Manuscript           Online article           Confidential

# Massively parallel quantification of CRISPR editing in cells by TRAP-seq enables better design of Cas9, ABE, CBE gRNAs of high efficiency and accuracy

**Authors**

Xi Xiang [1-4] *, Kunli Qu [1, 5] *, Xue Liang [1, 5] *, Xiaoguang Pan [1, 5] *, Jun Wang [1-3], Peng Han [1,2,5], Zhanying Dong [1], Lijun Liu [1,5], Jiayan Zhong [6], Tao Ma [6], Yiqing Wang [1], Jiaying Yu[1,2], Xiaoying Zhao [1,2], Siyuan Li [1,2], Zhe Xu [1,2], Jinbao Wang [6], Xiuqing Zhang [2,3], Hui Jiang [6], Fengping Xu [1,3], Lijin Zou [7], Huajing Teng [8], Xin Liu [3], Xun Xu [3, 9], Jian Wang [3], Huanming Yang [3, 10], Lars Bolund [1,3,4,5], George M. Church [11], Lin Lin [1, 4,12,‡] & Yonglun Luo [1,3,4,5,12,‡,#]

**Affiliations**

[1]Lars Bolund Institute of Regenerative Medicine, BGI-Qingdao, Qingdao 266555, China.

[2]BGI Education Center, University of Chinese Academy of Sciences, Shenzhen 518083, China.

[3]BGI-Shenzhen, Shenzhen 518083, China.

[4]Department of Biomedicine, Aarhus University, Aarhus 8000, Denmark.

[5]Qingdao-Europe Advanced Institute for Life Sciences, BGI-Shenzhen, Qingdao 266555, China

[6]MGI, BGI-Shenzhen, Shenzhen 518083, China.

[7]The First Affiliated Hospital of Nanchang University, Nanchang, Jiangxi 330006, P.R. China

[8]Key Laboratory of Carcinogenesis and Translational Research, Department of Radiation Oncology, Peking University Cancer Hospital & Institute, Beijing, China

[9]Guangdong Provincial Key Laboratory of Genome Read and Write, BGI-Shenzhen, Shenzhen, 518120，China

[10]Guangdong Provincial Academician Workstation of BGI Synthetic Genomics, BGI-Shenzhen, Shenzhen, 518120，China

[11] Department of Genetics, Blavatnik Institute, Harvard Medical School, Boston, MA 02115, USA

[12] Steno Diabetes Center Aarhus, Aarhus University, Aarhus 8200, Denmark

# lead contact

*These authors contribute equally to the study and should be regarded as co-first authors.

‡Correspondence addressed to: Lin Lin: lin.lin@biomed.au.dk or Yonglun Luo: alun@biomed.au.dk

**Abstract**

The CRISPR RNA-guided endonucleases Cas9, and Cas9-derived adenine/cytosine base editors (ABE/CBE), have been used in both research and therapeutic applications. However, broader use of this gene editing toolbox is hampered by the great variability of efficiency among different target sites. Here we present TRAP-seq, a versatile and scalable approach in which the CRISPR gRNA expression cassette and the corresponding surrogate site are captured by **T**argeted **R**eporter **A**nchored **P**ositional **Seq**uencing in cells. TRAP-seq can faithfully recapitulate the CRISPR gene editing outcomes introduced to the corresponding endogenous genome site and most importantly enables massively parallel quantification of CRISPR gene editing in cells. We demonstrate the utility of this technology for high-throughput quantification of SpCas9 editing efficiency and indel outcomes for 12,000 gRNAs in human embryonic kidney cells. Using this approach, we also showed that TRAP-seq enables high throughput quantification of both ABE and CBE efficiency at 12,000 sites in cells. This rich amount of ABE/CBE outcome data enable us to reveal several novel nucleotide features (e.g. preference of flanking bases, nucleotide motifs, STOP recoding types) affecting base editing efficiency, as well as designing improved machine learning-based prediction tools for designing SpCas9, ABE and CBE gRNAs of high efficiency and accuracy (>70%). We have integrated all the 12,000 CRISPR gene editing outcomes for SpCas9, ABE and CBE into a CRISPR-centered portal: The Human CRISPR Atlas. This study extends our knowledge on CRISPR gene and base editing, and will facilitate the application and development of CRISPR in both research and therapy.

**Keywords**

Gene editing, A-to-G base editing, C-to-T base editing, Genome engineering, System biology

## INTRODUCTION

56

57 Clustered Regularly Interspaced Short Palindromic Repeats (CRISPR) and CRISPR-associated
58 protein 9 (Cas9) are essential adaptive immune components in most bacteria. The system has
59 successfully been harnessed for programmable RNA-guided genome editing in prokaryotes, humans
60 and many other living organisms [1-5]. The *Streptococcus pyogenes* Cas9 (SpCas9) is the most
61 extensively studied and broadly applied Cas9 protein, amongst other Cas9 orthologs (e.g. SaCas9,
62 StCas9, NmCas9) [6-9] and Cas proteins (e.g. Cas12a, Cas13) [10, 11]. Guided by a programmable
63 small RNA molecule (also known as gRNA), the SpCas9 protein introduces a double-stranded DNA
64 break (DSB) to the DNA target site, which constitutes a complementary protospacer sequences and
65 a canonical protospacer adjacent motif (PAM) [2]. The classical CRISPR gene editing is achieved by
66 reparation of the DSBs in living organisms by the endogenous DNA repair mechanisms,
67 predominantly by the NHEJ and MMEJ pathways in mammalian cells. This process generates indels
68 (deletions or insertions) to the repaired site [12]. It is thus essential to have data from CRISPR editing
69 in cells to develop accurate prediction rule sets of CRISPR activity.

70

71 The simplicity of the CRISPR system, the flexibility for modifying the Cas9 protein, and the
72 increasing efforts from CRISPR scientists and pharmaceutical companies have extensively broadened
73 the CRISPR-Cas9-based gene editing toolkits. We are now enabled to epigenetically perturb
74 endogenous gene expression [13, 14], fluorescently label endogenous DNA elements[15] and site-
75 specifically edit single nucleotides [16-20]. The CRISPR base editors, which comprise two major
76 classes: adenine base editors (ABE) and cytosine base editors (CBE), have increasingly evolved as
77 attracting tools for gene editing. These base editors are created by fusing a catalytically dead Cas9
78 (dCas9) or Cas9 nickase (nCas9) to either an adenine deaminase or a cytidine deaminase [18, 19].
79 Without introducing double stranded DNA breaks, the ABE and CBE base editors, respectively, can
80 efficiently create an A to G (or T to C on the complementary strand) and C to T (or G to A on the
81 complementary strand) substitution within a small editing window of the target site [16, 21-23].
82 Albeit all these fantastic developments and applications of the CRISPR-Cas9 gene editing theme,
83 there is still an urgent need of methods and high throughput data on the Cas9-induced DBS repair
84 outcomes, as well as ABE and CBE efficiencies, to ensure a successful CRISPR gene editing outcome.
85 Such cataloged data of Cas9 and base editor efficiencies will allow the selection of experimentally
86 validated gRNAs, as well as for developing better rules for *in silico* Cas9, ABE, and CBE gRNAs
87 design.

88

89　Quantification of gRNA activity at the endogenous sites in cells is limited by scale. *In vitro*

90　approaches (in a test-tube) can overcome the scale but fails to recapitulate the effects of genome and

91　epigenome architectures and cellular DNA repair mechanisms on CRISPR editing [24, 25]. Methods

92　based on integrating synthetically barcoded DNA constructions were developed for large-scale

93　measuring of Cas9-induced DSB repair outcomes of gRNAs in cells [26-28]. Currently, we lack

94　large-scale ABE and CBE editing data for developing better rules for designing base editing gRNAs.

95　In this study, we developed an assay system for massively parallel quantification of a large-scale

96　CRISPR gRNAs activities in human cells. We optimized the design and procedures for generation

97　and in-cell CRISPR editing of synthetically barcoded DNA constructs. Each construct contains a

98　unique gRNA expression cassette and the corresponding surrogate target site. Using this method,

99　Targeted Reporter Anchored Positional Sequencing (hereafter referred as **TRAP-seq**), we

100　demonstrated the applicability of TRAP-seq for massively parallel quantification of the SpCas9-

101　induced DSB repair outcomes, ABE and CBE efficiency and profiles for 12,000 gRNAs in human

102　embryonic kidney cells.

103

104

105    **RESULTS**
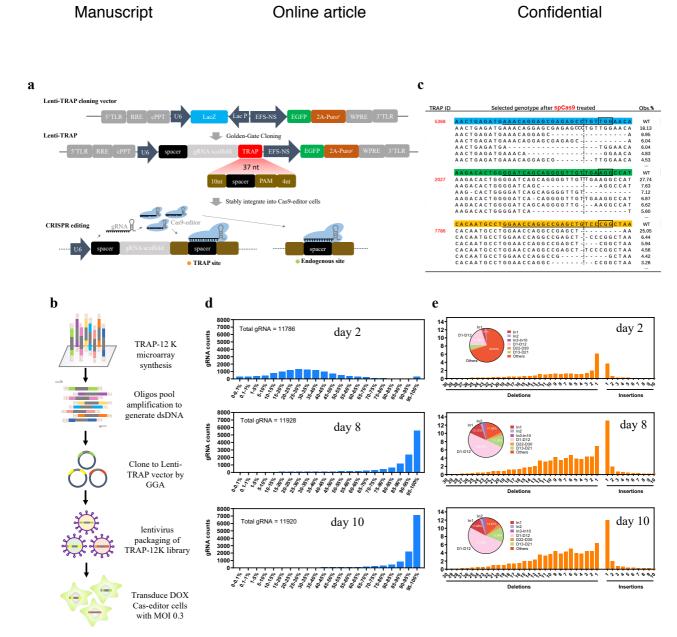
106    **Design and functional validation of the lentivirus-based TRAP-seq system**

107    To streamline vector cloning, gRNA expression and delivery into cells, we firstly designed a

108    lentivirus-based system with four main features: (1) A human U6 promoter; (2) Golden-Gate

109    Assembly (GGA) based cloning with a *lac Z* marker for precise and efficient insertion of an

110    expression cassette; (3) A green fluorescent protein (GFP) marker for measuring viral titer and real-

111    time tracking of viral delivery; (4) A puromycin selection gene for enrichment of stably transduced

112    cells (**Fig. 1a and S1**). Essentially, this lentivirus system allows conventional GGA-based insertion

113    of a synthetic DNA containing a gRNA spacer, scaffold and the corresponding surrogate target site

114    after the U6 promoter. As current microarray-based method can only faithfully synthesize oligo pools

115    of max 170 bp, we optimized the DNA design to contain a 102bp gRNA expression cassette (20bp

116    spacer + 82bp scaffold) and a 37bp surrogate target site, flanked by a 31bp GGA cloning site and

117    PCR handles (**Fig. 1a and S1**). We and several other groups previously demonstrated that such a

118    surrogate target site can faithfully recapitulate the endogenous editing efficiency and indel profile [27,

119    29]. To further validate the 37 bp surrogate target site, we firstly generated HEK293T cells expressing

120    a doxycycline (Dox)-inducible SpCas9 [1], an adenine base editor (ABE7.10) [19] or a cytosine base

121    editor (CBE, BE4-Gam) [30]. Next, we performed ICE-based analysis of three different sites (*AAVS1,*

122    *INHCB, TYMP*) in the HEK293T-SpCas9, HEK293T-ABE and HEK293T-CBE cells. The results

123    validated that the CRISPR editing efficiency and outcomes from the surrogate sites were closely

124    correlated (*$r^2$ = 0.96 – 0.99*) with those from the endogenous genome sites (**Fig. S2**). For

125    simplification, we hereafter named the system as Targeted Reporter Anchored Positional Sequencing

126    (TRAP-seq), the 170bp synthetic oligo/DNA as TRAP oligo/DNA, and the 37bp surrogate target site

127    as TRAP site.

Manuscript      Online article      Confidential



128

**Figure 1 High Throughput Quantification of SpCas9 efficiency in cells by TRAP-seq**

**a**. Schematic illustration of the TRAP-seq system. TLR, long terminal repeat; RRE, Rev Response Element; cPPT, central polypurine tract; U6, human U6 promoter; EFS-NS, short EFS promoter derived from EF1a. Figure not drown to scale.

**b**. Schematic illustration of oligo pool synthesis, PCR amplification, gold-gate assembly, lentivirus packaging, and transduction of the 12K TRAP-seq library.

**c**. Representative quantification of top 5 indel types for 3 TRAP sites. Dash line indicates the DSB site. Results for the 12,000 TRAP sites can be found at the CRISPR atlas database.

**d**. Bar plots of SpCas9 editing efficiency of all TRAP sites measured by targeted amplicon sequencing. Results are shown for Dox-induced HEK293T-SpCas9 cells from 2, 8 and 10 days after transduction. Corresponding results for Dox-free HEK293T-SpCas9 are shown in Fig. S10.

