

1 **Systematic Development of Reprogrammed Modular Integrases Enables Precise Genomic**
2 **Integration of Large DNA Sequences**

3 Friedrich Fauser^{1,2,3}, Sebastian Arangundy-Franklin^{1,2}, Jessica E Davis¹, Lifeng Liu¹, Nicola J
4 Schmidt¹, Luis Rodriguez¹, Nicholas A Scarlott¹, Rakshaa Mureli¹, Danny F Xia¹, Sarah J
5 Hinkley¹, Bhakti N Kadam¹, Nga Nguyen¹, Stephen Lam¹, Bryan Bourgeois¹, Emily Tait¹,
6 Mohammad Qasim¹, Vishvesha Vaidya¹, Adeline Chen¹, Andrew Nguyen¹, Patrick Li¹, David E
7 Paschon¹, Gregory D Davis¹ and Jeffrey C Miller^{1,3}

8
9 ¹Sangamo Therapeutics, Inc., Richmond, CA, USA

10 ²These authors contributed equally

11 ³E-mail: ffauser@sangamo.com and jmiller@sangamo.com

12
13 **Summary**

14 Despite recent progress in the ability to manipulate the genomes of eukaryotic cells¹⁻³, there
15 is still no effective and practical method to precisely integrate large synthetic DNA constructs into
16 desired chromosomal sites using a programmable integrase. Serine integrases can perform the
17 necessary molecular steps⁴, but only if their natural target site is first installed into the recipient
18 genome by other methods. A more elegant approach would be to directly reprogram the serine
19 integrase itself to target a desired site endogenous to the genome that is different from the natural
20 recognition site of the integrase⁵. Here, we describe the development of a platform of Modular
21 Integrases (the MINT platform), a versatile protein-guided genome editing method that can
22 facilitate site-directed targeted integration of synthetic DNA into chromosomal sites. Through a
23 combination of structural modeling, directed evolution, and screening in human cells we have
24 reprogrammed the specificity of the serine integrase Bxb1. We then utilized these reprogrammed
25 Bxb1 variants to enable precise integration of kilobase-sized constructs into multiple endogenous
26 locations within the human genome with up to 35% efficiency and promising genome-wide
27 specificity. We demonstrate the therapeutic potential of the MINT platform by retargeting Bxb1 to
28 the human TRAC and AAVS1 loci where wild-type Bxb1 has no detectable activity.

29 Introduction

30 The ability to integrate synthetic DNA constructs into desired chromosomal locations in
31 eukaryotic genomes would have broad implications for the development of genomic medicines,
32 synthetic biology, agriculture, and basic research. Initially, there was excitement that homology-
33 directed-repair stimulated by engineered nucleases would be able to accomplish this goal, but
34 numerous limitations were encountered when this was assessed⁶⁻⁹. Furthermore, despite the
35 significant progress in programmable DNA binding using zinc fingers (ZFs)¹⁰⁻¹², transcription
36 activator-like effectors (TALEs)^{13,14}, and RNA-guided proteins such as Cas9¹⁵, simply tethering
37 such engineered DNA binding domains to the catalytic domains of recombinases or transposons
38 utilizing a cut-and-paste mechanism has not yielded reagents capable of driving high levels of
39 integration in human cells¹⁶⁻²³. We initially explored such systems and observed high levels of
40 indels (Extended Data Fig. 1a) which may partially explain the difficulties others have
41 encountered.

42 Greater success was achieved by inserting ZFs into the tyrosine recombinase CRE to
43 perform inversions²⁴, however fusions to other types of integrases have not proved nearly as
44 successful²⁵⁻²⁹. There has been some recent progress using CRISPR-associated transposases^{30,31}
45 and other nucleic acid-guided systems³², but such systems require simultaneous delivery of
46 multiple components to mammalian cells and thus far have only achieved modest levels of
47 activity³³.

48 Large serine recombinases (LSRs - also known as serine integrases) have long been
49 proposed as the ideal tool for genome engineering^{34,35} due to their unique properties among other
50 integration systems (**Figure 1a**). In their natural context, an integrase dimer binds to a phage attP
51 site, while another dimer binds an attB site in the bacterial host's genome. The assembled and fully
52 active tetramer then forms a covalent bond to the DNA via its active-site serine residues producing
53 a temporary two base-pair overhang. Subunit rotation facilitates strand exchange and ligation,
54 without leaving free DNA ends or nicks, thus obviating the need to engage the host cell's DNA
55 repair mechanisms. Importantly, in the absence of a phage-encoded reversibility factor, integration
56 proceeds unidirectionally and irreversibly. Nevertheless, the use of LSRs as genome editing
57 reagents has been hampered by the difficulty in their engineering and deployment towards non-
58 cognate targets.

59 The best results to date have been achieved with systems that use a phage-derived serine
60 integrase such as Bxb1 combined with prime editing to simultaneously introduce the natural attB
61 target site into mammalian cells^{36,37}. However, the requirement to first use prime editing to
62 integrate the target site makes the process more complex for therapeutic development, requiring
63 the simultaneous delivery of numerous genetic components to the cell. A more straightforward
64 approach would involve directly reprogramming the target preference of the serine integrase itself
65 while retaining the desirable properties of natural LSRs.

66 Despite the utility of reprogramming an LSR such as Bxb1³⁸, there are several technical
67 challenges to overcome. First, while TALEs and CRISPR systems can easily be targeted to desired
68 DNA sequences using simple targeting rules, no one has been able to divine such a “simple code”
69 that governs Bxb1 retargeting. Directed evolution systems have engineered proteins such as ZFs
70 that lack a simple DNA targeting code, but the modularity and compactness of the ZF repeat
71 enables strategies that engineer one zinc finger at a time^{39,40}. Directed evolution can also be used
72 on non-modular proteins such as Cre, but this required over a hundred sequential rounds of directed
73 evolution for a given target site^{41,42}. A further challenge for reprogramming Bxb1 is that no
74 structure of Bxb1 bound to its target site is currently available.

75 In this study, we have used a combination of structural modeling and experimental
76 characterization to map critical protein-DNA interactions between Bxb1 and its target site. The
77 only known structures of an LSR suggest that the zinc ribbon domain (ZD) and the recombinase
78 domain (RD) are separated by a flexible polypeptide linker and recognize their portions of the attB
79 and attP targets in a modular fashion⁴³. On the assumption that Bxb1 recognizes its target site in a
80 similar way, we developed a strategy for LSR engineering that uses directed evolution to
81 reprogram the key specificity-determining amino acid residues of the RD and ZD domains in
82 parallel and then combine the successful variants into fully reprogrammed LSR variants. We used
83 this strategy to successfully target multiple sites endogenous to the human genome with up to 35%
84 efficiency and to target two therapeutically relevant loci with up to 1% efficiency.

85 **DNA-Bxb1 interaction mapping**

86 To compensate for the lack of a known structure of Bxb1 bound to its DNA target site, we
87 performed an initial experiment designed to map interactions between residues in the Bxb1 RD
88 and ZD domains and key regions of its natural DNA target site. Based on a sequence alignment of

89 Bxb1 and the regions of the LI integrase known to interact with its DNA target^{43,44}, we identified
90 89 residues in Bxb1 that are likely to interact with DNA (Extended Data Fig. 2). We performed
91 scanning mutagenesis on these 89 residues, combined those variants with plasmids bearing
92 symmetric changes to a single position of each half-site of both the attB and attP sequences that
93 are specified by wild-type Bxb1³⁵ and delivered them to human K562 cells. This experiment
94 identified two clusters of residues which, when mutated, could alter specificity at portions of the
95 DNA target site. A first cluster of residues at positions 231, 233, and 237 was able to alter
96 specificity at positions -10 and -9 of the DNA target site. A second cluster of residues at positions
97 314, 315, 316, 318, 323, and 325 was able to alter specificity at positions -18 through -13 of the
98 DNA target site. The residue at position 158 was also able to alter specificity at positions -7 and -
99 6 of the DNA target site. Saturation mutagenesis of these residues was able to achieve further
100 specificity shifts. Mapping the Bxb1 mutations that shifted targeting specificity onto a
101 Rosettafold⁴⁵ structural model of Bxb1 indicated two regions in the RD domain and one region in
102 the ZD domain where focused protein engineering should be able to alter the DNA targeting
103 specificity (**Figure 1b**). We refer to these regions as the helix, loop, and hairpin respectively based
104 on the structure of each region in our Bxb1 structural model (**Figure 1c,d**). Notably, the central
105 dinucleotide of the attachment sites is not directly recognized by Bxb1 but needs to match between
106 the attP and attB site for integration to occur.

107 **Directed evolution of Bxb1 DNA interaction domains**

108 In order to fully reprogram the specificity of Bxb1, screening of a large combinatorial
109 library of Bxb1 variants is needed. To achieve this, we developed a directed evolution system based
110 on a previously described method⁴⁶ whereby a “stuffer” sequence flanked by divergent attB and
111 attP target sequences is placed within the coding sequence of an antibiotic resistance gene such
112 that recombination of these test sequences by a Bxb1 variant expressed in the same bacterial cell
113 restores the open reading frame of the antibiotic resistance gene and allows the bacterial cell to
114 survive an antibiotic challenge (**Figure 2a**). In our system we placed the Bxb1 variant library and
115 the DNA target sites on separate plasmids to enable more rapid testing of the same library of Bxb1
116 variants against different DNA target sites (**Figure 2b**, Extended Data Fig. 3). This system also
117 provides a convenient means of assessing the DNA targeting specificity of selected Bxb1 variants
118 by using a single Bxb1 variant tested against a library of different target sites (**Figure 2c**).

119 To inform the design of our Bxb1 variant libraries, we analyzed the structure of the LI
120 integrase and observed that residues in the structure that correspond to residues 231-237 in Bxb1
121 form an alpha helix that docks with its target DNA in a manner reminiscent of the interaction
122 between a zinc finger (ZF) and its target trinucleotide⁴⁷. Thus, we also adopted the same residue
123 randomization scheme used with ZFs⁴⁰ and kept residue 235 fixed as a leucine since this seems
124 analogous to the leucine that is often at +4 of a ZF recognition helix. This resulted in a library of
125 Bxb1 helix variants with residues 231-234, and 236-237 randomized (**Figure 2b**).

126 Since ZFs can target 3 bp DNA sequences we hypothesized that the Bxb1 helix might also
127 be able to specify the DNA bases at positions -11, -10, and -9 of the target site. Thus, we performed
128 64 separate selections using our Bxb1 helix variant library against all 64 possible DNA triplets at
129 positions -11, -10 and -9. Upon deep sequencing the plasmids contained in bacteria that survived
130 a single round of antibiotic challenge, we observed enrichment of four-residue peptide motifs for
131 37 out of 64 selections. For these successful selections we identified the individual helix sequences
132 that best represented the enriched four-residue motifs and characterized the DNA targeting
133 preferences of Bxb1 variants bearing each selected helix individually. Many selected helices
134 demonstrated dramatic changes in target specificity at positions -11, -10, and -9 of the DNA target
135 site. A comparison of the target preference for the wild-type helix SATALKR and 19 selected helix
136 sequences is shown in **Figure 2d**.

137 Next, we used our directed evolution system to select Bxb1 variants with mutations at the
138 hairpin region (residues between 314 and 325). Since the entire hairpin region would not be
139 efficiently randomized by a typical library generated in E coli (in the range of 10e9 variants), we
140 adopted a structure-based scheme, where we fully randomized residues 314, 316, 318, 321, 323,
141 and 325 and partially randomized residue 322 (**Figure 2b**). The residues we fully randomized face
142 the DNA major groove in our structural model while position 322 is likely to influence hairpin
143 structure if it is a proline. We used this Bxb1 variant library in selections with target sites where
144 positions -19 to -12 matched potential Bxb1 target sites in the human genome. Several of our
145 selected hairpins have obvious target preference differences in comparison to the wild-type
146 hairpin, and the positions where the selected hairpins show the strongest sequence preferences
147 differ across selected hairpins (**Figure 2e**).

148 Finally, we performed selections using a library of Bxb1 variants where positions 154-159
149 (the loop region) were randomized (**Figure 2b**). We performed 16 separate selections using this
150 library against all possible DNA dimer sequences at positions -7 and -6 and found that,
151 surprisingly, a single residue change of S157G could shift the sequence preferences of wild-type
152 Bxb1 towards nearly every single dinucleotide at positions -7 and -6 except for CG (**Figure 2f**).
153 Other selected loop variants were able to show improved target specificity relative to the wild-type
154 loop, but the sequences that can be targeted specifically in the loop variants we have characterized
155 tend to be limited to sequences with A or T at position -7 and C or T at position -6 (**Figure 2f**).

156 **Identification of Bxb1 pseudo-sites in the human genome**

157 We next wanted to test the activity of our engineered Bxb1 variants at chromosomal target
158 sites in human cells. To avoid the likely issue of having no detectable integration with our initial
159 fully reprogrammed Bxb1 variants, we decided to first target sites in the human genome where
160 wild-type Bxb1 had detectable activity. In order to identify such Bxb1 pseudo-sites in human K562
161 cells, we performed both a computational search based on published Bxb1 specificity data³⁵ and
162 an experimental assessment utilizing anchored multiplex PCR⁴⁸ to map the location of integrated
163 donor constructs. To identify as many active sites as possible, we designed the computational
164 search to identify potential pseudo-sites with any sequence at the central dinucleotide and we
165 designed the experimental approach to use a pool of 16 different donor constructs each bearing a
166 different central dinucleotide (Extended Data Fig. 4a). Combining the results of both approaches,
167 we were able to identify 23 sites with at least 0.1% integration in human K562 cells and we selected
168 five sites with between 39% and 61% homology to the natural Bxb1 attB site where wild-type
169 Bxb1 achieves between ~0.20% and ~2.45% integration as test targets for our engineered Bxb1
170 variants (Extended Data Fig. 4b). Notably, with both computational and experimental approaches
171 we were only able to identify active attB pseudo-sites, and no attP pseudo-sites, mirroring the
172 direction of natural Bxb1-directed integration into its host genome. This is consistent with our
173 finding that “landing pad” cell lines with attB sites pre-integrated into the genome support higher
174 levels of integration than landing pad cells lines with attP sites pre-integrated (Extended Data Fig.
175 1b).

176 **Directed evolution yields Bxb1 variants with increased activity at human pseudo-sites**

177 As an initial performance validation in human cells of Bxb1 variants derived from our
178 directed evolution system, we tested selected Bxb1 helix variants for integration at a human
179 pseudo-site on chromosome 3 (Extended Data Fig. 5a). First, we performed a helix selection
180 against one half-site of this target site and tested selected variants as mixtures with wild-type Bxb1
181 since integration also requires binding to the other half-site and the attP site on the donor. This
182 selection yielded multiple families of helix sequences and we tested a total of eight selected helices
183 for their ability to carry out site-specific integration in K562 cells. Helix sequences with an S or T
184 at position 234 were active at the chromosome 3 target even in the absence of wild-type Bxb1.
185 Similar sequences were enriched in selections with other target sites, so we concluded that these
186 represented helix sequences with poor specificity, and we computationally removed similar
187 sequences from other selections. In contrast, selected Bxb1 variants with an N at position 234 such
188 as AGGNLKR were only active at the chromosome 3 target when mixed with wild-type Bxb1 and
189 the DNA targeting specificity characterization of this helix indicated a dramatic shift in specificity
190 towards a T at position -10 (**Figure 2d**) consistent with the GTC sequence at positions -11 to -9 of
191 the half-site we targeted with this helix selection (Extended Data Fig. 5a).

192 Next, we wanted to investigate whether helix and hairpin variants can be combined to
193 further engineer Bxb1. Initially we performed a selection using our Bxb1 variant hairpin library
194 against the same half-site of the chromosome 3 target that was used for the helix selections. We
195 then combined the most promising selected hairpin LARGRRKWARYR with both the AGGNLKR
196 helix and the more active, but less specific WSSSLKR helix. Both combinations retained full
197 activity at the chromosome 3 site when combined with wild-type Bxb1 but show a substantial
198 decrease without wild-type Bxb1. Thus, WSSSLKR in combination with LARGRRKWARYR
199 now requires wild-type Bxb1 to target the donor and/or other half-site with full activity indicating
200 improved specificity vs. the WSSSLKR helix variant alone. Since both AGGNLKR and
201 LARGRRKWARYR were able to target the intended chromosome 3 half-site more specifically
202 than the wild-type Bxb1 helix and hairpin, we reasoned the combination of this helix and hairpin
203 should be even more specific for the chromosome 3 site and characterized this combination of
204 helix and hairpin using an unbiased genome-wide specificity assay. This fully engineered Bxb1
205 variant appeared to have fewer integration sites within the human genome than wild-type Bxb1
206 and we observed that the intended site on chromosome 3 had higher levels of integrations than any
207 other site in the human genome (Extended Data Fig. 5b, Extended Data Fig. 6). Furthermore, if

208 we had performed the experiment using only a donor with a central dinucleotide that matched the
209 intended target, then only the intended target would have been identified in this experiment.

210 We then chose two additional endogenous targets identified by the computational genome
211 scan and two endogenous targets identified by the experimental genome-wide screen to serve as
212 DNA target sites for our directed evolution system. We used these targets as an initial test of a
213 screening system that utilizes a large plasmid library of artificial target sequences that will be
214 required to target sites where wild-type Bxb1 has no detectable activity (**Figure 3a**). This system
215 allowed us to rapidly identify evolved helix and hairpin variants that, when combined, are active
216 against their desired half-sites (**Figure 3b,c**, Extended Data Fig. 5d-e), demonstrating the
217 modularity of our MINT platform. We were then able to achieve substantial improvements in
218 activity relative to wild-type Bxb1 at 6 half-sites from these four human pseudo-sites (Extended
219 Data Fig. 5f). Four of these half-sites comprise full pseudo-sites in the *MACO1* and *GYS1* genes
220 and activity at these sites could be increased further by pairing the two variants that worked best
221 at the relevant left and right half-sites, achieving ~35% targeted integration at both loci (**Figure**
222 **3d**). For *MACO1* we observed 0.85% TI with wild-type Bxb1 and achieved a 41x increase in
223 activity with engineered Bxb1 variants, while we detected 0.23% TI with wild-type Bxb1 at *GYS1*
224 and achieved a 163x increase with engineered Bxb1 variants.

