

1 **Peptide fusion improves prime editing efficiency**

2

3 Minja Velimirovic^{1,2}, Larissa C. Zanetti^{1,3}, Max W. Shen⁴, James D. Fife¹, Lin Lin⁵, Minsun Cha¹,
4 Ersin Akinci^{1,6}, Danielle Barnum^{1,7}, Tian Yu¹, Richard I. Sherwood¹

5

6

7 ¹ Division of Genetics, Department of Medicine, Brigham and Women's Hospital and Harvard Medical
8 School, Boston, MA 02115

9 ² Centre Hospitalier Universitaire de Québec Research Center—Université Laval, Québec, Québec G1V
10 4G2, Canada

11 ³ Hospital Israelita Albert Einstein, São Paulo, SP 05652-900, Brazil

12 ⁴ Merkin Institute of Transformative Technologies in Healthcare, Broad Institute of Harvard and MIT,
13 Cambridge, MA 02142, USA

14 ⁵ Hubrecht Institute, 3584 CT Utrecht, the Netherlands

15 ⁶ Department of Agricultural Biotechnology, Faculty of Agriculture, Akdeniz University, Antalya, 07070,
16 Turkey

17 ⁷Vrije Universiteit Amsterdam, Medical School of V, De Boelelaan 1105, 1081 HV Amsterdam,
18 Netherlands

19 Corresponding author: R.I.S.

20

21 **Abstract**

22 Prime editing enables search-and-replace genome editing but is limited by low editing efficiency.
23 We present a high-throughput approach, PepSEq, to measure how fusion of 12,000 85-amino
24 acid peptides derived from human DNA repair-related proteins influences prime editing
25 efficiency. We show that peptide fusion can enhance prime editing, prime-enhancing peptides
26 combine productively, and a top dual peptide-prime editor increases prime editing significantly in
27 multiple cell lines across dozens of target sites.

28

29

30

31

32 Prime editing is a CRISPR-based ‘*search-and-replace*’ technology that mediates targeted
33 insertions, deletions, and all possible base-to-base conversions in mammalian cells in the
34 absence of double-stranded breaks (DSBs) or donor DNA templates¹. The prime editing
35 enzyme (PE2) consists of SpCas9-nickase fused to an engineered reverse transcriptase (RT).
36 PE2 is recruited to a target site by a prime editing guide RNA (pegRNA) which, in addition to a
37 standard genome-targeting spacer and SpCas9-binding hairpin, contains a 3’ sequence that
38 acts as a template for the fused RT to synthesize a programmed DNA sequence on one of the
39 nicked DNA strands. When cellular DNA repair machinery repairs the broken strand, this RT-
40 extended flap competes with the unedited flap, and the edited sequence sometimes replaces
41 the original sequence in the genome^{1,2}.

42 Because of its versatility, prime editing has enormous potential in improving understanding of
43 the effects of genetic changes on cellular and organismal function. However, prime editing is
44 limited by low efficiency. While editing efficiency is dependent on the experimental system, a
45 survey of lentiviral PE2 efficiency at thousands of sites found that PE2 rarely leads to installed
46 edits in >20% of alleles³. Analysis of features associated with prime editing efficiency at these
47 thousands of loci found that the strongest correlate is DeepSpCas9 score^{3,4}, suggesting that
48 prime editing is limited by the interaction strength between the SpCas9-pegRNA complex and
49 the target locus. Optimization of pegRNA features³, induction of a distal nick on the opposite
50 strand (designated PE3)¹, and pairing overlapping pegRNAs⁵ have been found to improve prime
51 editing efficiency, yet low efficiency remains an issue in deployment of prime editing.

52 Here, we screened a library of 12,000 85-amino acid peptides derived from DNA repair proteins
53 to identify peptides that improve prime editing efficiency when appended to the N-terminus of
54 PE2. Peptide and protein fusion is a well-established method of modulating genome editing
55 outcomes⁶⁻⁸. While scalable, sensitive protein fusion screening remains challenging, high-
56 throughput oligonucleotide library synthesis enables screening of highly diverse peptide fusion
57 constructs. Reasoning that peptides derived from DNA repair-related proteins may encode
58 domains capable of altering prime editing efficiency, we designed a library of 85-amino acid
59 peptides comprising complete 2X tiling of 417 DNA repair-related proteins^{9,10} and 29
60 housekeeping genes as controls (Supplementary Table 1). We also included 5,458 DNA repair-
61 related mutant peptides with all possible S->E and T->E phosphomimetic substitutions. This
62 library of 12,000 oligos was cloned N-terminal to a 33-amino acid XTEN linker followed by PE2
63 in a vector allowing Tol2 transposon-mediated genomic integration (Fig. 1a)¹¹.

64 To enable quantitative evaluation of peptide-PE2 editing efficiency in high-throughput, we
65 devised the Peptide Self-Editing sequencing assay (PepSEq). We designed a self-targeting
66 pegRNA that introduces a 6-nt mutation (CCTCTG->GAATTC) in the peptide-adjacent linker
67 sequence (sgPE-linker). Following Tol2-mediated genomic integration of a single peptide-PE2
68 library member per cell¹², cells are treated in pooled format with sgPE-linker. To evaluate prime
69 editing efficiency in pooled format, we perform paired-end nextgen sequencing (NGS), mapping
70 peptide-PE2 identity and genotypic outcome for each self-targeted allele (Fig. 1a). We
71 performed initial PepSEq screens in mouse embryonic stem cells (mESCs) because they are a
72 non-transformed cell line not known to possess DNA repair defects, in contrast to other common
73 models such as HEK293T, which is known to lack mismatch repair capacity¹³.