**e**. Bar plots of indel profiles (1-20 bp deletion, 1-10 bp insertion) for all TRAP sites introduced by SpCas9 in the Dox-induced HEK293T-SpCas9 cells from at 2, 8, and 10 days post transduction. Pie chat quantified the proportion of major indel groups: 1bp insertion (ins), 2bp insertion, 3-10 bp insertion, 1-12 bp deletion, 13-21bp deletion and 22-30 bp

143 deletion. Other indels and wild-type reads are presented together as "others". Corresponding results for Dox-free

144 HEK293T-SpCas9 are shown in Fig. S11 and S12.

145

146 **Generation of 12K TRAP-seq lentiviral library**

147 We next generated a 12K TRAP-seq library comprising 12,000 TRAP oligos by microarray synthesis

148 (**Fig. 1b**). The library targets 3,834 human protein-coding genes (**Table S1**) [31]. The gRNA spacers

149 were selected from the iSTOP database [32]. Out optimized workflow (also seen in methods) for PCR

150 amplification of the 12K TRAP-seq oligos and cloning into the lentivirus-based TRAP-seq vector

151 system is illustrated in **Fig. S3a.** A serial of optimizations in PCR conditions, GGA reactions and

152 lentiviral packaging were carried out to avoid unspecific amplification (**Fig. S3b, c**), maximize

153 successful ligation (**Fig. S4**) and properly quantify viral titer by FACS (**Fig. S5**), respectively.

154 To analyze the coverage of each TRAP oligo in our 12K TRAP-seq library, as well as to assess if the

155 whole procedure of vector cloning and lentivirus packaging/transduction affected the overall TRAP

156 representation, we performed targeted amplicon sequencing of the TRAP DNA in the 12K TRAP-

157 seq oligo library, GGA plasmid DNA and wildtype HEK293T cells transduced with the 12K TRAP-

158 seq lentivirus with a multiplexity of infection (MOI) of 0.3. With a constant sequencing depth (>

159 1,000X), all 12,000 TRAP oligos were detected in the 12K TRAP-seq library and the majority of

160 TRAP oligos (> 90%) were evenly distributed (**Fig. S6a**). Most importantly, over 99% of the TRAP

161 oligos were present in the 12K TRAP-seq plasmids and lentivirus preparation with high correlation

162 of representation for each TRAP oligo ($r^2 = 0.86$-$0.91$, **Fig. S6b**), suggesting that our optimized PCR,

163 GGA, lentivirus packaging and transduction methods faithfully retained the complexity of the 12K

164 TRAP-seq library without causing dramatic over/under-representation of the TRAP oligos.

165

166 **Quantification of SpCas9 editing at 12,000 sites by TRAP-seq**

167 To demonstrate applicability of the 12K TRAP-seq lentivirus library, we firstly investigated

168 massively parallel quantification of SpCas9 editing activity at all 12,000 TRAP sites. As

169 schematically shown in **Fig. S7**, we transduced the Dox inducible HEK293T-SpCas9 cells with the

170 12K TRAP-seq lentivirus (MOI = 0.3 and transduction coverage = 4,690 cells per TRAP). Puromycin

171 selection and Dox addition started two days after transduction to achieve maximum transduction and

172 gene editing efficiency (**Fig. S8**). To enable comparison and identification of CRISPR-introduced

173 indels, we also transduced wildtype HEK293T cells with the 12K TRAP-seq lentivirus with same

174 MOI and transduction coverage. We harvested genomic DNA from the transduced cells at three time

175 points: 2, 8, and 10 days after transduction (**Fig. S7 and Fig. S8**). Targeted PCRs were performed

176   with a pair of universal primers specifically amplifying the TRAP DNA, followed by targeted deep

177   sequencing with a DNA Nanoball sequencing technology [33]. With a constant sequencing depth

178   (**Fig S9**), the proportional representation of each TRAP correlated better in the Dox-free HEK293T-

179   SpCas9 cells (**r = 0.95**) than that in the Dox-induced HEK293T-SpCas9 cells (**R = 0.88**). Similar to

180   CRISPR knockout screening pool libraries [34, 35], our results suggested that there existed similar

181   cell fitness-related enrichment and depletion of the gRNAs in the 12K TRAP-seq library.

182

183   To measure the SpCas9 editing outcome, we firstly filtered out indels commonly found in both WT

184   and SpCas9 HEK293T cells (also see methods), which were introduced by oligo synthesis or PCRs.

185   We next analyzed the editing frequency and indel profiles for all 12,000 TRAP sites (**Fig. 1c, Table**

186   **S2,** also see CRISPR Atlas below). Although the SpCas9 expression was Dox inducible, significant

187   editing efficiencies were detected for all gRNAs in the HEK293T-SpCas9 cells at 2, 8 and 10 days

188   after transduction in Dox-free medium (**Fig. S10**), suggesting a substantial leakiness of SpCas9

189   expression. As expected, significantly higher editing efficiencies were achieved for all 12,000 gRNAs

190   in Dox-addition HEK293T-SpCas9 cells at 8 and 10 days after transduction (**Fig. 1d**). These results

191   support the notion that SpCas9 expression level and cultivation time affect gene editing efficiency

192   [36]. We and others had demonstrated that the indel outcomes introduced by SpCas9 comprises

193   mainly small deletions and insertions [37-39]. The distribution of indel profiles (deletion of 1-30 bp

194   and insertion of 1-10 bp) of the 12K TRAP sites were thus assessed in the transduced HEK293T-

195   SpCas9 cells. Two days after transduction, deletion or insertion of 1 bp were the two most frequent

196   indel types in cells (**Fig. 1e, S11**). Following increased editing time (Dox-free groups, **Fig. S10**) and

197   SpCas9 expression (Dox-induced groups, **Fig. 1e**), the frequency of other indel types rose

198   significantly and 1 bp insertion was the most dominant indel type which is in agreement with previous

199   findings (**Fig. 1e, S12**) [37]. With the indel outcome, we were capable of analyzing the mutation

200   consequence of all indels on protein translation. More than 70% of the total indels led to out-of-frame

201   genotypes (**Fig. S13**). In conclusion, we demonstrated that TRAP-seq is a simple method for

202   massively parallel quantification of gRNA editing outcomes in cells.

203

204   **Characterization of nucleotide features affecting SpCas9 efficiency and indel outcomes**

205   Development of more accurate rules for *in silico* CRISPR design relies heavily on datasets of gRNA

206   activity and indel from a large number of gRNAs. The rich gRNA activity and indels profile data

207   generated by TRAP-seq above were valuable for further improving the performance of CRISPR

208    design. We sought to investigate if the gRNA activity and indel outcomes measured by the TRAP-

209    seq can mirror previous findings about the effects of nucleotide features on CRISPR activity.

210    Nucleotide features such as secondary structure [24] and GC content [40] of the guide sequences

211    affect CRISPR editing efficiency. We analyzed the correlation between gRNA activity and GC

212    content and secondary structure (deltaG energy) in the gRNA spacer. Our TRAP-seq results further

213    confirmed that the gRNA spacer GC content (**Fig. 2a, S14**) and secondary structure (**Fig. 2b, S15**)

214    affected SpCas9 gene editing efficiency in cells. The optimal GC content and deltaG energy is [50-

215    70%] GC and [-5; -1] KJ/mol, respectively. Consistent with the previous finding [41], our TRAP-seq

216    results revealed that the SpCas9 disfavors motifs of "TT" and "GCC" at the N17-N20 region (**Fig.**

217    **2c, S16**). Recent reports have discovered that indel profiles for a given gRNA is predictable [27, 28].

218    The SpCas9 predominantly generates blunt-end double-stranded DNA breaks (DSB) between the

219    N17 and N18 nucleotide preceding the protospacer adjacent motif (PAM), which are most frequently

220    repaired by the NHEJ and MMEJ pathways in mammalian cells [42]. We compared the indel profiles

221    of approximately 12,000 gRNAs revealed by TRAP-seq to the predicted indel profiles by inDelphi,

222    a machine learning program for SpCas9 indel prediction [27]. Our results show that the overall indel

223    profiles (in the Dox-free or Dox-induced cells at day 8 and 10) are highly correlated with the predicted

224    ones by inDelphi (**median r = 0.51-0.65, Fig. 2d and Fig. S17**). Of note, the overall correlation

225    between the TRAP indels in transduced cell at day 2 and indel profiles predicted by inDelphi was

226    much lower (**median r = 0.31, Fig. S17**), suggesting that the indel outcome also depends on the

227    experimental conditions (e.g. Cas9 expression level, editing time etc.).

228

229    Among all indels, the 1bp insertion between N17 and N18 was the most abundant type (**Fig. 1e**).

230    Previous studies had discovered that 1bp insertion is not random [27]. We therefor asked whether the

231    nucleotides of 1bp insertion among our 12,000 TRAP sites followed the same principle. First, we

232    quantified the frequency of inserted adenine (A), thymine (T), cytosine (C), guanine (G) among all

233    1bp insertions. The results show that the 1bp insertion favor T and disfavor G (**Fig. 2e**). For a given

234    TRAP site, however, there is a strong preference of one nucleotide type (also seen **the CRISPR Atlas**

235    **resource**). Next, we divided all 12,000 TRAP sites into four groups based on the N17 or N18

236    nucleotide and quantified the frequency of inserted bases among all 1bp insertions. Our results show

237    that the N17 nucleotide strongly defines the inserted base (**Fig. 2f, S18**), but less extensively affected

238    by the N18 nucleotide (**Fig. S19**). Lastly, we divided all the 12,000 TRAP sites into 16 groups based

239    on N17N18 sequence motifs and compared the indel frequency of 1bp insertion versus deletions.

240    Despite a constantly higher frequency of deletions over insertions, motifs of GN (N = A, T, C, or G)

241    and MC (M = A or C) at N17N18 favor deletions over 1bp insertion, as compared to TN and MG

242    motifs (**Fig. 2g, S20**). Taken together, we show that high throughput TRAP-seq enables the

243    identification and validation of features affecting SpCas9 editing efficiency and indel outcomes. The

244    results corroborate previous findings that SpCas9 editing outcomes are predictable in cells.

245



246

247    **Figure 2 Characterization of features affecting SpCas9 efficiency and indel outcomes in cells**

248    **a**. Box plot between GC content (with an interval of 10%) and SpCas9 gRNA efficiency measured in Dox-induced cells

249    from Day 10. Results for other groups are shown in Fig. S14.

250    **b.** Box plot between deltaG energy (with an interval of 2) and SpCas9 gRNA efficiency measured in Dox-induced cells

251    from Day 10. Results for other groups are shown in Fig. S15.

252    **c.** Comparison of SpCas9 efficiency between gRNAs harboring the GCC or TT motif in N17-N20 seed region and

253    gRNAs without these two motifs. Data are shown for Dox-induced cells from Day 10. Results for other groups are

254    shown in Fig. S16. "****", p value less than 0.0001.

255    **d**. Correlation between TRAP-seq indels and indels predicted by inDelphi from 11,910 sites. Data are shown for Dox-

256    induced cells from Day 10. Results for other groups are shown in Fig. S17.

257    **e**. Pie chart of the proportion of 1bp insertion among four bases (A, T, C, G)

258    **f**. Correlation between the inserted 1 base and the nucleotide at N17 position. Data are shown for Dox-induced cells

259    from Day 10. Results for other groups are shown in Fig. S18.

260    **g**. Effects of N17N18 dinucleotide motifs on the indel frequency of 1bp insertion and deletions (1-30 bp). The gRNAs

261    are divided into 16 groups based on the N17N18 motifs. For each gRNA, the total indel frequencies of 1-30bp deletions

262    and 1bp insertion were analyzed. "n" indicates the number of gRNAs included for each group.

263    **h.** Comparison of SpCas9 gRNA efficacy predictions in a regression schema for various datasets and prediction models.

264    **i**. Top 20 features that weighted the most for the GNL machine learning model. Results are shown as the SHAP

265    (SHapley Additive exPlanations) values. The 30mere comprises 4bp upstream, 20bp protospacer, 3 bp PAM, and 3 bp

266    downstream sequences. Machine learning was based on gRNA efficiency data from Dox-induced cells at Day 10.