225 **Retargeting Bxb1 to clinically relevant sites**

226 Encouraged by our success improving Bxb1 performance at chromosomal pseudo-sites,
227 we proceeded to the more difficult challenge of targeting clinically relevant regions of the human
228 genome, such as the well-established safe harbor AAVS1 locus⁴⁹ and the T-cell receptor α constant
229 (TRAC) locus (**Figure 4a**), in which wild-type Bxb1 has no detectable activity. We first screened
230 Bxb1 variants against a library of potential target sites within these regions (Extended Data Fig.
231 7). This resulted in pairs of fully engineered Bxb1 variants (**Figure 4b**) for the TRAC and AAVS1
232 loci shown in **Figure 4a**. We achieved up to ~1% TI at the TRAC locus using our standard PCR-
233 based NGS assay (**Figure 4c**), with no detectable indels (Extended Data Figure 8). Similarly, we
234 achieved up to ~1% TI at the AAVS1 locus using a digital PCR-based assay (**Figure 4d**, Extended
235 Data Figure 9) and no detectable indels by standard PCR-based NGS assay. As in earlier
236 experiments, we utilized DNA donors with wild-type attP target sites with central dinucleotides
237 matched to the intended genomic target site and included wild-type Bxb1 together with the Bxb1

238 variants selected for the left and right half-sites of the genomic target site. Notably, we did not
239 observe any integration activity above background levels when we omitted the Bxb1 variant
240 targeted to either the left or right half-site, indicating full retargeting of Bxb1 to TRAC and AAVS1.

241 **Alternative delivery strategies and performance in human T cells**

242 We established our MINT platform in human K562 cells using routine plasmid delivery
243 protocols for the donor molecule. However, plasmid DNA is not an ideal delivery modality for
244 most relevant cell types. We envision that for therapeutic applications Bxb1 will be delivered as
245 mRNA while the donor may be delivered via AAV, minicircle, or comparable methods. To
246 demonstrate that donor delivery is compatible with clinically relevant delivery strategies, we
247 explored alternative delivery modalities for the genetic cargo. We tested both single-stranded AAV
248 (ssAAV) and partially double-stranded AAV (dsAAV) donors and observed that integration was
249 more efficient in our K562 landing pad cell line when making the attP site double-stranded by co-
250 delivering a complementary oligonucleotide (**Figure 5a**). We also used a self-complementary AAV
251 (scAAV) which resulted in up to ~25% TI (**Figure 5b**).

252 **Discussion and future directions**

253 In this work we have demonstrated the direct reprogramming of a site-specific LSR which
254 has been a long-standing challenge for genome engineering¹⁶. We achieved this through a
255 combination of structural modeling, directed evolution, and screening in human cells. We were
256 able to integrate up to ~3 kb of synthetic DNA into endogenous loci in the human genome with
257 efficiencies of up to 35%. We further demonstrate the therapeutic potential of our MINT platform
258 by retargeting Bxb1 to the human TRAC and AAVS1 loci where we achieved up to 1% integration.
259 In the future, we intend to deploy similar engineering techniques as outlined in this study to further
260 increase the activity of these reagents. We envision that our MINT platform will also accelerate
261 efforts in other research areas where Bxb1 was successfully used, such as for metabolite pathway
262 assembly⁵⁰ and various agrobiotechnology applications in model plants⁵¹ and crops^{52,53}. We expect
263 the modularity and efficiency of our approach will make it an ideal choice for applications where
264 error-free integration of large synthetic DNA constructs is required.

265 Techniques that use a cut-and-paste mechanism instead of a DNA replication dependent
266 method have no inherent limit to the size of the construct that can be integrated. This is not only

267 beneficial for the integration of synthetic DNA into a safe harbor site, as demonstrated in this study,
268 but also for targeting the first intron of a mutated gene to integrate the correct copy of the
269 corresponding coding sequence linked to a splice acceptor. Furthermore, targeting endogenous
270 attP-like sites in addition to attB-like sites would support other genomic rearrangements such as
271 deletions and inversions. Additionally, the ability to target two different endogenous loci in close
272 proximity would enable recombinase mediated cassette exchange (RCME) at endogenous loci
273 which could replace an entire genomic locus with a linear synthetic DNA donor construct. This
274 would enable unprecedented flexibility in types of genomic alterations that can be generated⁵⁴. We
275 were able to successfully target the ~4 kb human AAVS1 locus and a ~2 kb portion of the human
276 TRAC locus which demonstrates a high targeting density that would support RCME at an
277 endogenous locus.

278 In this study we established a method for retargeting Bxb1 that will likely enable the
279 reprogramming of other serine integrases that have recently been described^{37,55}. Since the RD and
280 ZD domains of Bxb1 behave in a modular fashion it is likely that RD and ZD domains from
281 different LSRs can be combined to further expand the ability to target desired genomic loci.
282 Ultimately, our goal is to create an archive of pre-characterized modular protein-guided integrases.
283 Such a pre-characterized archive of integrase variants would allow the *in silico* design of variants
284 to target any sequence of interest without the need for custom directed evolution selections. The
285 LSR engineering approach we demonstrated in this study serves as a blueprint for retargeting LSRs
286 to other gene-sized loci for therapeutic applications and beyond.

287 **Data availability**

288 Amino acid sequences and DNA sequences of constructs used in this study are provided
289 upon request. Illumina sequencing data underlying all experiments will be deposited in the NCBI
290 Sequence Read Archive before peer reviewed publication. Source data are provided with this
291 paper.

292 **Code availability**

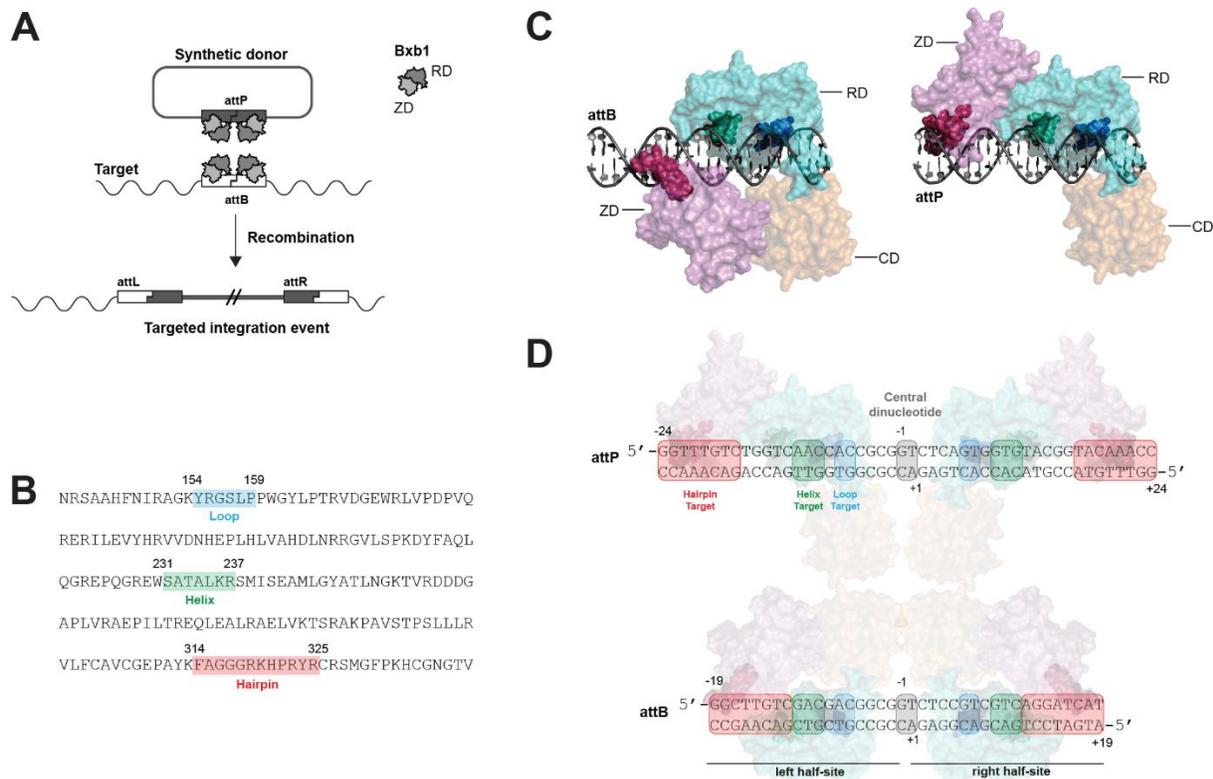
293 Custom computer code is available upon request although comparable analysis can be
294 performed with publicly available software.

295 **Acknowledgements**

296 We thank Gustavo de Alencastro for assistance with AAV construct design, as well as
297 Sumita Bhardwaj, Tammy Chen, Alicia Goodwin, Emma Petrouski, Sanjna Sridhar, and Hung
298 Tran for assistance with AAV production. We also thank Ed Rebar, Jon Melnick, Deepak Patil, and
299 Charles Paine for contributions to the early development of ZF-targeted Serine Recombinases. We
300 also thank Sandy Macrae, Adrian Woolfson, and Jason Fontenot for their relentless support of this
301 project.

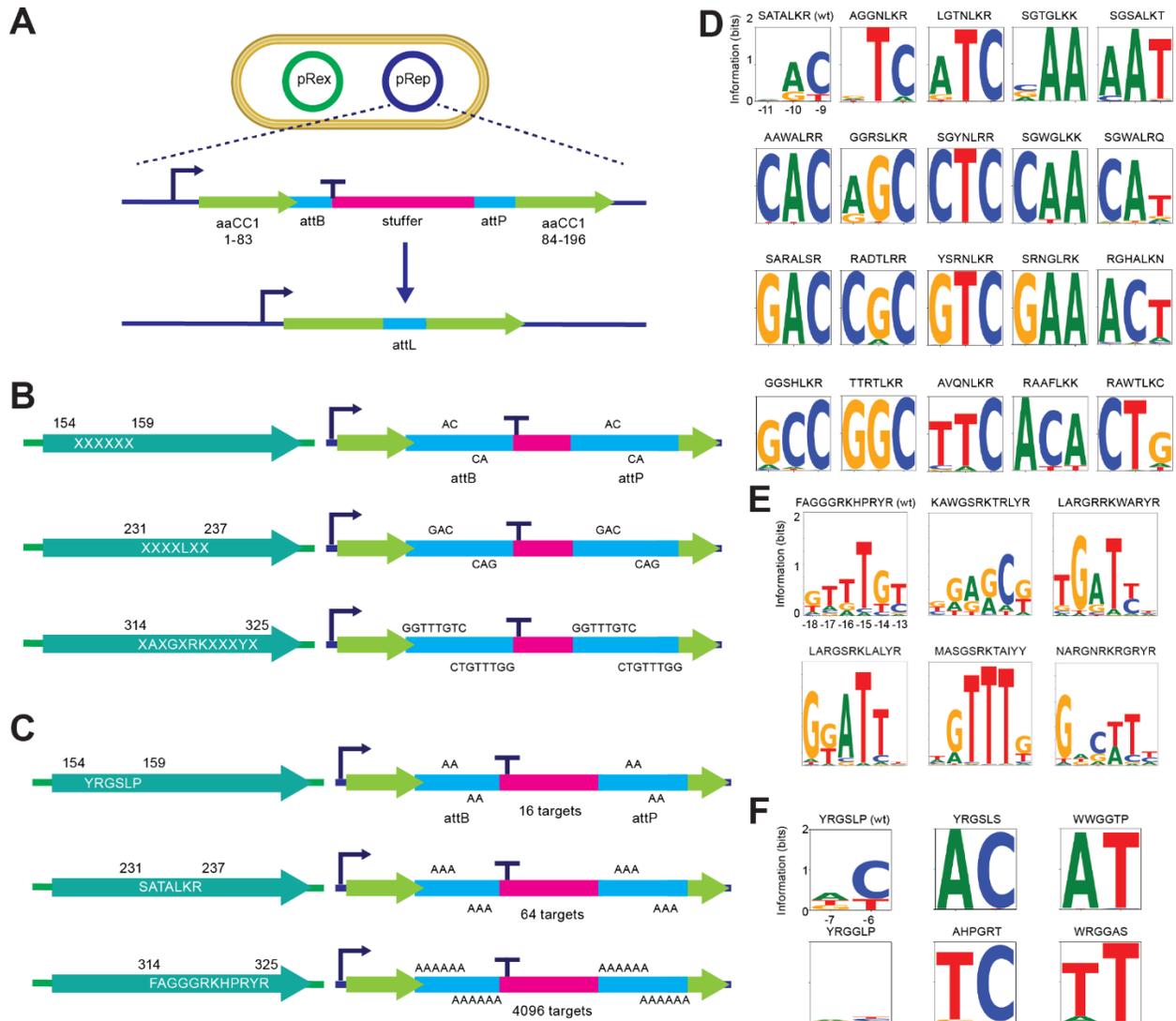
302 **Competing interests**

303 All authors contributed to this work as full-time employees of Sangamo Therapeutics.
304 Sangamo Therapeutics has filed patent applications regarding Integrase systems described in this
305 study, listing F.F., S.A.-F., N.A.S., and J.C.M. as inventors.



306

307 **Figure 1. Bxb1-DNA interaction mapping.** **a.** Schematic of Bxb1-induced targeted integration
308 of a synthetic circular donor into a genomic target. Attachment sites are bound by Bxb1 dimers
309 and non-reversible recombination is facilitated by a Bxb1 tetramer. **b.** Amino acid sequence of
310 wild-type Bxb1 with residues investigated in this study highlighted. **c.** Structural model of Bxb1
311 bound to attP and attB attachment sites. Bxb1-DNA interaction mapping identified three
312 specificity-determining regions of Bxb1 that can be reprogrammed: the loop and helix regions
313 (residues 154-159 and 231-237 shown in blue and green respectively) in the RD domain of Bxb1
314 and the hairpin region (residues 314-325 shown in red) within the ZD domain of Bxb1. **d.** Sequence
315 of the natural Bxb1 attP and attB attachment sites with the portions recognized by the hairpin,
316 helix, and loop regions of Bxb1 highlighted. The central dinucleotide of the attachment sites is not
317 directly recognized by Bxb1 but needs to match between the attP and attB site for integration to
318 occur.



319
 320 **Figure 2. Bxb1 domain engineering using directed evolution.** **a.** Schematic of our two-plasmid
 321 directed evolution system; pRex encodes an integrase gene while pRep contains an antibiotic
 322 marker disrupted by a stuffer sequence flanked by modified attB and attP sequences such that
 323 recombination by an active integrase excises the stuffer sequence and restores the open reading
 324 frame of the antibiotic resistance marker. This system enables either screening of libraries of
 325 integrase variants against a single DNA target, or screening of a library of DNA targets against a
 326 single integrase. **b.** Libraries where the loop, helix or hairpin submotif has been randomized are
 327 transformed along with pRep plasmids where the corresponding positions of the attB and attP
 328 target sites have been changed **c.** Specificity assay for individual selected integrase variants, where
 329 DNA target libraries are tested against a single integrase variant, each target site in the library
 330 contains the same modification at both half-sites of the attB and attP sites. **d.** Example DNA
 331 specificity plots of selected helix variants. **e.** Example DNA specificity plots of selected hairpin
 332 variants. **f.** Example DNA specificity plots of selected loop variants.