74 We performed three biological replicates of peptide-PE2 integration, each followed by two
75 biological replicates of sgPE-linker addition, collecting >30M NGS reads for each replicate and
76 10M reads prior to sgPE-linker treatment. We developed a computational pipeline to filter and
77 analyze the NGS data (Online Methods), removing peptide-PE2s with <100 total reads in a
78 given replicate from analysis. In total, 16-28% of all alleles were prime edited, 0.3-0.7% were
79 indels, and nearly all remaining reads were unedited (Fig. 1b-c, Supplementary Fig. 1,
80 Supplementary Table 2). Due to the low frequency of indels and other alleles, we focused

81 analysis on prime editing efficiency for each peptide-PE2. Replicates that shared peptide-PE2
82 integration had strong consistency in prime editing efficiency ($R = 0.42-0.66$) while those with
83 distinct peptide-PE2 integration had negligible consistency ($R = 0.01-0.06$) (Supplementary Fig.
84 1). It is not surprising to obtain poor replicate consistency in a high-throughput screen in which
85 the majority of library members are expected to be inert, so we analyzed the three combined
86 replicates with independent peptide-PE2 integration using a beta-binomial model to identify 105
87 top candidate prime editing-enhancing peptide-PE2s (Online Methods).

88 To obtain higher-resolution data on this set of candidate peptide-PE2s, we cloned a sub-library
89 with these 105 peptide-PE2s and 10 control peptide-PE2s, performing PepSeq in five biological
90 replicate peptide-PE2 integrations in mESCs each with $>1M$ NGS reads. Replicate consistency
91 was much higher ($R = 0.25-0.52$, Supplementary Fig. 2, Supplementary Table 3), and the 105
92 candidate peptide-PE2s as a group gave 15% higher prime editing than control peptide-PE2s
93 ($P < 0.0001$, Fig. 1d), indicating that peptide fusion can improve prime editing efficiency. This
94 screen identified 44 peptide-PE2s that significantly improve prime editing efficiency ($FDR =$
95 0.05), increasing prime editing efficiency up to 70% (Fig. 1e-f). The proteins from which these
96 peptides are derived are not robustly enriched in any particular DNA repair pathway, and none
97 encompass known functional domains that appear related to hypothesized prime editing
98 mechanisms¹. This result additionally indicates that the 12,000-peptide PepSeq screen was
99 able to flag true hits in spite of noise, a finding supported by the fact that the 44 peptide-PE2s
100 that significantly increase prime editing in the smaller screen have appreciable replicate
101 consistency in the 12,000-peptide screens ($R = 0.17-0.39$, Supplementary Fig. 2).

102 To gain insight into how these peptides function, we next asked whether peptides that increase
103 prime editing combine productively. We constructed a dual peptide-PE2 library in which nine top
104 candidate peptides and one control peptide were combined in all 100 possible combinations
105 separated by an eight amino acid linker (Fig. 2a), and we performed 10 biological replicate
106 PepSeq screens in mESCs and two replicates in HCT-116 colorectal carcinoma cells. These
107 replicates were highly concordant within and between cell lines (mESC median $R = 0.62$, HCT-
108 116 $R = 0.47$, mESC vs. HCT-116 median $R = 0.32$, Supplementary Fig. 3, Supplementary
109 Table 4), and 79 of the 81 candidate dual peptide-PE2s gave significantly higher prime editing
110 efficiency than the control-control pair in mESCs (Fig. 2b, Supplementary Fig. 3).

111 To explore relationships between peptides, we asked whether dual-peptide-PE2 activation could
112 be predicted by a linear model assuming consistent peptide-specific influences on prime editing.
113 We find high consistency among observed prime editing and expected prime editing under
114 additive assumptions by linear estimates ($r = 0.92$) (Fig. 2c). The high accuracy of the linear
115 model suggests that each peptide has independent (not redundant or synergistic) effects on
116 prime editing, either through interacting with distinct pathways or providing a fixed advantage in
117 protein stability or DNA binding.

118 Eight of the nine dual peptide-PE2s with highest prime editing included an N-terminal peptide
119 from NFATC2IP (NFATC2IPp1), and the dual-peptide with highest median prime editing in
120 mESC and HCT-116 (median 1.77X control-control in mESC, 1.88X in HCT-116) pairs
121 NFATC2IPp1 with a phosphomimetic peptide from IGF1 (IGF1pm1) (Fig. 2d). These two
122 peptides induce the strongest increases in prime editing in the single-peptide-PE2 screen (Fig.
123 1c) and in the linear model of dual-peptide-PE2 screen (Supplementary Table 4), providing
124 rationale to pursue IGFpm1-NFATC2IPp1-PE2 (IN-PE2) as an improved prime editor.