267

## An improved machine learning model to predict SpCas9 efficiency

269    To further streamline the prediction of SpCas9 efficiency and the identification of nucleotide features

270    important for gRNA activity, we randomly selected 80% of the 12K TRAP-seq gRNA efficiency and

271    trained the GNL-Scorer [43] - a machine learning algorism that we previously developed based on

272    the Bayesian Ridge regression (BRR) model and 2485 features. Our results showed that the GNL-

273    scorer trained with the TRAP-seq dataset (GNL-Scorer (Trap)) gave an accuracy prediction score of

274    over 70% (**Fig. S21**). To benchmark the performance of the TRAP-seq dataset and the GNL

275    prediction algorithm, we compared our dataset and GNL-Scorer (Trap) with seven previously

276    published datasets (from HEL, NB4, TF1, MOLMB, A375, Hela, HEK293T, HCT116 cells) and five

277    prediction tools (DeepCas9, Azimuth-2.0, TUSCAN, CRISPRater, SSC). Our results showed that the

278    GNL-Scorer (Trap) achieved the best accuracy score in five datasets (second best for the remaining

279    3 datasets) and have the best generalized prediction outcome across all test datasets (**Fig. 2h**). Using

280    the SHapley Additive exPlanations (SHAP) algorithm for explaining the feature output, our results

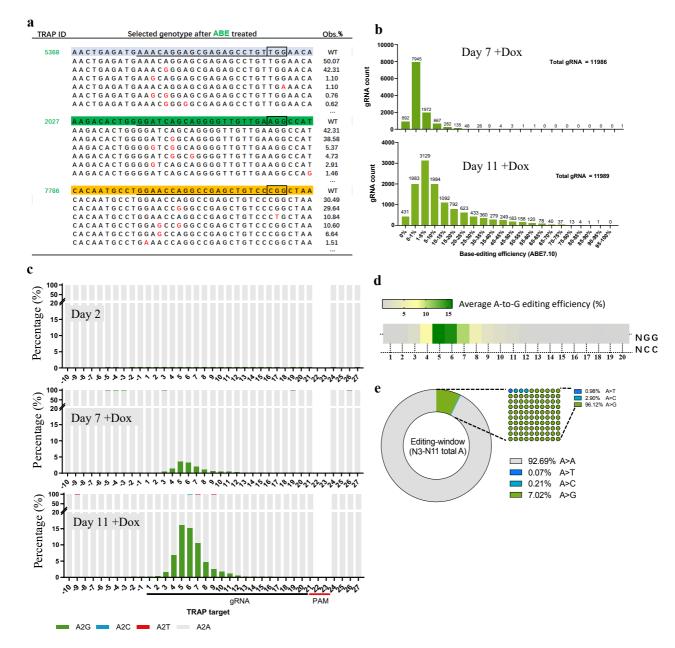281    further revealed features (such as melting temperature, GC content delta G energy, sequences motifs

282    etc.) that are important for the performance of our prediction model and SpCas9 efficiency (**Fig. 2i**).

283    Our results taken together suggest the TRAP-seq dataset based GNL-Scorer performs generally well

284    for gRNA knockout efficiency prediction that the 12K SpCas9 gRNA efficiency dataset revealed by

285    TRAP-seq enable better understanding of features affecting gRNA efficiency in cells, improve the

286    design of gRNAs of high knockout efficiency. This improved GNL-scorer algorithm for predicting

287    SpCas9 efficiency has been deposited to and available at public domain GitHub.

288

289    **Quantification of CRISPR-mediated adenine base editing at 12,000 sites by TRAP-seq**

290    Unlike SpCas9 gene editing, we still lack large-scale data of CRISPR adenine base editing (ABE)

291    efficiency. Such valuable data would enable us to develop better *in silico* ABE gRNA design tools.

292    Since the TRAP site could confidentially recapitulate the ABE editing outcome of the corresponding

293    endogenous site (**Fig. S2**), we sought to investigate the ABE editing outcomes in all 12,000 TRAP

294    sites using the 12K TRAP-seq library. Although all the 12,000 gRNAs were not specifically designed

295    for ABE editing, we reasoned that this "randomly" selected gRNA library would enable us to

296    unbiasedly identify rules affecting ABE editing efficiency. To do this, we firstly transduced

297    HEK293T-ABE cells with the 12K TRAP-seq lentivirus (MOI=0.3), and performed targeted

298    amplicon sequencing of the TRAP DNAs from cells at 2, 7, and 11 days after transduction and

299    cultured in Dox-free or Dox-addition medium (**Fig. S22**). For additional controls, we also performed

300    similar experiments with wild-type HEK293T cells, with constant transduction coverage (4,690 cells

301    per TRAP) and sequencing depth (approximately 1,000 reads per TRAP) (**Fig. S23**).

## Figure 3 Quantification of ABE efficiency at 12,000 sites by TRAP-seq

**a.** Representation of top 5 adenine editing outcomes for three TRAP sites. Full ABE results for all 12,000 sites can be found at the CRISPR atlas resource.

**b.** Quantification of overall ABE efficiency for all gRNAs from Dox-induced HEK293T-ABE cells 7- and 11-days post transduction. Other groups are presented in Fig. S23.

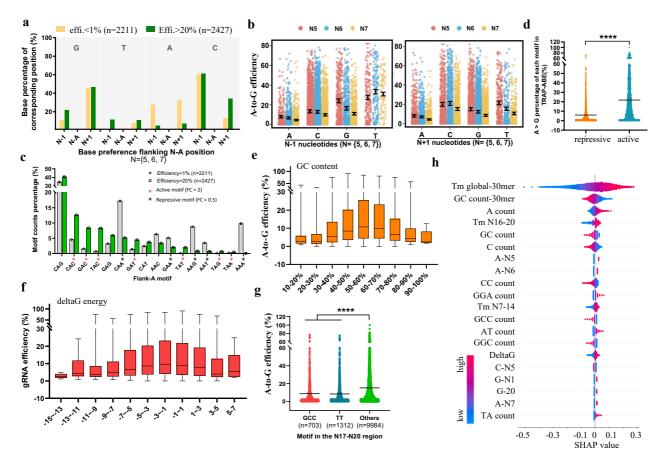**c**. Quantification of overall percentage of A-to-G, A-to-T, A-to-C, and A-to-A (unedited) events across the 37bp region of all TRAP sites. Results are from Dox-induced HEK293T-ABE cells at 2, 7- and 11-days post transduction.

**d.** Heatmap quantification of the overall A-to-G editing efficiency within the 20nt protospacer region.

**e.** Summary of substitution of Adenines within the ABE editing window N3-N11.

313　Next, we quantified the edited adenine events for each 37bp TRAP site in the 12K TRAP-seq library

314　(**Fig. 3a, Table S3**). Our results showed that substantial editing events (within an editing window

315　from N3 to N11) appeared 7 days after transduction in the HEK293T-ABE with Dox induction. The

316　editing efficiency increased 4-5 folds when extending the cultivation time to 11 days (**Fig. 3b, S24**).

317　Most importantly, the ABE editing window remained constant between N3 and N11 (**Fig. 3c, S25**),

318　supporting the notion that the ABE base editor is highly conserved with respect to its editing region

319　[19]. Our TRAP-seq results also validated that the highest ABE editing efficiency was observed for

320　adenines located at N5, N6 and N7 (**Fig. 3d**). Lastly, quantification of all 42,790 edited adenine sites

321　revealed that the ABE base editor conservatively generated A-to-G substitution (96.12%), and a small

322　proportion of A-to-C (2.9%) and A-to-T (0.98%) substitutions (**Fig. 3e**). Although additional

323　experiments will be required to test the TRAP-seq library in more cell lines, these initial studies

324　suggest that the TRAP-seq is a highly valuable method for massively parallel quantification of

325　CRISPR adenine base editing efficiency in cells.

326

327　**Characterization of nucleotide features affecting ABE efficiency**

328　We also sought to characterize features that affect ABE efficiency. To enable comparisons, we first

329　selected two groups of gRNAs based on ABE efficiency: (1) high efficiency ABE gRNAs (n = 2,331,

330　at least one edited adenine site had an efficiency over 20% with the protospacer region N1-N20) and

331　(2) low efficiency ABE gRNAs (n = 2,589, the efficiency of any edited adenine site within the

332　protospacer N1-N20 lower than 1%). Next, we compared the base percentage between the low and

333　high efficiency gRNAs across the 37bp TRAP region for the low and high efficiency ABE gRNAs.

334　Not surprisingly, high efficiency gRNAs showed overrepresentation of Adenine within the editing

335　window N3-N7 (**Fig. S26**). Interestingly, high efficiency gRNAs favored Guanine over Thymine in

336　the seed region (N17 to N20), the distal protospacer region (N-1 to N4) and the N21 base of the PAM.

337　The presence of Cytosine in N20 was disfavored for high efficiency gRNAs.

338

## Figure 4 Characterization of nucleotide features affecting ABE efficiency

**a.** Proportion of frank-A (A bases located at the core editing window N5-N7) bases (A, T, C, G) between high (edited A efficiency > 20% for at least one A base within N5-N7) and low (edited A efficiency < 1% for any A base within N5-N7) efficiency ABE gRNAs. Results are based on the data from the Dox-induced HEK293T-ABE cells from day 11.

**b.** A-to-G editing efficiency for A bases located at N5-N7, grouped based on the flanking bases. Results are based on the data from the Dox-induced HEK293T-ABE cells from day 11.

**c.** Comparison of the frequency of N**A**N (N = A, T, C, G) trinucleotide motifs between high and low efficiency ABE gRNAs (bold A referred to the deaminated bases within the N5-N7 core editing window). Red and black asterisks indicate the active and repressive motifs based on a cutoff of two-fold difference.

**d.** Scatter plot of edited A efficiency between sites within the active and repressive motifs. The number of "n" indicates number of sites. "****", p value less than 0.0001.

**e.** Box plot of overall A-to-G editing efficiency for all TRAP gRNAs according to the gRNA spacer GC content (with 10% interval). Results are based on the data from the Dox-induced HEK293T-ABE cells from day 11. Complementing results for other groups can be found in Fig. S26.

**f.** Box plot of overall A-to-G editing efficiency for all TRAP gRNAs according to the gRNA spacer deltaG energy (with an interval of 2). Results are based on the data from the Dox-induced HEK293T-ABE cells from day 11. Complementing results for other groups can be found in Fig. S27.

**g**. Dot plot of overall ABE efficiency between gRNAs have the GCC or TT motifs within the seed region N17-20 and gRNAs without these two motifs. "n" indicates the number of gRNAs within each group. Results were based on the

15

358  data from the Dox-induced HEK293T-ABE cells from day 11. Complementing results for other groups can be found in

359  Fig. S28. "****", p value less than 0.0001.

360  **h**. Top 20 features that weighted the most based on the GNL machine learning model that affect the overall ABE

361  efficiency in cells. Results are shown as the SHAP (SHapley Additive exPlanations) values. The 30mere comprises 4bp

362  upstream, 20bp protospacer, 3 bp PAM, and 3 bp downstream sequences. Complementary SHAP results for each edited

363  Adenine site within the N3-N11 window were shown in Fig. S30. Results are based on the ABE editing data from the

364  Dox-induced HEK293T-ABE cells from day 11.

365

366  We next sought to investigate the effect of intrinsic nucleotide preference on base editing efficiency.

367  To identify the intrinsic nucleotide preference for ABE7.10, we focused on the ABE efficiency at N5,

368  N6 and N7, which were the three highest ABE sites (**Fig. 3d**). Using a similar strategy to enable

369  comparison, we first selected two groups of gRNAs: N5-N7 high efficiency (n = 2427, at least one

370  edited adenine site had an efficiency over 20% within the core editing window N5-N7) and N5-N7

371  low efficiency (n = 2211, the efficiency of any edited adenine site within core editing window N5-

372  N7 was lower than 1%). Next, we analyzed the preference of flanking bases (A, T, C, G) at N5-N7

373  for the low and high efficiency gRNAs. Our results revealed that the high efficiency ABE gRNAs

374  favored upstream keto (K) bases (G, T) and downstream pyrimidine (Y) bases (T, C) (**Fig. 4a**). The

375  presence of flanking adenine was strikingly overrepresented in the low efficiency ABE gRNAs (Fig.

376  **4a**). To validate this observation, we analyze the correlation between A-to-G editing efficiency of all

377  12,000 TRAP sites in N5-N7 and their flanking bases. Consistent with the previous finding, the

378  average ABE editing efficiency is higher if the edited adenine is flanked by Keto bases upstream and

379  pyrimidine bases downstream (**Fig. 4b**), as compared to edited sites flanked by amino bases (A, C)

380  upstream and purine bases (A, G) downstream. The overall ABE efficiency is much lower (approx.

381  2 to 7 folds) if the flanked base is adenine. Based on these observations, we further compared the

382  frequency of tri-nucleotide flank-A motifs between the low and high ABE efficiency gRNAs. Using

383  a cutoff of two folds, we categorized seven (B**A**C, K**A**T, T**A**R) and five (**A**AD, S**A**A) motifs (bold **A**

384  refers to the deaminated adenine) as active and repressive flank-A motifs, respectively (**Fig. 4c**). To

385  further validate that, we assessed A-to-G editing efficiency between all active or repressive flank-A

386  motifs within our 12K TRAP sites. Our results showed that the A-to-G editing efficiency is

387  significantly higher (p < 0.0001, fold change = 4) in the active motifs as compared to the repressive

388  ones (**Fig. 4d, S27**).

389

390    Since the ABE editor shares principles of gRNA-guided DNA binding with SpCas9 nickase, we

391    reasoned that many of the features (such as GC content, gRNA secondary structure (deltaG energy),

392    N17-N20 motifs) that were known to influence SpCas9 editing efficiency should also affect ABE

393    efficiency. To address that, we performed Pearson correlation analysis between the ABE efficiency

394    and GC content and the deltaG energy of the guide sequence. Our results demonstrated that both GC

395    content (**Fig. 4e, S28**) and secondary structure affect ABE efficiency (**Fig. 4f, S29**). We also

396    investigated and validated that the presence of TT and GCC motifs at the N17-N20 seed region

397    negatively affects ABE efficiency (**Fig. 4e, S30**).