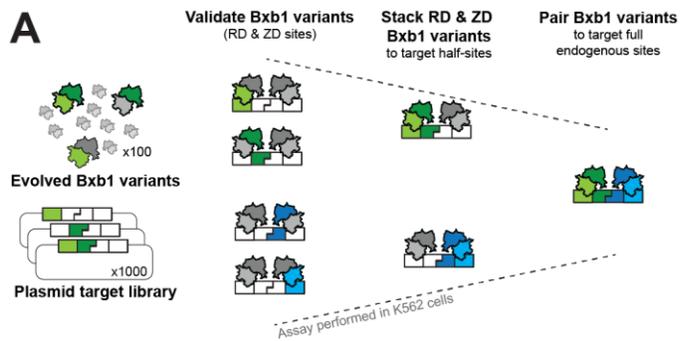


Figure 3. Performance of engineered Bxb1 variants at pseudo-sites in the human genome. a. Schematic of plasmid-based system for testing evolved Bxb1 variants against artificial DNA targets. Our LSR engineering strategy divides each endogenous target site into four “quarter-sites” where the left and right half-sites are each further divided into the portion recognized by the RD domain and the portion recognized by the ZD domain of Bxb1. Individual selected Bxb1 variants are screened against a library of plasmid targets that includes their intended quarter-sites. Successful RD and ZD variants are then combined and tested against the same library of plasmid targets that also includes the relevant half-sites. A single plasmid target library can contain quarter-site and half-site targets for numerous full endogenous target sites. Left and right site candidates derived from this assay can then be tested as pairs against the endogenous target site. See Extended Data Fig. 7b,c for additional details. **b.** Sequence of two Bxb1 pseudo-sites in the human genome. Both sites were identified experimentally using wild-type Bxb1. **c.** Bxb1 peptide sequences of evolved Bxb1 variants that showed improved performance against their corresponding half-sites. **d.** Results from a PCR-based assay demonstrating improved performance of evolved Bxb1 variants against their chromosomal endogenous

B

Pseudo attB site: MACO1
 5' -GCCCCTTCTCCTACAGAGCAAGCAGCAGGGTAAATTC
 CGGGGAAGAGGATGCTCTCGTTCGTCGCCATTAAAGA-5'

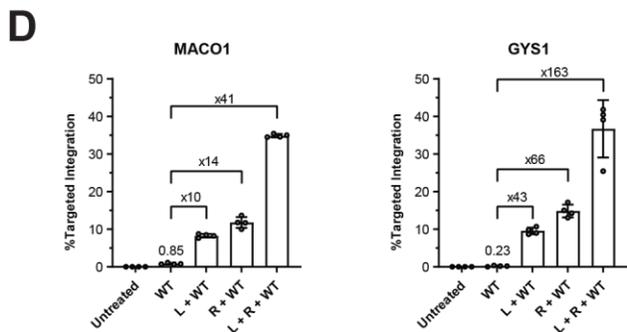
Pseudo attB site: GYS1
 5' -GGGATTCCCATAAACCGTGCACTCAGCTGCGGGAAGGCA
 CCCTAAGGGTATTGGCACGTGAGTCGACGCCCTTCCGT-5'

Hairpin		Loop Helix		Hairpin	
ZD	RD	RD	ZD	RD	ZD
left half-site			right half-site		

C

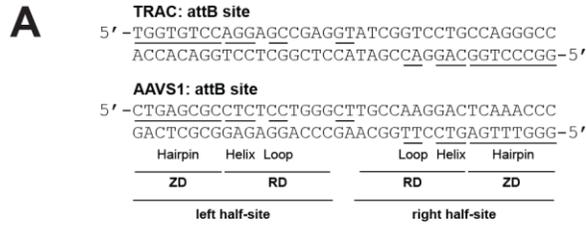
Bxb1 peptide sequences derived from directed evolution

	Helix	Hairpin
Wild-type Bxb1	SATALKR	FAGGGRKHPRYR
MACO1-L	SATALKR	MAGGHRKQALYR
MACO1-R	RAWSLKR	MAGGPRKKGRYR
GYS1-L	HGWSLKV	LARGSRKLALYR
GYS1-R	HGCTLKR	NARGNRKGRYR



366

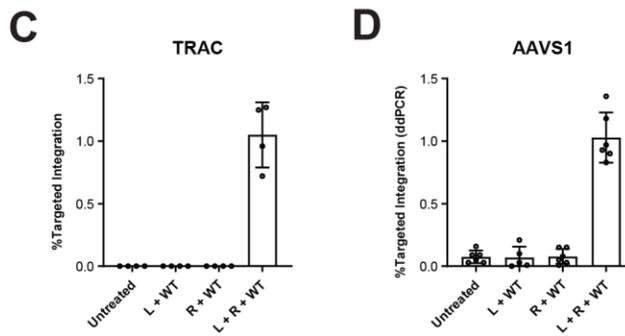
367 targets in K562 cells. The presence of a wild-type Bxb1 expression construct is necessary to bind
 368 the wild-type attP sequence on the donor plasmid.



B

	Bxb1 peptide sequence from directed evolution		
	Loop	Helix	Hairpin
Wild-type Bxb1	YRGSLP	SATALKR	FAGGGRKHPRYR
TRAC-L*	YRGGLP	YGSALKQ	LARGPRKRAGYK
TRAC-R*	YRGGLP	SQWALKC	RAWGKRKYAYYQ
AAVS1-L	YRGGLP	YPWSLRR	KAWGSRKTRLR
AAVS1-R	YRGGLP	AGGNLKR	MARGGRKSAIYY

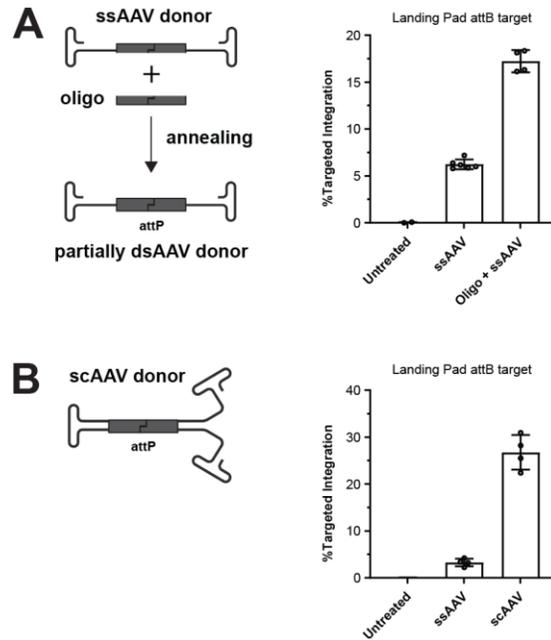
*TRAC-L and TRAC-R variants include an additional D257K mutation



369

370

371 **Figure 4. Retargeting Bxb1 to the human TRAC and AAVS1 loci. a.** Sequence of the pseudo-
 372 sites recognized by the TRAC and AAVS1 Bxb1 variants shown in panel b. **b.** Bxb1 peptide
 373 sequences of evolved Bxb1 variants for each of the two TRAC and AAVS1 half-sites generated by
 374 the strategy shown in **Figure 3a.** **c.** Results from our standard PCR-based NGS assay
 375 demonstrating targeted integration at the depicted TRAC site. See Extended Data Figure 8 for
 376 additional details. **d.** Results from a digital PCR-based assay demonstrating targeted integration at
 377 the depicted AAVS1 site. See Extended Data Fig. 9 for additional details.



378

379 **Figure 5. Alternative donor delivery strategies. a.** Utilization of a single-stranded AAV (ssAAV)
 380 donor for Bxb1-mediated targeted integration. We tested ssAAV as a donor against an attB landing
 381 pad in K562 cells and noticed measurable integration levels that can be increased through co-
 382 delivery of an oligonucleotide that is complementary to the attP sequence, therefore making the
 383 ssAAV donor partially double-stranded. **b.** Utilization of a self-complementary AAV (scAAV)
 384 donor for Bxb1-mediated targeted integration.

385 **Methods**

386 **Cloning of expression constructs and donors used in mammalian cells**

387 Most constructs were cloned using NEBuilder® HiFi DNA Assembly (NEB, Catalog
388 #E2621X) or Q5® Site-Directed Mutagenesis (NEB, Catalog #E0554S), or synthesized and cloned
389 by Twist Biosciences utilizing their clonal gene services. DNA sequences for all constructs can be
390 found in **supporting sequence information**. All constructs were sequence confirmed using Sanger
391 sequencing services provided by Elim Biopharm Inc., and whole plasmid sequences were verified
392 using either the Nextera XT DNA library prep kit (Illumina, Catalog #FC-131-1096) or whole
393 plasmid sequencing services (Plasmidsaurus Inc.; Elim Biopharma Inc.).

394 **High throughput Bxb1 variants assembly**

395 The Bxb1 variant gene fragments (bases 460 – 1087, corresponding amino acids 144 – 362
396 with loop, helix and hairpin regions included) were synthesized as eBlocks (Integrated DNA
397 Technologies) and assembled to full length fusion expression cassette through 2-step PCR. In the
398 first step PCR, the eBlocks were amplified with overlapping gene fragments DF148 and DF164
399 with AccuPrime *Pfx* SuperMix (Invitrogen, Catalog #12344040) and the following thermocycler
400 conditions: initial melt of 95 °C for 3 min; 15 cycles of 95 °C for 30 s, 68 ° for 30 s, and 68 °C for
401 2 min 30 s; followed by a final extension at 68 °C for 5 min; hold at 4 °C. The full-length expression
402 cassette PCR products were then amplified with forward primer 5'-
403 GCAGAGCTCTCTGGCTAACTAGAG-3' and reverse primer 5'-
404 TTTTTTTTTTTTTTTTTTTTTTTTTTTTTTTTTTTTGGAGGGGTCACAGGGATGCCACCCGTAG
405 ATC-3' and the following thermocycler conditions: initial melt of 95 °C for 3 min; 30 cycles of 95
406 °C for 30 s, 68 ° for 30 s, and 68 °C for 2 min 40 s; followed by a final extension at 68 °C for 5
407 min; hold at 4 °C. The final PCR products were examined on a gel and purified with AMPure XP
408 beads (Bechman# A63881). The purified PCR products were used for mRNA production.

	sequence (5'->3')
DF148	gcagagctctctggctaactagagaaccactgcttactggcttatcgaaattaatacactcactatagggagagccaa gctgactagcgtttaaactaagctgatccactagtccagtggtggaattgccaccatgcccgcctggtgttataag gcttagtcgagtgacggatgccacgactagccctgaacgccaaactggaaagtgccagcagctctgcgcacagcgcg ggtgggatgtggtgggtgtggtgaagactggacgtatcaggcgcagtgaccctttgacaggaaacggcgccca acctcgtcgtggtcgtcttcgaggaacagccggttgacgtgattgtagcgtaccgggtggatcactacaagatcaa ttcggcactgacagcagctgttcattgggccgaggaccataagaagtgggtgtctcagcgacggaagcacatttgac

	accacaacaccattcgctgcagtagtaatcgcgctgatggggactgtcgcaaatggagttggaagctatcaaagaac ggaatcgcagtgacgcacttcaatatcagggctgggaag
DF164	cgccttctgtgaagagcaagttctcgatttctcggatgatgcagaacggctgaaaaggtttgggtggccggttctgacag tgctgtggagctggcggaggtaatgcagagctttagatctcacgtctctgatcggtagtccagcgtaccgcgcaggat ctcccaacgagaagctcttgacgctagaatagcagcttggcggccagacaggaagaattggagggcttggagcac ggcctagcgggtgggagtgccgggaaacaggtcaacggtttgggtgattgggtggcgcgaacaagacaccgcagcaaa aaatacgtggctcagaagcatgaacgtacgccttactttgacgtgcgagggcactgacacgcaccatagattcggag atctcaagagtacgagcagcacttgcggttgggatccgtgtgagcgcctgcatacaggtatgtccggtagcggttctg gctcccatcatcaccatcacggctctggcccaagaagaagaggatctgataactcgagctctagaatcaac ctctggattacaaaattgtgaaagattgactggtattcttaactatgttgctccttttacgctatgtggatacgtgcttaatgc ctttgatcatgctattgcttccgtagtggcttctcttctcctctgtataaatcctggttctgtctctttatgaggagttgtgg cccgtgtcaggaacgtggcgtggtgtgactgtgttctgacgcaacccccactggttggggcattgccaccactg tcagctccttccgggacttctgcttccccctccctattgccacggcggaaactcatcgccgctgccttggccgctgctgg acaggggctcggctgttgggactgacaattccgtggtgtgctgggggaagctgacgtccttccatggctgctgcctgt ggtgccacctggattctgcgcccacgtccttctgctacgtcccttcggccctcaatccagcggaccttctcccgcggc ctgctgccggtctgcggcctctccgcgtcttcgcttcgcccctcagacgagtcggatctcccttggggccgctccccg cctgtctagatggcccgtttaaaccgctgatcagcctcactgtgcttctagtggcag

409

410 **mRNA production**

411 The mRNA was prepared with purified PCR products via the mMACHINE
412 T7 Ultra kit (ThermoFisher, Catalog #AM1345) following the manufacturer’s instructions.
413 Synthesized mRNA was purified with RNA cleanXP beads (Beckman, Catalog #A63987) and
414 quantified with Quant-IT kit (Invitrogen, Catalog #Q10213).

415 **AAV production**

416 AAV genome plasmids containing the relevant attP sequences for Bxb1-mediated
417 recombination, flanked by inverted repeat terminal sequences (ITR) for AAV packaging, were
418 cloned using NEBuilder® HiFi DNA Assembly (NEB, Catalog #E2621X). ITRs were mutated to
419 produce scAAV⁵⁶. Bovine growth hormone (bGH) polyA termination sequence was inserted
420 upstream of the 3’ ITR to aid transcription termination as well as AAV titration. All constructs were
421 sequence confirmed using Sanger sequencing services provided by Sequetech Corporation. The
422 whole plasmid sequences were verified using Nextera XT DNA library prep kit (Illumina, Catalog
423 #FC-131-1096) or whole plasmid sequencing services provided by Plasmidsaurus Inc. to ensure
424 ITRs are intact prior to AAV production.

425 The recombinant AAV vectors (rAAV) were produced by triple transfection of suspension
426 human embryonic kidney (HEK) cells in 400 ml flasks, then purified by cesium chloride density
427 gradient centrifugation and dialysis. The rAAV genome was packaged into AAV-DJ (synthetic
428 chimeric capsid of AAV 2/8/9) and titrated using bGH PolyA qPCR assay on the Quantstudio™ 3
429 Real-time PCR system.

430 **Directed evolution system**

431 **Preparation of Integrase Mutant libraries.** As a starting point for engineering of Bxb1
432 variants in *E. coli*, we cloned a codon optimized ORF into plasmid pRex, which contains a pBR322
433 origin of replication as well as the L-rhamnose inducible pRhaBAD promoter. Libraries were
434 prepared using inverse PCR and primers encoding degenerate bases using an NNK degeneracy
435 scheme. In order to diversify the “loop” submotif, primers were designed to target residues 154-
436 159. In order to diversify the “helix” submotif, primers were designed to target residues 231-234,
437 and residues 236-237. In order to diversify the “Beta-hairpin” submotif, primers were designed to
438 target residues 314, 316, 318, 321, 323, and 325 using an NNK randomization scheme, while
439 residue 322 was randomized using an SSK randomization scheme. Briefly, iPCR reactions were
440 set up in 50 uL volumes, using 1x KOD ONE Master mix, 0.5uM forward primer, 0.5 uM reverse
441 primer, and 10 ng of template plasmid DNA. Following purification of the PCR products via silica
442 column (Qiagen PCR cleanup kit), and elution in 50 uL of buffer EB, compatible overhangs for
443 ligation were generated by digestion of PCR products using 5 units each of BsaI (NEB, R3733S)
444 and DpnI (NEB, R0176S) in 60 uL of 1x CutSmart buffer. The resulting digest products were then
445 ligated at 10 ng/uL in 1x T4 Ligase buffer, at 8 U/uL T4 Ligase, for 1 hour at room temperature.
446 Following ligation, DNA products were purified using the Qiagen PCR cleanup kit and eluted in
447 50 uL of buffer EB. The resulting ligated products were used to transform 50 uL of
448 electrocompetent cells (NEB 10Beta, C3020K; or ThermoFisher OneShot Top10, C404052) using
449 a BMX electroporator and a 96-well cuvette (BTX, 45-0450-M) using the manufacturer’s protocol
450 and allowed to recover for 1 hour at 37 °C in 1ml total volume of SOC media. The resulting
451 transformations yielded libraries in the order of 4e8 to 5e9 CFUs. After recovery, cells were
452 transferred to 1 L of LB media containing 34 ug/mL chloramphenicol and grown overnight with
453 shaking at 37 °C. The resulting culture was harvested, and plasmids were purified using a Qiagen
454 Plasmid Plus Giga Kit.

455 **Preparation of selection cassette plasmids.** Selection plasmid pRep was based on a p15A
456 origin of replication and contained a Spectinomycin resistance cassette. The selection cassette
457 contained a Tet promoter upstream of an open reading frame encoding a gentamicin resistance
458 gene (aacC1) disrupted by a “stuffer sequence”. The stuffer sequence was flanked with modified
459 attB sequence at its 5’ end, and with a modified attP sequence at its 3’ end. The attB and attP
460 sequences were modified such that recombination by an active integrase would leave behind an
461 attL sequence encoding an in-frame peptide insertion, giving rise to an active aacC1. Each
462 selection cassette was generated by the cloning of two ~1050 bp DNA fragments into a linearized
463 pRep backbone plasmid using the NEB HiFi assembly master mix, according to the manufacturer’s
464 instructions. Sequence-confirmed plasmids were either prepared by Mini- or Midi- Kits (Qiagen)
465 according to the manufacturer’s instructions.

466 **Selection of Bxb1 variants by aacC1 reassembly.** Typical selection reactions used 1ug of
467 Bxb1 library, and 1ug selection cassette plasmid. OneShot Top10 electrocompetent cells were
468 mixed with library and selection plasmids before being transferred to a 96-well cuvette, and
469 electroporated as above. After recovery, cells were transferred to 4 mL of LB containing 34 ug/mL
470 chloramphenicol, 50 ug/mL spectinomycin and 0.2% L-rhamnose, then the cultures were shaken
471 at 37 °C for 16 h to allow recombination to take place. After recombination, cultures were
472 harvested, media removed, and the cells were resuspended and transferred to 50 mL Super Broth
473 containing 34 ug/mL chloramphenicol and 20 ug/mL gentamicin. After overnight culture, 1 mL of
474 each culture was used to extract plasmid DNA using the Qiagen Miniprep kit. For selections using
475 agar plates, cells were resuspended in 1 mL of Superbroth and plated on 25x25 cm LB agar plates
476 containing 34 ug/mL chloramphenicol and 20 ug/mL gentamicin. After overnight incubation at 37
477 °C, colonies were scraped from plates, resuspended in 10 mL LB broth, and a 500 uL volume was
478 removed and used in plasmid purification using Qiagen Miniprep kits.

479 **NGS sequencing of Bxb1 variant submotifs post-selection.** Samples were prepared for
480 sequencing using Illumina MiSeq or NextSeq using an integrase-adaptor hybrid primer pair
481 (PCR1) followed by an Illumina-adaptor specific primer pair (PCR2). Briefly, 50uL reactions were
482 prepared using 1x Q5 Hotstart Master mix (NEB M0494S), 0.5uM Forward primer, 0.5uM Reverse
483 primer, and 50ng of plasmid DNA. PCR was performed using the following cycling conditions,
484 98 °C 30 s, then 15 cycles of 98 °C 5 s, 68 °C 7 s, 72 °C 10 s, followed by a final extension step

485 at 72 °C for 30s. The resulting PCR product was used as template for PCR2. Briefly, 50 uL
486 reactions were prepared using 1x KOD ONE polymerase master mix (Millipore-Sigma), 0.5 uM
487 Adapter primer 1, 0.5 uM Adapter primer 2, and 1 uL of PCR1 as template. PCR was performed
488 using the following cycling conditions 98 °C 30 s, then 13 cycles of 98 °C 5s, 60 °C 7 s, 72 °C 5
489 s followed by a final extension step at 72 °C for 30 s. The resulting PCR products were column
490 purified using a Qiagen PCR cleanup kit, and samples were sequenced using standard Illumina
491 Kits for MiSeq or NextSeq.