125 To ask whether IN-PE2 increases prime editing efficiency across a larger collection of target
126 sites, we designed a lentiviral library comprising 100 pegRNA-target pairs spanning a range of
127 edit types and predicted editing efficiencies³ (Methods). After stable integration of this library in
128 three human and mouse cell lines (mESC, HEK293T, U2OS), we treated cells with either IN-
129 PE2 or a control PE2 containing the 5' linker sequence but lacking an N-terminal peptide

130 (CTRL-PE2). Among the targets with sufficient library representation and editing, IN-PE2
131 yielded significantly higher prime editing than CTRL-PE2 in all three cell lines (median 1.63X in
132 mESC at 19 sites, 1.31X in HEK293T at 27 sites, 1.23X in U2OS at 12 sites) Fig. 2e,
133 Supplementary Fig. 4, Supplementary Table 5). The results indicate that IN-PE2 leads to a
134 consistent increase in prime editing efficiency across a variety of targets (Supplementary Fig. 4).
135 We performed pegRNA-target library experiments with six additional peptide-PE2s with one to
136 three top candidate peptides fused to PE2, finding IN-PE2 to display the most robust prime
137 editing of these seven peptide-PE2s (Supplementary Fig. 4). The consistent increase in prime
138 editing efficiency in cell lines without known DNA repair defects (mESC, U2OS) and those with
139 known deficiencies in mismatch repair (HEK293T¹³, HCT-116¹⁴) suggests that IN-PE2 is unlikely
140 to function through interaction with mismatch repair machinery.

141 We next asked whether IN-PE2 increases prime editing at endogenous genomic loci. We
142 transduced HEK293T and U2OS with a pool of six lentiviral pegRNAs encoding missense
143 variants in exon 8 of the NF2 gene, each in two biological replicates. We subsequently treated
144 cells with CTRL-PE2 or IN-PE2 and performed genomic DNA NGS to determine editing
145 outcomes. IN-PE2 treatment led to significantly increased prime editing efficiency in both cell
146 lines across the six loci (median 1.17X in HEK293T, 1.15X in U2OS, $p < 0.01$ in each cell line)
147 (Fig. 2f, Supplementary Fig. 5). There was variability in the magnitude of increased prime
148 editing across sites, but all six loci had increased prime editing in each cell line, suggesting that
149 IN-PE2 consistently increases prime editing efficiency at native genomic targets.

150 To gain insight into the mechanism by which IN-PE2 increases prime editing efficiency, we
151 constructed IN-GFP-PE2 and CTRL-GFP-PE2 fusions to address whether IN-PE2 increases the
152 amount of protein in cells. We found that cells possess 1.58X the amount of IN-GFP-PE2 as
153 CTRL-GFP-PE2 (N=5, $p < 0.0001$, Fig. 2g, Supplementary Fig. 6). Cycloheximide timecourse
154 experiments show that IN-GFP-PE2 and CTRL-GFP-PE2 are degraded at a similar rate
155 (Supplementary Fig. 6), altogether suggesting that the NI peptides increase either transcription
156 or translation of the PE2 enzyme and offering a plausible explanation for the increased activity
157 of IN-PE2.

158 In summary, through screening 12,000 peptide-PE2 fusion proteins using PepSeq, a sensitive,
159 NGS-based self-editing platform, we identify a prime editor that consistently increases editing
160 efficiency across dozens of targets in four human and mouse cell lines. Because prime editing
161 applications are currently limited by editing efficiency, we anticipate that IN-PE2 will be a
162 valuable tool in elucidating how DNA sequence influences genome function.

163

164

165 **Acknowledgments**

166 The authors thank Mandana Arbab for technical assistance. The authors acknowledge funding
167 from 1R01HG008754 (R.I.S.), 1R21HG010391 (R.I.S., C.A.C.), American Cancer Society
168 (R.I.S.), American Heart Association (R.I.S.), Qatar Biomedical Research Institute (R.I.S.), and
169 the São Paulo Research Foundation- FAPESP nº 2019/13813-6 and 2017/25009-1 (L.C.Z.).
170 Figures created with BioRender.

171

172 **Author contributions**

173 Conceptualization, Methodology, Writing – Original Draft and Writing – Reviewing and Editing:
174 R.I.S., M.V., L.L., M.W.S., J.D.F. Investigation and Validation: R.I.S., M.V., L.C.Z., M.C., E.A.,
175 D.B., M.W.S., J.D.F., T.Y. ; Software, Formal Analysis and Visualization: R.I.S., M.W.S., J.D.F.,
176 T.Y., M.V. ; Funding Acquisition and Supervision: R.I.S.

177

178 **Competing interests**

179 The authors report no competing interests.

180

181 **Figures**

182

183 **Figure 1. The high-throughput Peptide Self-Editing sequencing assay (PepSeq) identifies**
184 **peptides capable of increasing prime editing efficiency.**

- 185 (a) In PepSeq, a library of peptides from human DNA repair-related genes is cloned N-
186 terminal to SpCas9-PE2 separated by a linker and integrated into cells at one copy per
187 cell. Cells are subsequently treated with a pegRNA targeting the linker sequence that
188 installs a fixed edit. Paired-end genomic DNA NGS of the peptide sequence and the
189 editing site allows calculation of prime editing outcomes in high throughput.
190 (b) Observed prime editing outcome frequencies for a 12,000-peptide PepSeq screen. Box
191 plot indicates median and interquartile range, and whiskers indicate extrema.
192 (c) Overall distribution of prime editing outcome frequencies across all peptides and
193 replicates.
194 (d) Comparison of prime edited allele fraction for 105 DNA repair-related peptides vs. 10
195 housekeeping control peptides from a 115-peptide PepSeq screen.
196 (e) Volcano plot showing control-normalized prime editing fold change (x-axis) vs. vs. $-\log_{10}$
197 p-value (y-axis) from 115-peptide PepSeq screen.
198 (f) Comparison of control-normalized prime edited allele fraction for nine top peptides and
199 all control peptides from a 115-peptide PEPSeq screen.