398

399    **An improved machine learning model to predict ABE efficiency**

400    We sought to apply our GNL-scorer machine learning model [43] to systematically identify features

401    of importance for ABE efficiency (GNL-scorer_ABE) and ABE efficiency prediction rules, as well

402    as developed a new machine learning-based tool for predicting ABE efficiency. We randomly

403    selected 80% (the remaining 20% used for model evaluation) of the ABE efficiency data and trained

404    the Bayesian Ridge regression (BRR)-based GNL-scorer model with 2485 features (**Fig. S31a**). The

405    ABE efficiency prediction was performed for both A-to-G edited site efficiency in N3-N11 and the

406    overall ABE efficiency. Our results showed that the accuracy of predicting ABE is above 60% for

407    the core ABE editing window (**Fig. S31b**), and the accuracy of predicting the over ABE efficiency is

408    approximately 70% (**Fig. S21**). The SHAP algorithm-based feature outputs further demonstrated that

409    our machine learning results consistently revealed that features such as the global melting temperature,

410    GC content, deltaG energy, and nucleotide compositions (such as the presence of Adenine in N5-N7,

411    Guanine in N20) greatly affect the ABE efficiency (**Fig. 4f and S32**). Collectively, we demonstrate

412    that the rich ABE editing efficiency data revealed by TRAP-seq enable us to systematically define

413    factors influencing ABE efficiency and improve ABE gRNA design for future studies.

414

415    **Quantification of CBE-mediated recoding efficiency at 11,979 sites by TRAP-seq**

416    After demonstrating that the TRAP-seq method is versatile for massively parallel quantification of

417    SpCas9 and ABE efficiency, we sought to test the performance of TRAP-seq for high throughput

418    quantification of CBE efficiency. As mentioned earlier, all the gRNA spacers of our 12K TRAP-seq

419    library were retrieved from the iSTOP database. This allows us to address the STOP recoding

420    efficiency of all 11,979 gRNA and 3,832 genes by CBE in cells (**Fig. 5a**). First, based the optimized

421    and constant conditions of lentivirus library transduction, we performed 10 parallel TRAP-seq-library
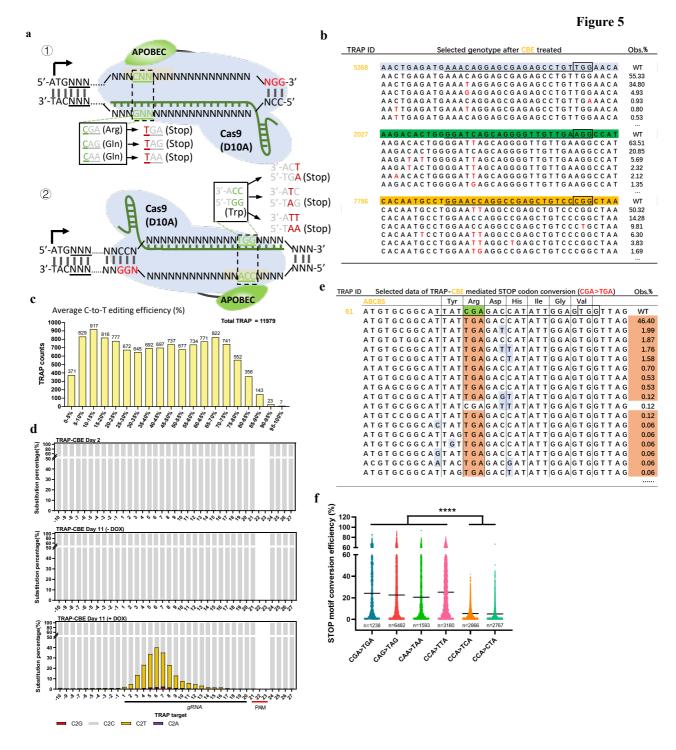
422   based CRISPR CBE editing experiments in the Dox-inducible HEK293T-CBE cells (**Fig. S33**).

423   Constant sequencing depths (1,000X coverage) and TRAP representation (**r = 0.96-0.97**) were

424   achieved for Dox-free and Dox-addition HEK293T-CBE cells 2 and 11 days after transduction by

425   targeted amplicon sequencing of the TRAP region (**Fig. S34**). Next, to assess CBE editing, we

426   quantified the efficiency of C-to-T edit, as well as C-to-G and C-to-A edits, of all C bases within the

427   37bp of all 11,979 TRAP sites (**Fig. 5b**, **Table S4**, full results were shown in the CRISPR Atlas

428   database). As seen from the overall (**Fig. 5c, S35**) and C-to-T site (**Fig. 5d**) CBE efficiency plots,

429   quantification of the average efficiency of all edited Cs for the 11979 gRNAs revealed that there was

430   an even distribution of gRNAs editing efficiency from 5% to 75%. There were 371 gRNAs with very

431   low CBE efficiency (0-5%) (**Fig. 5c**). The C-to-T CBE efficiency was primarily detected in the Dox-

432   induced HEK293T-CBE cells at 11 days after transduction, indicating there was very minor leakiness

433   of CBE editor expression (**Fig. 5c, S35**). Of particular note, compared to ABE, the editing window

434   of CBE was broader (N1 to N16) but the highest cytosine editing efficiency was similarly found at

435   N6 (**Fig. 5d, S36a**). Quantification of all edited Cs revealed that the majority (93.18%) were C-to-T

436   edits, however, low frequency of unspecific C-to-G (4.54%) and C-to-A (2.27%) edits were also

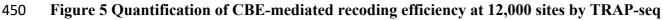437   observed for the CBE (BE4-gam) base editor (**Fig. S36b**).

438

439   We next investigated the stop-codon recoding efficiency by CBE in cells. With the current setups in

440   HEK293T-CBE cells, the median STOP efficiency of all 11,979 gRNA was approximately 22% (**Fig.**

441   **5e, S37, S38a**). A total of 3,481 genes (90%) were successfully knocked out by stop-codon recoding

442   with CBE (**Fig. 38b**). We also sought to analyze the efficiency of the 6 types of recoding into stop

443   codons (**Fig. 5a**). As seen in **Fig. 5f**, the recoding efficiencies of C**C**A-to-C**T**A, C**C**A-to-T**C**A and

444   C**C**A-to-T**T**A were significantly higher (**p value < 0.0001**) than **C**GA-to-TGA, **C**AG-to-TAG and

445   **C**AA-to-TAA (bold **C** refers to the deaminated cytosines). Taken together, we here demonstrated that

446   the TRAP-seq technology enables massively parallel quantification of CBE-mediated recoding

447   capacity in cells.

448

**Figure 5**

449

**Figure 5 Quantification of CBE-mediated recoding efficiency at 12,000 sites by TRAP-seq**

**a.** Schematic illustration of the 6 stop-codon recoding schemes by CBE. Cas9 (D10A) is the Cas9 nickase used in the

testing CBE editor. The stop-codon recoding scheme is draw based on sense and anti-sense strands. "ATG", translation

start site; "NGG", PAM of SpCas9.

**b.** Representation of top 5 cytosine base editing outcomes for three TRAP sites. Full CBE frequency results can be

found in the CRISPR Atlas database.

456    **c.** Quantification of overall CBE efficiency for all gRNAs in Dox-induced HEK293T-CBE cells at 11 days. Other

457    groups are presented in Fig. S33.

458    **d**. Quantification of overall percentage of C-to-T, C-to-G, C-to-A, and C-to-C (unedited) events across the 37bp region

459    of all TRAP sites. Results are based on HEK293T-CBE cells from 2 days post transduction, and transduced HEK293T-

460    CBE cells from 11 days cultured in Dox-free or Dox-addition medium.

461    **e.** Representation of CBE-mediated recoding efficiency at one TRAP site measured by TRAP-seq. Complementary

462    results referred to Fig. S35 and CRISPR atlas database.

463    **f**. Comparison of C-to-T efficiency between sites within the 6 stop-codon recoding types. "n" indicates the number of

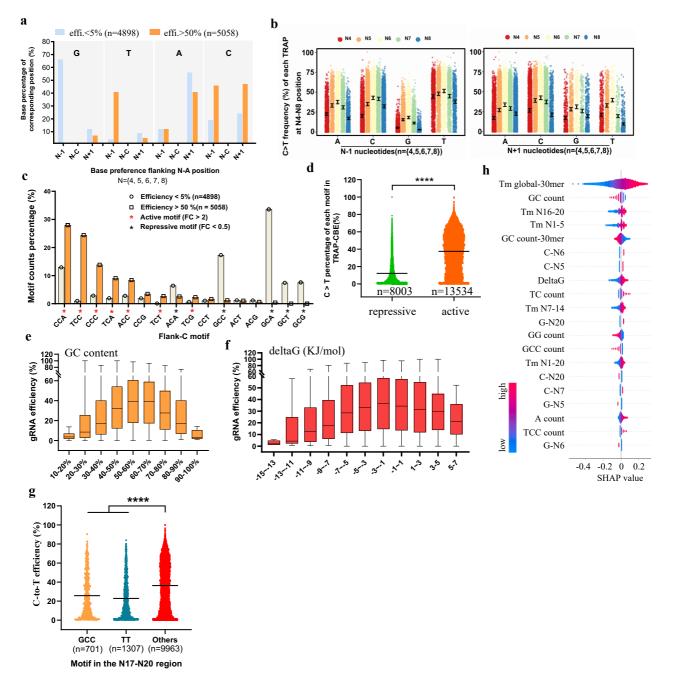464    sites within each group. "****", p value less than 0.0001.

465

466    **Characterization of nucleotide features affecting CBE efficiency**

467    The large-scale CBE efficiency data from 56,887 edited Cs and 11,979 gRNAs enable us to

468    characterize the features that affect CBE efficiency in cells. To simplify the analysis and

469    characterization, we first selected two opposing types of CBE gRNAs based on their efficiencies: low

470    efficiency (n = 1,844, the efficiency of any edited C within the protospacer is < 5%) and high

471    efficiency (n = 1,731, At least one edited C efficiency within the protospacer is > 60%). Next, we

472    compared the base preference across the 37bp TRAP regions between the two types of gRNAs. Our

473    results clearly show that the high efficiency CBE gRNAs favor the presence of Cytosine (N4-N8),

474    but disfavor Guanine (N3-N7) and Adenine (N5-N6) within the core editing window (**Fig. S39**). The

475    presence of Thymine within the seed region (N8-N20) of protospacers was underrepresented in the

476    high efficiency CBE gRNAs, which was in contrast to the N3-N6 region. Similar to our findings in

477    ABE, high efficiency CBE gRNAs favor the presence of Guanine at proximal PAM region (N19-

478    N21) and disfavored Cytosine at N20.

479

480    Since the core editing window of CBE is N4-N8 (**Fig. S36**), we focused on the deaminated Cytosines

481    within N4-N8 when further analyzing the effect of flanking bases on CBE efficiency. To enable

482    comparison, we selected C-to-T edited sites of low (< 5%, N = 4,898) and high (> 50%, N = 5,058)

483    efficiency within the N4-N8 region. Our results show that the presence of Thymine upstream is

484    strikingly overrepresented in the highly efficient C-to-T editing, whereas the presence of Guanine

485    upstream is only present in the low efficiency CBE sites (**Fig. 6a**). In addition, the highly efficient

486    CBE sites are less frequently flanked by Adenine downstream and more frequently flanked by

487    Cytosine, as compared to low efficiency CBE sites (**Fig. 6a**). To validate this observation, we

488    calculated the C-to-T editing efficiency at N4-N8 for all 12,000 TRAP sites. Consistent with previous

489    observations, the overall efficiency of CBE sites flanked by Thymine upstream is approximately two-

490    fold higher than with other flanking bases. Of note, sites flanked by Guanine upstream (as well as

491    downstream, but to less extent) show much lower CBE efficiency (fold changes = 2 - 12 folds) (**Fig.**

492    **6b**). This provides highly valuable knowledge for designing CBE gRNAs with better editing outcome.



493

494    **Figure 6 Characterization of nucleotide features affecting CBE efficiency in cells by TRAP-**

495    **seq**

496    **a.** Proportion of frank-C (N4-N8) bases (A, T, C, G) between high (at least one edited C base within N4-N8 was >50%

497    edited) and low (edited C efficiency for any C base within N4-N8 was less than 1%) efficiency CBE gRNAs. Results

498    are based on the data from the Dox-induced HEK293T-CBE cells from day 11.

499　　**b.** C-to-T editing efficiency for C bases located at N4-N8, grouped based on the flanking bases. Results are based on the
500　　data from the Dox-induced HEK293T-CBE cells from day 11.
501　　**c.** Comparison of the presence of the 16 NCN trinucleotide motifs between high and low efficiency CBE gRNAs (bold
502　　C referred to the deaminated cytosine within the N4-N8 core editing window). Red and black asterisks indicate the
503　　active and repressive motifs based on a cutoff of two-fold difference.
504　　**d.** Scatter plot of edited C efficiency between sites within the active and repressive motifs. The number of "n" indicates
505　　number of sites. "****", p value less than 0.0001.
506　　**e.** Box plot of overall C-to-T editing efficiency for all TRAP gRNAs according to the gRNA spacer GC content (with
507　　10% interval). Results are based on the data from the Dox-induced HEK293T-CBE cells from day 11. Complementing
508　　results for other groups can be found in Fig. S36.
509　　**f.** Box plot of overall C-to-T editing efficiency for all TRAP gRNAs according to the gRNA spacer deltaG energy
510　　(with an interval of 2). Results are based on the data from the Dox-induced HEK293T-CBE cells from day 11.
511　　Complementing results for other groups can be found in Fig. S38.
512　　**g**. Scatter plot of overall CBE efficiency between gRNAs have the GCC or TT motifs within the seed region N17-20
513　　and gRNAs without these two motifs. "n" indicates the number of gRNAs within each group. Results were based on the
514　　data from the Dox-induced HEK293T-CBE cells from day 11. Complementing results for other groups can be found in
515　　Fig. S41.
516　　**h**. Top 20 features that weighted the most based on the GNL machine learning model, which affect the overall CBE
517　　efficiency in cells. Results are shown as the SHAP (SHapley Additive exPlanations) values. The 30mere comprises 4bp
518　　upstream, 20bp protospacer, 3 bp PAM, and 3 bp downstream sequences. Complementary SHAP results for each edited
519　　cytosine site within the protospacer N1-N20 were shown in Fig. S42. Results are based on the CBE editing data from
520　　the Dox-induced HEK293T-CBE cells from day 11.