492 **Molecular specificity assays.** We built a plasmid reporter to assay the molecular specificity
493 of novel integrase submotifs. Three versions of the pRex plasmids were used as starting point to
494 avoid background activity of the wild-type integrase during cloning. By placing an extra adenine
495 residue within the loop, helix, or hairpin submotif, we created pRex variants where Bxb1 was
496 inactivated by a frameshift-inducing mutation but could be rescued by subsequent mutagenesis.
497 Recombination cassettes from the pRep selection plasmids were amplified by PCR and cloned into
498 the pRex plasmids by Gibson assembly upstream of the Bxb1 ORF. To assay the specificity of
499 selected loops, a single reverse primer and different forward primers which encoded the new loop
500 residues were used in iPCR using the small library of loop target plasmids (16 targets). A similar
501 procedure was used to generate helix mutants using the small library of helix target plasmids (64
502 targets).

503 To generate a library of ZD hairpin targets, we amplified and cloned an oligo pool (Twist
504 Biosciences) which comprised all possible attB and attP sites where positions -18 to -13 have been
505 randomized into our pRex plasmid, downstream of the integrase ORF. Such a library ensures that
506 the ZD hairpin targets on attBL and attPL are identical, and the targets on attBR and attPR are
507 inverted repeats of the left side targets. Mutant hairpins are generated by two different iPCR
508 primers, which together encode the novel hairpin sequence.

509 iPCR reactions were carried out in 50 uL volumes, with 1 x KOD ONE master mix, 0.5uM
510 forward primer, 0.5uM reverse primer, and 20 ng of a plasmid pool containing the relevant
511 recombination cassettes. Thermocycling was carried out using an initial denaturation step of 98 °C
512 for 30s, followed by 35 cycles of 98 °C 10s, 60 °C 5s, 68 °C 30s. PCR amplicons were purified
513 using AMPure XP beads (Beckman) and eluted in 40 uL EB buffer (Qiagen). 30 uL of the purified
514 PCR amplicons were then digested and ligated in a 50 uL one-pot reaction containing 1x T4 ligase

515 buffer (NEB), 20U DpnI (NEB), 20 U BsaI-HFv2 (NEB), and 400 U of Salt-T4 DNA Ligase
516 (NEB). All reactions were then incubated at 37 °C 30 minutes, 20 °C 30 minutes, 37 °C for 30
517 minutes. 20 uL of thawed chemically competent NEB 5alpha cells were added 2 uL of the
518 digested/ligated amplicons, and cells were transformed according to the manufacturer's
519 instructions. After recovery, cells were added to 800 uL of media containing 34ug/mL
520 chloramphenicol and 0.2% L-rhamnose (w/v). The resulting culture was incubated for 16 hours at
521 37 °C with shaking in 96-well deep-well plates, then harvested by centrifugation. Plasmid DNA
522 was extracted using a Qiaprep 96 Turbo Miniprep kit.

523 To assess the specificity of each clone, primers which flank the resulting attL sequence and
524 contain Illumina adaptor sequences were used to amplify the products of recombination. Briefly,
525 PCR reactions were carried out in 20uL volumes, using 1 x Hotstart Taq master mix (NEB,
526 M0496S), attB_MiSeq forward primer, and attP_Miseq reverse primer, using 1uL of purified
527 plasmid as template. Thermocycling was carried out as follows 98 °C 30s; 25 cycles of [98 °C 5s,
528 53 °C 10s, 72 °C 5s] Final extension 72 °C 30seconds. A second PCR to install sequencing
529 barcodes was carried out using specific barcoding primers, with the recombination products from
530 each clone being represented by a unique combination of forward and reverse barcoding primers.
531 PCR reactions were carried out in 20 uL volumes, using 1 x KOD ONE master mix, 0.5 uM
532 forward primer, 0.5 uM reverse primer, and 1 uL of the previous PCR product. Thermocycling was
533 carried out as follows 98C 30s; 12 cycles of [98 °C 5s, 60 °C 5s, 68 °C 3s]. The resulting PCR
534 products were column purified using Qiagen PCR cleanup kits, then sequenced using standard
535 Illumina kits for NextSeq or MiSeq.

536 **General mammalian cell culture condition**

537 K562 cells (ATCC, CCL243) were cultured using RPMI-1640 growth medium
538 supplemented with 10% FBS (Fetal Bovine Serum) and 1x PSG (Penicillin-Streptomycin-
539 Gentamycin, Gibco, 10378-016) and maintained at 37 °C with 5% CO₂.

540 **K562 tissue culture nucleofection protocol and genomic DNA preparation**

541 Expression constructs were routinely dosed as plasmid DNA (pDNA) in K562 cells. K562
542 cells were electroporated with pDNA using the SF cell line 96-well Nucleofector kit (Lonza,
543 Catalog#V4SC-2960) or SF Cell Line 384-well Nucleofector Kit (Lonza, Catalog#V5SC-2010),

544 using manufacturer's protocol. Prior to electroporation, K562 cells were centrifuged at ~300 x g
545 for 5 min, and washed with 1x DPBS (Corning, Catalog#21-031-CV). For 96-well nucleofection,
546 cells were resuspended at 2e5 cells per 12 µl of supplemented SF cell line 96-well Nucleofector
547 solution. 12 µl of cells were mixed with 8 µl of pDNA and transferred to the Lonza Nucleocuvette
548 plate. Nucleofector program 96-FF-120 was used to electroporate K562 cells with the pDNA mix
549 on the Amaxa Nucleofector 96-well Shuttle System (Lonza). After electroporation, cells were
550 incubated for 10 min at room temperature and transferred to a 96-well tissue culture plate
551 containing 180 µl of complete medium (prewarmed to 37 °C).

552 For 384-well nucleofection, cells were resuspended at 1e5 cells per 14 µl of supplemented
553 SF cell line 384-well Nucleofector solution. 14 µl of cells were mixed with 6 µl of pDNA and
554 transferred to the Lonza Nucleocuvette plate. Nucleofector program FF/120/DA was used to
555 electroporate K562 cells with the pDNA mix on the Amaxa HT Nucleofector System (Lonza,
556 AAU-1001). After electroporation, cells were incubated for 10 min at room temperature and
557 transferred to a 384-well tissue culture plate containing 60 µl of complete medium (prewarmed to
558 37 °C). K562 cells were incubated for ~72 h and then harvested for quantification of editing events.

559 **Digital Droplet PCR quantification of targeted integration at AAVS1 target site**

560 For AAVS1 target integration ddPCR quantification, 2 targeted integration probes
561 (targeting attL and attR) and 2 reference probes were designed. The target integration probes were
562 designed targeting either attL or attR region which would not be detected in non-transfected cells
563 or plasmid donors. The reference probes were designed within a 10 kb region of the AAVS1 attB
564 target site to mitigate the risk of copy number difference due to abnormal karyotype in K562 cells.
565 The probes and primers were designed with PrimerPlus3 (<https://www.primer3plus.com/>) based
566 on the recommendation by BioRad manual and were synthesized at IDT (Integrated DNA
567 Technologies, Inc.). The primers and probes were tested in duplex format at different annealing
568 temperatures and with synthetic eblock mixes, non-transfected cell lysate and NTC. The attL and
569 reference2 probes and primers and 57.1°C were selected for AAVS1 target site ddPCR
570 quantification.

571 attL probe: /56-FAM/CGCCTCTCC/ZEN/TGGGCTCTCAGTGGTGTACG/3IABkFQ/

572 attL-for: gcatgagatggtgacgag

573 attL-rev: ggccggtgacatattcctc

574 Reference2 probe: /5HEX/CGGATCCCG/ZEN/CGCCCAACTCAAGATTGG/3IABkFQ/

575 Ref-for: agcacaccttgatcttcacc

576 Ref-rev: agtctctgtcccgattttgg

577 DNA was extracted with QuickExtract DNA Extraction Solution (Lucigen#QE09050). For each
578 reaction, 50 µl of QuickExtract DNA solution was added to approximately 0.5-1 million pelleted
579 cells, followed by mixing and incubation at 65 °C for 15 min and heat inactivation at 98 °C for
580 5 min. The cell lysates were mixed by vortexing for 15 seconds before ddPCR. Each ddPCR
581 reaction was prepared and analyzed with a QX200 ddPCR system (Bio-Rad) and ddPCR Supermix
582 for Probes (No dUTP) (Bio-Rad, Catalog #1863024) per Bio-Rad's standard recommendations.
583 All reactions were mixed to 22 µl including 10 U of HindIII-HF (NEB, Catalog #R3104L) and up
584 to 2 µl of QuickExtract lysates. Forward primer, reverse primer and probe were at a 3.6:3.6:1 ratio.
585 Droplets were generated in the droplet generator per Bio-Rad's protocol. Thermocycler conditions:
586 95 °C for 10 min; 40 cycles of 95 °C for 30 s and 57.1 °C for 60 s; 98 °C for 10 min; and hold at
587 8 °C. QX Manager Software 2.1 Standard Edition (Bio-Rad) was used for QC and the analysis.
588 The thresholds were set manually at 3000 for channel1/FAM and 1000 for channel2/HEX. All final
589 data was exported into Microsoft Excel for further analysis. The target integration ratio was
590 calculated by the equation: Targeted Integration (%) = 100*CattL/Cref2 (C: volumetric
591 concentration (copies/µl)).

592 **AAV transduction**

593 For evaluating rAAVs. K562 cells were electroporated with 800ng of Bxb1 pDNA using
594 the SF cell line 96-well Nucleofector kit (Lonza, Catalog#V4SC-2960), using manufacturer's
595 protocol. K562 cells were centrifuged at ~300 x g for 5 min, and washed with 1X PBS (Corning,
596 Catalog#21-031-CV). For 96-well nucleofection, cells were resuspended at 2E5 cells per 12 µl of
597 supplemented SF cell line 96-well Nucleofector solution. 12 µl of cells were mixed with 8 µl of
598 Bxb1 pDNA and transferred to the Lonza Nucleocuvette plate. Nucleofector program 96-FF-120
599 was used to electroporate K562 cells with the pDNA mix on the Amaxa Nucleofector 96-well
600 Shuttle System (Lonza). After electroporation, cells were incubated for 10 min at room temperature

601 and transferred to a 96-well tissue culture plate containing 180 μ l of complete medium (prewarmed
602 to 37 °C). 30mins post-electroporation, rAAV constructs were dosed at MOI (multiplicity of
603 infections) at 500,000 vg/cell. rAAV donor only control wells were included in parallel. K562 cells
604 were incubated for ~72 h and then harvested for quantification of editing and circularization
605 events. A PCR-based NGS assay was used to measure targeted integration events.

606 **K562 landing pad cell line generation**

607 K562 cells were electroporated with a pair of zinc-finger nuclease (ZFN) mRNA and an
608 ultramer with the attB sequence using the SF cell line 96-well Nucleofector kit (Lonza,
609 Catalog#V4SC-2960), using manufacturer's protocol. The ZFNs generate a double-strand break
610 in the genome that facilitates integration of the ultramer via homology directed repair through the
611 corresponding homologous ends. K562 cells were incubated for ~72 h at 37 °C with 5% CO₂. One
612 third of the cells were harvested for quantification of bulk integration of the ultramer using a PCR-
613 based NGS assay. 2/3rd of the cells were maintained for diluting to singles. Samples showing ~10%
614 integration were selected and diluted to singles in a 96-well plate and incubated for ~1.5 weeks.
615 At the end of ~1.5 weeks, the plates were examined under a microscope for the growth of single
616 clones. Cells from wells showing single clones were transferred to a 24-well tissue culture plate
617 and transferred to 37 °C with 5% CO₂ for 3 days or until 75-80% confluency is reached. 50% of
618 the cells were harvested for quantification of ultramer integration using the PCR-based NGS assay.
619 Cells were maintained in fresh medium until the NGS assay was completed. Since AAVS1 has 3
620 alleles, a sample showing 34.39% integration at a single allele was selected for further expansion.
621 The other two wild-type alleles showed a 6bp deletion, but this did not disrupt the performance of
622 the cell line.

623 Ultramer: 5'-AGGAGACTAGGAAGGAGGAGGCCTAAGGATGGGGCTTTTCGGCC
624 GGCTTGTCGACGACGGCGGTCTCCGTCGTCAGGATCATCCGGCAGATAAAAGTACCC
625 AGAACCAGAGCCACATTAACCGGCC-3'

626 **PCR-based NGS assay for targeted integration and indel quantification**

627 3 days post transfection, cells were spun down at ~500 x g for 5 min. Supernatant was
628 discarded, cells were washed in PBS (Corning, Catalog #21-031-CV), and cells were resuspended
629 in 50 μ l of QuickExtract DNA Extraction Solution (Lucigen, Catalog #QE09050). Genomic DNA

630 was extracted by treating the cells to the following protocol: 65 °C for 15 min, 98 °C for 8 min.
631 Target sites were amplified from the genomic DNA using Accuprime HiFi reagents (Invitrogen,
632 Catalog #12346094) and the following PCR conditions: initial melt of 95 °C for 5 min; 30 cycles
633 of 95 °C for 30 s, 55 °C for 30 s and 68 °C for 40 s; and a final extension at 68 °C for 10 min.
634 Primers containing adapters (forward primer adapter: ACACGACGCTCTTCCGATCT; reverse
635 primer adapter: GACGTGTGCTCTTCCGATCT), targeting specific target sites were used at a
636 final concentration of 0.1 μM. Sequences for the primers used can be found in **supporting**
637 **sequence information**. The PCR productions obtained were then subjected to a second PCR to
638 add Illumina barcodes to the PCR fragments generated in the first PCR. We used Phusion High-
639 Fidelity PCR MasterMix with HF Buffer (NEB, Catalog #M0531L) for the second PCR and used
640 the following PCR conditions, initial melt of 98 °C for 30 s; 12 cycles of 98 °C for 10 s, 60 °C for
641 30 s and 72 °C for 40 s; and a final extension at 72 °C for 10 min. PCR libraries generated from
642 the second PCR were pooled and purified using QIAquick PCR purification kit (Qiagen, Catalog
643 #28106). Samples were diluted to a final concentration of ~2 nM after they were quantified using
644 the Qubit dsDNA HS Assay kit (Invitrogen, Catalog #Q33231). The libraries were then run on
645 either an Illumina MiSeq using a standard 300-cycle kit or an Illumina NextSeq 500 or an Illumina
646 NextSeq 2000 using a mid-output 300-cycle kit using standard protocol.

647 **Pooled screening of attB sites in K562 cells**

648 Extended Data Fig. 7a outlines the experimental screening of Bxb1 variants against
649 plasmid libraries of artificial target sequences in human K562 cells. Target libraries were designed
650 as shown in Extended Data Fig. 7b,c and cloned using oligo pools (Twist Bioscience). The
651 resulting target libraries were then co-transfected with a universal donor and individual Bxb1
652 variants. Activity was measured using PCR-based NGS assay.

653 **Genome-wide specificity assay**

654 Our genome-wide specificity assay was adapted from the GUIDE-seq protocol⁵⁷, with
655 modifications designed to measure Bxb1-induced off-target editing events. K562 cells were
656 transfected using conditions similar as described above. Transfections involving integrations with
657 multiple donors with different core dinucleotides were performed separately per donor plasmid
658 and then cells were pooled and expanded for 1 week before being spun down for genomic DNA

659 extraction. Without the DpnI site addition to donor plasmids, the cells were grown out for 3-4
660 weeks prior to DNA extraction and the DpnI digestion step below was not followed.

661 Genomic DNA was extracted from K562 cells using Qiagen DNeasy Blood & Tissue kits
662 (Catalog #69504) following the Purification of Total DNA from Animal Blood or Cells Spin-
663 Column Protocol for cultured cells. The optional RNase A incubation was followed for all samples.
664 DNA was eluted in 60 μ L Elution Buffer and quantified using the Qubit fluorometer and the Qubit
665 dsDNA HS Assay Kit (Invitrogen, Catalog #Q33231) following the recommended protocol.

666 Next, 10 μ M adapters were prepared by annealing the MiSeq common oligo to each GS_i5
667 oligo in a 96-well plate format to make a barcoded Y adapter plate. Annealing was performed with
668 1X oligo annealing buffer (10 mM Tris HCL pH 7.5, 50 mM NaCl, and 0.1 mM EDTA) by
669 following the below thermocycling method: initial melt of 95 $^{\circ}$ C for 2 min, step-down from 80 $^{\circ}$ C
670 to 4 $^{\circ}$ C with -1 $^{\circ}$ C per cycle and 1 min incubation at each temperature, hold at 4 $^{\circ}$ C until further
671 use. Adapters were stored at -20 $^{\circ}$ C, and before use were thawed on ice. 400 ng (133,000 haploid
672 human genomes) genomic DNA was brought up to 50 μ L using 1X IDTE pH 7.5 (IDT, Catalog
673 #11-05-01-05) in each tube of a Covaris 8 microTUBE-130 AFA Fiber H Slit Strip V2 (Covaris,
674 Catalog #520239). Samples were sonicated on a Covaris ME220 using the following settings on a
675 ME220 Rack 8 AFA-TUBE TPX Strip (Covaris, Catalog #PN500609) using the waveguide
676 (Covaris PN 500526): Power 0.0 W, Temperature 19.7 $^{\circ}$ C, Duration(s) 65.0, Peak Power 40.0,
677 Duty %Factor 10.0, Cycles/Burst 10000, Avg. Power 4.0.

