1 **High throughput PRIME editing screens identify functional DNA variants in the human**

2 **genome**

3

4 **Authors:** Xingjie Ren[1,*], Han Yang[1,*], Jovia L. Nierenberg[2], Yifan Sun[1], Jiawen Chen[3], Cooper

5 Beaman[1], Thu Pham[4], Mai Nobuhara[4], Maya Asami Takagi[1], Vivek Narayan[1], Yun Li[3,5,,6], Elad

6 Ziv[1,7], Yin Shen[1,8,9,#]

7

8 **Affiliations**

9 [1] Institute for Human Genetics, University of California, San Francisco, San Francisco, CA, USA.

10 [2] Department of Epidemiology and Biostatistics, University of California, San Francisco, San

11 Francisco, CA, USA.

12 [3] Department of Biostatistics, University of North Carolina at Chapel Hill, Chapel Hill, NC 27599,

13 USA.

14 [4] Pharmaceutical Sciences and Pharmacogenomics Graduate Program, University of California,

15 San Francisco, San Francisco, CA, USA.

16 [5] Department of Genetics, University of North Carolina, Chapel Hill, NC, USA.

17 [6] Department of Computer Science, University of North Carolina, Chapel Hill, NC, USA.

18 [7] Division of General Internal Medicine, Department of Medicine, and Helen Diller Family

19 Comprehensive Cancer Center, University of California, San Francisco, San Francisco, CA, USA.

20 [8] Department of Neurology, University of California, San Francisco, San Francisco, CA, USA.

21 [9] Weill Institute for Neurosciences, University of California, San Francisco, San Francisco, CA,

22 USA.

23 * These authors contributed equally to the work.

24 [#]Corresponding authors: Yin Shen Yin.Shen@ucsf.edu

**Abstract**

Despite tremendous progress in detecting DNA variants associated with human disease, interpreting their functional impact in a high-throughput and base-pair resolution manner remains challenging. Here, we develop a novel pooled prime editing screen method, PRIME, which can be applied to characterize thousands of coding and non-coding variants in a single experiment with high reproducibility. To showcase its applications, we first identified essential nucleotides for a 716 bp *MYC* enhancer via PRIME-mediated saturation mutagenesis. Next, we applied PRIME to functionally characterize 1,304 non-coding variants associated with breast cancer and 3,699 variants from ClinVar. We discovered that 103 non-coding variants and 156 variants of uncertain significance are functional via affecting cell fitness. Collectively, we demonstrate PRIME capable of characterizing genetic variants at base-pair resolution and scale, advancing accurate genome annotation for disease risk prediction, diagnosis, and therapeutic target identification.

**Main**

Advances in genome sequencing have led to the identification of hundreds of millions of genetic variants in the human population, with a fraction conferring risk for common illnesses such as diabetes, neurological disorders, and cancers[1]. A major barrier to understanding the genetic underpinnings of these complex diseases is the paucity of functional annotation for disease-associated variants, especially because such variants are predominantly located within non-coding regions. Growing evidence suggests that non-coding risk variants may contribute to disease pathogenesis by disrupting gene regulation. Even protein-coding variants discovered from individuals with disease are frequently classified as Variants of Uncertain Significance (VUS). Therefore, more precise and higher throughput functional characterization methods for elucidating disease-associated variant function at base-pair resolution, and multiplexed across genomic loci, are necessary to realize the potential of personalized medicine.

The development of genome editing technologies has enabled us to perturb and assess DNA sequences in desired regions at a large scale. However, there are still fundamental barriers to utilizing these methods for precision genome annotation. For example, CRISPRa, CRISPRi, CRISPR deletion, and CRISPR indel have been applied in genetic screening strategies for characterizing both genes and *cis*-regulatory regions[2], but have failed to pinpoint casual variants for diseases. Traditional methods of characterizing DNA variants (SNPs) by knock-in via homologous recombination are inefficient and low throughput. Base editors also have limitations, introducing specific mutations (C→T, A→G, T→C, G→A) with varied target efficiencies[3]. Thus, there is still a significant deficit in methods for effectively characterizing the roles of putative disease-causing variants in human health and diseases. Robust high throughput methods making desired edits at base-pair resolution are urgently needed to achieve a better understanding of the genetic underpinnings of disease.

Prime editing (PE), a versatile and precise genetic engineering method, has been developed to introduce any type of edit, including point mutation, insertion, and deletion[4]. In particular, PE2, employs the *Streptococcus pyogenes* Cas9 (SpCas9) H840A nickase and Moloney murine leukemia virus (M-MLV) reverse transcriptase. The spacer in the prime editing guide RNA (pegRNA) directs the Cas9 nickase and M-MLV complex to the target site, while the RT template sequence provides the desired editing information. Thus, both targeting and editing information can be easily programmed in the same pegRNA to perform single nucleotide substitution, insertion or deletion. PE3, a newer iteration of PE, can further increase editing efficiency by promoting the replacement of non-edited strands using an additional single-guide (sgRNA) for nicking[5]. Prime editors' capacity for precision genome editing suggests the possibility

71    of high throughput variant-level genome manipulation. Recently, PE screens were used to identify

72    VUS at the *NPC1* locus based on a lysosome functional assay by transfection of pegRNAs and

73    targeted sequencing of this region[6]. Although transient transfection of PE machinery followed by

74    targeted sequencing of the edited locus enables the identification of editing events, its scope is

75    limited to just that locus, and thus, scaling up for massively parallel assessment of multiple loci is

76    not feasible. Besides increased throughput, improved control of transgene copy number, stable

77    expression of PE machinery, and direct loci comparison are also desired.

78         Here, we enable high throughput pooled screens of thousands of DNA variants in the

79    human genome by lentiviral delivery of PE, namely PRIME. We demonstrate the utility of PRIME

80    for three different applications, including the saturation mutagenesis analysis of a 716 bp

81    enhancer, the functional characterization of 1,304 breast cancer-associated variants, and the

82    evaluation of 3,699 clinical variants' impact on cell fitness. Our results establish the

83    generalizability of PRIME for precisely characterizing genetic variants in the human genome.
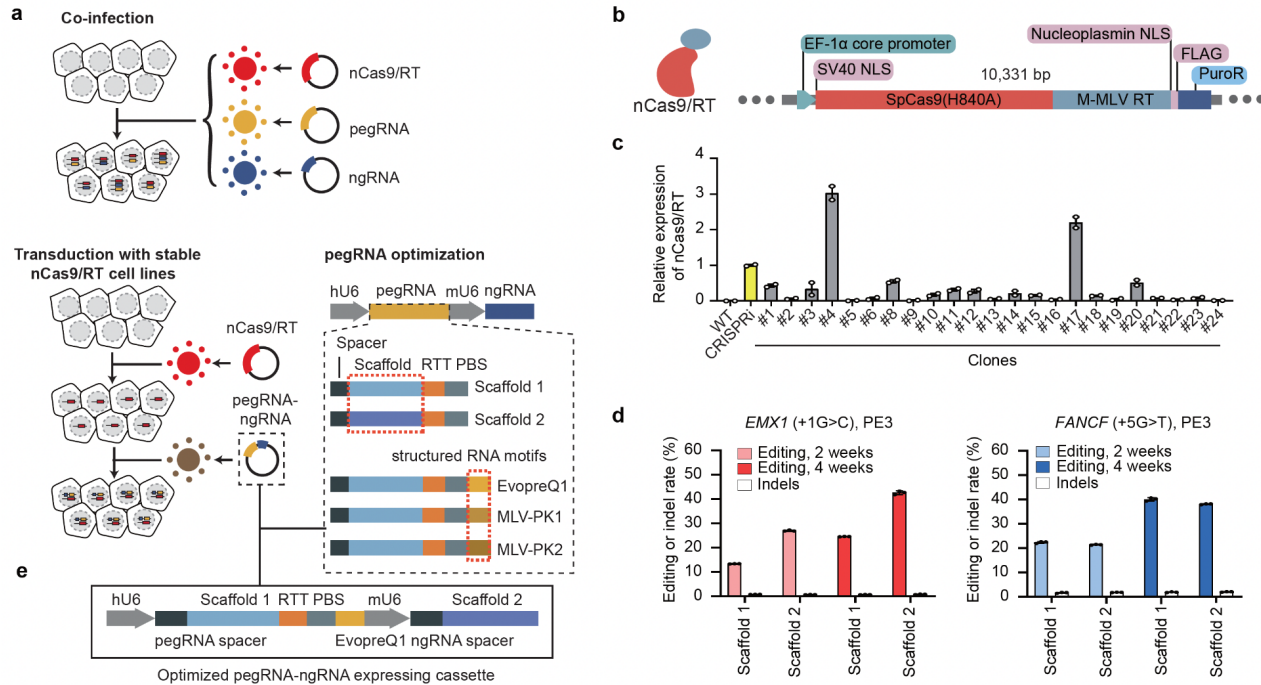
84