521

522　　Since the flanking bases played an important role on CBE efficiency, we reasoned that there exists a
523　　preference of tri-nucleotide flank-C motifs (NCN; N=A, T, C, G; bold C refers to the deaminated
524　　cytosine located at N4, N5, N6, N7 or N8) for active or repressive CBE. To identify these motifs, we
525　　compared the frequency of all 16 NCN motifs between the low (N=4898, efficiency < 5%) and high
526　　(N=5058, efficiency > 50%) efficiency cytosine editing sites in N4-N8. Based on a two-fold
527　　difference, we identified seven (TCN, CCM and ACC; N=A, T, C, G; M = A, C) and five (ACA and
528　　GCN) as active and repressive motifs, respectively (**Fig. 6c**). To validate this, we further compared
529　　the C-to-T editing efficiency between the active and repressive motifs for all 21537 edited C sites in
530　　our 12K TRAP-seq library. The C-to-T editing efficiency of the active motifs were significantly
531　　higher (fold change = 3, p < 0.0001) than those within the repressive motifs (**Fig. 6d, S40**). We also
532　　sought to investigate if the CBE efficiency shared features with SpCas9 editing efficiency using our
533　　CBE 12K TRAP-seq data. Thus, we analyzed the correlation between CBE efficiency and the GC
534　　content, the gRNA spacer secondary structure, as well as the proximal PAM motifs. Our results show

535   that, similar to SpCas9 and ABE, the CBE efficiency is affected by the gRNA spacer GC content

536   (**Fig. 6e, S41**) and secondary structure (**Fig. 6f, S42**). The CBE efficiency of gRNAs is significantly

537   (**p < 0.0001**) lower with TT or GGC motifs at the proximal (17-N20) PAM region (**Fig. 6g, S43**).

538

539   **An improved machine learning model to predict CBE efficiency**

540   We further took advantage of our GNL-scorer machine learning model to development a prediction

541   tool and systematically evaluate the effect of 2485 features on CBE efficiency. Based on randomly

542   selecting 80% of 12K gRNAs CBE efficiency for model training and 20% for model evaluation, our

543   results showed that the accuracy prediction score of CBE efficiency by the GNL machine learning

544   model reaches approximately 80% (**Fig. S21**). Apart from predicting the overall CBE efficiency, our

545   machine learning-based prediction tool provides highly precise predicting outcome of the site-

546   specific CBE efficiency within the editing window (**Fig. S44**). Consistently, our machine learning

547   results further showed that features such as melting temperatures, GC content, deltaG energy,

548   nucleotide composition (e.g. the presence of cytosine at N5-N7, nucleotide counts, motifs) greatly

549   affect CBE efficiency (**Fig. 6h and Fig S45**). Finally, the machine learning-based CBE efficiency

550   prediction algorism (GNL-scorer_CBE) has also been deposited to the GitHub to facilitate the design

551   of CBE gRNAs of high efficiency. Taken together, we hereby demonstrate that high-throughput

552   quantification of CBE efficiency by TRAP-seq enables the better understanding and design of highly

553   efficient CBE gRNAs in cells.

554

555   **The CRISPR Atlas**

556   As part of this work, a human CRISPR atlas database (http://www.crispratlas.com/crispr) has been

557   launched to present and integrate all the SpCas9, ABE and CBE efficiency and editing outcomes for

558   12,000 gRNAs in HEK293T cells. This CRISPR atlas is presented with a gene-centered and gRNA-

559   centered summary of the overall efficiency of gRNAs and editing outcomes. For SpCas9, total

560   efficiency, indel (1-30bp deletions, 1-10bp insertions) profiles were presented for each gRNA. For

561   ABE and CBE, the overall gRNA efficiency and graphical presentation of base substitution efficiency

562   across the 37bp TRAP region were shown for all 12,000 TRAP sites. We believe that the CRISPR

563   atlas database generated by this study will complement the existing CRISPR resources in gRNA

564   design [44], efficiency prediction [45, 46] and indel prediction [27], thus streamlining the application

565   of CRISPR in functional studies.

566

**DISCUSSION**

In conclusion, the work described here demonstrates the broad applicability of the TRAP-seq system for massively parallel quantification of SpCas9, ABE and CBE efficiency in human cells. Recent studies published by other groups have demonstrated that the surrogate target sites can well mimic the SpCas9 indel outcomes at the corresponding endogenous sites, and thus predict the SpCas9 indel profile for a given gRNA in cells [26, 28, 40]. Consistent with that, we demonstrate corroborating findings with our TRAP-seq method and further expand the data of SpCas9 knockout efficiency and indel profiles with 12,000 gRNAs. This will aid the improvement of CRISPR gRNA design for gene knockout purposes with machine learning models [40, 45]. With such a large amount of SpCas9 efficiency data, it is possible to systematically identify both previously known as well as novel features that affect CRISPR gene editing efficiency. Importantly, according to our knowledge, this is the first time that both ABE and CBE efficiencies are measured at such a large scale in cells. Based on the 12K ABE and CBE TRAP-seq data, our analyses identify several novel features (such as the preference of flanking bases, active/repressive tri-nucleotide motifs) that strongly influence ABE and CBE efficiency in cells, respectively. We believe that incorporating the nucleotide features of importance for ABE and CBE efficiency from this study will improve the performance of in silico base editing designers such as BE-Designer [47] and Beditor [48]. However, we acknowledge that there might be a difference in the DNA repair machinery between different cell types and organisms, which will potentially affect the SpCas9, ABE and CBE efficiency and outcome. Additional experiments will be required to test the SpCas9, ABE and CBE efficiency with TRAP-seq in more cell lines in the future.

The concept of using surrogate target sites to capture the gene editing outcomes is highly attractive. We and other groups have generated dual-fluorescence-based surrogate systems for rapid evaluation of ZFNs, TALENs and CRISPR-Cas9 activity in cells [29, 49, 50]. The DSBs generated by CRISPR-Cas9 were predominantly repaired by the NHEJ and MMEJ pathways, which will lead to the introduction of small indels at the DBS site. However, large deletions or chromosomal rearrangements have also been reported in CRISPR editing as outcomes of repaired mediated by e.g. HDR or SSA in cells [51, 52]. The TRAP-seq system developed in this study is based on a 37bp surrogate target site. Thus, SpCas9 editing outcomes such as large deletions or chromosomal arrangements will not be captured by our method. However, for ABE and CBE, the editing outcomes would not be affected by such a size-related problem.

599

600 Earlier, we have discovered that chromatin accessibility at the editing sites affect CRISPR gene
601 editing efficiency [24]. Since the TRAP-seq library were randomly inserted in the genome of the
602 targeted cells, the chromatin accessibility state of the surrogate site might be different from the
603 endogenous target site. It would be interesting to apply the TRAP-seq system to systematically
604 analyze the epigenetic factors (e.g. DNA methylation, chromatic accessibility) on ABE, CBE
605 efficiency in future studies.

606

607 In this study, we demonstrate the TRAP-seq system with applications in massively parallel
608 quantification of editing efficiency for SpCas9 and base editors (ABE and CBE) derived from the
609 SpCas9. The current CRISPR-based gene editing toolbox has been greatly expanded with the
610 engineered SpCas9 (e.g. xCas9, eSpCas9, SpCas9-HF1), the SpCas9 orthologs (e.g. SaCas9, StCas9,
611 NmCas9) and other Cas proteins (e.g. Cas12a) [53-55]. However, features affecting the editing
612 efficiency and indel outcomes are still rarely explored for most of these Cas proteins, which will limit
613 the applications of this great toolbox. We believe that the TRAP-seq will become an important
614 technology for the whole CRISPR gene editing society to better understand how CRISPR gene editing
615 works in cells. The CRISPR atlas database generated by this study will become a CRISPR-centered
616 portal, in which we provide experimentally validated gRNAs for CRISPR gene editing. Taken
617 together, the TRAP-seq technology, the SpCas9/ABE/CBE efficiency of 12,000 gRNAs, and the
618 CRISPR atlas database will enable us to better functionally understand how CRISPR works in cells
619 and improve CRISPR in both research, therapeutic and drug discovery applications.

620

621

622 **MATERIALs and METHODS**

623 **Vector construction**

624 The empty pLenti-TRAP-seq vector backbone (shown in **Fig. S1**) was generated by a serial of cloning.

625 Briefly, we replaced the SpCas9 open reading frame (ORF) in pLentiCRISPRv2-puro (Addgene

626 plasmid # 98290) plasmid with an enhance green fluorescence protein (EGFP) ORF based on *XbaI*

627 and *BmHI* digestion and T4 ligation. Next, we replaced the gRNA cassette in the EGFP-inserted

628 pLentiCRISPRv2-puro with a synthetic Golden-Gate Assembly cassette, and hence generated a

629 lentivirus-based vector (hereafter referred as pLenti-TRAP-seq) allowing the insertion of TRAP DNA

630 to the Golden-Gate cloning site by GGA. The full sequencing of the pLenti-TRAP-seq vector can be

631 downloaded from our CRISPR atlas database ([www.crispratlas.com/crispr](www.crispratlas.com/crispr)). Original plasmid stock

632 can be acquired from the corresponding authors' lab.

633

634 The doxycycline inducible SpCas9, ABE, and CBE vectors were generated by subcloning. Briefly,

635 based on a PiggyBac transposon system (full GenBank vector sequences can be found at the CRISPR

636 atlas website), which consists of an all-in-one expression system: (1) An expression cassette of a TRE

637 promote-driven protein expression cassette with multiple cloning sites (MCS). (2) An expression

638 cassette of a consecutive promoter-driven Tetracycline-Controlled Transcriptional Activation and

639 hygromycin. The coding sequences of SpCas9 (Addgene plasmid ID # 41815), ABE 7.10 (Addgene

640 plasmid ID # 102919), and CBE (Addgene plasmid ID # 100806) were PCR amplified and inserted

641 to the MCS of the PiggyBac transposon system. All vectors were validated by Sanger sequencing.

642

643 **TRAP 12K oligos design and microarray synthesis.**

644 A typical TRAP oligo consists of the BsmBI recognition site "**cgtctc**" with 4 bp specific nucleotides

645 "acca" upstream, following the GGA cloning linker "aCACC", one bp "g" for initiating transcription,

646 then the 20 bp gRNA sequences of "gN20", 82bp gRNA scaffold sequence, 37 bp surrogate target

647 sequences (10bp upstream sequences, 23 bp protospacer and PAM sequences, 4 bp downstream

648 sequence), the downstream linker "GTTTg" and another BsmBI binding site and its downstream

649 flanking sequences "acgg". An example of the typical TRAP oligo sequence was shown below:

650 "acca**cgtctc**aCACCg⟦GTCCCCTCCACCCCACAGTG⟧GTTTTAGAGCTAGAAATAGCAAGTTA

651 AAATAAGGCTAGTCCGTTATCAACTTGAAAAAGTGGCACCGAGTCGGTGCTTTTTT<u>ACT</u>

652 <u>TTTATCTGTCCCCTCCACCCCACAGTGGGGCCAC</u>GTTTg**gagacg**acgg". Sequences in the

653   black frame is the 20 bp gRNA spacer. The underline sequence is the 37 bp surrogate target site,

654   which is termed as "TRAP target" shortly in this study.

655

656   For the 12K TRAP oligo design, we used bioinformatic tools to automatically generate the 12K TRAP

657   oligo pools. Briefly, 1) we selected approximately 7,000 genes from the a drugable gene database

658   (http://dgidb.org); 2) Discard all the exons which the DNA length was less than 23 bp with filtering;

659   3) Select the first three coding exons of each gene. If the exons number is less than 3, keep all the

660   exons; 4) Extract all the possible gRNA sequences (including the PAM sequence "NGG") in these

661   filtered exons sequence, analyzes and predictd the off-target sites of each gRNA using FlashFry

662   version 1.80 (https://github.com/mckennalab/FlashFry), discarded gRNAs with potential off-target

663   of 0-3 bp mismatches in human genome; 5) Rank each gRNA based on the number of off-target site

664   in an ascending order; 6) Map and extract the 10 bp upstream and 4 bp downstream flanking sequence

665   of each selected gRNA,  construct the TRAP target sequence as 10 bp upstream + 23 bp gRNA

666   (include PAM) + 4 bp downstream = 37 bp; 7) Filter out TRAPs with BsmBI recognition site, because

667   of GGA cloning; 8) Compared all the selected gRNAs with the database of CRISPR-iSTOP [56]; 9)

668   Construct the full length sequence of each TRAP, which is 170 bp; In total, the first 12K TRAP-seq

669   oligos contain 3832 genes and 12000 TRAPs were contained in the final TRAP-library. The 12K

670   oligo pools was synthesized in Genscript® (Nanjing, China).