200

201 **Figure 2. A dual-peptide-PE2 displays improved prime editing efficiency across dozens**
202 **of loci in four cell lines.**

- 203 (a) Construct design for a dual-peptide PepSeq library with all pairs of 10 peptides.
204 (b) Comparison of control-normalized prime edited fraction for 81 dual-peptide pairs.
205 (c) Comparison of control-normalized dual-peptide prime edited fraction predicted by a
206 linear model vs. observed median prime edited fraction.
207 (d) Comparison of prime edited fraction for IN-PE2 vs. control-control-PE2 in mESC and
208 HCT-116 dual-peptide screen replicates.
209 (e) Comparison of median control-normalized prime edited fraction for IN-PE2 in the 100-
210 target library across three cell lines.
211 (f) Comparison of prime edited fraction for IN-PE2 vs. CTRL-PE2 at a set of six
212 endogenous sites in HEK293T and U2OS.
213 (g) Flow cytometric GFP fluorescence intensity for two representative replicates of IN-GFP-
214 PE2 vs. CTRL-GFP-PE2 in mESCs.
215 (h) Comparison of IN-GFP-PE2 vs. CTRL-GFP-PE2 control-normalized flow cytometric GFP
216 fluorescence levels. N = 5.

217

218

219 **Online Methods**

220 **Peptide library design**

221 We designed 85-amino acid peptides covering all annotated human DNA repair
222 proteins^{16,17}, tiling by starting each peptide 45-amino acids after the prior peptide using a codon-
223 optimized library design¹⁵. We also included mutant peptides with all possible S-->E and T-->E
224 phosphomimetic substitutions. 147 wild-type peptides targeting 29 housekeeping genes were
225 also included as controls. Unique 9-nt sequences were inserted in phosphomimetic peptides to
226 facilitate sequence mapping for downstream analysis. The sequence design was performed with
227 “seqinr” and “Biostrings” packages in R.

228 **Cell Culture**

229 All cell lines were obtained from ATCC and were cultured in: McCoy’s 5A media (Thermo
230 Fisher) + 10% FBS (Thermo Fisher) (U2OS, HCT-116); DMEM (Thermo Fisher) + 10% FBS
231 (HEK293); mESCs were maintained on gelatin-coated plates feeder-free in mESC media
232 composed of Knockout DMEM (Life Technologies) supplemented with 15% defined foetal
233 bovine serum (FBS) (HyClone), 0.1 mM nonessential amino acids (NEAA) (Life Technologies),
234 Glutamax (GM) (Life Technologies), 0.55 mM 2-mercaptoethanol (b -ME) (Sigma), 1X ESGRO
235 LIF (Millipore), 5 nM GSK-3 inhibitor XV and 500 nM UO126. Cells were regularly tested for
236 mycoplasma.

237

238 **Peptide library cloning and screening**

239 The SpCas9-PE2-encoding sequence from pCMV-PE2¹ (Addgene Plasmid #132775) was
240 subcloned into the p2T-CAG-SpCas9-BlastR plasmid¹⁸ (Addgene Plasmid #107190) to create
241 p2T-CAG-SpCas9PE2-5pLinker-BlastR (Addgene Plasmid #173066)

242 Specified oligonucleotide libraries were synthesized by Twist Bioscience (12,000-peptide) or
243 IDT (115-peptide and dual-peptide) and were cloned into the NheI site of p2T-CAG-
244 SpCas9PE2-5pLinker-BlastR through amplification with Q5® High-Fidelity DNA Polymerase
245 (New England Biolabs) using primers Cas9NTLib_GA_fw and Cas9NTLib_GA_rv (see
246 Supplementary Table) followed by ligation using the NEBuilder HiFi DNA Assembly Kit (NEB) for 1
247 h at 50 °C. Assembled plasmids were purified by isopropanol precipitation with GlycoBlue
248 Coprecipitant (Thermo Fisher) and reconstituted in TE and transformed into NEB® 10-beta
249 Electrocompetent *E. coli* (NEB). Following recovery, the library was grown in liquid culture in LB
250 medium overnight at 37 °C and isolated by ZymoPURE™ II Plasmid Maxiprep Kit (Zymo
251 Research). Library integrity was verified by restriction digest with Agel (New England Biolabs)
252 for 1 h at 37 °C, and library diversity was validated by Sanger sequencing sampling.

253 Mouse ESC cells were plated at ~20-25% confluence onto 25-cm plates the day before
254 transfection so that they reach ~50-75% confluency on the day of transfection. For stable Tol2
255 transposon plasmid integration, cells were transfected using Lipofectamine 3000 (Thermo
256 Fisher) following standard protocols, and equimolar amounts of Tol2 transposase plasmid and
257 transposon-containing plasmid. To generate lines with stable Tol2-mediated genomic
258 integration, selection with the appropriate selection agent at an empirically defined
259 concentration (blastidicin, hygromycin, or puromycin) was performed starting 24 h after
260 transfection and continuing for >1 week.

261

262 In cases where sequential plasmid integration was performed such as integrating
263 pegRNA/target library and then Cas9, the same Lipofectamine 3000 transfection protocol with
264 Tol2 transposase plasmid was performed each time, and >1 week of appropriate drug selection
265 was performed after each transfection.

266 **Deep sequencing, library preparation**

267 Genomic DNA was extracted from harvested cells with the PureLink Genomic DNA Purification
268 Mini Kit (Invitrogen). For library experiments, sequences including the peptide and the prime
269 editing site were PCR amplified using Q5® High-Fidelity DNA Polymerase (NEB) and primers
270 as specified (Supplementary Table 6). For each replicate, the first PCR included a total of 10-20
271 µg of genomic DNA. To determine the number of cycles required to complete the exponential
272 phase of amplification we first performed qPCR, followed by PCR using primers that included
273 both Illumina adaptor and barcode sequences (Supplementary Table 6). For measuring PE2
274 efficiencies at endogenous sites, the independent first PCR was performed in a 200ul reaction
275 volume that contained 1000ng of the initial genomic DNA template per sample. The second
276 PCR to attach the Illumina adaptor and barcode sequences was then performed using purified
277 product from the first PCR. After bead purification, pooled samples were sequenced using
278 NextSeq (Illumina).