**Optimization of PE efficiency in mammalian cells delivered by lentivirus**
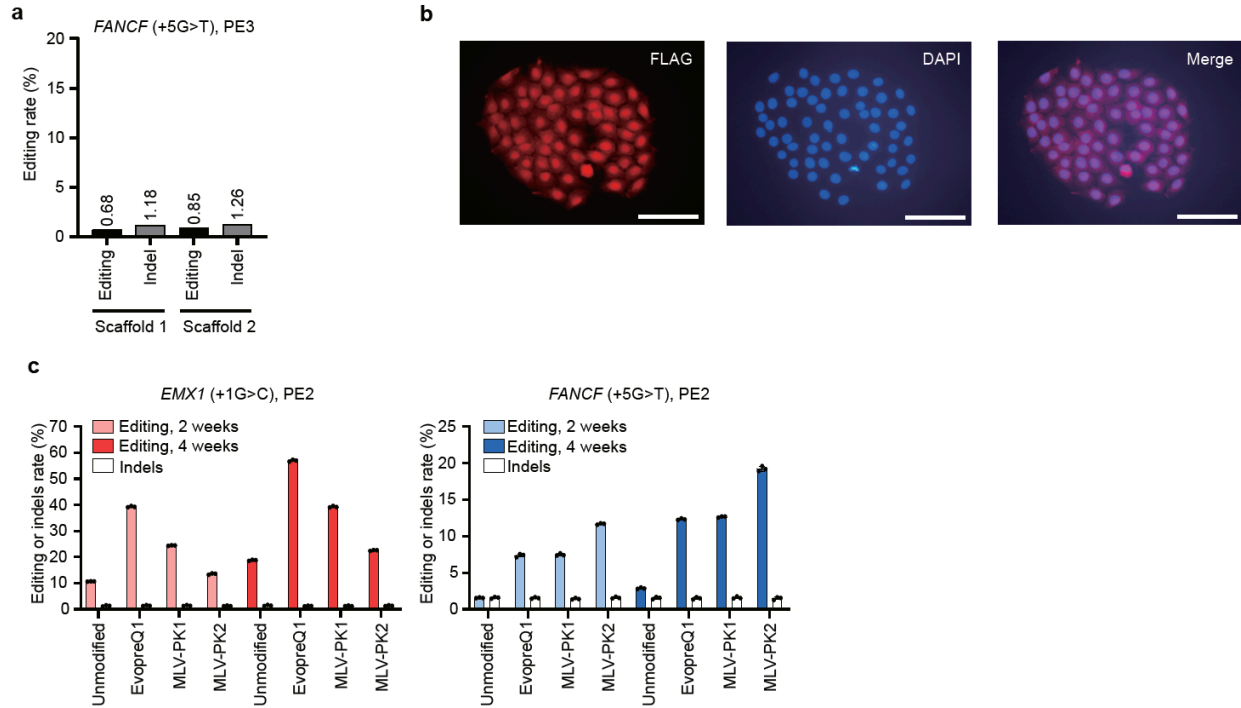
86         To enable PE screens with delivery by lentivirus, we initially installed PE3 by infecting

87    MCF7 cells using three different viruses: 1) virus expressing Cas9 (H840A) nickase (nCas9) and

88    Moloney murine leukemia virus (M-MLV) reverse transcriptase; 2) virus expressing pegRNA; 3)

89    virus expressing nick sgRNA (ngRNA). Unfortunately, this strategy yielded less than 1% PE

90    efficiency with a relatively high indel rate. This is because of the low efficiency of coinfecting three

91    different viruses in the same cell (**Fig. 1a**, **Supplementary Fig. 1a**).

92         Packaging all PE3 components within the same virus is challenging. To increase PE

93    efficiency and facilitate a pooled screening approach with a lentiviral library, we infected MCF7

94    cells with lentivirus containing an nCas9 and M-MLV reverse transcriptase stable expression

95    cassette (**Fig. 1b**). After puromycin selection, we isolated multiple clones and selected one with

96    the highest nCas9 expression (**Fig. 1c**, RT-qPCR, clone #4, **Supplementary Fig. 1b**) for

97    subsequent experiments. The stable expression of nCas9/M-MLV allows for high efficiency

98    pegRNA/ngRNA packaging and lentiviral delivery, with greater editing efficiency than the co-

99    infection method (**Fig. 1d**). To further improve PE efficiency, we assessed editing efficiency using

100   three different structured RNA motifs (EvopreQ1, MLV-PK1, and MLV-PK2) at the 3' terminus of

101   the pegRNA[7-9]. Cells treated with pegRNAs containing scaffold structure RNA motifs exhibited

102   consistently higher editing efficiencies at both the *EMX1* and *FANCF* locus compared to using PE

103   without structured RNA motifs (**Supplementary Fig. 1c**), so we added evopreQ1 to the pegRNA

104 design for all pooled screens. Scaffold 1[5] and 2[10] had no significant effects on PE efficiency,

105 suggesting the feasibility of dual pegRNA and ngRNA delivery from the same viral particle (**Fig.**

106 **1d**). All PE experiments in clonal MCF7 cells (MCF7-nCas9/RT) exhibited relatively low indel rates

107 (0.7% to 1.95%). Thus, we used MCF7-nCas9/RT cells and lentiviral delivery of both the pegRNA

108 with scaffold 1 and ngRNA with scaffold 2 in the same construct (**Fig 1e**).



109

110 **Figure 1. Optimizing PE efficiency in mammalian cells using lentiviral delivery.**

111 (a) The different strategies tested for optimizing PE efficiency in MCF7 cell lines. Top: co-infecting

112 three different viruses to deliver PE machinery. Bottom: dual pegRNA/ngRNA viral infection of

113 clonal MCF7 line stably expressing nickase Cas9 (nCas9) and Moloney murine leukemia virus

114 reverse transcriptase (M-MLV RT). Two scaffolds and three different structured RNA motifs tested

115 are also shown. (b) Lentiviral construct for generating nCas9/RT expressing MCF7 clones. PuroR,

116 Puromycin resistance gene. M-MLV RT, Moloney murine leukemia virus reverse transcriptase.

117 (c) RT-qPCR analysis showing the relative expression of nCas9/RT in different clones, normalized

118 to the dCas9 expression of an established CRISPRi iPSC line (Yellow). Error bars represent the

119 s.e.m. (d) The editing efficiency and indel rate for *EMX1* and *FANCF* loci at 2 weeks and 4 weeks

120 after PE installation using two different RNA scaffolds. Error bars represent the s.d. (e) Improved

121 vector for expression of pegRNA and ngRNA for PRIME. RTT: reverse transcription template,

122 PBS: primer binding site.

123

**Supplementary Figure 1. Optimizing PE efficiency in MCF7 cell line.**

(a) Prime editing efficiency and indel rate by co-infection of pegRNA, ngRNA and nCas9/RT expressing lentiviruses in MCF7 cells. (b) Immunofluorescent staining showing the localization of nCas9/RT (red, FLAG tagged) in the nucleus (blue, DAPI) in MCF7-nCas9/RT cells. Scale bars, 1000 μm. (c) Editing efficiency and indel rate by PE using three different structured RNA motifs to the 3' terminus of pegRNAs at 2 and 4 weeks post infection in MCF7-nCas9/RT cells. Error bars represent the s.d.

131

**PRIME enables nucleotide-resolution analyses of enhancer function**

Enhancers can modulate cell type-specific gene expression and are highly enriched with disease-associated variants. Knowledge of the endogenous function for each nucleotide in enhancers should reveal crucial transcription factors that govern enhancer activation and facilitate the development of better models for gene regulatory networks and the prediction of disease-associated non-coding variant regulatory effects. To test whether PRIME can quantify the impact of each base in an enhancer, we focused on an MCF7-specific *MYC* enhancer identified from a CRISPRi screen[11]. This enhancer is located 405 kb downstream of *MYC* and displays enhancer signatures, including open chromatin, H3K27ac, and H3K4me1 signals, in addition to forming a chromatin loop with the *MYC* promoter (**Fig. 2a**). Deletion of this enhancer caused an 85% downregulation of *MYC* expression in MCF7 cells confirming its enhancer activity for *MYC*

143    (**Supplementary Fig. 2a**). Since *MYC* downregulation is correlated with MCF7 cell survival[12], we

144    performed a PE-enabled high throughput saturation mutagenesis screen of this *MYC* enhancer