671

672   **TRAP-12K plasmid library preparation**

673   First, the TRAP 12K oligos were cleaved and harvested from the microarray and diluted to 1 ng/$\mu$L.

674   Next, we performed PCR amplifications using the primers: TRAP-oligo (BsmBI GGA)-F: 5'-

675   TACAGCTaccacgtctcaCACC-3'; TRAP-oligo (BsmBI GGA)-R: 5'-AGCACAAccgtcgtctccAAAC-

676   3'.

677

678   The PCR reaction was carried out using PrimeSTAR HS DNA Polymerase (Takara, Japan) following

679   the manufacturer's instruction. Briefly, each PCR reaction contained 1 $\mu$L oligo template, 0.2 $\mu$L

680   PrimeSTAR polymerase, 1.6 $\mu$L dNTP mixture, 4 $\mu$L PrimeSTAR buffer, 1 $\mu$L forward primer (10

681   uM) and 1 $\mu$L reverse primer (10 uM) and ddH2O to a final volume of 20 $\mu$L.

682

683   The thermocycle program was 98°C 2min, (98°C/10s，55°C*/10s，72°C/30s) with 21 cycles, then

684   72°C for 7min and 4°C hold. To avoid amplification bias of oligos introduced by PCR, we conducted

685   gradient thermocycles and performed PCR products gray-intensity analysis to determine the optimal

686   PCR cycles of 21. The best thermocycles should be in the middle of an amplification curve. In this

687   study, the *PCR cycles* was 21 for oligos amplification. But for PCR amplification of TRAP from

688   cells integrated with TRAP lentivirus, the *PCR cycle* was 25. The final TRAP PCR product length

689   was 184 bp. We performed **72 ×** parallel PCR reactions for 12K oligos amplification, then these PCR

690   products were pooled and gel purified by 2% agarose gel.  1 $\mu$g purified PCR product were quantified

691   with PCR-free next generation sequencing (MGI Tech).

692

693   The PCR products of TRAP oligos were then used for Golden Gate Assembly (GGA) to generate the

694   TRAP 12K plasmids library. For each GGA reaction, the reaction mixture contained 100 ng pLenti-

695   TRAP-seq vector, 10 ng purified 12K TRAP oligos-PCR products, 1 $\mu$L T4 ligase (NEB), 2 $\mu$L T4

696   ligase buffer (NEB), 1 $\mu$L BsmBI restriction enzyme (ThermoFisher Scientific, FastDigestion) and

697   ddH2O to a final volume of 20 $\mu$L. The GGA reactions were performed at 37°C 5 min and 22°C 10

698   min for 10 cycles, then 37°C 30 min and 75°C 15 min. **36 ×** parallel GGA reactions were performed

699   and the ligation products were pooled into one tube.

700

701   Transformation was then carried out using home-made chemically competent DH5a cells. For each

702   reaction, 10 $\mu$L GGA ligation product was transformed in to 50 $\mu$L competent cells and all the

703   transformed cells were spread on one LB plate (15 cm dish in diameter) with Xgal, IPTG and Amp

704   selection. High ligation efficiency was determined by the presence of very few blue colonies (also

705   see **Fig. S4**). To ensure that there is sufficient coverage of each TRAP in the 12K TRAP-seq library,

706   **42 ×** parallel transformations were performed and all the bacterial colonies were scraped off and

707   pooled together for plasmids midi-prep. For NGS-based quality quantification of TRAP coverage,

708   midi-prep plasmids were used as DNA templates for TRAP PCR amplifications, followed by gel

709   purification and NGS sequencing.

710

711   **TRAP-12K lentivirus packaging.**

712   HEK293T cells were used for lentivirus package. All the cells were cultured in Dulbecco's modified

713   Eagle's medium (DMEM) (LONZA) supplemented with 10 % fetal bovine serum (FBS) (Gibco), 1%

714   GlutaMAX (Gibco), and penicillin/streptomycin (100 units penicillin and 0.1 mg streptomycin/mL)

715   (The culture medium was named as D10 shortly) in a 37 °C incubator with 5% $CO_2$ atmosphere and

716    maximum humidity. Cells were passaged every 2-3 days when the confluence was approximately 80-
717    90%.

718

719    For lentivirus packaging: **Day 1**: Wild-type HEK293T cells were seeded to a 10 cm culture dish, 4 ×
720    $10^6$ cells per dish (10 dishes in total); **Day 2**: Transfection. Briefly, we refreshed the medium with 7
721    mL fresh culture medium to 1 hour before transfection (be gently, as the HEK293T cells are easy to
722    be detached from the bottom of dish); Next, we performed transfection with the PEI 40000
723    transfection method. For 10 cm dish transfection, the DNA/PEI mixture contains 13 $\mu$g pLenti-
724    TRAPseq 12K vectors, 3 $\mu$g pRSV-REV, 3.75 $\mu$g pMD.2G, 13 $\mu$g pMDGP-Lg/p-RRE, 100 $\mu$L PEI
725    40000 solution (1 $\mu$g/ $\mu$L in sterilized ddH2O) and supplemented by serum-free opti MEM without
726    phenol red (Invitrogen) to a final volume of 1 mL. The transfection mixture was pipetted up and down
727    several times gently, then kept at room temperature (RT) for 20 min, then added into cells in a
728    dropwise manner and mix by swirling gently. **Day 3**: Changed to fresh medium; **Day 4**: Harvest and
729    filter all the culture medium of the 10 cm dish through a 0.45 $\mu$m filter, pool the filtered media into
730    one bottle. Each 10 cm dish generated approximately 7~8 mL lentivirus crude. Add polybrene
731    solution (Sigma-Aldrich) in to the crude virus to a final concentration of 8 $\mu$g/mL. Aliquot the crude
732    virus into 15 mL tubes (5 mL/tube) and store in -80 °C freezer.

733

734    **Lentivirus titer quantification by flow cytometry (FCM).**
735    As the pLenti-TRAP-seq vector expresses a EGFP gene, the functional titer of our lentivirus prep
736    was assayed by FCM as described previously [57]. Briefly, 1) **DAY 1**: split and seed HEK293T cells
737    to 24-well plate, 5 × $10^4$ cells per well. Generally, 18 wells were used to perform the titter detection,
738    a gradient volume of the crude lentivirus was added into the cells and each volume was tested by
739    replicates. In this experiment, the crude virus gradients were 2.5 $\mu$L, 5 $\mu$L, 10 $\mu$L, 20 $\mu$L, 40 $\mu$L, 80
740    $\mu$L, 160 $\mu$L and 320 $\mu$L for each well. Another 2 wells of cells were used for cell counting before
741    transduction; 2) **DAY 2**: Conduct lentivirus transduction when cells reach up to 60~80% confluence.
742    Before transduction, detach the last two wells of cells using 0.05% EDTA-Trypsin to determine the
743    total number of cells in one well ($N_{(initial)}$). Then change the remaining wells with fresh culture
744    medium containing 8 $\mu$g/mL polybrene, then add the gradient volume of crude virus into each well
745    and swirling gently to mix; 3) **DAY 3**: Change to fresh medium without polybrene; 4) **DAY 4**: Harvest
746    all the cells and wash them twice in PBS. Fix the cells in 4% formalin solution at RT for 20 min, then
747    spin down the cell pellet at 2,000 rpm for 5 min. Discard the supernatant and re-suspend the cell pellet

748    carefully in 600 $\mu$L PBS, and conduct FCM analysis immediately. FCM was performed using a BD

749    LSRFortessa™ cell analyzer with at least 30,000 events collected for each sample in replicates.

750

751    The FCM output data was analyzed by the software Flowjo vX.0.7. Percentage of GFP-positive cells

752    was calculated as: $\mathcal{Y}\% = N_{\text{(GFP-positive cells)}} / N_{\text{(total cells)}} \times 100\%$. Calculate the GFP percentage of all

753    samples. For accurate titter determination, there should be a linear relationship between the GFP

754    positive percentages and crude volume. The titter (Transducing Units (TU/mL) calculation according

755    to this formula: $\text{TU/mL} = (N_{\text{(initial)}} \times \mathcal{Y}\% \times 1000) / V$. V represents the crude volume ($\mu$L) used for

756    initial transduction.

757

758    **Generation of Doxycycline-inducible spCas9/ABE7.10/CBE stable cell lines**

759    TRE-spCas9, TRE-ABE7.10 and TRE-CBE stable cells were generated by PiggyBac transposon

760    systems. For stable cell lines establishment, HEK293T cells were transfected with pPB-TRE-spCas9-

761    Hygromycin (or pPB-ABE7.10-hygromycin, pPB-CBE-hygromycin) vector and pCMV-hybase with

762    a 9:1 ratio. Briefly, $1 \times 10^5$ HEK293T cells were seeded in 24-well plate and transfections were

763    conducted 24 h later using lipofectamine 2000 reagent following the manufacturer's instruction.

764    Briefly, 450 ng pPB-TRE-spCas9-Hygromycin (or pPB-ABE7.10-hygromycin, pPB-CBE-

765    hygromycin) vectors and 50 ng pCMV-hybase were mixed in 25 $\mu$L optiMEM (tube A), then 1.5 $\mu$L

766    lipofectamine 2000 reagent was added in another 25 $\mu$L optiMEM and mix gently (tube B). Incubate

767    tube A and B at RT for 5 min, then add solution A into B gently and allow the mixture incubating at

768    RT for 15 min. Add the AB mixture into cells evenly in a dropwise manner. Cells transfected with

769    pUC19 were acted as negative control. Culture medium was changed to selection medium with 50

770    $\mu$g/mL hygromycin 48h after transfection. Completion of selection took approximately 5-7 days until

771    the negative cells were all dead in the un-transfected cells. The cells were allowed to grow in 50

772    $\mu$g/mL hygromycin containing D10 medium for 3-5 days for further expansion. PCR-based

773    genotyping were carry out using the primers: ***spCas9-iden-F***: gacacctacgatgatgatctcg; ***spCas9-iden-***

774    ***R***: tggtgctcatcatagcgcttga; ***ABE7.10-iden-F***: 5'-cagtactcgtgctcaacaatcg-3'; ***ABE7.10-iden-R***: 5'-

775    ggcgttgcgaacaccgaataca-3'; ***BE4-iden-F***: 5'-ttcttcgatccgagagagctcc-3'; ***BE4-iden-R***: 5'-

776    ctgcaccttgtgttcggacag-3'.

777

778    For functional tests of the spCas9, ABE7.10 and CBE4 expression cells, individual TRAP constructs

779    packaged in lentivirus particles were transduced into the cells. Transduced cells were harvested for

780    indel analysis 6 days after transduction. Indel analysis were carried out for both the TRAP site and

781    the endogenous genome target sites.

782

783    **12K TRAP-seq library lentivirus transduction**

784    HEK293T-SpCas9, -ABE7.10 and -CBE4 cells were cultured in D10 medium with 50 $\mu$g/mL

785    hygromycin throughout the whole experiment. For 12K TRAP-seq library transduction, we followed

786    the procedures showed in **Fig. S7, S21, S31**. Briefly, 1) Day -1: $2.5 \times 10^6$ cells per 10 cm dish were

787    seeded, and 12 dishes in total. For each group, one dish was used for cell number determination before

788    transduction and one dish for drug-resistance (puromycin) test control and the remaining 10 dishes

789    were used for the 12K TRAP-seq lentivirus library transduction; 2) Day 0: We determined the

790    approximate cell number per dish by cell countering. This was used to determine the volume of crude

791    lentivirus used for transduction using a multiplicity of infection (MOI) of 0.3. The low MOI (0.3)

792    ensures that most infected cells receive only 1 copy of the lentivirus construct with high probability

793    [34]. The calculation formula is: V = N × 0.3 / TU. V = volume of lentivirus crude used for infection

794    (mL); N = cell number in the dish before infection; TU = the titter of lentivirus crude (IFU/mL). In

795    this study, take the TRE-ABE7.10 group for instance, there were $1.875 \times 10^7$ cells in one dish, the

796    TRAP 12K lentivirus crude titter = $3.8 \times 10^6$ IFU/mL. Thus V = $1.875 \times 10^7 \times 0.3 / 3.8 \times 10^6$ =

797    1.48 mL. The 12K TRAP-seq transduction coverage per dish is $1.875 \times 10^7 \times 0.3 / 12000 = 469 \times$.