279 **Library data processing**

280 Designed library peptides were identified in sequenced reads by exact string matching to
281 the first eight nucleotides of the peptide sequence, which were unique across the library.
282 Sequenced target sites were aligned to the designed reference using Needleman-Wunsch with
283 match score 1, mismatch cost -1, gap open cost -5, and gap extend cost 0. Reads with mean
284 PHRED quality score below 30 were filtered. Mismatches at nucleotides with less than PHRED
285 quality score 30 were filtered. Indels with less than three matching nucleotides on both sides
286 with at least PHRED quality score 30 were filtered.

287

288 **Identifying peptide hits**

289 We excluded peptides with less than 100 reads in any experiment. We used a beta-
290 binomial model to infer peptide editing effects from replicate data while adjusting for sampling
291 noise from limited sequencing reads. We model a peptide i with parameters α_i, β_i used to
292 sample a peptide effect $p_{ij} \sim \text{Beta}(\alpha_i, \beta_i)$ for an experiment or batch j . Samples from
293 experiment or batch j are taken for sequencing, yielding a binomial distribution over the number
294 of edited reads $y_{ij} \sim \text{Bin}(n_j, p_{ij})$ for read depth n_j . Given k samples of y_{ij} over k biological
295 replicates or batches, we infer the maximum likelihood estimate (MLE) of α_i, β_i for peptide i . As
296 our beta-binomial model is conjugate, the MLE of α_i, β_i can be found analytically by solving the
297 system of equations:

298

$$\begin{aligned} 0 &= -k\Gamma(\alpha_i) - \Gamma(\alpha_i + \beta_i) + \sum_j \Gamma(y_{ij} + \alpha_i) - \Gamma(n_j + \alpha_i + \beta_i) \\ 0 &= -k\Gamma(\beta_i) - \Gamma(\alpha_i + \beta_i) + \sum_j \Gamma(n_j - y_{ij} + \alpha_i) - \Gamma(n_j + \alpha_i + \beta_i) \end{aligned}$$

299

300 Where $\Gamma(\cdot)$ is the Gamma function. We solved this using Sympy¹. When solutions could not be
301 found due to numerical instability, we used a fast approximation that solves the MLE of α_i, β_i by
302 matching the observed mean and variance, motivated by viewing the beta-binomial distribution
303 as an overdispersed binomial distribution. The additional variance over a binomial distribution is
304 related to the sum $\alpha_i + \beta_i$.

$$\begin{aligned} \text{obs var} &= \text{var}_j(y_{ij}/n_j) \\ \text{expvar} &= \frac{\text{mean}_j(y_{ij}/n_j) * (1 - \text{mean}_j(y_{ij}/n_j))}{\text{mean}_j(n_j)} \\ \alpha_i + \beta_i &= \frac{\text{mean}(n_j) - 1}{(\text{obsvar}/\text{expvar}) - 1} - 1 \\ \frac{\alpha_i}{\alpha_i + \beta_i} &= \text{mean}_j(y_{ij}/n_j) \end{aligned}$$

306
307
308

309 To interpret α_i, β_i , we convert them into the mean $\frac{\alpha_i}{\alpha_i + \beta_i}$ and variance $\frac{\alpha_i \beta_i}{(\alpha_i + \beta_i)^2 (\alpha_i + \beta_i + 1)}$ of a
310 Beta distribution.

311 We selected peptides for follow-up evaluation using several metrics. To increase
312 confidence in hits, we preferred peptides present in higher numbers of replicates. We prioritized
313 peptides based on the probability of observing a higher edited read count under its inferred
314 peptide effect parameters compared to edited read counts sampled from inferred control peptide
315 effect parameters under our beta-binomial model, which prefers higher MLE mean and lower
316 MLE variance. We also selected peptides with high MLE mean effect even if their variance was
317 high.

318

319 **100 target site library design.**

320 An oligonucleotide pool containing 100 target sequences was synthesized by IDT. Each
321 oligonucleotide contained the following elements 3' to 5': 19nt PE1 stub, 4nt barcode, ~40nt
322 variable target, 6nt poly A terminator, ~30nt PBS/template, 86nt hairpin, 20nt spacer, 20nt U6
323 stub. The barcode stuffer allowed individual pegRNA and target sequence pairs to be identified
324 after deep sequencing. To test the effect of PBS and RT template length on PE2 efficiency, we
325 prepared pegRNAs with 8 different combinations of edit types. Types of mutations:

- 326 · 3 x 1-nt substitution
 - 327 ○ PAM NGG-->NCG
 - 328 ○ PAM NGG-->NGT
 - 329 ○ Seed 1 transversion—nt nearest PAM, AàT, TàA, CàG, GàC
- 330 · 3 x >1-nt substitution
 - 331 ○ PAM NGG-->NCT
 - 332 ○ Seed 2-3 transversion (AàT, TàA, CàG, GàC) + PAM NGGàNTC
 - 333 [discontinuous, maintain 2 intervening nt]