678 Sheared DNA was purified using 1 volume Ampure XP beads (Beckman Coulter, Catalog
679 #A63880). After beads were added, the solution was mixed and incubated for 5 minutes at room
680 temperature. The mixture was then incubated on a magnet for 5 minutes before the supernatant
681 was removed. 150 μ L freshly made 70% ethanol was then used to wash the beads twice, allowing
682 the solution and beads to sit for 30 seconds each time. After the second wash, the beads were dried
683 for 6 minutes before adding 15 μ L IDTE pH 7.5 and mixing off the magnet. After 2 minutes the
684 mixture was placed on a magnet and incubated for another 2 minutes. 14.5 μ L of the supernatant
685 was collected for the next step.

686 The reaction was brought up to 50 μ L with the addition of CutSmart (final concentration
687 1X) and 1 μ L DpnI (NEB, Catalog #R0176S) and incubated for 1 hour at 37 $^{\circ}$ C. DNA was purified
688 using 0.8x Ampure XP beads using the same bead clean-up protocol as before. Next, the following

689 End repair and A-tailing mixture was added to each 14.5 uL DNA mixture, while the reaction was
690 kept on ice: 0.5 uL 10 mM dNTP mix (Invitrogen, Catalog #18427013), 2.5 uL 10X T4 DNA
691 Ligase Buffer (Enzymatics, Catalog #B6030), 2 uL End-repair mix (Enzymatics, Catalog #Y9140-
692 LC-L), 2 uL 10X Platinum Taq Polymerase PCR Rxn Buffer (-Mg² free) (Invitrogen, Catalog
693 #10966034), 0.5 uL Taq DNA Polymerase Recombinant (5u/ uL) (Invitrogen, Catalog
694 #10342020), and 0.5 uL dsH₂O to a total of 22.5 uL. The solution was mixed and incubated on a
695 thermocycler with the following program: 12 °C for 15 min, 37 °C for 15 min, 72 °C for 15 min,
696 hold at 4 °C until further use.

697 Next, 2 uL T4 DNA Ligase (Enzymatics, Catalog #L6030-LC-L) and 1 uL of one of the 10
698 uM annealed Y adapters (chosen from a 96-well plate of GS_i5 adapters) was added per end-
699 repaired and A-tailed reaction. Each reaction was mixed and incubated on a thermocycler with the
700 following program: 16 °C for 30 min, 22 °C for 30 min, hold at 4 °C until further use. DNA was
701 then purified using 0.9 volumes Ampure XP beads using the same bead clean-up protocol but using
702 23 uL IDTE pH 7.5 to resuspend the DNA-bead mixture for elution and collecting 22 uL
703 supernatant after incubation on the magnet.

704 Next, ligated and sheared DNA fragments were amplified using a primer specific to all
705 adapters (P5_1) and a primer specific to the sequence of interest (GSP1 +/-). To each tube on ice,
706 the following reagents were added: 22 uL DNA from previous step, 3 uL 10X Platinum Taq
707 Polymerase PCR Rxn Buffer (-Mg² free) (Invitrogen, Catalog #10966034), 1.2 uL 50 mM MgCl₂
708 (Invitrogen, Catalog# 10966034), 0.6 uL 10 mM dNTP mix (Invitrogen, Catalog #18427013), 0.5
709 uL 10 μM P5_1 primer, 1 uL 10 μM GSP1+/-, 1.5 uL 0.5 M TMAC (Sigma Aldrich, Catalog
710 #T3411), and 0.3 uL Platinum Taq DNA polymerase (5 U/μL) (Invitrogen, Catalog #10966034) to
711 a total of 30.1 uL. Each reaction was mixed and incubated with the following thermocycler
712 conditions: initial melt of 95 °C for 2 min; 15 cycles of 95 °C for 30 s, 70 °C (-1 °C/cycle) for 2
713 min, and 72 °C for 30 s; followed by 10 cycles of 95 °C for 30 s, 55 °C for 1 min, and 72 °C for
714 30 s; followed by a final extension at 72 °C for 5 min; hold at 4 °C until further use.

715 Amplified DNA was purified using the bead clean-up protocol but with 1.2 volumes of
716 Ampure XP beads and using 21 uL IDTE pH 7.5 to resuspend the DNA-bead mixture for elution
717 and collecting 20.4 uL supernatant after incubation on the magnet.

718 Next, the second round of PCR amplification was performed to add plate barcodes (GS_i7
719 sequences). To each tube on ice, the following reagents were added: 1.5 ul 10 μM Plate adapter
720 (GS_i7), 20.4 ul DNA from previous step, 3 ul 10X Platinum Taq Polymerase PCR Rxn Buffer (-
721 Mg²⁺ free) (Invitrogen, Catalog #10966034), 1.2 ul 50 mM MgCl₂ (Invitrogen, Catalog
722 #10966034), 0.6 ul 10 mM dNTP mix (Invitrogen, Catalog #18427013), 0.5 ul 10 μM P5_2 primer,
723 1 ul 10 μM GSP2+/-, 1.5 ul 0.5 M TMAC (Sigma Aldrich, Catalog #T3411), and 0.3 ul Platinum
724 Taq DNA polymerase (5 U/μL) (Invitrogen, Catalog #10966034) to a total of 30 ul. Each reaction
725 was mixed and incubated with the following thermocycler conditions: initial melt of 95 °C for 5
726 min; 15 cycles of 95 °C for 30 s, 70 °C (-1 °C/cycle) for 2 min, and 72 °C for 30 s; followed by 10
727 cycles of 95 °C for 30 s, 55 °C for 1 min, and 72 °C for 30 s; followed by a final extension at 72
728 °C for 5 min; hold at 4 °C until further use.

729 Afterwards, 25 uL of each barcoded reaction was combined into one pool and 0.7 volumes
730 of Ampure XP beads were added. Samples were mixed 10 times and incubated for 5 minutes at
731 room temperature. They were then added to the magnet for 5 minutes then the supernatant was
732 discarded. 2 1x volumes of freshly made 70% ethanol were added, incubating for 30 seconds each
733 time before discarding the supernatant. After the last wash, the beads were air-dried for 6-8
734 minutes. 75 uL IDTE pH 7.5 was added and the tubes were removed from the magnet and mixed,
735 then incubated for 2 minutes. The reaction was separated on the magnet for 2-4 minutes and then
736 the supernatant was collected into a new tube for NGS library sample submission. Pooled eluate
737 was quantified using the Qubit and the Qubit dsDNA HS kit following the recommended protocol.

738 Final products were sequenced on a MiSeq or NextSeq 2000 with paired-end 150 bp reads
739 with the cycle settings: 148-10-22-148 for MiSeq or 151-10-22-151 for NextSeq. Samples were
740 sequenced to obtain at least 3,000,000-fold coverage per sample. For MiSeq reactions, 3 uL of 100
741 uM custom sequencing primer Index1 (5'-3':
742 ATCACCGACTGCCCATAGAGAGGACTCCAGTCAC) was added to MiSeq Reagent cartridge
743 position 13 and 3 uL of 100 uM custom sequencing primer Read2 (5'-3':
744 GTGACTGGAGTCCTCTCTATGGGCAGTCGGTGAT) was added to MiSeq Reagent cartridge
745 position 14.

746 For NextSeq reactions, 1.98 ul of 100 uM custom sequencing primer Read2 was added to
747 600 ul Illumina HP21 primer mix for a 0.3 uM final concentration and 3.98 ul of 100 uM custom

748 sequencing primer Index1 was added to 600 ul BP14 primer mix for a 0.6 uM final concentration.
749 Then 550 µl of custom primer mix was added to custom 1 well or custom 2 well on the reagent
750 cartridge separately.

751 **Computational analysis of directed evolution experiments**

752 Data processing was performed with custom python scripts and logo plots were generated
753 using Logomaker⁵⁸. Sequence reads that don't match sequences encoded by the directed evolution
754 library were filtered out as were sequence reads that were likely artifacts due to being rare sequence
755 reads with only one or two differences from much more common sequence reads. The remaining
756 sequence reads were translated into peptide sequence and the rest of the analysis only considered
757 the peptide sequence of the randomized region. For hairpin selections, the partially randomized
758 position 322 could be either a G, A, R, or P and the motif analysis was conducted separately based
759 on the identity of position 322. The clearest signal was 4 residue motifs within the 6 fully
760 randomized positions. For example, the helix library with a XXXXLXX randomization scheme
761 had a strongly enriched motif of XGGNLXR where X is any amino acid. To account for potential
762 biases in the starting library, motifs were scored based on the FDR corrected p-value of a given
763 motif occurring by chance given the amino acid frequency at each position when considered
764 independently. For helix selections, enriched 4-residue motifs that were observed in selections with
765 a wide variety of different DNA target sequences were considered to be non-specific and were
766 filtered out. Compatible four residue motifs (e.g. XGXNLKR and XXGNLKR) were then
767 combined to create motifs. Several peptides that best matched the most enriched 4 residue motifs
768 in each combined motif were then chosen for additional characterization.

769 **Computational analysis of pooled screening of attB sites in K562 cells**

770 A custom python script was used to count sequence tags corresponding to different
771 recombined targets. Counts were normalized to the total number of sequence reads obtained for a
772 given sample. Normalized sequence tag counts for a given Bxb1 variant + wild-type Bxb1 at a
773 given site were then compared to the normalized sequence tag counts for wild-type Bxb1 alone at
774 for the same sequence tag.

775 **Computational analysis of chromosomal targeted integration events**

776 Sequence data was first processed using our indel analysis software pipeline. The output
777 from this analysis was further processed using a custom python script that identified aligned
778 sequence reads that contained the sequence from the right attP half-site from the donor. Sequence
779 reads that contained this integration tag, but were not the expected length were scored as “TI +
780 indel”. Sequence reads that contained this integration tag and were the expected length were scored
781 as “perfect TI”. Sequence reads that did not contain the integration sequence tag were scored as
782 either wild-type amplicon or non-TI indel based on the output of the indel analysis software.

783 **Computational identification of Bxb1 pseudo-sites**

784 Raw data from Bessen et al. was processed to produce a position weight matrix for the
785 Bxb1 attB and attP target sites. A custom python script scanned the human genome for potential
786 target sites that matched the strongly preferred G nucleotide at position -4 and at position +4.
787 Sequences that met these criteria were then scored against the position weight matrix. The left and
788 right half-sites of the natural Bxb1 attB site are quite different from each other and likely represent
789 different binding modes of the Bxb1 hairpin region. Thus, for attB sites, each potential site was
790 scored against a position weight matrix representing an inverted repeat of the left half-site of the
791 natural attB sequence, and inverted repeat of the right half-site of the attB sequence, the composite
792 attB left and right half-sites on the top strand of DNA, or the composite attB left and right half-
793 sites on the bottom strand of DNA. The top 12 scores for each category were experimentally
794 characterized for a total of 48 potential attB pseudo-sites in the human genome. The natural attB
795 site is much more symmetric so the position weight matrix for the left and right half-sites was
796 averaged together and this averaged position weight matrix was used to score both sites of potential
797 attP pseudo-sites in the human genome. The top 48 scoring potential attP pseudo-sites were also
798 characterized experimentally.

799 **Computational analysis of genome-wide specificity assay**

800 NGS reads were demultiplexed, adapter trimmed, and filtered for a minimum quality
801 threshold of 14 over all bases. Samples then underwent analysis for plasmid integration site
802 detection. NGS samples that were analyzed for plasmid integration site detection were processed
803 to remove remaining contaminant unintegrated plasmid reads due to incomplete DpnI digestion or
804 fragment removal, aligned to the hg38 genome, and potential integration sites were summarized.
805 First, reads that contained both attP sequence 5' and 3' of the dinucleotide were removed from

806 analysis, corresponding to unintegrated donor plasmid reads. Then all sequence up to the start of
807 the dinucleotide (up to and including the 5' attP sequence) was removed, leaving the remaining
808 sequence to align to the hg38 genome using Bowtie2. Alignments with a MAPQ less than 23 were
809 removed from the analysis. Next, common read1 start locations, which correspond to unique
810 genomic shear locations and ligation events generated in the protocol, were used to deduplicate
811 common reads. Unique read1 start and dinucleotide positions were then summed per dinucleotide
812 position to generate a list of deduplicated reads per potential integration event. Next, all reads per
813 AMP-seq reaction were summed per alignment position in the genome and per alignment
814 orientation (top and bottom strands). Then, positions were combined into a single potential
815 integration location per AMP-seq reaction and alignment orientation if they fell within a 50 bp
816 window of one another, all reads per this grouping were summed and the coordinate with the most
817 reads was kept per group. Lastly, potential integration locations across top and bottom strand
818 alignments and between both AMP-seq "plus" and "minus" reactions per original transfected
819 sample were combined. Common potential integration locations were merged into one potential
820 integration location if within 50 bp of one another (this would encompass alignments separated by
821 a dinucleotide as a result of sequencing upstream and downstream of integration in both AMP-seq
822 reactions). The final list of potential integration loci were inspected for expected integration
823 genotypes (2 merged locations, in opposite alignment orientation, separated by a dinucleotide that
824 corresponds to the donor plasmid dinucleotide used in the assay).

825 References

- 826 1 Anzalone, A. V. *et al.* Search-and-replace genome editing without double-strand breaks or donor
827 DNA. *Nature* **576**, 149-157 (2019). <https://doi.org:10.1038/s41586-019-1711-4>
- 828 2 Gaudelli, N. M. *et al.* Programmable base editing of A*T to G*C in genomic DNA without DNA
829 cleavage. *Nature* **551**, 464-471 (2017). <https://doi.org:10.1038/nature24644>
- 830 3 Komor, A. C., Kim, Y. B., Packer, M. S., Zuris, J. A. & Liu, D. R. Programmable editing of a target
831 base in genomic DNA without double-stranded DNA cleavage. *Nature* **533**, 420-424 (2016).
832 <https://doi.org:10.1038/nature17946>
- 833 4 Stark, W. M. Making serine integrases work for us. *Curr Opin Microbiol* **38**, 130-136 (2017).
834 <https://doi.org:10.1016/j.mib.2017.04.006>
- 835 5 Forum: CRISPR roundtable with Doudna and Liu. *Nat Biotechnol* **38**, 943 (2020).
836 <https://doi.org:10.1038/s41587-020-0619-8>
- 837 6 Urnov, F. D. *et al.* Highly efficient endogenous human gene correction using designed zinc-finger
838 nucleases. *Nature* **435**, 646-651 (2005). <https://doi.org:10.1038/nature03556>
- 839 7 Perez, C. *et al.* Factors affecting double-strand break-induced homologous recombination in
840 mammalian cells. *Biotechniques* **39**, 109-115 (2005). <https://doi.org:10.2144/05391GT01>
- 841 8 Kosicki, M., Tomberg, K. & Bradley, A. Repair of double-strand breaks induced by CRISPR-Cas9
842 leads to large deletions and complex rearrangements. *Nat Biotechnol* **36**, 765-771 (2018).
843 <https://doi.org:10.1038/nbt.4192>
- 844 9 Haapaniemi, E., Botla, S., Persson, J., Schmierer, B. & Taipale, J. CRISPR-Cas9 genome editing
845 induces a p53-mediated DNA damage response. *Nat Med* **24**, 927-930 (2018).
846 <https://doi.org:10.1038/s41591-018-0049-z>
- 847 10 Ichikawa, D. M. *et al.* A universal deep-learning model for zinc finger design enables transcription
848 factor reprogramming. *Nat Biotechnol* **41**, 1117-1129 (2023). [https://doi.org:10.1038/s41587-022-](https://doi.org:10.1038/s41587-022-01624-4)
849 [01624-4](https://doi.org:10.1038/s41587-022-01624-4)
- 850 11 Miller, J. C. *et al.* Enhancing gene editing specificity by attenuating DNA cleavage kinetics. *Nat*
851 *Biotechnol* **37**, 945-952 (2019). <https://doi.org:10.1038/s41587-019-0186-z>
- 852 12 Paschon, D. E. *et al.* Diversifying the structure of zinc finger nucleases for high-precision genome
853 editing. *Nat Commun* **10**, 1133 (2019). <https://doi.org:10.1038/s41467-019-08867-x>
- 854 13 Boch, J. *et al.* Breaking the code of DNA binding specificity of TAL-type III effectors. *Science* **326**,
855 1509-1512 (2009). <https://doi.org:10.1126/science.1178811>
- 856 14 Miller, J. C. *et al.* A TALE nuclease architecture for efficient genome editing. *Nat Biotechnol* **29**,
857 143-148 (2011). <https://doi.org:10.1038/nbt.1755>
- 858 15 Jinek, M. *et al.* A programmable dual-RNA-guided DNA endonuclease in adaptive bacterial
859 immunity. *Science* **337**, 816-821 (2012). <https://doi.org:10.1126/science.1225829>
- 860 16 Akopian, A., He, J., Boocock, M. R. & Stark, W. M. Chimeric recombinases with designed DNA
861 sequence recognition. *Proc Natl Acad Sci U S A* **100**, 8688-8691 (2003).
862 <https://doi.org:10.1073/pnas.1533177100>
- 863 17 Chaikind, B., Bessen, J. L., Thompson, D. B., Hu, J. H. & Liu, D. R. A programmable Cas9-serine
864 recombinase fusion protein that operates on DNA sequences in mammalian cells. *Nucleic Acids*
865 *Res* **44**, 9758-9770 (2016). <https://doi.org:10.1093/nar/gkw707>
- 866 18 Gaj, T., Mercer, A. C., Sirk, S. J., Smith, H. L. & Barbas, C. F., 3rd. A comprehensive approach to
867 zinc-finger recombinase customization enables genomic targeting in human cells. *Nucleic Acids*
868 *Res* **41**, 3937-3946 (2013). <https://doi.org:10.1093/nar/gkt071>
- 869 19 Gordley, R. M., Gersbach, C. A. & Barbas, C. F., 3rd. Synthesis of programmable integrases. *Proc*
870 *Natl Acad Sci U S A* **106**, 5053-5058 (2009). <https://doi.org:10.1073/pnas.0812502106>
- 871 20 Mercer, A. C., Gaj, T., Fuller, R. P. & Barbas, C. F., 3rd. Chimeric TALE recombinases with
872 programmable DNA sequence specificity. *Nucleic Acids Res* **40**, 11163-11172 (2012).
873 <https://doi.org:10.1093/nar/gks875>