798    As we performed 10 replicates for each group, the overall coverage would reach to about 4690 ×. In

799    this study, $V_{spCas9}$ = 1.26 mL, $V_{ABE7.10}$ = 1.48 mL, $V_{wt}$= 1.37 mL for each dish. For transduction, we

800    added aforementioned volume of crude virus to each group in a dropwise manner and mix by swirling

801    gently. The infected cells were cultured in a 37 ℃ incubator; 3) Day 1: 24 hours after transduction,

802    split the transduced cells of each dish to 3 dishes equally; 4) Day 2: For the 3 dishes of split (30 dishes

803    in total, 3 divided into sub-groups), sub-group 1 (10 dishes) were harvested and labeled as the Day 2

804    after 12K TRAP-seq transduction. All cells from this sub-group were pooled into one tube and stored

805    in -20 ℃ freezer for genomic DNA extraction; The sub-group 2 (10 dishes) was changed to fresh D10

806    medium contains 50 $\mu$g/mL hygromycin + 1 $\mu$g/mL puromycin (Dox-free group); The sub-group 3

807    (10 dishes) was changed to D10 medium contains 50 $\mu$g/mL hygromycin + 1 $\mu$g/mL puromycin + 1

808    $\mu$g/mL doxycycline (Dox-induction group). For the WT HEK293T cells (Group 3) screening,

809    hygromycin but not puromycin should be excluded from the culture medium; 5) The transduced cells

810    were spitted every 2~3 days when cell confluence reaches up to 90%. At the indicated time points in

811    Fig. S7, 21, 31, cells were harvested and stored in -20 ℃ freezer for further genomic DNA extraction.

812

**PCR amplicons of TRAPs from cells**

813

814    Genomic DNA was extracted using the phenol-chloroform method. The genomic DNA were digested

815    with RNase A (OMEGA) to remove RNA contamination (In this study, 50 $\mu$g RNase A worked well

816    to digest the RNA contamination in 100 ~ 200 $\mu$g genomic DNA after incubating in 37 ℃ for 30 min).

817    Then the genomic DNA was purified and subjected to PCR for amplification of the TRAP DNA. The

818    PCR primers were: **TRAP-NGS-F1**: 5'-GGACTATCATATGCTTACCGTA-3' and **TRAP-NGS-**

819    **R1**: 5'-ACTCCTTTCAAGACCTAGCTAG-3'. The PCR product length is 252 bp. In this study, 5

820    $\mu$g genomic DNA was used as temperate in one PCR reaction which contained approximately 7.6 $\times$

821    $10^5$ copies of TRAP construct (assuming $1 \times 10^6$ cells contain 6.6 $\mu$g genomic DNA), which covered

822    about 63 $\times$ coverage of the 12K TRAP-seq library. In total, 32 $\times$ parallel PCR reactions were

823    performed to achieve approximately 2,016 $\times$ coverage of each TRAP construct. For each PCR

824    reaction, briefly, 50 $\mu$L PCR reaction system consists of 5 $\mu$g genomic DNA, 0.5 $\mu$L PrimeSTAR

825    polymerase, 4 $\mu$L dNTP mixture, 10 $\mu$L PrimeSTAR buffer, 2.5 $\mu$L forward primer (10 uM) and 2.5

826    $\mu$L reverse primer (10 uM) and supplemented with ddH2O to a final volume of 50 $\mu$L. The

827    thermocycle program was 98°C 2min, (98°C for 10s，55°C for 10s，72°C for 30s) with 25 cycles,

828    then 72°C for 7min and 4°C hold. Then purify all the PCR products by 2% gel, pool the products

829    together and conduct deep amplicon sequencing.

830

**Deep amplicon sequencing**

831

832    MGISEQ-500 (MGI of BGI in China) was used to perform the amplicons deep sequencing following

833    the standard operation protocol. First, PCR-free library was prepared using MGIeasy FS PCR-free

834    DNA library Prep kit following the manufacturer's instruction. Briefly, measure the concentration of

835    purified PCR products using Qubit 4 ™ fluorometer (Invitrogen) and dilute the concentration of each

836    sample to 10 ng/ $\mu$L. 10 $\mu$L diluted PCR product was mixed with an A-Tailing reaction which

837    contained A-Tailing enzyme and buffer, incubated at 37°C for 30 minutes then 65°C for 15 min to

838    inactive the enzyme. Then the A-Tailed sample was mixed with PCR Free index adapters (MGI.), T4

839    DNA Ligase and T4 ligase buffer to add index adapter at both 3' and 5' ends of PCR products. The

840    reaction was incubated at 23°C for 30 min and then purified with XP beads. Then denature the PCR

841    products to be single-strand DNA (ssDNA) by incubating at 95 ℃ for 3 min and keep on 4 ℃ for the

842    subsequent step. Transform the ssDNA to be circles using cyclase (MGI) at 37 ℃ for 30 min and

843    then digested to remove linear DNA using Exo enzyme at 37 ℃ for 30 min. Purify the products again

844    by XP beads and assay the concentration of library by Qubit 4 ™ fluorometer. The amplicons libraries

845    were subjected to deep sequencing on the MGISEQ-2000 platform. In this study, for each lane 4

846    samples (6 ng each) were pooled together for deep sequencing. To avoid sequencing bias induced by

847    base unbalance of TRAP sequence, 12 ng whole-genome DNA library (balance library) was mixed

848    with the 4 PCR samples in a final concentration of 1.5 ng/ $\mu$L and sequenced in one lane. All the

849    samples were subjected to pair-ended 150 bp deep-sequencing on MGISEQ-500 platform.

850

851    **Data analysis**

852    In order to evaluate the sequencing quality of amplicons and filter the low-quality sequencing data,

853    the default parameters of Fastqc-0.11.3 (http://www.bioinformatics.babraham.ac.uk/projects/fastqc/)

854    and fastp-0.19.6 (https://github.com/OpenGene/fastp) were used to carry out the filtration procedure

855    and generate the clean dataset of each sample. The clean sequencing segments of pair-ended TRAP

856    segments were merged using FLASh-1.2.1 (http://ccb.jhu.edu/software/FLASH/index.shtml) to

857    obtain full-length TRAP constructs. The expression characteristics of all the sequences were analyzed

858    by python 3.6, and most of the BsmBI linker fragments changed in orientation (GTTTGGAG->

859    GTTTGAAT). Therefore, in order to obtain the amplified fragment reads of each TRAP reference

860    sequence, the TRAP sequence BsmBI Linker was removed from the reference sequence. The BWA-

861    MEM algorithm of bwa(http://bio-bwa.sourceforge.net/) was used for local alignment, and the reads

862    of all samples were divided into 12,000 independent libraries. Due to the existence of sequencing or

863    synthesis introduced errors, each library was then filtered. In order to simplify the filtering process,

864    the filtration strategy varies from TRAP-ABE7.10, TRAP-CBE4.0-gam to TRAP-SpCas9. For

865    ABE7.10 and CBE4.0-gam, they mainly cause single-base variation, rarely introduce insertion and

866    deletion, the trap sequence length remains unchanged before and after editing. Therefore, we filter

867    the sequence of each library by locking the intermediate 37bp sequence starting with gRNA + scaffold

868    fragment and ending with GTTT. While TRAP-SpCas9 mainly cause insertion and deletion, the

869    length of trap sequence change around 37bp. Therefore, we adopt three steps to filter the sequence of

870    each library. The first step is to obtain the sequence containing gRNA + scaffold fragment as dataset1.

871    The second step is to obtain the sequence containing GTTTGAAT in dataset1 as dataset2. The third

872    step is to extract the intermediate trap sequence from dataset2, which removed the length limit. In

873    order to eliminate the interference of background noise before analyzing editing efficiency, all

874    mutations or indels found in WT HEK293T cells group were removed from the Dox group in advance.

875

876    For the TRAP-ABE7.10 and TRAP-CBE4.0-gam, the total editing efficiency for each trap is

877    calculated according to the following formula:

878        $Total\ editing\ efficiency$

879    $$= \frac{The\ number\ of\ reads\ whose\ trap\ sequence\ inconsistent\ with\ reference}{The\ total\ reads\ number\ of\ each\ trap}\%$$

880    , and the substitution percentage in the 37bp editing window is calculated according to the following

881    formula:

882        $Substitution\ percentage\ of\ ABE$

883    $$= \frac{The\ number\ of\ reads\ with\ actual\ A/T/C/G/N\ base\ in\ specified\ position}{The\ total\ number\ of\ reads\ with\ theoretical\ A\ base\ in\ specified\ position}\%$$

884        $Substitution\ percentage\ of\ CBE$

885    $$= \frac{The\ number\ of\ reads\ with\ actual\ A/T/C/G/N\ base\ in\ specified\ position}{The\ total\ number\ of\ reads\ with\ theoretical\ C\ base\ in\ specified\ position}\%$$

886

887    For the TRAP-SpCas9 system, the total editing efficiency for each trap is calculated according to the

888    following formula:

889        $Total\ editing\ efficiency$

890    $$= \frac{The\ number\ of\ reads\ whose\ trap\ length\ is\ not\ equal\ to\ 37bp}{The\ total\ reads\ number\ of\ each\ trap}\%$$

891    , and the average fraction of indels from 30bp deletion to 10bp insertion is calculated according to

892    the following formula:

893        $Average\ fraction\ of\ indels$

894    $$= \frac{The\ number\ of\ reads\ whose\ trap\ length\ range\ from\ 7-47bp}{The\ total\ reads\ number\ of\ all\ 12k\ trap}\%$$

895

896    Example of selecting low and high efficency gRNA for ABE and CBE. At least one site in this

897    sequence has more than 20% efficient for each TRAP, it is considered that the whole sequence has

898    more than 20% analytical value. Firstly, for the range of N1-N20, divide all TRAP library into a

899    group with an efficiency of greater than 20% and others with an efficiency of 1%, compare the base

900    distribution of the two groups. Then for the range of N5-N7, divide all TRAP library into the above

901    two groups, calculating the base mutation preference of *N-1 & N + 1* sites, motif preference and the

902    editing efficiency of N17-N20 sites including GCC, TT and other motif, respectively. For the CBE

903    system, the statistical method is basically the same as that of ABE. Due to the high editing efficiency

904    of CBE, in order to correct statistical deviations, CBE divides the efficiency greater than 50% and
905    less than 5% into two groups when calculating the N1-N20 base preference and base mutation
906    preference of *N-1 & N+1* site from N4-N8 . For the editing window is wider, the efficiency greater
907    than 50% and less than 5% is divided into two groups during the motif preference in N4-N8 window
908    and the editing efficiency of N17-N20 sites including GCC, TT and other motif . Python-3.6 and R
909    scripts were used for efficiency and motif analysis of all the TRAP samples. All visualizations use
910    GraphPad Prism8.2 and R package ggplot2.

911

912    **GNL machine learning featurization**
913    The feature set applied in our model construction contains 2485 features, which includes following
914    five categories. (*i*) 604 features of "one-hot" encoding of the nucleotide. There are two subsets in this
915    category: position-dependent and position-independent. And each category applies to the one
916    nucleotide and pairwise nucleotide. Such as "_nuc_pd_Order2" consisted of e.g. AA_1/AT_1/AG_1,
917    and "_nuc_pd_Order1" consisted of e.g. A_1/T_1/G_1/C_1. (*ii*) 3 GC features, which consists of GC
918    count, GC count < 10, GC count > 10. (iii) 16 features of the two nucleotides flanking the NGG PAM
919    in the 5' and 3'. (iv) Five thermodynamic features. We calculated five thermodynamic features ausing
920    the "Tm_staluc function" in Biopython package. All these features above were derived from the
921    30mer of target sequences. (v) 1856 features of three nucleotides with "position-dependent" and
922    "position-independent", such as ACG/AGG and ACG_1/ACT_2. Note that, all these features were
923    encoded by 30mer context sequence.(vi) Free energy (DeltaG), which was calculated by the local
924    version of "mfold" (http://unafold.rna.albany.edu/?q=mfold/download-mfold). We used the binary
925    programme "quikfold" in the same bin file of "mfold" to calculate multiple input sequences in the
926    same time. All the parameters were set as default, to make sure the outputs of each 20mer sequences
927    be the same as the webpage.

928

929    **Model training: Comparison with other machine learning models**
930    We trained model to predict the cleavage efficiency of each site in the editing windows among ABE
931    and CBE editing system. For the whole sequence sites of gRNA, three type of editing system was
932    uniformly trained by the same BRR model. To select the optimized predictive model, we initially
933    compared the predictive performance among eight models, they are, Bayesian Ridge regression
934    (BRR), gradient boosted regression tree (GBRT), decision tree (DT), L1-regression (L1-reg), L2-
935    regression (L2-reg), linear regression (LR), neural network (NN), random forest (RF). All these

936   algorithms were trained under the same features vector spaces. The mean performance of each

937   algorithm was conducted by 10-fold cross validation. Because all of them can be used as regression,

938   the performance was evaluated by SCC (Spearman Correlation coefficient). The model with best

939   performance, highest SCC value and lowest S.D., was selected. Finally, BRR was outperform than

940   other counterparts. All the models applied for training using the scikit-learn package in python.

941   During training, optimal hyper-parameter was chose using the inner 10-fold cross validation, using

942   the grid search. After that, $\alpha$ and $\lambda$ were both set as 1.e-6.

943

944   **Bayesian Ridge Regression**

945   Bayesian Ridge regression is changed from Bayesian linear regression by adding the prior of

946   coefficient "$\omega$" as spherical Gaussian and the priors over lambda are chosen to be gamma distributions,

947   which is similar to the classical Ridge regression [58]. Bayesian linear regression is briefly shown as

948   (1), and the coefficients of w is hypothesized as the spherical distribution to find a maximum

949   posteriori estimation of $\omega$ as (2) shows.

$$P(y|X, \omega, \alpha) = N(y|X\omega, \alpha) \tag{1}$$

950   Where $\alpha$ is treated as a random variable that is to be estimated from the data as gamma distribution.

951

$$P(\omega|\lambda) = N(\omega|0, \lambda^{-1}I_p) \tag{2}$$

952

953   Where $\lambda$ is also treated as a random variable that is to be estimated from the data, and also be

954   hypothesized as gamma distribution.