145    in MCF7 cells dependent on the cell survival phenotype (**Fig. 2b**).
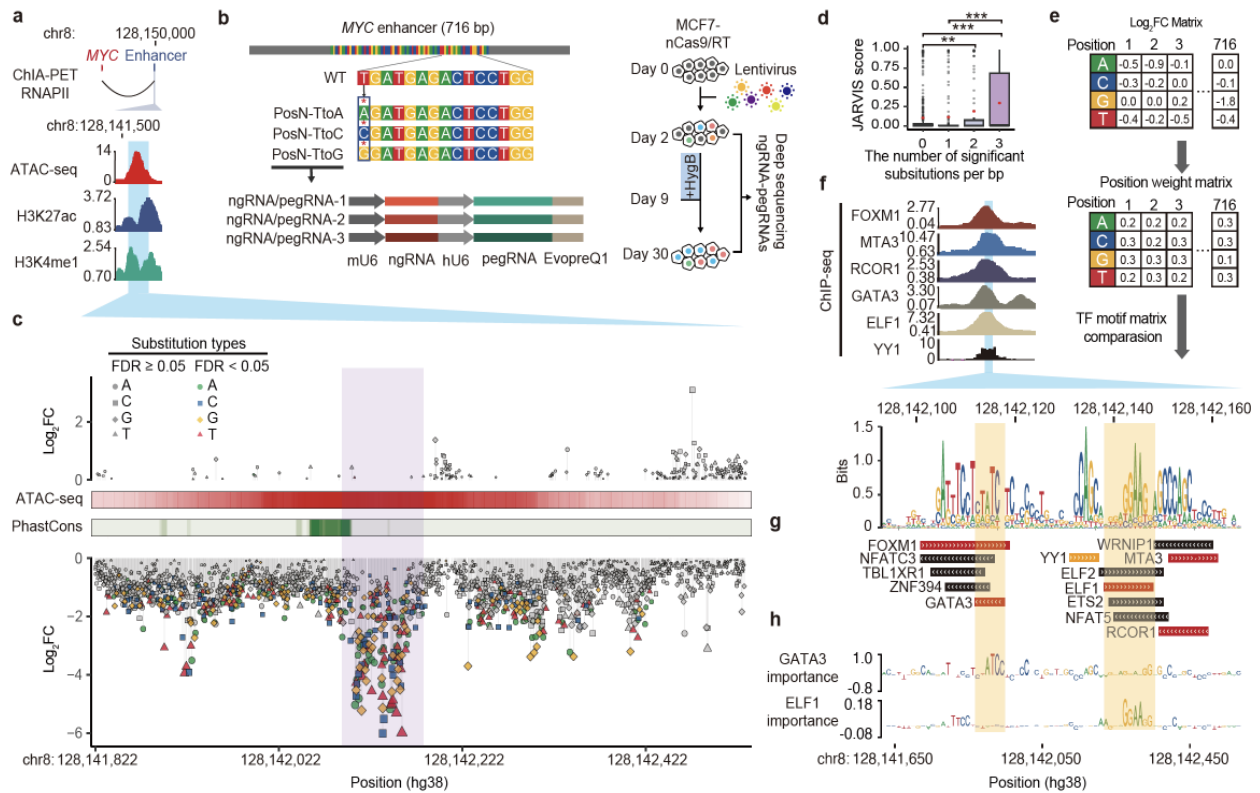
146    To dissect the enhancer's function at base-pair resolution, we designed a library of 6,252

147    pairs of pegRNA/ngRNA to generate 2,148 single nucleotide substitutions within the 716 bp *MYC*

148    enhancer region (**Supplementary Table 1**). Specifically, we changed the original base into three

149    other nucleotides, and each event was independently evaluated three times in the same screen

150    (**Fig. 2b**). We also included 94 positive control pegRNA/ngRNA pairs, which introduced stop

151    codons (iSTOPs) in *MYC*, and 400 negative control pegRNA/ngRNA pairs. 246 of the negative

152    controls were non-human genome targeting, and 154 targeted the *AAVS1* safe harbor locus

153    (**Supplementary Table. 1**). We then infected MCF7-nCas9/RT cells with lentiviral libraries

154    expressing these pegRNA/ngRNA pairs (**Supplementary Fig. 2b**). Two days after infection,

155    virus-transduced cells were hygromycin selected for one week and expanded in regular media for

156    another 3 weeks. We collected cells at 2 and 30 days post-infection, amplified the integrated

157    pegRNA/ngRNA pairs, and determined the relative depletion or enrichment of each

158    pegRNA/ngRNA between these two time points by deep sequencing (**Fig. 2b**). We performed this

159    screen 3 times (**Supplementary Fig. 2c**) and used negative controls, including non-human

160    targeting and AAVS1 targeting paired pegRNA/ngRNAs for data normalization. Fold changes

161    (FC) for each pegRNA/ngRNA pair between day 2 and day 30 samples post-infection were

162    calculated using the MAGeCK pipeline[13] (**Supplementary Table 1**). As expected, 78% (73/94) of

163    iSTOPs were depleted ($log_2FC < 0$) 30 days post-infection. iSTOP depletion rates were negatively

164    correlated with their distance from the transcription start site (TSS) of *MYC*, consistent with the

165    observation that gene knockout is more efficient when perturbations are introduced at the 5'

166    terminus[14] (**Supplementary Fig. 2d**). In addition, two iSTOPs (amino acid position 350 and 355)

167    targeting the region between the nuclear localization signal (NLS) and the carboxy-terminal

168    domain (CTD) domain were also significantly depleted (**Supplementary Fig. 2d**). The N-terminus

169    of MYC contains its core transcription transactivation domain which binds multiple partners[15]. It is

170    possible that those two iSTOPs created a truncated MYC still capable of binding to cofactors, but

171    unable to bind MYC DNA targets, interfering with the functions of wild type MYC and its cofactors.

172    To investigate the effects of each nucleotide on enhancer function, we defined sensitive

173    base pairs (SBP) as nucleotides that affect cell fitness when substituted at least once (FDR <

174    0.05, $|log_2FC| > 1$). 334 of the 716 (46.6%) tested base pairs were SBP with $log_2FC < -1$

175    (**Supplementary Table 1**), indicating that mutations at those locations reduce enhancer activity

176    and cell fitness. 23.1% (77/334) of SBPs were depleted at day 30 with all three substitutions (FDR

177     < 0.05, log$_2$FC < -1). Additionally, none of the tested sequences were significantly enriched at day

178     30 with increased cell growth phenotype, indicating that perturbation of these sequences

179     exclusively attenuated enhancer activity (**Fig. 2c**).

180         Deep learning models have been developed to prioritize non-coding regions and predict

181     their relevance to human disease. Encouragingly, SBPs with two or more significant substitutions

182     (n = 172) were predicted to be more deleterious than SBPs with only one significant substitution

183     (n = 162) or non-SBPs (n = 382) by JARVIS[16] (**Fig. 2d**). This demonstrates the success of PRIME

184     in validating computationally predicted functional sequences. We further established a continuous

185     bin density analysis, detecting variation in SBP density along the enhancer to define SBP-

186     enriched regions (**Supplementary Fig. 2e and f**). We identified the core enhancer region in the

187     enhancer with a high density of SBPs, based on the slope value of the cumulative curve of SBPs

188     with three significant substitutions, as a larger slope value indicates a higher density of SBPs in

189     the region. The core enhancer region was defined by a minimal slope cut-off of 0.43 (Z score-

190     derived $P$ < 0.05). The core enhancer region (chr8:128,142,093-128,142,181, hg38) colocalized

191     with an open chromatin summit. This region contains SBPs with the most extensive fold changes

192     when mutated, indicating its strong effect on enhancer activity. (**Fig. 2c**, highlighted in purple).

193     Notably, the enhancer's core sequence was located next to a highly conserved region (**Fig. 2c**).

194     This is not surprising because enhancers undergo rapid evolutionary changes compared to

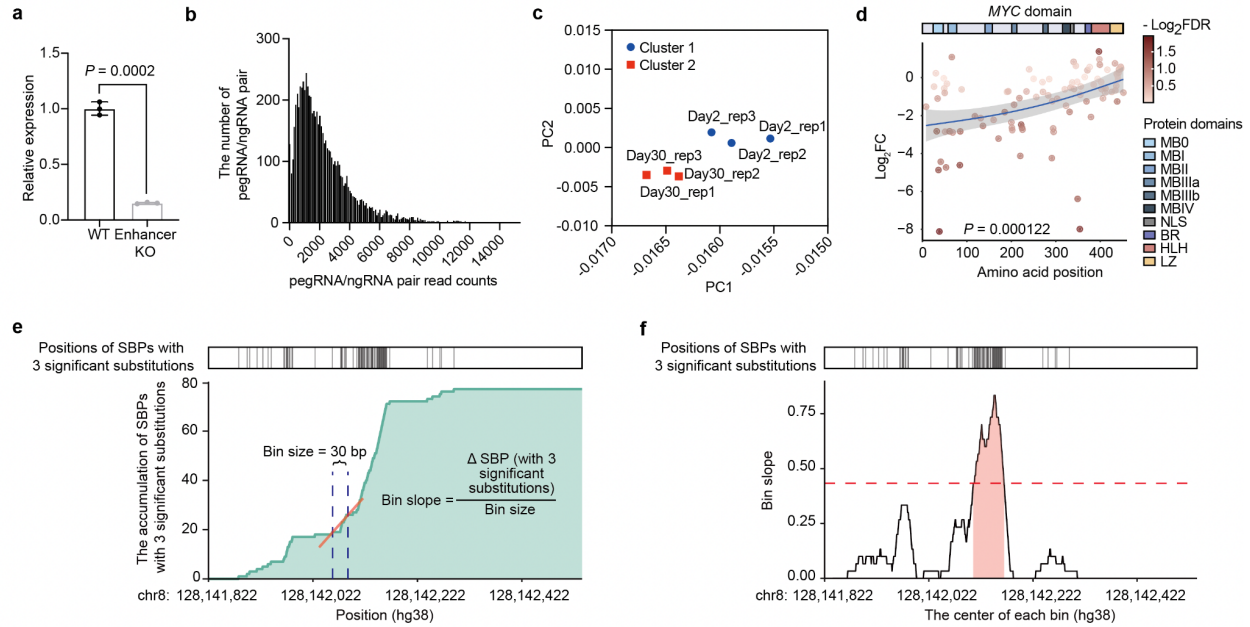195     protein-coding sequences[17].