874 21 Prorocic, M. M. *et al.* Zinc-finger recombinase activities in vitro. *Nucleic Acids Res* **39**, 9316-9328
875 (2011). <https://doi.org:10.1093/nar/gkr652>

876 22 Standage-Beier, K. *et al.* RNA-Guided Recombinase-Cas9 Fusion Targets Genomic DNA Deletion
877 and Integration. *CRISPR J* **2**, 209-222 (2019). <https://doi.org:10.1089/crispr.2019.0013>

878 23 Voziyanova, E., Li, F., Shah, R. & Voziyanov, Y. Genome targeting by hybrid FIp-TAL
879 recombinases. *Sci Rep* **10**, 17479 (2020). <https://doi.org:10.1038/s41598-020-74474-2>

880 24 Mukhametzyanova, L. *et al.* Activation of recombinases at specific DNA loci by zinc-finger domain
881 insertions. *Nat Biotechnol* (2024). <https://doi.org:10.1038/s41587-023-02121-y>

882 25 Feng, X., Bednarz, A. L. & Colloms, S. D. Precise targeted integration by a chimaeric transposase
883 zinc-finger fusion protein. *Nucleic Acids Res* **38**, 1204-1216 (2010).
884 <https://doi.org:10.1093/nar/gkp1068>

885 26 Owens, J. B. *et al.* Transcription activator like effector (TALE)-directed piggyBac transposition in
886 human cells. *Nucleic Acids Res* **41**, 9197-9207 (2013). <https://doi.org:10.1093/nar/gkt677>

887 27 Owens, J. B. *et al.* Chimeric piggyBac transposases for genomic targeting in human cells. *Nucleic*
888 *Acids Res* **40**, 6978-6991 (2012). <https://doi.org:10.1093/nar/gks309>

889 28 Ye, L. *et al.* TAL effectors mediate high-efficiency transposition of the piggyBac transposon in
890 silkworm *Bombyx mori* L. *Sci Rep* **5**, 17172 (2015). <https://doi.org:10.1038/srep17172>

891 29 Yant, S. R., Huang, Y., Akache, B. & Kay, M. A. Site-directed transposon integration in human
892 cells. *Nucleic Acids Res* **35**, e50 (2007). <https://doi.org:10.1093/nar/gkm089>

893 30 Klompe, S. E., Vo, P. L. H., Halpin-Healy, T. S. & Sternberg, S. H. Transposon-encoded CRISPR-
894 Cas systems direct RNA-guided DNA integration. *Nature* **571**, 219-225 (2019).
895 <https://doi.org:10.1038/s41586-019-1323-z>

896 31 Strecker, J. *et al.* RNA-guided DNA insertion with CRISPR-associated transposases. *Science* **365**,
897 48-53 (2019). <https://doi.org:10.1126/science.aax9181>

898 32 Durrant, M. G. *et al.* Bridge RNAs direct modular and programmable recombination of target and
899 donor DNA. *bioRxiv*, 2024.2001.2024.577089 (2024). <https://doi.org:10.1101/2024.01.24.577089>

900 33 Lampe, G. D. *et al.* Targeted DNA integration in human cells without double-strand breaks using
901 CRISPR-associated transposases. *Nat Biotechnol* **42**, 87-98 (2024).
902 <https://doi.org:10.1038/s41587-023-01748-1>

903 34 Merrick, C. A., Zhao, J. & Rosser, S. J. Serine Integrases: Advancing Synthetic Biology. *ACS Synth*
904 *Biol* **7**, 299-310 (2018). <https://doi.org:10.1021/acssynbio.7b00308>

905 35 Bessen, J. L. *et al.* High-resolution specificity profiling and off-target prediction for site-specific
906 DNA recombinases. *Nat Commun* **10**, 1937 (2019). <https://doi.org:10.1038/s41467-019-09987-0>

907 36 Anzalone, A. V. *et al.* Programmable deletion, replacement, integration and inversion of large DNA
908 sequences with twin prime editing. *Nat Biotechnol* **40**, 731-740 (2022).
909 <https://doi.org:10.1038/s41587-021-01133-w>

910 37 Yarnall, M. T. N. *et al.* Drag-and-drop genome insertion of large sequences without double-strand
911 DNA cleavage using CRISPR-directed integrases. *Nat Biotechnol* **41**, 500-512 (2023).
912 <https://doi.org:10.1038/s41587-022-01527-4>

913 38 Xu, Z. *et al.* Accuracy and efficiency define Bxb1 integrase as the best of fifteen candidate serine
914 recombinases for the integration of DNA into the human genome. *BMC Biotechnol* **13**, 87 (2013).
915 <https://doi.org:10.1186/1472-6750-13-87>

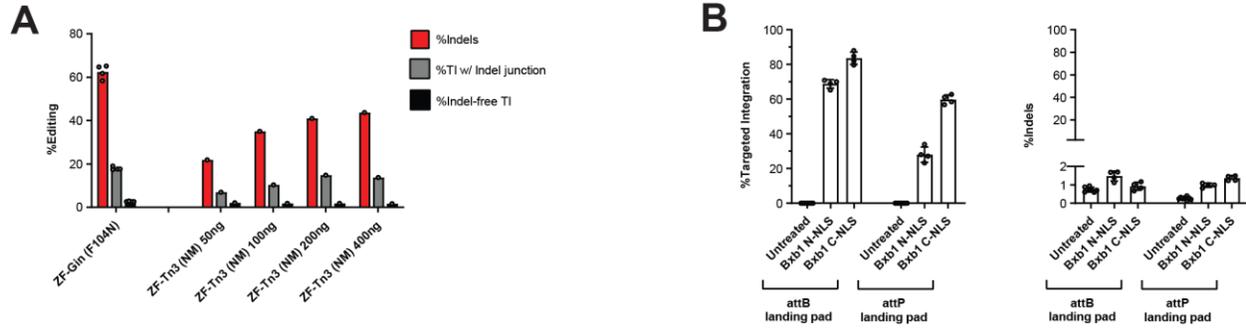
916 39 Rebar, E. J. & Pabo, C. O. Zinc finger phage: affinity selection of fingers with new DNA-binding
917 specificities. *Science* **263**, 671-673 (1994). <https://doi.org:10.1126/science.8303274>

918 40 Greisman, H. A. & Pabo, C. O. A general strategy for selecting high-affinity zinc finger proteins
919 for diverse DNA target sites. *Science* **275**, 657-661 (1997).
920 <https://doi.org:10.1126/science.275.5300.657>

921 41 Karpinski, J. *et al.* Directed evolution of a recombinase that excises the provirus of most HIV-1
922 primary isolates with high specificity. *Nat Biotechnol* **34**, 401-409 (2016).
923 <https://doi.org:10.1038/nbt.3467>

924 42 Lansing, F. *et al.* Correction of a Factor VIII genomic inversion with designer-recombinases. *Nat*
925 *Commun* **13**, 422 (2022). <https://doi.org:10.1038/s41467-022-28080-7>
926 43 Rutherford, K., Yuan, P., Perry, K., Sharp, R. & Van Duyne, G. D. Attachment site recognition and
927 regulation of directionality by the serine integrases. *Nucleic Acids Res* **41**, 8341-8356 (2013).
928 <https://doi.org:10.1093/nar/gkt580>
929 44 Li, H., Sharp, R., Rutherford, K., Gupta, K. & Van Duyne, G. D. Serine Integrase attP Binding and
930 Specificity. *J Mol Biol* **430**, 4401-4418 (2018). <https://doi.org:10.1016/j.jmb.2018.09.007>
931 45 Baek, M. *et al.* Accurate prediction of protein structures and interactions using a three-track neural
932 network. *Science* **373**, 871-876 (2021). <https://doi.org:10.1126/science.abj8754>
933 46 Gersbach, C. A., Gaj, T., Gordley, R. M. & Barbas, C. F., 3rd. Directed evolution of recombinase
934 specificity by split gene reassembly. *Nucleic Acids Res* **38**, 4198-4206 (2010).
935 <https://doi.org:10.1093/nar/gkq125>
936 47 Pavletich, N. P. & Pabo, C. O. Zinc finger-DNA recognition: crystal structure of a Zif268-DNA
937 complex at 2.1 Å. *Science* **252**, 809-817 (1991). <https://doi.org:10.1126/science.2028256>
938 48 Zheng, Z. *et al.* Anchored multiplex PCR for targeted next-generation sequencing. *Nat Med* **20**,
939 1479-1484 (2014). <https://doi.org:10.1038/nm.3729>
940 49 Kotin, R. M., Linden, R. M. & Berns, K. I. Characterization of a preferred site on human
941 chromosome 19q for integration of adeno-associated virus DNA by non-homologous
942 recombination. *EMBO J* **11**, 5071-5078 (1992). [https://doi.org:10.1002/j.1460-](https://doi.org:10.1002/j.1460-2075.1992.tb05614.x)
943 [2075.1992.tb05614.x](https://doi.org:10.1002/j.1460-2075.1992.tb05614.x)
944 50 Wang, X. *et al.* Bxb1 integrase serves as a highly efficient DNA recombinase in rapid metabolite
945 pathway assembly. *Acta Biochim Biophys Sin (Shanghai)* **49**, 44-50 (2017).
946 <https://doi.org:10.1093/abbs/gmw115>
947 51 Thomson, J. G. *et al.* The Bxb1 recombination system demonstrates heritable transmission of site-
948 specific excision in Arabidopsis. *BMC Biotechnol* **12**, 9 (2012). [https://doi.org:10.1186/1472-](https://doi.org:10.1186/1472-6750-12-9)
949 [12-9](https://doi.org:10.1186/1472-6750-12-9)
950 52 Blechl, A., Lin, J., Shao, M., Thilmony, R. & Thomson, J. The Bxb1 Recombinase Mediates Site-
951 Specific Deletion in Transgenic Wheat. *Plant Molecular Biology Reporter* **30**, 1357-1366 (2012).
952 <https://doi.org:10.1007/s11105-012-0454-2>
953 53 Jiang, L. *et al.* Target lines for recombinase-mediated gene stacking in soybean. *Theor Appl Genet*
954 **135**, 1163-1175 (2022). <https://doi.org:10.1007/s00122-021-04015-6>
955 54 Inniss, M. C. *et al.* A novel Bxb1 integrase RMCE system for high fidelity site-specific integration
956 of mAb expression cassette in CHO Cells. *Biotechnol Bioeng* **114**, 1837-1846 (2017).
957 <https://doi.org:10.1002/bit.26268>
958 55 Durrant, M. G. *et al.* Systematic discovery of recombinases for efficient integration of large DNA
959 sequences into the human genome. *Nat Biotechnol* **41**, 488-499 (2023).
960 <https://doi.org:10.1038/s41587-022-01494-w>
961 56 McCarty, D. M. *et al.* Adeno-associated virus terminal repeat (TR) mutant generates self-
962 complementary vectors to overcome the rate-limiting step to transduction in vivo. *Gene Ther* **10**,
963 2112-2118 (2003). <https://doi.org:10.1038/sj.gt.3302134>
964 57 Tsai, S. Q. *et al.* GUIDE-seq enables genome-wide profiling of off-target cleavage by CRISPR-Cas
965 nucleases. *Nat Biotechnol* **33**, 187-197 (2015). <https://doi.org:10.1038/nbt.3117>
966 58 Tareen, A. & Kinney, J. B. Logomaker: beautiful sequence logos in Python. *Bioinformatics* **36**,
967 2272-2274 (2020). <https://doi.org:10.1093/bioinformatics/btz921>

968



969

970 **Extended Data Figure 1. Performance of Zinc Finger Serine Recombinase fusion proteins**
 971 **and wild-type Bxb1. a.** Performance of Zinc Finger (ZF)-targeted Serine Recombinases Gin and
 972 Tn3. The data shown in this graph is derived from separate experiments and is representative for
 973 similar work that has been performed over a period of more than two years. In contrast to Bxb1,
 974 both Gin and Tn3 fusion proteins can result in high levels of indels causing low product purity and
 975 hindering further improvement of target integration frequencies. We also observed indels within
 976 the assayed TI junction site. Data is derived from a PCR-based NGS. **b.** Performance of wild-type
 977 Bxb1 against natural attB and attP target sequences in human cells. We first established K562
 978 landing pad cell lines by installing the natural attB or attP sequence in the human AAVS1 locus.
 979 We noticed improved performance of Bxb1 with a C-terminal NLS compared the construct with a
 980 N-terminal NLS. This guided future Bxb1 designs where all evolved variants presented in this
 981 study have a C-terminal NLS. We also noticed higher targeted integration into the attB landing
 982 pad. Notably, no or only minimal levels of indels were observed within the landing pad target
 983 sequences. Data is derived from a PCR-based NGS.

```

LI integrase      1 MKAAYIRVSTQEQVENYSIOAQTEKLTALCRSKDWDVYDTFIDGGYSGS-----NMNRP 55
M+A + IR+S      S + Q E   LC + WDV   D SG+      RP
Bxb1 integrase   1 MRALVVIRLSRVTDATT-SPERQLESCQQLCAQRGWDVVGVGAEDLDVSGAVDPFDRKRRP 59

LI integrase     56 ALNEMLS-KLHEIDAVVVYRLDRLRSRQKDTITLIEEYFLKNNVEFV--SLSETLDTSSP 112
L L+ +   D +V YR+DRL+RS +   L+ ++ +++ + V +   DT++P
Bxb1 integrase   60 NLARWLAFEEQPFVIVAYRVDRLRSIRHLQQLV--HWAEDHKKLVVSATEAHFDITTP 117

LI integrase     113 FGRAMIGILSVFAQLERETIRDRMVMGKIKRIEAGLPLTTAKGRITFGYDVID----TKLY 168
F +I ++  AQ+E E I++R      I AG   +   +GY      +L
Bxb1 integrase   118 FAAVVIALMGTVAQMELEAIKERNRSAAHFNIRAGKYRGSLS--PPWGYLPTRVDGEWRLV 175

LI integrase     169 INEEEAQQLRLIYDIFEEEQ-SITFLQKRLKKGKGF-----KVRTYNYRY 210
+ + +++ +Y   +   +   +   L + G               + +
Bxb1 integrase   176 PDFVQREERILEVYHRVVDNHEPLHLVAHDLNRRGVLSPKDYFAQLQGREPQGREWSATAL 235

LI integrase     211 NNWLTNDLYCGYVSYKDKVHVKG-----IHEPIISEEQFYRVQEIFSRMGK-NPNMNKE 263
+ ++  GY +   K               EPI++ EQ  ++   + + P ++
Bxb1 integrase   236 KRSMISEAMLGYATLNGKTVRDDGAPLVRAEPILTREQLEALRAELVKTSRAPAVSTP 295

LI integrase     264 SASLLNVLVCSKCGLGFVHRRKDTVSRGKKYHYRYYSCKTYKHTHELEKCGNKIWRADK 323
SLL ++ C+ CG               + G + H R Y C++      CGN   +
Bxb1 integrase   296 --SLLLRVLFCAVCG-----EPAYKFAGGGRKHPR-YRCRSMGFPK---HCGNGTVAMAE 344

LI integrase     324 LEELIIDRVNNSYFASRNID--KEDELDSLNEKLEKIEHAKKRLFDLYINGSYE-----V 376
+   ++V +   + ++-               +L +A+  L L + +Y
Bxb1 integrase   345 WDAFCEEQVLDLLGDAERLE-KVWVAGSDSAVELAEVNALVDLTSLIGSPAYRAGSPQR 403

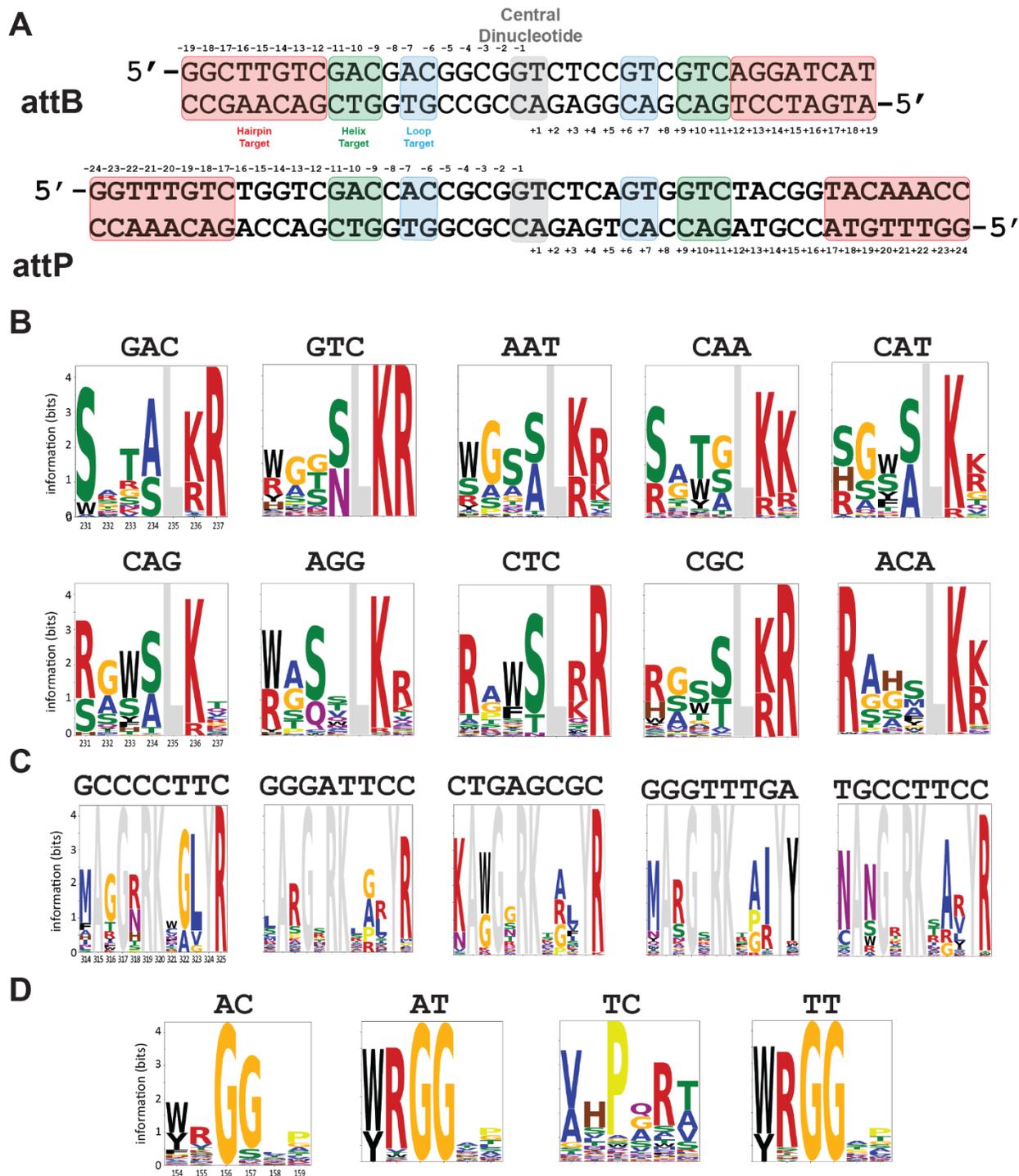
LI integrase     377 SELDSMMNDIDAQINYEAQIEANEELKKNKIQENLADLATVDFNSLEFREKQLYLKSLS 436
LD+ + + A+   E               + + Q               + K +L+S+
Bxb1 integrase   404 EALDARIAALAARQELEGLEARPSGWREWRETGQRFGDWWR-----EQDTAAKNTWLRSM 458

LI integrase     437 INKIYIDG---EQVTIEWL----- 452
++ D      TI++
Bxb1 integrase   459 NVRLTFDVRGGLTRIIDFGDLQEYEQHLRLGVSVERLHTGM 499