955

956   **Model explanation**

957   In addition to train the model with high performance, we also interested in the model importance for

958   our final model. We used **SHAP (SHapley Additive exPlanations) algorithm [59], which** is a

959   unified approach to explain the output of any machine learning model. Importantly, we excluded the

960   necessary site when training the site model for each editing system. E.g. We drop A1 when training

961   the N1 site of ABE system. So, the importance of the left features can be ranked by the SHAP value,

962   and the top 20 important features of each model were shown for each editor.

963

964   **DATA AVAILABILITY**

965    NGS data: CNBG accession number: TBD

966    Code for machine learning: https://github.com/TerminatorJ/CRISPR-TRAP-seq.git

967    CRISPR atlas: www.crispratlas.com/crispr

968

979

980    **AUTHOR CONTRIBUTION**

981    L.L. and Y.L. conceived the idea. L.B, G.C, L.L and Y.L oversaw the whole study. X.X, K.Q, L.L.

982    and Y.L. and designed the study. X.X, K.Q, X.L, and X.P performed most of the experiments and

983    analyses. All authors have contributed to the execution of the experiments and studies. L.L. and Y.L.

984    drafted the manuscript. All authors discussed the results and contributed to the final manuscript.

985

986    **CONFLICT OF INTEREST**

987    The authors declare no conflict of interest.

988

989    **REFERENCEs**

990    1.    Mali, P., et al., *RNA-guided human genome engineering via Cas9.* Science, 2013. **339**(6121):
991          p. 823-6.
992    2.    Cong, L., et al., *Multiplex genome engineering using CRISPR/Cas systems.* Science, 2013.
993          **339**(6121): p. 819-23.
994    3.    Jinek, M., et al., *A programmable dual-RNA-guided DNA endonuclease in adaptive bacterial
995          immunity.* Science, 2012. **337**(6096): p. 816-21.
996    4.    Farboud, B., et al., *Enhanced Genome Editing with Cas9 Ribonucleoprotein in Diverse Cells
997          and Organisms.* J Vis Exp, 2018(135).
998    5.    Hwang, W.Y., et al., *Efficient genome editing in zebrafish using a CRISPR-Cas system.* Nat
999          Biotechnol, 2013. **31**(3): p. 227-9.

1000   6.    Friedland, A.E., et al., *Characterization of Staphylococcus aureus Cas9: a smaller Cas9 for all-*
1001         *in-one adeno-associated virus delivery and paired nickase applications.* Genome Biol, 2015.
1002         **16**: p. 257.
1003   7.    Muller, M., et al., *Streptococcus thermophilus CRISPR-Cas9 Systems Enable Specific Editing*
1004         *of the Human Genome.* Mol Ther, 2016. **24**(3): p. 636-44.
1005   8.    Hou, Z., et al., *Efficient genome engineering in human pluripotent stem cells using Cas9 from*
1006         *Neisseria meningitidis.* Proc Natl Acad Sci U S A, 2013. **110**(39): p. 15644-9.
1007   9.    Kleinstiver, B.P., et al., *Broadening the targeting range of Staphylococcus aureus CRISPR-*
1008         *Cas9 by modifying PAM recognition.* Nat Biotechnol, 2015. **33**(12): p. 1293-1298.
1009   10.   Zetsche, B., et al., *Cpf1 is a single RNA-guided endonuclease of a class 2 CRISPR-Cas system.*
1010         Cell, 2015. **163**(3): p. 759-71.
1011   11.   Cox, D.B.T., et al., *RNA editing with CRISPR-Cas13.* Science, 2017. **358**(6366): p. 1019-1027.
1012   12.   Yeh, C.D., C.D. Richardson, and J.E. Corn, *Advances in genome editing through control of DNA*
1013         *repair pathways.* Nat Cell Biol, 2019. **21**(12): p. 1468-1478.
1014   13.   Qi, L.S., et al., *Repurposing CRISPR as an RNA-guided platform for sequence-specific control*
1015         *of gene expression.* Cell, 2013. **152**(5): p. 1173-83.
1016   14.   Lin, L., et al., *Genome-wide determination of on-target and off-target characteristics for RNA-*
1017         *guided DNA methylation by dCas9 methyltransferases.* Gigascience, 2018. **7**(3): p. 1-19.
1018   15.   Chen, B., J. Guan, and B. Huang, *Imaging Specific Genomic DNA in Living Cells.* Annu Rev
1019         Biophys, 2016. **45**: p. 1-23.
1020   16.   Kang, B.C., et al., *Precision genome engineering through adenine base editing in plants.* Nat
1021         Plants, 2018. **4**(7): p. 427-431.
1022   17.   Zong, Y., et al., *Precise base editing in rice, wheat and maize with a Cas9-cytidine deaminase*
1023         *fusion.* Nat Biotechnol, 2017. **35**(5): p. 438-440.
1024   18.   Komor, A.C., et al., *Programmable editing of a target base in genomic DNA without double-*
1025         *stranded DNA cleavage.* Nature, 2016. **533**(7603): p. 420-4.
1026   19.   Gaudelli, N.M., et al., *Programmable base editing of A\*T to G\*C in genomic DNA without*
1027         *DNA cleavage.* Nature, 2017. **551**(7681): p. 464-471.
1028   20.   Koblan, L.W., et al., *Improving cytidine and adenine base editors by expression optimization*
1029         *and ancestral reconstruction.* Nat Biotechnol, 2018. **36**(9): p. 843-846.
1030   21.   Rees, H.A., et al., *Improving the DNA specificity and applicability of base editing through*
1031         *protein engineering and protein delivery.* Nat Commun, 2017. **8**: p. 15790.
1032   22.   Nishida, K., et al., *Targeted nucleotide editing using hybrid prokaryotic and vertebrate*
1033         *adaptive immune systems.* Science, 2016. **353**(6305).
1034   23.   Zhang, Y., et al., *Programmable base editing of zebrafish genome using a modified CRISPR-*
1035         *Cas9 system.* Nat Commun, 2017. **8**(1): p. 118.
1036   24.   Jensen, K.T., et al., *Chromatin accessibility and guide sequence secondary structure affect*
1037         *CRISPR-Cas9 gene editing efficiency.* FEBS Lett, 2017. **591**(13): p. 1892-1901.
1038   25.   Uusi-Makela, M.I.E., et al., *Chromatin accessibility is associated with CRISPR-Cas9 efficiency*
1039         *in the zebrafish (Danio rerio).* PLoS One, 2018. **13**(4): p. e0196238.
1040   26.   Kim, H.K., et al., *High-throughput analysis of the activities of xCas9, SpCas9-NG and SpCas9*
1041         *at matched and mismatched target sequences in human cells.* Nat Biomed Eng, 2020. **4**(1):
1042         p. 111-124.
1043   27.   Shen, M.W., et al., *Predictable and precise template-free CRISPR editing of pathogenic*
1044         *variants.* Nature, 2018. **563**(7733): p. 646-651.

28. Allen, F., et al., *Predicting the mutations generated by repair of Cas9-induced double-strand breaks.* Nat Biotechnol, 2018.

29. Zhou, Y., et al., *Enhanced genome editing in mammalian cells with a modified dual-fluorescent surrogate system.* Cell Mol Life Sci, 2016. **73**(13): p. 2543-63.

30. Komor, A.C., et al., *Improved base excision repair inhibition and bacteriophage Mu Gam protein yields C:G-to-T:A base editors with higher efficiency and product purity.* Sci Adv, 2017. **3**(8): p. eaao4774.

31. Sjostedt, E., et al., *An atlas of the protein-coding genes in the human, pig, and mouse brain.* Science, 2020. **367**(6482).

32. Billon, P., et al., *CRISPR-Mediated Base Editing Enables Efficient Disruption of Eukaryotic Genes through Induction of STOP Codons.* Mol Cell, 2017. **67**(6): p. 1068-1079 e4.

33. Xu, Y., et al., *A new massively parallel nanoball sequencing platform for whole exome research.* BMC Bioinformatics, 2019. **20**(1): p. 153.

34. Shalem, O., et al., *Genome-scale CRISPR-Cas9 knockout screening in human cells.* Science, 2014. **343**(6166): p. 84-87.

35. Morgens, D.W., et al., *Genome-scale measurement of off-target activity using Cas9 toxicity in high-throughput screens.* Nat Commun, 2017. **8**: p. 15178.

36. Tycko, J., V.E. Myer, and P.D. Hsu, *Methods for Optimizing CRISPR-Cas9 Genome Editing Specificity.* Mol Cell, 2016. **63**(3): p. 355-70.

37. Kosicki, M., et al., *Dynamics of Indel Profiles Induced by Various CRISPR/Cas9 Delivery Methods.* Prog Mol Biol Transl Sci, 2017. **152**: p. 49-67.

38. Lin, L. and Y. Luo, *Tracking CRISPR's Footprints.* Methods Mol Biol, 2019. **1961**: p. 13-28.

39. Moller, H.D., et al., *CRISPR-C: circularization of genes and chromosome by CRISPR in human cells.* Nucleic Acids Res, 2018. **46**(22): p. e131.

40. Wang, D., et al., *Optimized CRISPR guide RNA design for two high-fidelity Cas9 variants by deep learning.* Nat Commun, 2019. **10**(1): p. 4284.

41. Graf, R., et al., *sgRNA Sequence Motifs Blocking Efficient CRISPR/Cas9-Mediated Gene Editing.* Cell Rep, 2019. **26**(5): p. 1098-1103 e3.

42. Ata, H., et al., *Robust activation of microhomology-mediated end joining for precision gene editing applications.* PLoS Genet, 2018. **14**(9): p. e1007652.

43. Wang, J., et al., *GNL-Scorer: A generalized model for predicting CRISPR on-target activity by machine learning and featurization.* J Mol Cell Biol, 2020.

44. Concordet, J.P. and M. Haeussler, *CRISPOR: intuitive guide selection for CRISPR/Cas9 genome editing experiments and screens.* Nucleic Acids Res, 2018. **46**(W1): p. W242-W245.

45. Xue, L., et al., *Prediction of CRISPR sgRNA Activity Using a Deep Convolutional Neural Network.* J Chem Inf Model, 2019. **59**(1): p. 615-624.

46. Labuhn, M., et al., *Refined sgRNA efficacy prediction improves large- and small-scale CRISPR-Cas9 applications.* Nucleic Acids Res, 2018. **46**(3): p. 1375-1385.

47. Hwang, G.H., et al., *Web-based design and analysis tools for CRISPR base editing.* BMC Bioinformatics, 2018. **19**(1): p. 542.

48. Dandage, R., et al., *beditor: A Computational Workflow for Designing Libraries of Guide RNAs for CRISPR-Mediated Base Editing.* Genetics, 2019. **212**(2): p. 377-385.

49. Ramakrishna, S., et al., *Surrogate reporter-based enrichment of cells containing RNA-guided Cas9 nuclease-induced mutations.* Nat Commun, 2014. **5**: p. 3378.

1089  50.  Kim, H., et al., *Surrogate reporters for enrichment of cells with nuclease-induced mutations.*
1090       Nat Methods, 2011. **8**(11): p. 941-3.
1091  51.  Cullot, G., et al., *CRISPR-Cas9 genome editing induces megabase-scale chromosomal*
1092       *truncations.* Nat Commun, 2019. **10**(1): p. 1136.
1093  52.  Kosicki, M., K. Tomberg, and A. Bradley, *Repair of double-strand breaks induced by CRISPR-*
1094       *Cas9 leads to large deletions and complex rearrangements.* Nat Biotechnol, 2018. **36**(8): p.
1095       765-771.
1096  53.  Hu, J.H., et al., *Evolved Cas9 variants with broad PAM compatibility and high DNA specificity.*
1097       Nature, 2018. **556**(7699): p. 57-63.
1098  54.  Slaymaker, I.M., et al., *Rationally engineered Cas9 nucleases with improved specificity.*
1099       Science, 2016. **351**(6268): p. 84-8.
1100  55.  Kleinstiver, B.P., et al., *High-fidelity CRISPR-Cas9 nucleases with no detectable genome-wide*
1101       *off-target effects.* Nature, 2016. **529**(7587): p. 490-5.
1102  56.  Kuscu, C., et al., *CRISPR-STOP: gene silencing through base-editing-induced nonsense*
1103       *mutations.* Nat Methods, 2017. **14**(7): p. 710-712.
1104  57.  Kutner, R.H., X.Y. Zhang, and J. Reiser, *Production, concentration and titration of*
1105       *pseudotyped HIV-1-based lentiviral vectors.* Nat Protoc, 2009. **4**(4): p. 495-505.
1106  58.  Tipping, M.E., *Sparse Bayesian learning and the relevance vector machine.* Journal of
1107       machine learning research, 2001. **1**(Jun): p. 211-244.
1108  59.  Lundberg, S.M. and S.-I. Lee. *A unified approach to interpreting model predictions*. 2017.
1109
1110