196         Our functional data provide a unique opportunity to calculate and construct a position

197     weight matrix (PWM). Using fold changes for each nucleotide, we generated a functional PWM

198     (**Fig. 2e**). Comparing our functional PWM with curated transcription factors (TFs) motifs from the

199     JASPAR, HOCOMOCO, and SwissRegulon databases[18-20] identified 13 TFs with matched motif

200     PWMs (**Fig. 2g and h, Supplementary Table 2**). 5 predicted TFs (GATA3, ELF1, FOXM1, MTA3

201     and RCOR1) have already been shown to bind to the *MYC* enhancer based on ENCODE ChIP-

202     seq datasets[21], and YY1 is predicted to bind to this enhancer in MCF7 by Avocado through the

203     ENCODE project[22] (**Fig. 2f**). Furthermore, *GATA3* and *YY1* are essential cell survival genes in

204     MCF7[23], confirming the utility of PE-enabled saturation mutagenesis for interrogating enhancer

205     function at base pair resolution. Essential nucleotides for the ELF1 and GATA3 binding motifs

206     identified by our screens were consistent with those imputed by BPNet[24], further validating the

207     importance of quantitative roles of each nucleotide discovered by PRIME. Combined, we

208     demonstrated that PRIME is useful for elucidating nucleotide-resolution functional annotations of

209     non-coding cis-regulatory elements.

**Figure 2. Functional characterization of a *MYC* enhancer by saturation mutagenesis using PRIME.**

(a) (Top) The target enhancer is downstream of *MYC*. (Bottom) The enhancer region is highly enriched with ATAC-seq, H3K27ac, and H3K4me1 ChIP-seq signals. The blue area indicates the region selected for PRIME. (b) (Top) Diagram showing the design of saturation mutagenesis screening at the 716 bp enhancer. Each nucleotide was subjected to substitution with three nucleotides by PE. (Middle) Each substitution event was covered by three uniquely designed pegRNA/ngRNA pairs. (Bottom) The PRIME workflow. (c) Log$_2$(fold change) of each substitution at each base pair ordered by their genomic locations. Mutations with a significant effect on cell fitness are colored. ATAC-seq signals and conservation scores calculated by PhastCons are shown. (d) JARVIS scores for base pairs with different numbers of significant substitutions. Box plots indicate median, IQR, Q1 − 1.5 × IQR, and Q3 + 1.5 × IQR. Outliers are shown as gray dots. Mean values are shown as red dots. *P* values were calculated using a two-tailed two-sample t-test. (e) The creation of a functional PWM for identifying potential TF binding sites. (f) (Top) ChIP-seq signals of 6 TFs in MCF7. The blue region indicates the core enhancer region. (Bottom) The sequence logo plot for the core enhancer regions generated by the functional PWM from (e). (g) Matched TF binding sites. (h) (Top) Dense tracks showing BPNet model-derived nucleotide importance scores for GATA3 and ELF1 binding sites.

229

**Supplementary Figure 2. Characterize enhancer function and results of PRIME in MCF7 cells.** (a) CRISPR/Cas9 knockout of the *MYC* enhancer in MCF7 decreased *MYC* expression. *P* values were calculated using a two-tailed two-sample t-test. Error bars represent the s.e.m. (b) Distribution of pegRNA/ngRNA pair read counts in the cloned plasmid library. (c) PCA analysis demonstrates the high reproducibility of PRIME between biological replicates. (d) The correlation between locations of PE-induced stop codons and their effect sizes. The blue line and *P* value were calculated using generalized additive models. The shaded areas indicate 95% confidence intervals. (e) (Top) The position of SBPs with three significant substitutions. (Bottom) Cumulative distribution plot of SBPs with three significant substitutions along the *MYC* enhancer and the formula for calculating the slope of each continuous bin. (f) Line plot of slopes for each continuous bin along the *MYC* enhancer. The red dashed line is the cutoff for a significant slope, which is based on a slope with a Z score-derived *P* value equal to 0.05. The red region is the core enhancer region, derived from the bins' slopes greater than the cutoff (slope > 0.43).

243

**Characterization of breast cancer-associated variants**

Next, we tested the feasibility of characterizing >5,000 disease-associated DNA variants at various genomic loci, including non-coding variants from GWAS and variants detected fro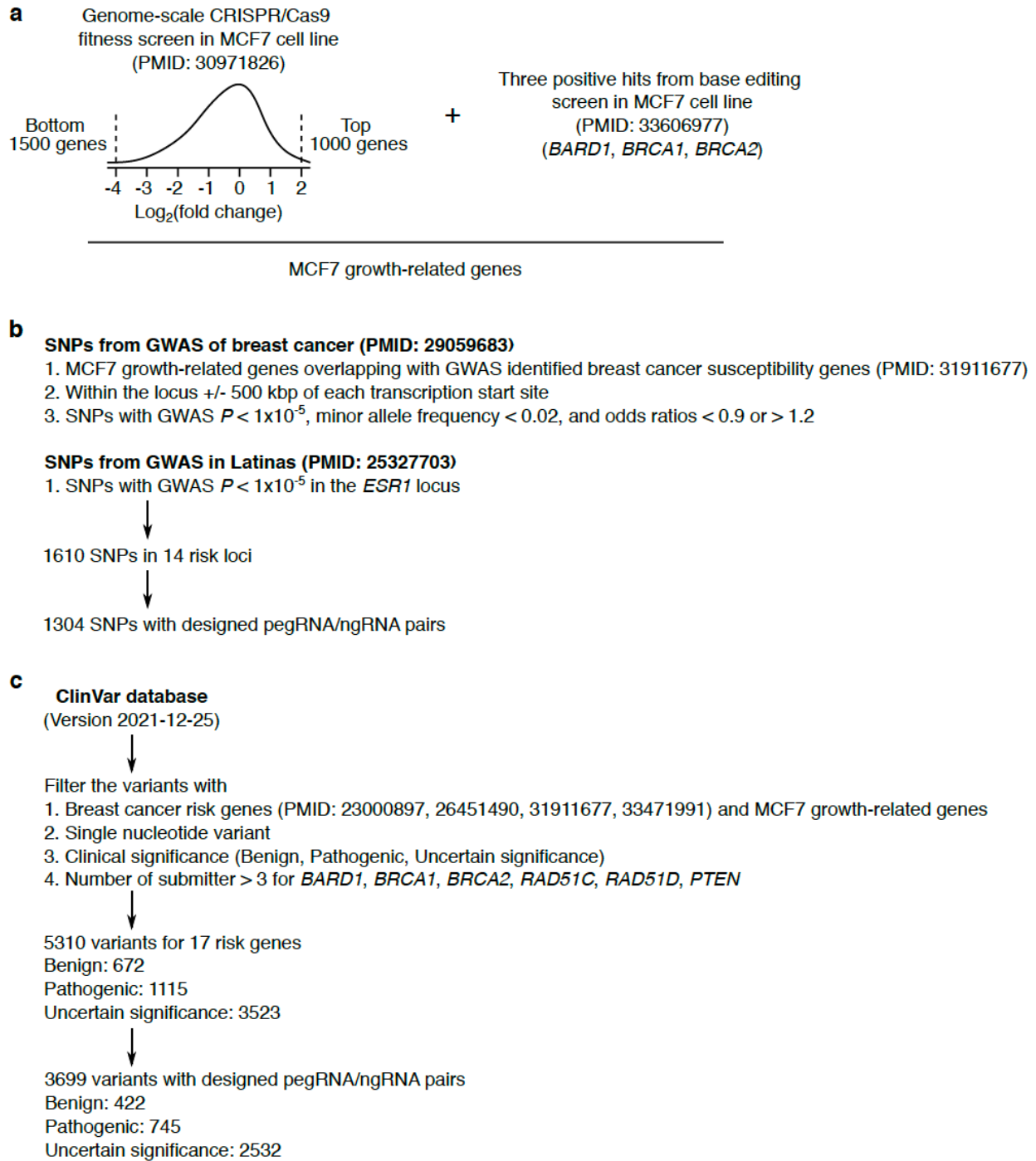m clinical samples. For GWAS-identified variants, we focused on breast cancer, the most common cancer in women in the U.S. To test the feasibility of characterizing DNA variants associated with breast cancer, we used the summary statistics from the largest GWAS to date, including samples of mostly European ancestry[25]. Candidate genes from a comprehensive fine mapping effort for

251    this GWAS[26] overlapping with growth phenotype genes prioritized by CRISPR screens[23, 27] were

252    selected. These include: *CCND1*, *PSMD6*, *MYC*, *UBA52*, *DYNC1I2*, *ESR1*, *MRPS18C*, *NOL7*,

253    *EWSR1*, *BRCA2*, and *GRHL2*, which were negatively selected in a CRISPR knockout screen,

254    and *CUX1*, *CASP8*, and *TNFSF10*, which are tumor suppressor genes and positively selected in

255    a CRISPR knockout screen (**Supplementary Fig. 3a**). We then selected 1,304 single nucleotide

256    polymorphisms (SNPs) (**Supplementary Fig. 3b** and **Supplementary Table 3**) within 500 kbp

257    upstream and downstream of these genes that were previously associated with breast cancer[25]

258    and had been implicated as possibly acting through these genes[26]. We also selected 3,699

259    variants from the ClinVar database (**Supplementary Fig. 3c)**, 2,840 of which were identified from

260    patients who were tested for hereditary breast cancer[28]. To systematically assess variants' impact

261    on cell fitness, we designed two libraries: one to introduce reference alleles (Ref library) and

262    another to introduce alternative alleles (Alt library) targeting the selected variants (**Fig. 3a**)

263    (**Supplementary Table 3**). 250 non-targeting pegRNA/ngRNA pairs were added as negative

264    controls, respectively. For the Alt library, 115 pegRNA/ngRNA pairs introducing stop codons

265    (iSTOPs) in 23 MCF7 growth-related genes were included as positive controls, while

266    pegRNA/ngRNA pairs introducing reference sequences were used for those loci in the Ref library.