```

984 -

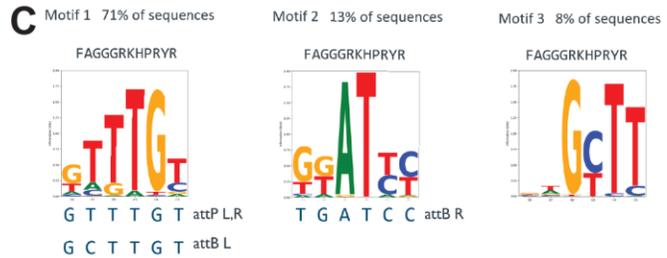
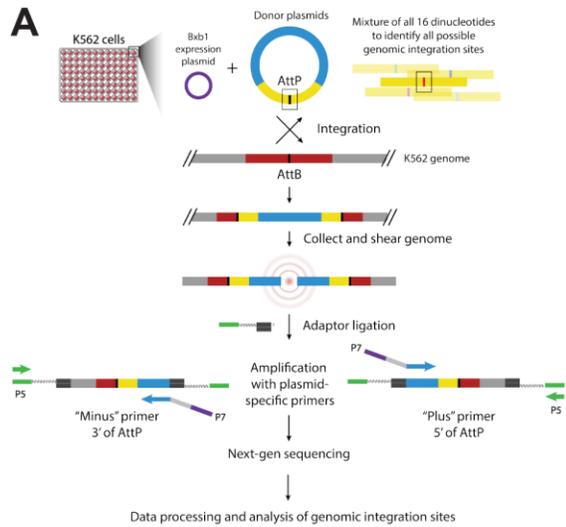
985 **Extended Data Figure 2. Sequence alignment of the large serine recombinase from the**
986 **Listeria innocua prophage and Bxb1.** The alignment was modified to reflect predicted secondary
987 structures. The region of Bxb1 that was probed with saturation mutagenesis is underlined.



988

989 **Extended Data Figure 3. a.** Schematic of attB and attP target sites with each base numbered and
 990 regions varied for Bxb1 hairpin, helix, and loop selections annotated. **b.** Example sequence logo
 991 plots summarizing enriched peptide motifs at amino acid residue positions 231-237 for helix
 992 selections with the corresponding helix DNA targets for each selection shown above each sequence
 993 logo. Position 235 was not randomized during selection and is shown in grey. Helices with an
 994 arginine (R) at position 237 were frequently observed in selections for a C at position -9 of the

995 DNA target, helices with a lysine (K) at position 237 were often observed in selections for a A at
996 position -9 of the DNA target, helices with an asparagine (N) at position 234 were often observed
997 in selections for a T at position -10 of the DNA target, helices with an alanine (A) or glycine (G)
998 at position 234 were often observed in selections for A at position -10 of the DNA target, and a
999 tryptophan (W) was often observed in selections with a C at position -11 of the DNA target. These
1000 correlations mimic interactions observed with engineered zinc fingers if one considers position
1001 233 to correspond to +2 of the zinc finger recognition helix, 234 to correspond to +3 of the zinc
1002 finger recognition helix, and 237 to correspond to +6 of the zinc finger recognition helix and if
1003 one considers the zinc finger DNA triplet to be the reverse complement of the Bxb1 helix target
1004 triplet. For example, from Ichikawa *et al.* (2023), the zinc finger helix QSGTLRR can target GCA,
1005 TKAYLLK can target TGA, QSSNLRT can target AAA, DPSALIR can target ATC, and
1006 RKWTLQQ can target AAG. This implies some similarities between how the Bxb1 helix and how
1007 the zinc finger recognition helix interacts with target DNA. **c.** Example logo plots summarizing
1008 the results of selections at amino acid positions 314-325 against different hairpin DNA targets.
1009 Residues that were not randomized are shown in gray. **d.** Example logo plots summarizing the
1010 results of selections at amino acid positions 154-159 against different loop DNA targets.



D

Pseudo-site	Sequence of indicated half-site	Hairpin target motif
chr18:23228790-23228837(L)	TGCTTGTGCCTACAGCC	motif 1
chr18:23228790-23228837(R)	AGTTTGTGAACCTCTGCA	motif 1
chr2:118918322-118918369(R)	GGTGTGTGGACTACAGTG	motif 1
chr10:95438013-95438060(R)	TGGATCTCACACCAGCC	motif 2
chr8:47341622-47341669(L)	TGGATCTCCACCCCTGTG	motif 2
chr12:104547644-104547691(R)	AAGATCTGACACCAGAG	motif 2
chr20:32457257-32457304(L)	CCTGCTTGTCTCAGCTGAG	motif 3
chr10:106891934-106891981(R)	GATGCTTACCAATTGAG	motif 3
chr20:32457257-32457304(R)	GCAGCTTCAGCCTCTGCC	motif 3

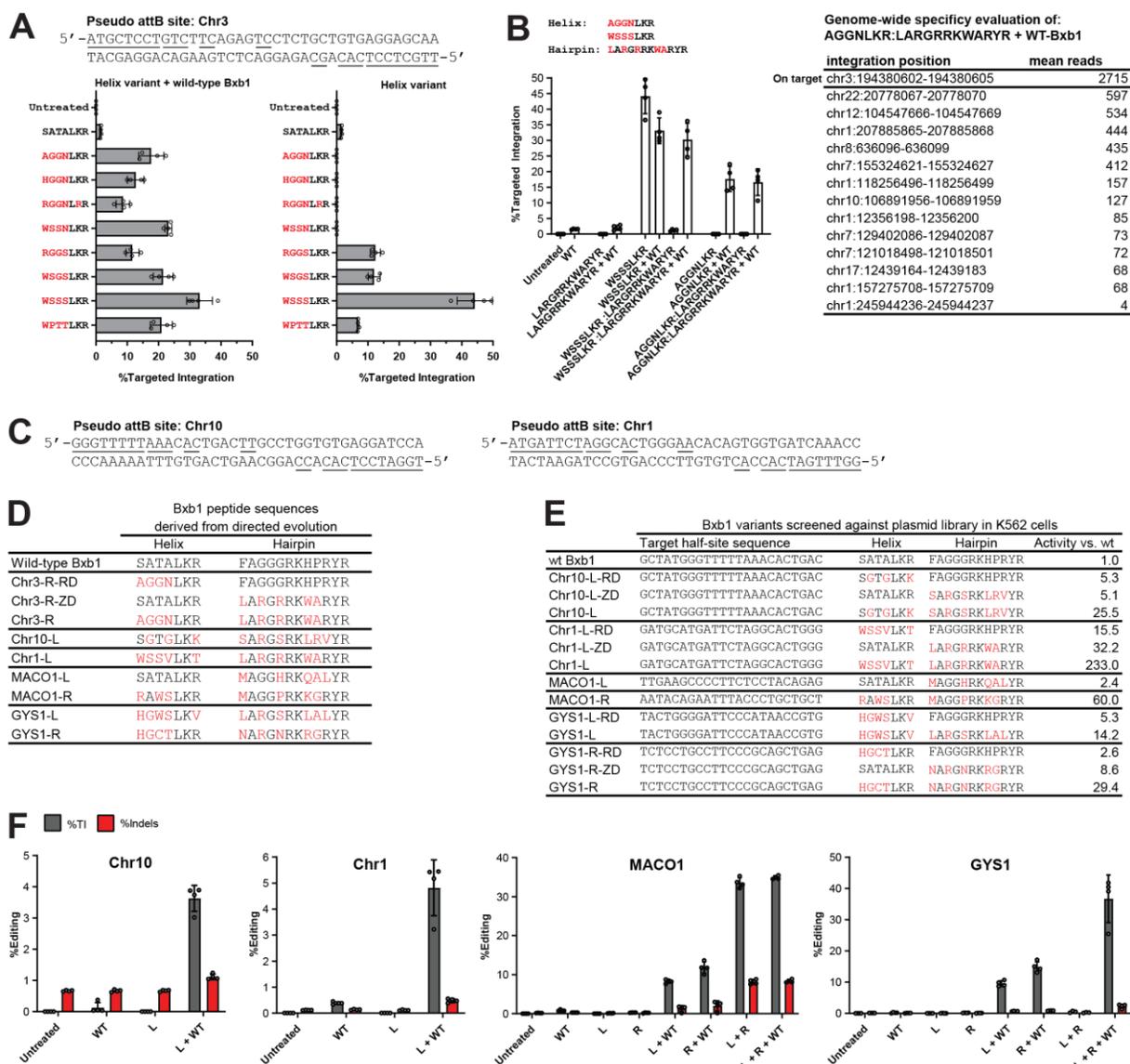
B

Genomic coordinates	Name in ED Fig. 5	Target sequence	%TI	Source
chr12:104547644-104547691		GTGATCTGACACCAGC <u>CA</u> CTCTGGTGTGCAAACTCT	6.92	Experimental
chr10:106891934-106891981		CCTTCTCTACCACAGCGT <u>TT</u> CTCAATTGGTAAGCATC	6.27	Experimental
chr3:194380580-194380627	Chr3	TTGCTCCTCACAGCAGAGG <u>ACT</u> CTGTAAGACAGGAGCAT	2.45	Computational
chr20:32457257-32457304		CCTGCTTGTCTCAGCTGAG <u>AG</u> CGGACAGGCTGAAGCTGC	2.11	Experimental
chr1:25477416-25477463	MACO1	GCCCTTCTCTACAGAG <u>CA</u> AGCAGCAGGTAATCT	1.59	Experimental
chr3:157294086-157294133		TGGTACTATACAGAG <u>AA</u> TGCTGTATGTTAAAGAGCT	1.49	Experimental
chr2:118918322-118918369		CTGATCTTGCAACTGA <u>AC</u> ACTGTAGTCCACACACC	1.29	Computational
chr8:47341622-47341669		TGGATCTCCACCCCTGT <u>TT</u> TGCTGTGTCTAAGTTTT	1.13	Experimental
chr6:87076238-87076285		CCAGCCTGGGCAACAGAG <u>CA</u> CTCTGTTGTCAAAAAA	1.09	Experimental
chr5:179235153-179235200		TTCTTTAAACCCTGAT <u>CA</u> TGCAAGTATAAGCTCA	0.85	Experimental
chr1:12564789-12564836		TTGCTTCTACCGTAGAT <u>G</u> CCCTCTGTCCACAAACC	0.64	Computational
chrX:50461699-50461746		GTTTTGTTCACTGTAGA <u>AT</u> CTCAGTGGTTAGCACAAAT	0.49	Computational
chr22:20778045-20778092		TTGCTCCTGACAGCAGTGGGAGCTATTGTCTAAGAGAT	0.43	Experimental
chr7:129402064-129402111		GAACTTACACACAGAGT <u>TT</u> CCAGTTGTTCAATCTCT	0.39	Experimental
chr3:37249472-37249519		CAATCTTGGCTCACTGCA <u>AC</u> CTCTGCAGTTGAAGCAAT	0.36	Experimental
chr18:33954336-33954383		CAGATCTATACGGCTGACT <u>AC</u> ACAGTGGTGAAGCAAT	0.34	Computational
chr1:114204100-114204147	Chr1	ATGATCTTAGGCACTGGG <u>AC</u> ACAGTGGTGAACAAACC	0.32	Computational
chr5:179764006-179764053		TGATCTCTGACTGCAGCAGT <u>CT</u> CAGCTGAGGAGCAGT	0.26	Experimental
chr19:48970994-48971041	GYS1	GGGATCCCAATAACCGT <u>GC</u> ACTCAGCTGCGGGAAGGCA	0.26	Experimental
chr1:207885846-207885893		TGTTTGGCCCCAAGTGCCT <u>CT</u> GCCACTGTCGACACACT	0.25	Computational
chrX:103694808-103694855		ATGTTCTCACCACAGCT <u>TA</u> CCCACTCTTCAAACTCA	0.25	Experimental
chr10:95438013-95438060	Chr10	GGGTTTTAAACACTGACT <u>TT</u> CGCTGGTGTGAGGATCCA	0.20	Computational
chr11:58774490-58774537		AGAAATAGTACACAGCT <u>TA</u> CCAGCATGTTAAAATCA	0.12	Experimental
chr18:23228790-23228837		TGCTTGTGCCTACAGCC <u>CT</u> TGCAAGTTCACAAATC	0.04	Computational
chr1:203596841-203596888		CTCTTGATGACTGCAGAGT <u>AT</u> TCATTGTTGACAAATC	0.04	Computational
chr4:78331572-78331619		GGGTTGTGGCAATGGAG <u>AT</u> CTCAGTGGTGAATAATCA	0.03	Computational

1011

1012 **Extended Data Figure 4. Identification of Bxb1 pseudo-sites in the human genome using an**
 1013 **unbiased genome-wide specificity assay. a.** Schematic of the unbiased-genome specificity assay
 1014 used in this study to experimentally identify pseudo-sites for wild-type Bxb1 in the human genome.
 1015 **b.** List of 26 validated Bxb1 pseudo-sites in the human genome that were either identified through
 1016 a computational search or experimentally using the assay shown in panel **a**. The central
 1017 dinucleotides in each site are underlined. Only the portion corresponding to an attB site are shown
 1018 in the table. Site-specific donors were used to validate targeted integration using a PCR-based NGS
 1019 assay. The 23 sites above the line have more than 0.1% TI **c.** Additional analysis of the Bxb1
 1020 hairpin specificity data for the wild-type Bxb1 hairpin in **Figure 2e** indicates the selected
 1021 sequences belong to at least three distinct DNA motifs. Plots of this data split into three separate
 1022 motifs are shown. The first motif is consistent with positions -18 to -13 of the hairpin targets in
 1023 both left and right half-sites of the natural attP site and the left half-site of the natural attB site. The
 1024 second motif is consistent with the right half-site of the natural attB site. **d.** Examples of half-sites

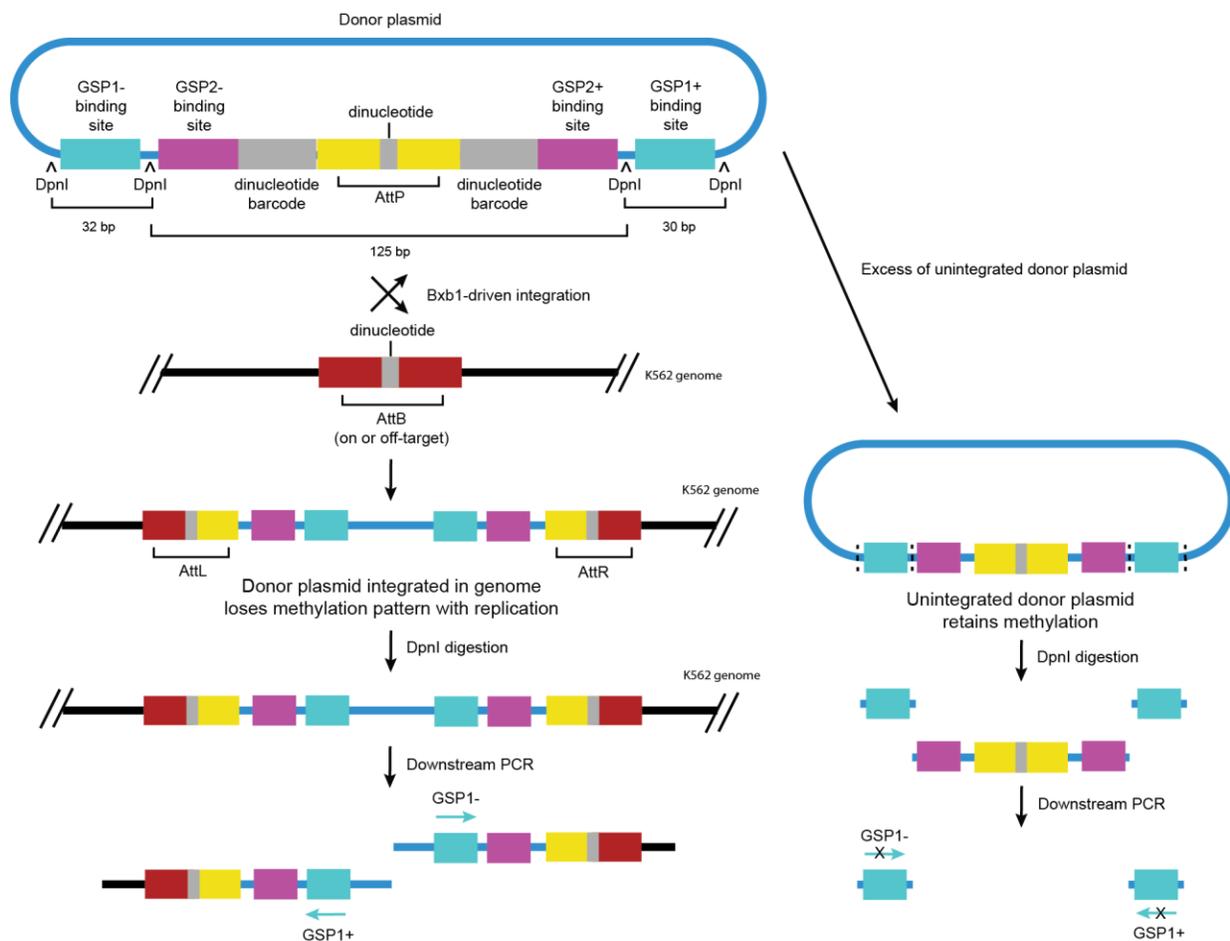
1025 from validated human pseudo-sites for wild-type Bxb1 that correspond to each of the three DNA
1026 sequence motifs are shown in panel c. Positions -18 to -13 of each half-site are underlined.