- 334 ○ 6-nt PAM + seed change to GAATTC
- 335 · 1 x 1-nt ins
- 336 ○ PAM NGG-->NGTG
- 337 · 1 x 1-nt del
- 338 ○ PAM NGG-->not NGG
- 339 - delete 1st G unless G after PAM
- 340 - Delete seed 1 unless 1st base of PAM is identical

341 To design a library of 100 pegRNA-target pairs we used 96 pegRNA-target pairs from Kim et al
342 high-throughput library data that vary in prime editing efficiency. In their library, for all sites with
343 DeepSpCas9 score >20, average PE efficiency is 11%, SD=9%. We chose 4 target sites 0-1 SD
344 below average (1%, 3%, 5%, 8% PE), 4 sites around average (11%, 14%, 17%, 20%), 2 sites
345 ~1 SD above average (30%, 40%), 2 near top (50%, 60%). Our library also includes 4
346 substitution mutations from Anzalone et al that showed highest prime editing activity. Oligo
347 library was cloned into pLenti-sgRNA-FE vector using NEBuilder HiFi DNA Assembly Kit (NEB).
348 Assembled plasmids were purified by isopropanol precipitation with GlycoBlue Coprecipitant
349 (Thermo Fisher) and reconstituted in TE and transformed into Endura™ Electrocompetent Cells
350 (Lucigen). After library diversity was verified, library mastermix was used to produce lentivirus.
351

352 **Production of lentivirus and cellular infection**

353 HEK293FT cells (15×10^6) were seeded on 150-mm cell culture dishes containing DMEM.
354 The next day, cells were transfected with pCMV-VSV-G (Addgene #8454), pRSV-Rev (Addgene
355 #12253), pMDLg/pRRE (Addgene #12251) and library, at a ratio of 1:2:3:4, using TransIT®-
356 Lenti Transfection Reagent (Mirus Bio). At 8h after transfection, cells were refreshed with
357 maintaining medium. At 24h and 48h after transfection, the lentivirus-containing supernatant
358 was collected, filtered through a 0.45- μ m pore filter (Corning), concentrated using Lenti-X™
359 Concentrator (TakaraBio), aliquoted and stored at -80°C .

360 In preparation for lentivirus transduction, cells (U2OS, HCT-116, mESC, HEK293FT) were
361 seeded on 100-mm dishes (at a density of 2×10^5 , 6.5×10^4 , 6.5×10^4 , 1×10^5 cells per
362 cm^2) and concentrated lenti was added to the media. The cells were then incubated overnight,
363 after which cells were refreshed with maintaining medium before adding blasticidin at 48h and
364 keeping it for minimum of next 5 d to remove untransduced cells. To preserve its diversity, the
365 cell library was maintained at a count of at least 1×10^7 cells throughout the study.

366 **Measurement of PE2 efficiencies at endogenous sites**

367
368
369 To validate the results of the high-throughput experiments, 6 individual pegRNA-encoding
370 plasmids targeting endogenous NF2 locus were constructed and used to produce lentiviral

371 particles. In preparation for transfection, HEK293T and U2OS cells were seeded on 10 cm
372 plates at a density of 4×10^4 cells per cm^2 and transduced with a lentivirus carrying pegRNA-
373 encoding plasmid. After cells were selected for successful lentiviral integration, they were
374 transfected using Lipofectamine 3000 with plasmid encoding IN-PE2 or PE2-control and
375 equimolar amounts of Tol2 transposase plasmid, according to the manufacturer's instructions.
376 After a week of selection for successful integration of constructs, cells were harvested for gDNA
377 extraction followed by library preparation for NGS. Primers used to sequence NF2 locus are
378 listed in Supplementary table 6.

379

380

381 **Code availability**

382 Custom code used to process and analyze peptide library data are available at
383 <https://github.com/maxwshen/prime-peptide>.

384

385

386

387

388

389

390

391

392

393

394

395

396 **References**

- 397 1. Anzalone, A. V. *et al.* Search-and-replace genome editing without double-strand breaks or
398 donor DNA. *Nature* **576**, 149–157 (2019).
- 399 2. Anzalone, A. V., Koblan, L. W. & Liu, D. R. Genome editing with CRISPR-Cas nucleases,
400 base editors, transposases and prime editors. *Nat. Biotechnol.* **38**, 824–844 (2020).
- 401 3. Kim, H. K. *et al.* Predicting the efficiency of prime editing guide RNAs in human cells. *Nat.*
402 *Biotechnol.* 1–9 (2020) doi:10.1038/s41587-020-0677-y.
- 403 4. Kim, H. K. *et al.* SpCas9 activity prediction by DeepSpCas9, a deep learning–based model
404 with high generalization performance. *Sci. Adv.* **5**, eaax9249 (2019).
- 405 5. Lin, Q. *et al.* High-efficiency prime editing with optimized, paired pegRNAs in plants. *Nat.*
406 *Biotechnol.* (2021) doi:10.1038/s41587-021-00868-w.
- 407 6. Zhang, X. *et al.* Increasing the efficiency and targeting range of cytidine base editors
408 through fusion of a single-stranded DNA-binding protein domain. *Nat. Cell Biol.* **22**, 740–750
409 (2020).
- 410 7. Komor, A. C., Kim, Y. B., Packer, M. S., Zuris, J. A. & Liu, D. R. Programmable editing of a
411 target base in genomic DNA without double-stranded DNA cleavage. *Nature* **533**, 420–424
412 (2016).
- 413 8. Pickar-Oliver, A. & Gersbach, C. A. The next generation of CRISPR–Cas technologies and
414 applications. *Nat. Rev. Mol. Cell Biol.* **20**, 490–507 (2019).
- 415 9. Owusu, M. *et al.* Mapping the Human Kinome in Response to DNA Damage. *Cell Rep.* **26**,
416 555-563.e6 (2019).
- 417 10. Wood, R. D., Mitchell, M. & Lindahl, T. Human DNA repair genes, 2005. *Mutat. Res.* **577**,
418 275–283 (2005).
- 419 11. Sherwood, R. I. *et al.* Discovery of directional and nondirectional pioneer transcription
420 factors by modeling DNase profile magnitude and shape. *Nat Biotechnol* **32**, 171–8 (2014).
- 421 12. Lin, L. *et al.* Comprehensive Mapping of Key Regulatory Networks that Drive Oncogene