267    The cloned plasmids were packaged into lentiviral libraries and transduced into MCF7-nCas9/RT

268    cells. Cells were collected 2 and 32 days post infection, and pegRNA/ngRNA pairs were amplified

269    and deep sequenced (**Fig. 3b**). PRIME replicates using either Ref or Alt library (n = 4) were

270    reproducible at the read count level (**Supplementary Fig. 4a**).

**a**　Genome-scale CRISPR/Cas9
fitness screen in MCF7 cell line
(PMID: 30971826)

Bottom
1500 genes

Top
1000 genes

+

Three positive hits from base editing
screen in MCF7 cell line
(PMID: 33606977)
(*BARD1, BRCA1, BRCA2*)

-4  -3  -2  -1  0  1  2
$\text{Log}_2$(fold change)

MCF7 growth-related genes

**b**

**SNPs from GWAS of breast cancer (PMID: 29059683)**
1. MCF7 growth-related genes overlapping with GWAS identified breast cancer susceptibility genes (PMID: 31911677)
2. Within the locus +/- 500 kbp of each transcription start site
3. SNPs with GWAS $P < 1 \times 10^{-5}$, minor allele frequency $< 0.02$, and odds ratios $< 0.9$ or $> 1.2$

**SNPs from GWAS in Latinas (PMID: 25327703)**
1. SNPs with GWAS $P < 1 \times 10^{-5}$ in the *ESR1* locus

1610 SNPs in 14 risk loci

1304 SNPs with designed pegRNA/ngRNA pairs

**c**　**ClinVar database**
(Version 2021-12-25)

Filter the variants with
1. Breast cancer risk genes (PMID: 23000897, 26451490, 31911677, 33471991) and MCF7 growth-related genes
2. Single nucleotide variant
3. Clinical significance (Benign, Pathogenic, Uncertain significance)
4. Number of submitter $> 3$ for *BARD1, BRCA1, BRCA2, RAD51C, RAD51D, PTEN*

5310 variants for 17 risk genes
Benign: 672
Pathogenic: 1115
Uncertain significance: 3523

3699 variants with designed pegRNA/ngRNA pairs
Benign: 422
Pathogenic: 745
Uncertain significance: 2532

271

272　**Supplementary Figure 3. Strategies for prioritizing genomic loci and clinical variants.** (a)

273　The MCF7 growth-related genes were selected from the CRISPR/Cas9 knockout screen and

274　base editing screen in MCF7 cells. (b) The strategy used for selecting breast cancer-related

275　SNPs. (c) The strategy used for selecting clinical variants.

276   From Alt library screens, 33.04% (38/115) of iSTOPs showed a significant cell fitness
277 effect (FDR < 0.05), which is comparable to the 31.8% positivity rate of iSTOPs for common
278 essential genes reported from the base editing screen in MCF7 cells[29]. Furthermore, the fold
279 changes for iSTOPs were highly correlated with those for sgRNAs from MCF7 CRISPR knockout
280 screens of the same genes[23] (**Supplementary Fig. 4b**). More pegRNA/ngRNA pairs were
281 depleted (FDR < 0.05, Alt screen n = 322 and Ref screen n = 337) than enriched (FDR < 0.05,
282 Alt screen n = 148 and Ref screen n = 209) (binomial test, $P = 4.78 \times 10^{-8}$ for Alt screen and $P =$
283 $6.85 \times 10^{-16}$ for Ref screen) for both Alt and Ref screens on day 32 compared to day 2
284 (**Supplementary Fig. 4c and d, Supplementary Table 4, 5**). Theoretically, when a designed
285 peg/ngRNA pair matches the wild type MCF7 genotypes, they should have no effect on cell
286 growth. Notably, however, certain pegRNAs matching the wild type MCF7 genotype, exhibited
287 significant effects on cell growth beyond what was predicted, while the proportion of significant
288 hits for each genotype group were independent of initial MCF7 genotypes (Chi-square test $P =$
289 0.9998 on the Ref library and $P = 0.999$ on the Alt library, Cochran-Mantel-Haenszel test $P =$
290 0.9665 for the Ref library and Alt library together). For example, in the Ref library, 11.2% (59 out
291 of 528) of pegRNAs at sites with a Ref/Ref MCF7 genotype exhibited significant depletion, similar
292 to the 10.2% (55 out of 540) at heterozygous sites and 7.9% (18 out of 227) at Alt/Alt genotype
293 sites (**Fig. 3c**). These changes at sites where alleles were not expected to change suggests the
294 presence of undesired consequences of constitutive nCas9 expression, similar to CRISPR
295 inhibition (CRISPRi) once editing machinery is recruited to target sites[30]. To test for potential
296 CRISPRi activity of nCas9 in PE, we compared the results between iSTOPs in the Alt library and
297 the corresponding pegRNA/ngRNA pairs in the Ref library. While pegRNAs in the Ref library
298 exhibited smaller effects on Day 32 compared to iSTOPs targeting the same loci, they were still
299 depleted on Day 32, confirming unintended consequences due to nCas9 occupancy at target
300 genomic loci (**Supplementary Fig. 4e**). Combined, we found that prolonged PE expression
301 exhibits undesired activity similar to CRISPRi, a crucial factor for consideration when analyzing
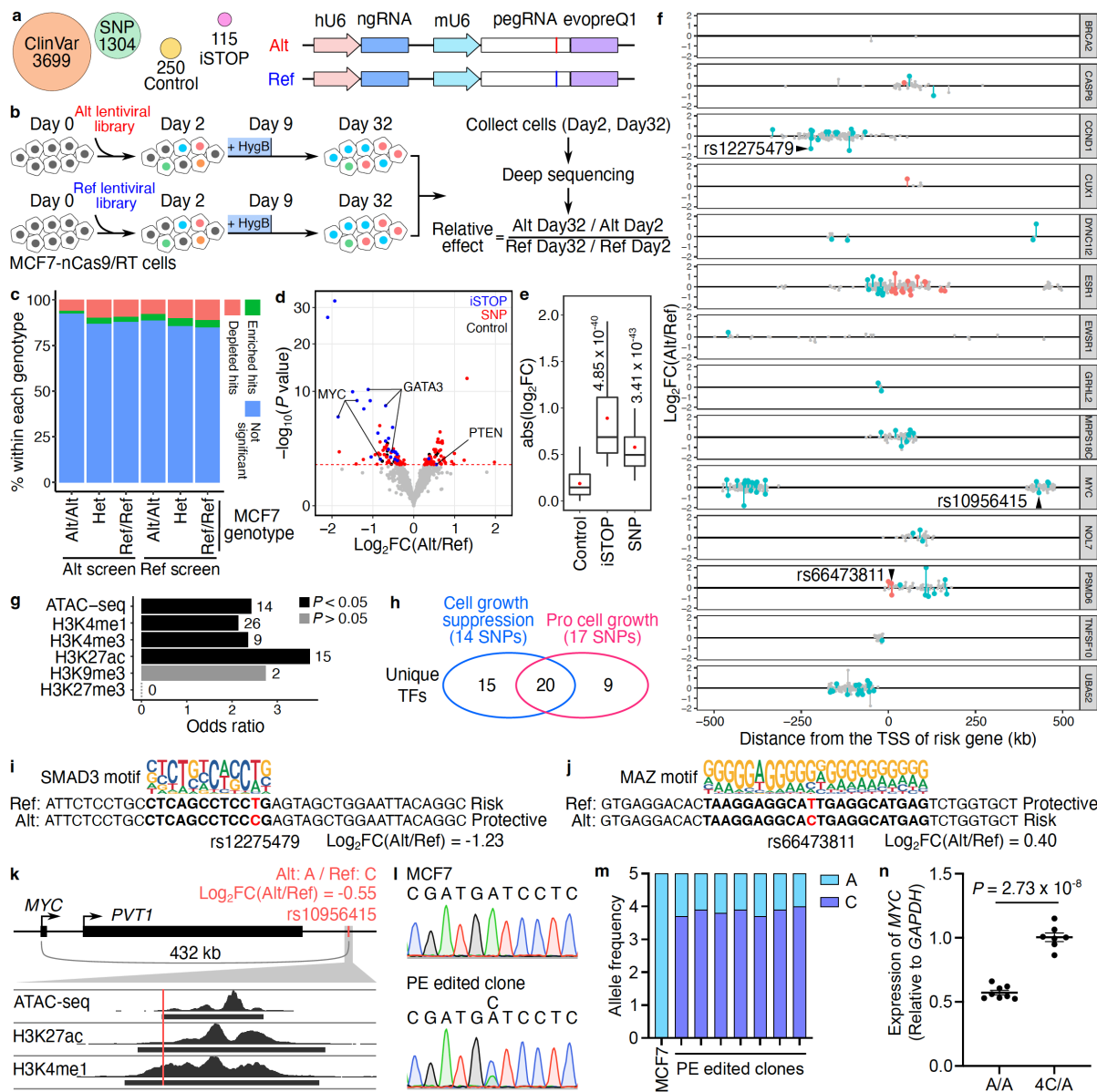302 lentivirus-mediated PE screens.