1027

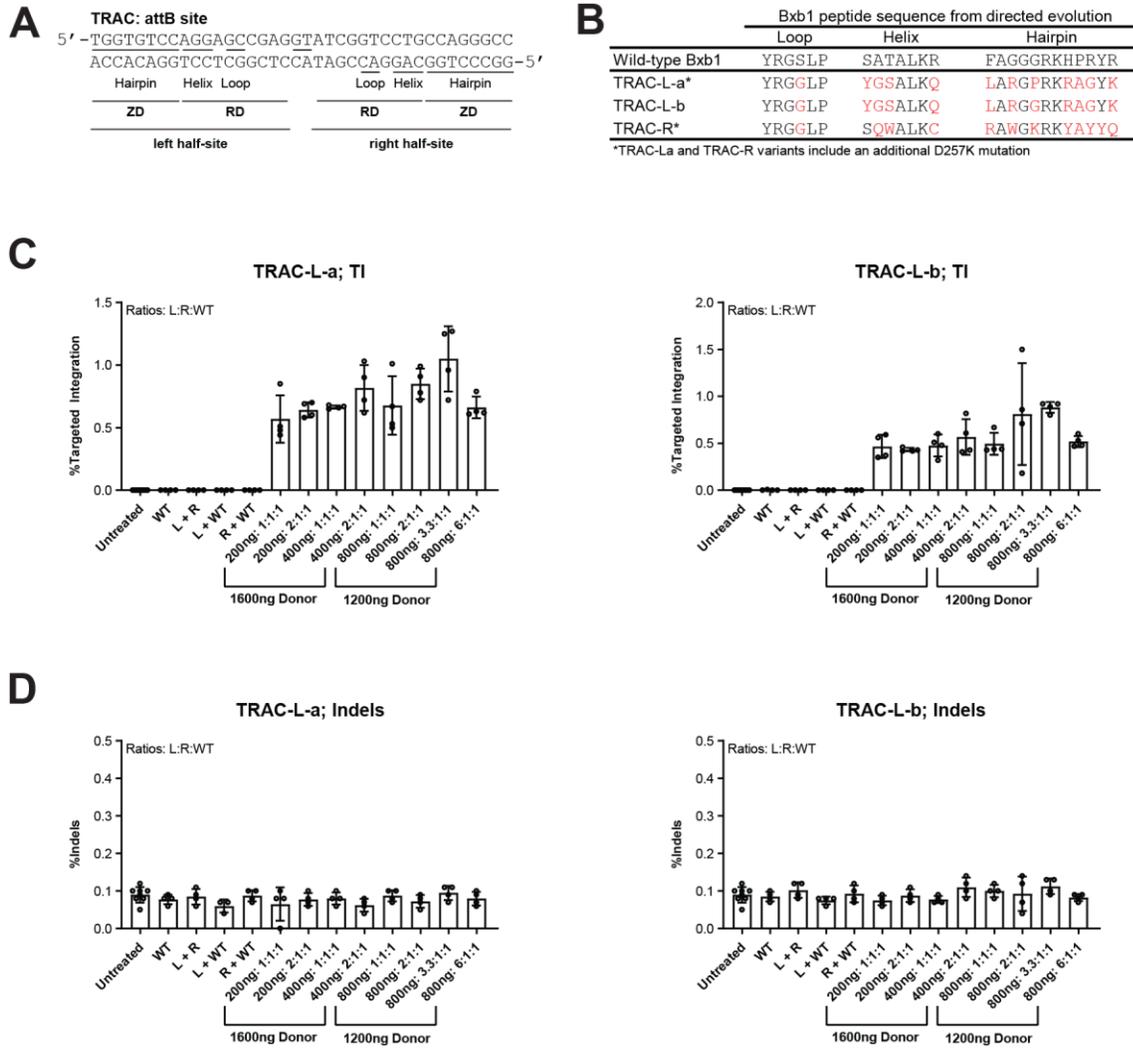
1028 **Extended Data Figure 5. Performance of Bxb1 variants against pseudo-sites in the human**
 1029 **genome.** **a.** Depicted is an attB pseudo-site on chromosome 3 and performance measurements of
 1030 various Bxb1 helix domain variants against this target site using a PCR-based NGS assay.
 1031 SATALKR is the wild-type helix peptide sequence and the target DNA sequence is shown as the
 1032 reverse complement of the sequence in Extended Data Figure 4b to make it easier to visualize the
 1033 region we are targeting in this experiment. **b.** On-target performance of a Bxb1 variant with
 1034 combined helix and hairpin variations, and data summary of a genome-wide specificity evaluation
 1035 using a modified version of the assay described in Extended Data Fig. 4 (see Extended Data Fig.
 1036 6). Note that the indicated Bxb1 variant has to be mixed with wild-type Bxb1 in order to be active
 1037 at the intended chr3 target site so the genome-wide specificity assay was performed with a mixture
 1038 of the indicated Bxb1 variant and wild-type Bxb1. The top 3 off-targets as well as two others are
 1039 also targets for wild-type Bxb1 (Extended Data Figure 4b) and are likely caused by the presence
 1040 of wild-type Bxb1 in the experiment. The experiment was also performed with a pool of donor
 1041 constructs containing all 16 possible central dinucleotide sequences. The intended target on
 1042 chromosome 3 is the only target on this list with a GA or TC central dinucleotide and thus the other

1043 sites presumably wouldn't have been detected if only a single donor with a GA or TC dinucleotide
1044 had been used in the genome-wide specificity experiment. **c.** Sequence of four additional Bxb1
1045 pseudo-sites in the human genome. **d.** Bxb1 peptide sequences of evolved Bxb1 variants that
1046 showed improved performance against the half-sites of pseudo-sites shown in **Figure 3b** and panel
1047 c. **e.** Screening data using synthetic DNA targets tested in human K562 cells that was used to
1048 identify the constructs shown in panel **d.** Activity is determined by the number of DNA sequence
1049 reads corresponding to recombined versions of each synthetic target; activity is normalized to the
1050 activity of wild-type Bxb1 alone against the same synthetic target site. **f.** Results from a PCR-
1051 based NGS assay demonstrating improved performance of evolved Bxb1 variants against their
1052 chromosomal endogenous targets in K562 cells. The presence of a wild-type Bxb1 expression
1053 construct is necessary to bind the wild-type attP sequence on the donor plasmid.

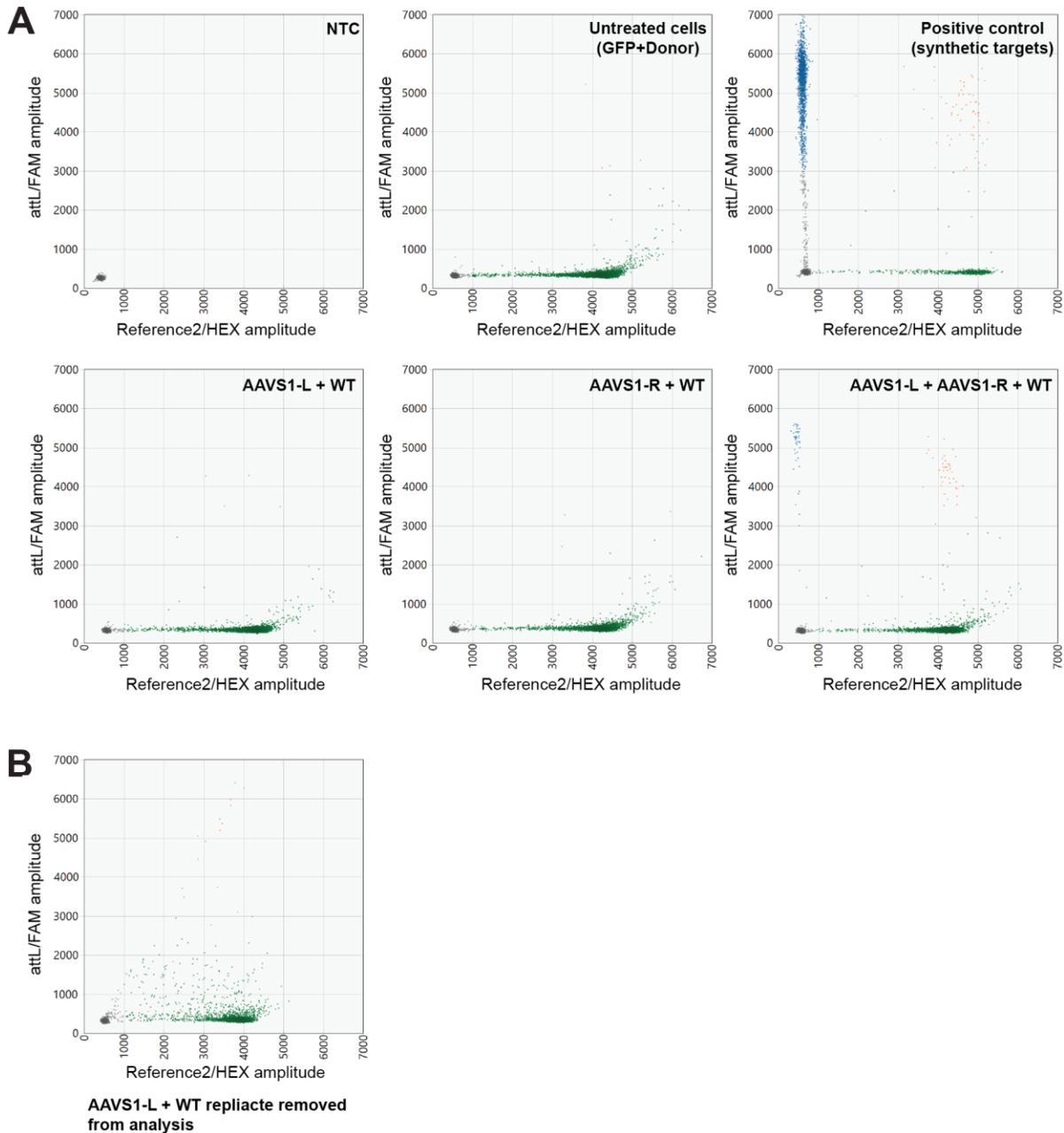


1054

1055 **Extended Data Figure 6. Improved genome-wide specificity assay.** Schematic of the modified
 1056 unbiased-genome specificity assay used in this study to experimentally identify integration site
 1057 shown in Extended Data Figure 5b. Strategically placing DpnI recognition sites in the donor
 1058 molecule supports the enzymatic removal of excess unintegrated donor plasmid, resulting in a
 1059 substantial reduction of donor plasmid-derived background signal.



1071
 1072 **Extended Data Figure 8. Bxb1 retargeting to the human TRAC locus.** a. Depicted is an attB-
 1073 like site in the human TRAC locus. b. Bxb1 peptide sequences of evolved Bxb1 variants that
 1074 showed improved performance against the half-sites of the TRAC site shown in panel a. c. Results
 1075 from a PCR-based NGS assay demonstrating targeted integration mediated by Bxb1 variants from
 1076 panel b at the TRAC site shown in panel a. A dose titration was performed where the total dose
 1077 was kept at 2000ng pDNA (Bxb1 expression constructs and donor combined). d. Indel analysis of
 1078 the experiment shown in panel c. All experiments were performed in human K562 cells.



1079

1080 **Extended Data Figure 9. ddPCR analysis of Bxb1 variants targeted to the human AAVS1**
 1081 **locus. a.** Example 2-D ddPCR scatterplots of differently treated samples. Each 2-D plot contains
 1082 1-4 clusters of droplets: (1) double-negative droplets containing no targeted DNA templates (grey
 1083 dots clustered at the left bottom in each plot); (2) reference-only droplets (green); (3) target
 1084 integration/attL-only droplets (blue); and (4) double-positive droplets containing both target
 1085 integration/attL and reference DNA templates (orange). **b.** The 2-D plot of an AAVS1-L + WT
 1086 replicate discarded from the analysis with noisy FAM signal likely due to shredded droplets.

Supporting Sequence Information

Bxb1 C-NLS MRALVVIRLSRVTDATTSPPERQLESCQQLCAQRGWDVVGVAEDLDVSGAVDPFDRKRRPNLARWLAFFEEQPFVIVAYRVDRLTRSIRHLQQLVHWAEDHKLVVS
ATEAHFDTTTTFPAVVIALMGTVAQMELEAIKERNRSAAHFNIRAGKYRGSLLPPWGYLPTRVDGEWRLVPPVQREIRILEVYHRVVDNHEPLHLVAHDLNRRGVLS
PKDYPAQLQGREPQGREWSATALKRSMISEAMLGYATLNGKTVRDDDGAPLVRAEPIILTREREALRAELVKT SRAKPAVSTPSLLLRLVLFCAVCGEPAYKFAAGG
RKHPRYRCRSMGFPHKCGNGTVAMAEDAFCEEQVLDLLGDAERLEKVVVAGSDSAVELAEVNAELVDLTSLIGSPAYRAGSPQREALDARIAALAAARQEELEGLE
ARPSGWEWRETGQRFQDWWREQDTAAKNTWLRSMNVRLTFDVRGGLTRTIDFGDLQEYEQHLRLGVSVERLHTGMSGSGSGSHHHHHHSGPKKRV**

MACO1

Donor fragment GGTTTGTCTGGTCAACCACCGCGCACTCAGTGGTGTACGGTACAAACCAAGCAGCTCGGAGCCTCTCTGTCACCTTGCTCTTTAGA
NGS fw primer ACACGACGCTCTCCGATCTNNNNCAATTAGTTGGCTGTATAAATTTGG
NGS rev primer GACGTGTGCTCTCCGATCTTCTAAAGAGCAAGTGACAGAGAGG

GYS1

Donor fragment GGTTTGTCTGGTCAACCACCGCGCACTCAGTGGTGTACGGTACAAACCAAGCCGAACTTAGCTTCCCTCATCGCCTA
NGS fw primer ACACGACGCTCTCCGATCTNNNNCGATGTGTCTCCATGAAGCA
NGS rev primer GACGTGTGCTCTCCGATCTTAGGCGATGAGGGAAGCTAAGT

Chr3

Donor fragment GGTTTGTCTGGTCAACCACCGCGCACTCAGTGGTGTACGGTACAAACCCCTGTGACTACATTTAGTGAGCAGGTGGAATGAACAA
NGS fw primer ACACGACGCTCTCCGATCTNNNNAGCCATTTCTCTCCATAGCAAAT
NGS rev primer GACGTGTGCTCTCCGATCTTTGTTTCAATCCACTGCTCACT

Chr1

Donor fragment GGTTTGTCTGGTCAACCACCGCGCACTCAGTGGTGTACGGTACAAACCAAGTACCAATCTGGAGGGATGAGGTGAGGGAGGATAGACAATA
NGS fw primer ACACGACGCTCTCCGATCTNNNNCAGCAGTCGATGTGGGAAC
NGS rev primer GACGTGTGCTCTCCGATCTTATTGTCTATCCCTCCCTCACCT

Chr10

Donor fragment GGTTTGTCTGGTCAACCACCGCGTCTCAGTGGTGTACGGTACAAACCAAGCATCTTTCAAATAGCACCTCATTTTATCTGAAGACCCAG
NGS fw primer ACACGACGCTCTCCGATCTNNNNAGCCAGAGTTAACCAAGCTAC
NGS rev primer GACGTGTGCTCTCCGATCTCTGGTCTTCCAGGATAAAATGAGG

TRAC

Donor fragment GGTTTGTCTGGTCAACCACCGCGTCTCAGTGGTGTACGGTACAAACCTCCCGTCATCCAGGTGCTCATATGCTGTAAGTTCCC
NGS fw primer ACACGACGCTCTCCGATCTNNNNGCCAGGTCATGCAACATGTAC
NGS rev primer GACGTGTGCTCTCCGATCTGGGAACCTACAGCATATGAGC

AAVS1

Donor fragment GGTTTGTCTGGTCAACCACCGCGTCTCAGTGGTGTACGGTACAAACCAAGCAGCCCGCCTTAGGGAAGCGGGACCCTGCTCTGGCGGAGGAATATGTC
NGS fw primer ACACGACGCTCTCCGATCTNNNNCATGAGATGGTGGACGAGGA
NGS rev primer GACGTGTGCTCTCCGATCTGACATATTTCTCCGCCAGAG