- 422 Expression. *Cell Rep.* **33**, 108426 (2020).
- 423 13. Trojan, J. *et al.* Functional analysis of hMLH1 variants and HNPCC-related mutations using
424 a human expression system. *Gastroenterology* **122**, 211–219 (2002).
- 425 14. Papadopoulos, N. *et al.* Mutation of a mutL homolog in hereditary colon cancer. *Science*
426 **263**, 1625 (1994).
- 427 15. Xu, G. J. *et al.* Systematic autoantigen analysis identifies a distinct subtype of scleroderma
428 with coincident cancer. *Proc. Natl. Acad. Sci.* **113**, E7526 (2016).
- 429 16. Knijnenburg, T. A. *et al.* Genomic and Molecular Landscape of DNA Damage Repair
430 Deficiency across The Cancer Genome Atlas. *Cell Rep.* **23**, 239-254.e6 (2018).
- 431 17. Friedberg, E. C. *et al.* *DNA Repair and Mutagenesis*. (ASM Press, 2005).
432 doi:10.1128/9781555816704.
- 433 18. Shen, M. W. *et al.* Predictable and precise template-free CRISPR editing of pathogenic
434 variants. *Nature* **563**, 646–651 (2018).

Figure 1

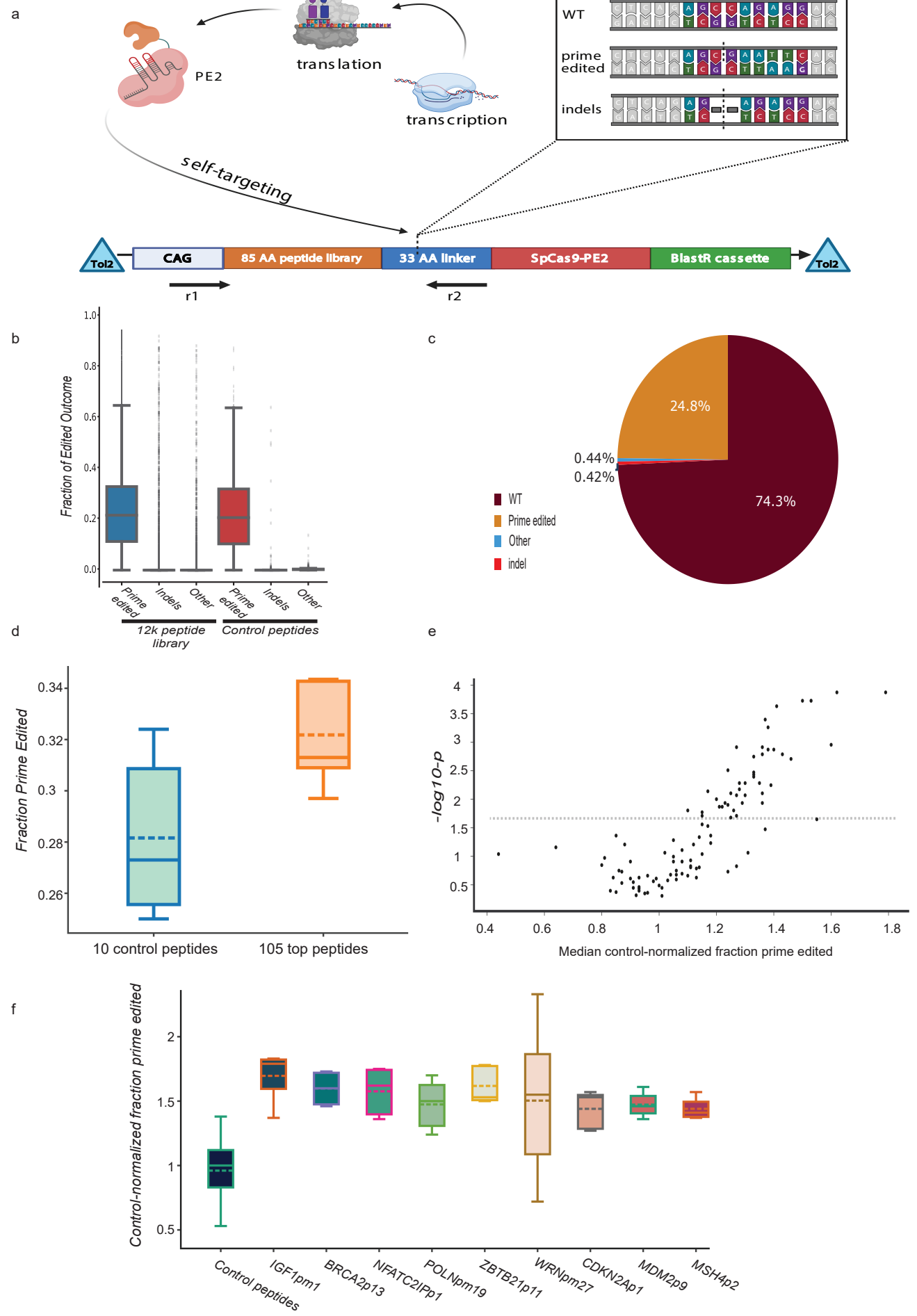
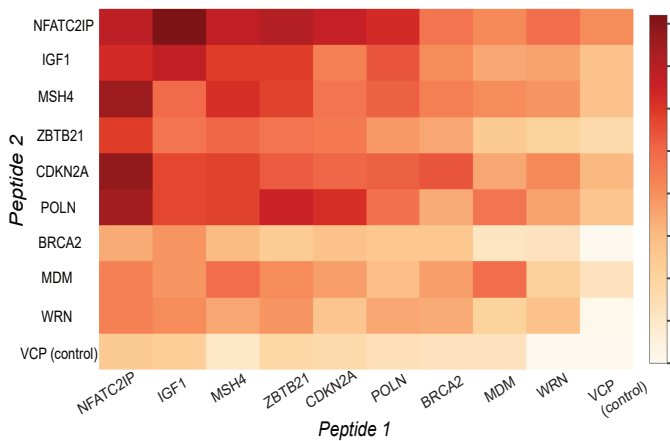