303   To correct for this undesired PE activity, we compared the ratio of FC for each
304 pegRNA/ngRNA pair from Alt and Ref screens by DESeq2[31]. We determined functional SNPs
305 based on their relative impact on cell growth between Alt and Ref PEs. In total, 56 SNPs with Ref
306 alleles and 47 SNPs with Alt alleles were identified to promote cell growth ($P < 0.05$, empirical
307 significance threshold to control type-I error at 5%, **Supplementary Fig. 4f**, **Fig. 3d**, and
308 **Supplementary Table 4)**. As expected, identified functional SNPs had smaller effect sizes than
309 stop codons and significantly larger effect sizes than negative control PEs (**Fig. 3e**). Additionally,

310    iSTOPs for genes promoting cell growth, such as *MYC* and *GATA3*, were depleted, while the

311    iSTOP for the cell growth suppressor *PTEN* was enriched, validating our analysis approach (**Fig.**
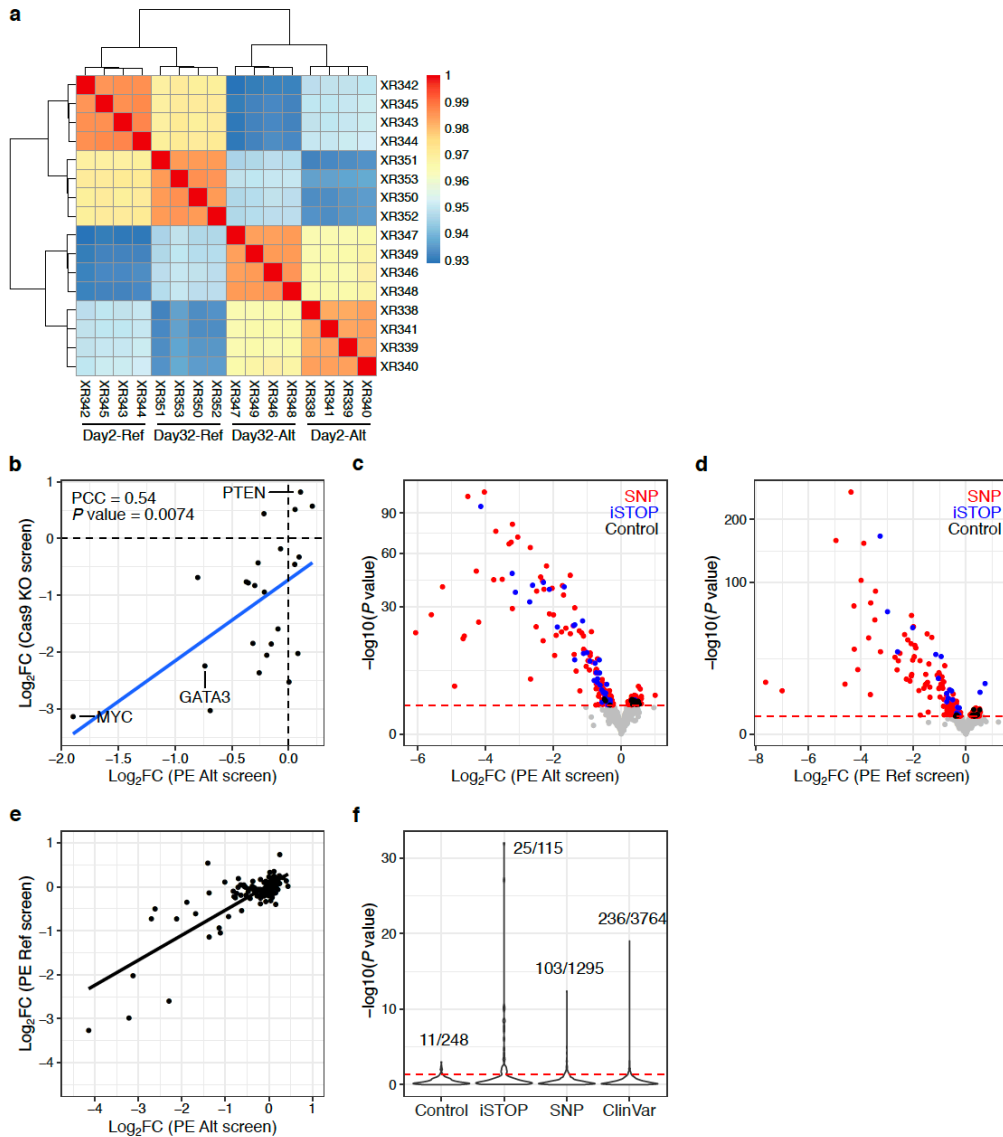
312    **3d**).

313         Since risk variants can either be the Ref or Alt allele, we further annotated functional SNPs

314    based on genetic annotation of breast cancer risk variants. Since most GWAS SNPs are likely

315    not causal, we expected that only a fraction of the 1,304 tested SNPs would exhibit a biological

316    effect. We calculated the mean likelihood of a variant being causal using CAVIAR and found that

317    the mean expectation for a variant being causal was ~8.9% when we made the assumption of

318    only one causal variant in each linkage disequilibrium (LD) clump. If we allowed for more than

319    one causal variant in each LD clump the mean probability of being causal for the variants was

320    ~13.0%. Compared to the reference allele, 50 risk SNPs' alternative alleles were pro-growth, and

321    53 risk SNPs' alternative alleles reduced cell growth (**Fig. 3f**). 18.45% (19/103) of the functionally

322    validated risk SNPs were located within the risk gene's body. The rest were located in distal

323    regions with an average distance of 185.8 kb from the risk gene's TSS (**Fig. 3f**). All tested loci

324    contained at least one SNP with a significant effect on cell growth, except for the *BRCA2* locus,

325    in which only 2 SNPs were tested. Finally, identified functional SNPs were significantly enriched

326    for active chromatin marks (two-tailed Fisher's exact test, $P < 0.05$), including ATAC-seq,

327    H3K27ac, H3K4me1, and H3K4me3 signals, relative to their corresponding genomic background

328    (1 Mbp surrounding selected cell growth genes) (**Fig. 3g**).

329         To explore potential mechanisms for functional SNPs' regulation of cell fitness changes,

330    we searched candidate TF binding motifs against the human motif database HOCOMOCO[19] using

331    40 bp regions centered on 103 identified functional SNPs. We retrieved 281 and 391 motifs (FDR

332    < 0.05 and TF expression > 1 FPKM) containing Alt and Ref alleles, respectively. After removing

333    redundant motifs for each SNP locus, we identified 90 TF binding sites for 35 unique TFs

334    associated with the cell growth suppression phenotype ($\log_2FC(Alt/Ref) < 0$) and 55 sites for 29

335    unique TFs associated with the pro cell growth phenotype ($\log_2FC(Alt/Ref) > 0$) (**Fig. 3h** and

336    **Supplementary Table 6**). In particular, the Alt allele (protective allele), rs12275479 (T>C) at the

337    *CCND1* locus disrupts the SMAD3 binding motif and is associated with reduced cell growth in our

338    screens, consistent with the TGFβ-SMAD3 axis decreased the number of mammosphere-

339    initiating cells in MCF7[32] (**Fig. 3f** and **i**). In another example, we found that a MAZ binding site of

340    MAZ is affected by the rs66473811 (T>C) Alt allele at the PSMD6 locus. MAZ is a transcription

341    factor that promotes breast cancer cell proliferation via driving tumor-specific expression of

342    *PPARγ1* gene and regulating *MYC* expression[33, 34] in line with that Alt allele being the risk allele

343    (**Fig. 3f** and **j**). To validate our PRIME results, we selected rs10956415 from the *MYC* locus, which

344     exhibited a moderate effect on cell growth in the screen ($Log_2FC(Alt/Ref) = -0.55$) (**Fig. 3f**).

345     rs10956415 is located in a candidate enhancer region 432 kb downstream of *MYC* (**Fig. 3k**).

