at least once

1052    in an assay of naïve adaptation (Yosef et al., 2013). Pre-spacers are separated depending on the

1053    presence or absence of an internal AAG sequence.

1054

1055    **Figure 13. The impact of an AAG within a pre-spacer on primed adaptation efficiency is**

1056    **dependent on position of the AAG relative to the PAM.** **(A)** Ratio of actual pre-spacer usage

1057    frequency to the frequency predicted by an exponential decay model for all pre-spacers with an

1058    internal AAG. Data are only shown for pre-spacers with an AAG PAM, on the non-target strand,

1059    <100 kb upstream of the *lacZ*+ (black datapoints) or *mhp*T- (blue datapoints) protospacer. Ratios

1060    are shown as a function of the position of the internal AAG sequence within the pre-spacer. **(B)**

1061    Frequency of pre-spacer usage for all pre-spacers on the stp2-mut1 and stp2-mut2 plasmids, for

1062    pre-spacers on the non-target strand with an AAG PAM. The single modified pre-spacer that

1063    differs between stp2-mut1 and stp2-mut2 is shown in orange. The frequency of pre-spacer usage

1064    from the stp2-mut1 and stp2-mut2 plasmids are shown on the x-axis and y-axis, respectively. **(C)**

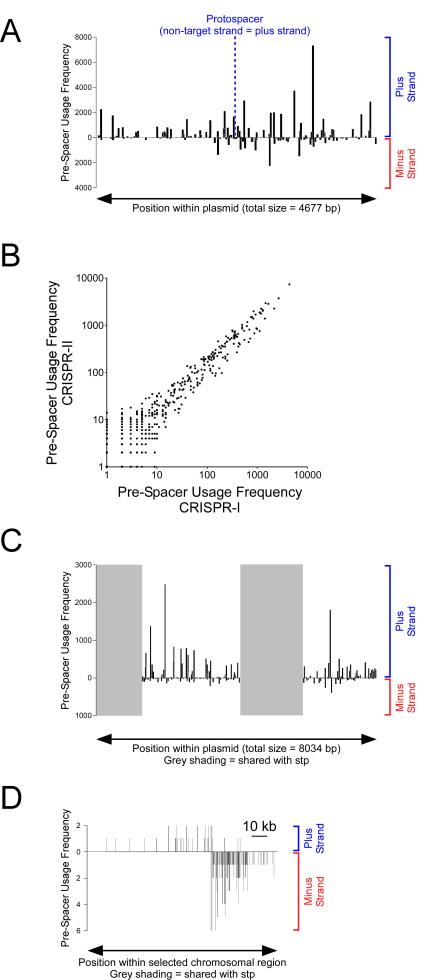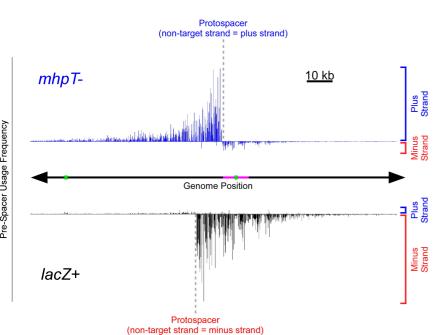1065    As for (B) but for the stp2 and stp2-mut3 plasmids.
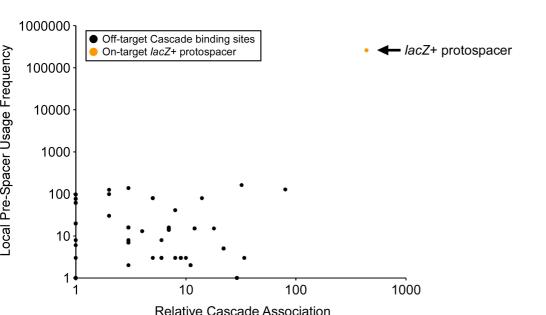
1066

1067    **Figure 14. Sequences at the PAM-distal end of the pre-spacer impact the efficiency of**

1068    **primed adaptation. (A)** DNA sequence logo for the 133 "enriched" pre-spacers without an

1069    internal AAG that are used >4-fold more frequently than expected based on the exponential

1070    decay model. Note that the y-axis maximum is set to 0.5 bits. **(B)** DNA sequence logo for the

1071    195 "de-enriched" pre-spacers without an internal AAG that are used >4-fold less frequently

1072    than expected based on the exponential decay model. Note that the y-axis maximum is set to 0.5

1073    bits. **(C)** Schematic showing the unmodified and modified pre-spacer sequences for the stp2,

1074    stp2-mut4, and stp2-mut5 plasmids. The PAM is indicated by a dashed box. Grey highlighting

1075    indicates changes relative to the unmodified pre-spacer. **(D)** Frequency of pre-spacer usage for

1076    all pre-spacers on the stp2 and stp2-mut3 plasmids, for pre-spacers on the non-target strand with

1077    an AAG PAM. The single modified pre-spacer that differs between stp2 and stp2-mut4 is shown

1078    in orange. The frequency of pre-spacer usage from the stp2 and stp2-mut4 plasmids are shown

1079    on the *x*-axis and *y*-axis, respectively. **(E)** As for (D) but for the stp2 and stp2-mut5 plasmids.

1080

Figure 1

Figure 2

Figure 3

Figure 4

Figure 5

A



*lacZ+*

B

*mhpT-*

C

*codA+*

Figure 6

Figure 7

A



B



| | | % of non-AAG pre-spacers |
|---|---|---|
| -1 slip | N**AAG**NNN... | 22.8 ± 8.2% |
| Canonical pre-spacer | **AAG**NNN... | N/A |
| +1 slip | **AG**NNN... | 36.7 ± 2.7% |
| +2 slip | **G**NNN... | 1.5 ± 0.2% |

C

Figure 8

Figure 9

A

```
32 nt, +1 slip         AAGN NNNNNNNNNNNNNNNNNNNNNNNNNNNN
32 nt,  no slip        AAG NNNNNNNNNNNNNNNNNNNNNNNNNNNN
Canonical length, no slip  AAG NNNNNNNNNNNNNNNNNNNNNNNNNNNN
34 nt,  no slip        AAG NNNNNNNNNNNNNNNNNNNNNNNNNNNNNN
34 nt, -1 slip       NAAG NNNNNNNNNNNNNNNNNNNNNNNNNNNNN
```

B

32 nt pre-spacers

C

34 nt pre-spacers

Figure 10

Figure 11

Figure 12

A

PAM

```
AAGCTTACTAATTATAATAGAAACCAGATACTAAA   Unmodified (stp2)
AAGCTTACTAATTATAATAGAAGCCAGATACTAAA   mut1
AAGCTTGCCGGCTGCAGCAGAAGCCAGATACTAAA   mut2
AAGCTTGCCGGCTGCAGCAGAAACCAGATACTAAA   mut3
```

-2 -1 1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18 19 20 21 22 23 24 25 26 27 28 29 30 31 32 33

Position within pre-spacer (nt)

B



C

Figure 13

Figure 14