Figure 2

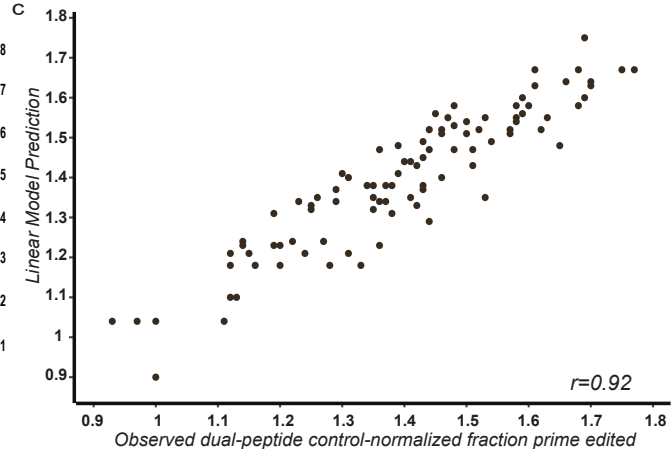
a



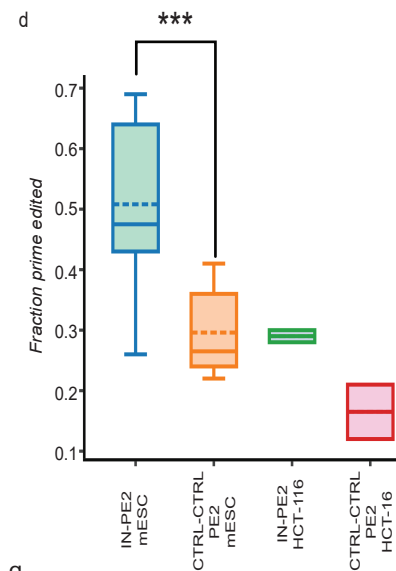
b



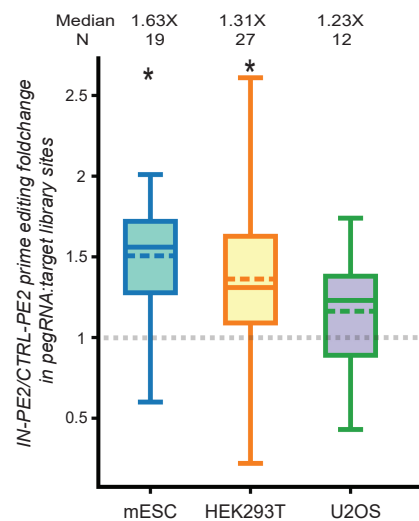
c



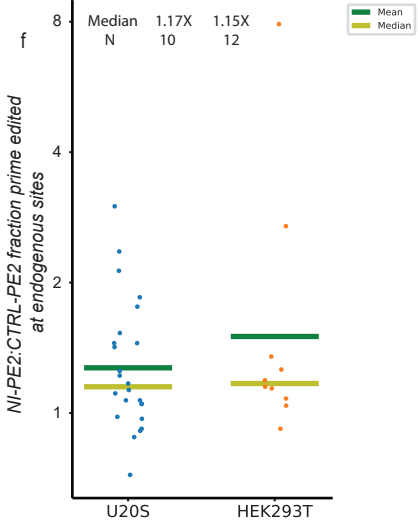
d



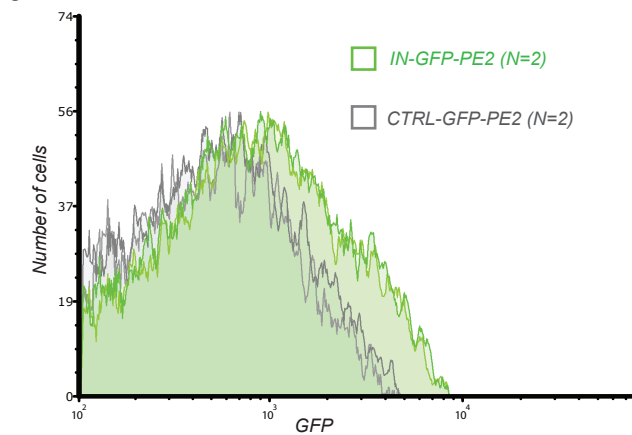
e



f



g



h