346     MCF7 cells are homozygous for the alternative allele (A) at the rs10956415 locus, which has a

347     copy number of five in this cell line[35] (**Fig. 3l** and **m**). Using prime editing, we converted 4 copies

348     of the alternative allele (A) to the reference allele (C) in seven independent clones, yielding a

349     43.2% average increase in *MYC* expression compared to unedited cells with 5 copies of A alleles

350     (**Fig. 3m** and **n**). Since *MYC* expression level positively correlates with MCF7 cell growth[36], the

351     *MYC* expression of the PE edited clones is consistent with the cell growth phenotype of

352     rs10956415 observed in the screening. Together, these results support the use of PRIME to

353     functionally characterize GWAS-identified variants.



354

355 **Figure 3. PRIME reveals functional SNPs associated with breast cancer.** (a) Alt and Ref

356 library design overview. In the design, we included breast cancer-associated variants (SNP),

357 clinical variants (ClinVar), introduced stop codons (iSTOP), and non-targeting controls. For each

358 variant, pegRNA/ngRNA pairs introducing either the Alt or Ref allele were designed. (b) Workflow

359 of PRIME with Alt and Ref libraries. MCF7-nCas9/RT cells were infected with either lentiviral

360 library. Cells were collected on days 2 and 32 post-infection. The abundance of pegRNA/ngRNA

361 pairs in the samples collected on days 2 and 32 were deep sequenced. The relative effect of each

362 variant was determined based on its relative impact on cell growth between Alt versus Ref alleles.

363 (c) The percentage of significant hits (FDR < 0.05) identified from Alt and Ref screens for Alt/Alt,

364 Het, and Ref/Ref genotypes in MCF7. (d) The functional SNPs (red) with either a positive or a

365 negative impact on cell growth were determined by their relative effect in the Alt versus Ref

366 screens. Blue dots represent significant iSTOPs, and black dots represent controls. The red

367 dashed line indicates 0.05 FDR. (e) Absolute effects of identified functional iSTOPs and SNPs

368 are higher than the effects of negative controls (*P* values were calculated by two-tailed two-

369 sample t-test). (f) The genomic distance of SNPs tested at each risk locus relative to each gene's

370 TSS. Red dots are functional SNPs within gene bodies, blue dots are functional SNPs in distal

371 regions, and gray dots are SNPs with non-significant effects. (g) Relative enrichment of genomic

372 features for identified functional SNPs (*P* values were calculated by two-tailed Fisher's exact test).

373 The numbers of SNPs overlapping each genomic feature are labeled next to each bar. (h) Venn

374 diagram showing the numbers of unique transcription factors (TFs) with differential binding sites

375 centered on functional SNPs. The numbers of SNPs that alter TF binding sites are also in the

376 parentheses. (i, j) Examples of functional SNPs disrupting TF binding sites. (i) The Alt protective

377 allele of rs12275749 (position shown in f) affects the SMAD3 binding site and (j) The Alt risk allele

378 of rs66473811 (position shown in f) is matched with the MAZ binding motif. (k) rs10956415 located

379 within a candidate enhancer region overlapping with ATAC-seq, H3K27ac and H3K4me1 peaks

380 in MCF7 cells. (l) Representative Sanger sequencing results for the rs10956415 locus in unedited

381 MCF7 cells and a PE edited clone. (m) Allele frequencies of alternative (A) and reference (C)

382 alleles of rs10956415 in unedited MCF7 cells and PE edited clones. (n) Relative *MYC* expression

383 in control clones and PE edited clones ($P = 2.73 \times 10^{-8}$, two-tailed two-sample t-test).

**Supplementary Figure 4. Quality control and primary analysis of disease variants.** (a) Heatmap with pairwise correlations and hierarchical clustering of read counts from PRIME. (b) Pearson correlations between the $\log_2$(fold change) of iSTOPs in the Alt library screen and the $\log_2$(fold change) of gRNAs in the CRISPR/Cas9 knockout screen for each target gene. (c) Volcano plot of the results from the Alt library screen. (d) Volcano plot of the results from the Ref library screen. (e) The $\log_2$(fold change) for each iSTOP from the Alt and Ref library screens. (f) Violin plot showing the 5% FDR cutoff used for the relative effect analysis comparing the Alt and Ref libraries. Numbers above peaks indicate the significant data points versus the total data points in each category when using 5% FDR. We used the 5% percentile of $P$ values from negative controls as the empirical significance threshold to achieve a false discovery rate (FDR) of 5% indicated by the red dashed line in d-f.

396

**PRIME can characterize clinical variants of uncertain significance**

398     Genetic variants detected in clinical samples provide a valuable resource for

399     understanding the etiologies of human diseases. However, many clinically discovered variants

400     are annotated as Variants of Uncertain Significance (VUS) due to unpredictable functional

401     consequences, even in well-characterized protein-coding genes. To assess the capacity of

402     PRIME to functionally annotate VUS using MCF7 growth phenotypes, we designed

403     pegRNA/ngRNA pairs for 2,532 VUS, 745 pathogenic variants, and 422 benign variants for 17

404     genes (**Supplementary Fig. 3c** and **Supplementary Table 3)**. 76.78% of the variants tested

405     were from breast cancer patients (**Supplementary Table 3**). By comparing the relative effect

406     sizes of each Alt and Ref allele pair, we identified 236 functional clinical variants affecting cell

407     growth in 15 genes, including 49 pathogenic variants, 156 VUS, and 31 benign variants (**Fig. 4a**

408     and **Supplementary Table 5**). The average effect sizes for pathogenic variants, VUS, and benign

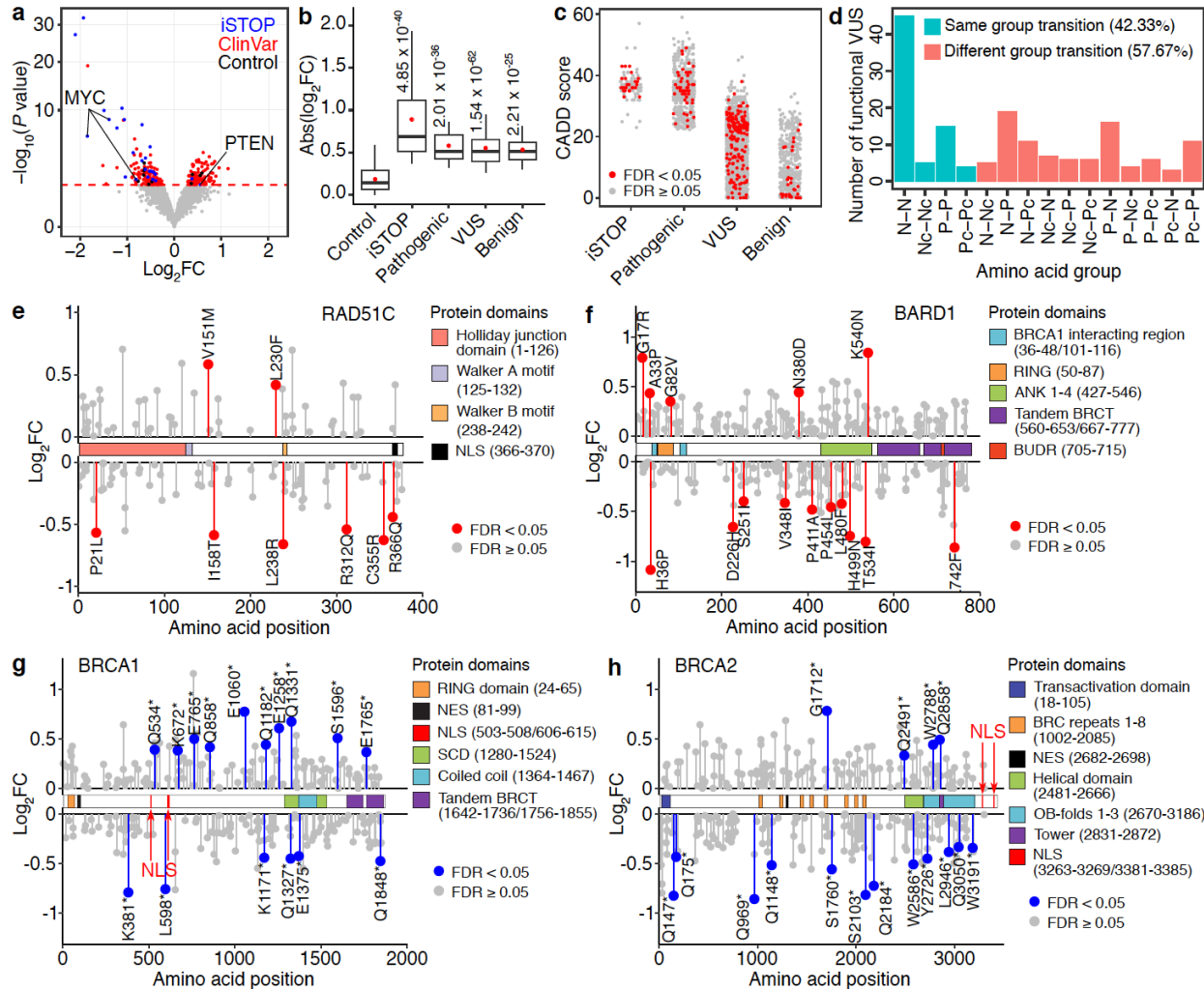409     variants were between that of negative controls and iSTOPs (**Fig. 4b**).

410     Several computational metrics have been used to assess the deleteriousness of variants[37,]

411     [38]. One such method is CADD, which integrates diverse genome annotations into a single,

412     quantitative score estimating the relative pathogenicity of human genetic variants[37]. iSTOPs and

413     pathogenic variants have similarly high CADD scores relative to other categories (**Fig. 4c**). The

414     CADD scores for the VUS and benign variants exhibit a broad distribution with median scores

415     much lower than those of iSTOPs and pathogenic variants. Interestingly, the CADD scores for

416     identified functional variants within the VUS or benign variant groups did not have higher CADD

417     scores as expected, indicating the limitation of solely relying on computational prediction for

418     variants annotation and underscoring the importance of validating clinical variants with functional

419     assays, even for those located in well-studied protein-coding genes. For example, one benign

420     variant in BARD1 (Arg378Ser) with a low CADD score (CADD = 4.317) would not be classified as

421     functional. However, this variant exhibited a significant cell growth suppression effect in MCF7

422     cells based on our screening results. BARD1 (Arg378Ser) can impair the nuclear localization of

423     the BRCA1/BARD1 complex, and synergistically promote tumor formation with BARD1

424     (Pro24Ser) *in vivo*[39]. Furthermore, most of the identified functional VUS were missense variants,

425     and about half of the significant VUS from our screens changed amino acid type within the same

426     group based on polarity (**Fig 4d**), complicating the determination of their molecular

427     consequences. Our results offer novel insights into the potential roles of clinical variants in

428     disease pathogenesis through their modulation of cell fitness, and provide annotations for VUS

429     and benign variants previously uncharacterized.

430     Functional and structural domains are integral contributors to protein function. 60% of the

431     functional VUS identified are located within an annotated protein domain in the UniProt

432     database[40], supporting their pathogenicity. For example, we identified 8 VUS in *RAD51C* (**Fig.**

433     **4e**), a cancer susceptibility gene and an essential gene for MCF7 survival. Two variants, one

434     (Pro21Leu) in the RAD51C functional domain (amino acid: 1-126) for Holliday junction processing

435     and the other (Arg366Gln) in the NLS region (amino acid: 366-370), were associated with reduced

436     cell growth by our screens (**Fig. 4e**). We also identified functional variants that were not located

437     in any annotated domain, including a functional RAD51C VUS (Arg312Gln) associated with a

438     phenotype of reduced MCF7 growth (**Fig. 4e**). Since Arg312Trp in RAD51C results in homologous

439     recombination deficiency and reduced colony formation phenotypes in MCF10A cells, and

440     abolishes RAD51C-RAD51D interaction[41], Arg312Gln may produce a similar pathogenic

441     consequence on protein function. When comparing the RAD51C sequence with other RAD51

442     family proteins, we observed functional VUS were located in both conserved and non-conserved

443     amino acids (**Supplementary Fig. 5a**), underscoring the challenge of predicting variant function

444     based solely on protein sequence conservation.

445     Protein-protein interaction (PPI) is another essential functional activity in many biological

446     processes. In this study, we also identified functional VUS located in protein binding regions with

447     the potential to affect PPI. For example, BARD1 interacts with BRCA1 through RING domains,

448     and BRCA1-BARD1's ubiquitin ligase activity is indispensable for DNA double-strand break

449     repair[42, 43]. We identified a functional VUS (His36Pro) in the BARD1 RING domain (**Fig. 4f**),

450     suggesting the structural consequences of this clinical variant affecting BARD1-BRCA1

451     heterodimer formation (**Supplementary Fig. 5b**). Consistent with these findings, AlphaFold

452     predicts that the His36Pro variant disrupts hydrogen bond formation between His36 in BARD1

453     and Asp96 in BRCA1 (**Supplementary Fig. 5c**).
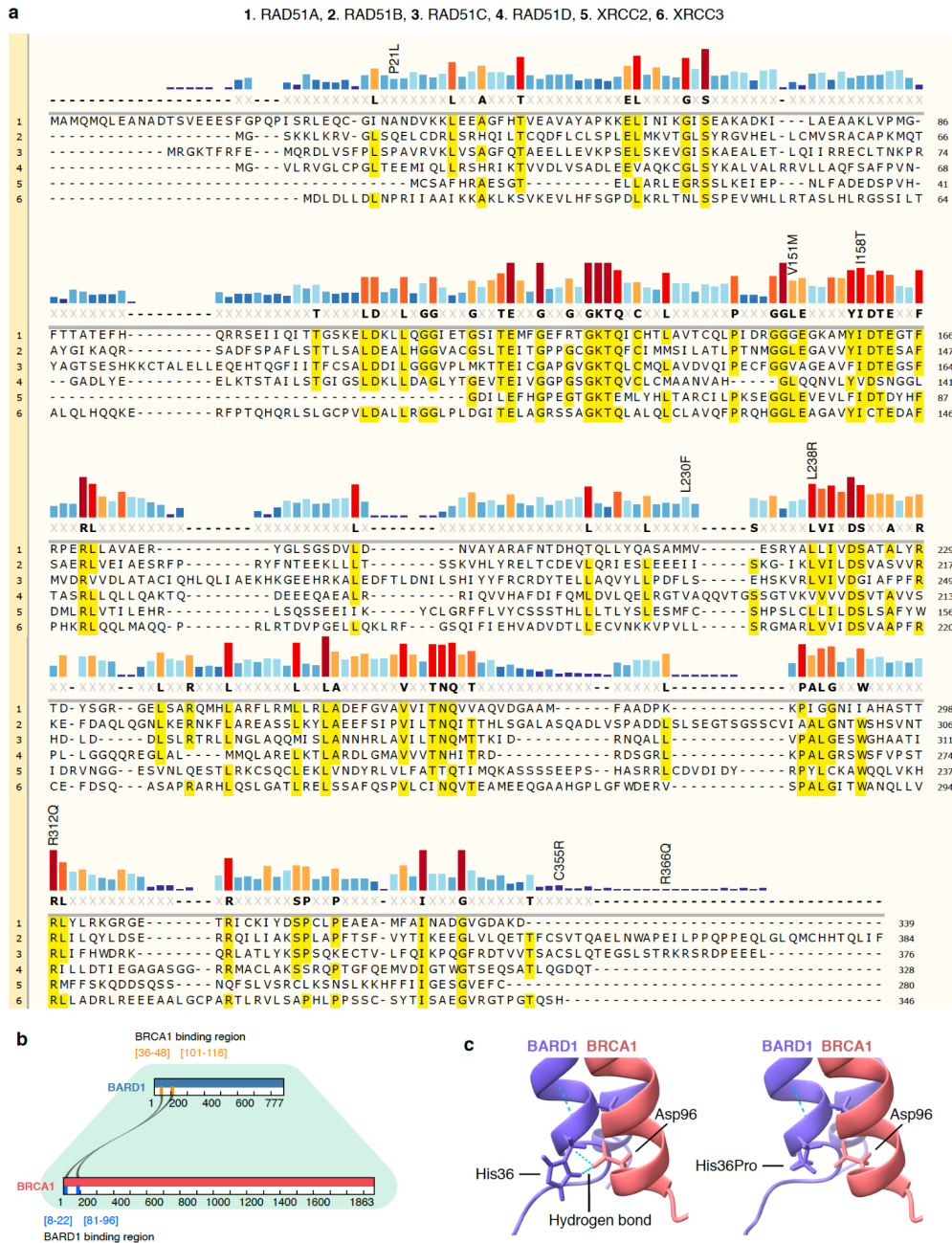
454     Nonsense mutations can generate new stop codons and truncated proteins. Although

455     most are annotated as pathogenic variants in ClinVar, the functional consequences of many

456     remain uncharacterized[28]. In our screens, 563 nonsense clinical variants were tested in 13 breast

457     cancer risk genes with 38 variants identified as positive hits in 7 genes. Remarkably, 39.47%

458     (15/38) exhibited unexpected phenotypes compared to the knockout phenotypes of cell death of

459     these genes. Specifically, a similar number of functional nonsense variants in *BRCA1* (n = 15)

460     and *BRCA2* (n = 16) (**Fig. 4g, h**) were identified; however, 60% (9/15) in *BRCA1* could promote

461     MCF7 cell growth compared to 25% (4/16) in *BRCA2*. After locating variants within BRCA1 and

462     BRCA2, we noticed that truncated proteins resulting from all gain-of-function nonsense variants

463     in BRCA1 still retained their NLS. These results were confirmed by a different nonsense mutation

464    at Q858, located downstream of the NLS in BRCA1, which resulted in truncated BRCA1 with NLS

465    and increased cell growth of MCF7[29]. However, for all of the functional variants identified in

466    BRCA2, their NLSs were located at the c-terminus[44] and were thus removed from the truncated

467    proteins, leading to the loss of BRCA2 nuclear localization. Collectively, these results demonstrate

468    the capability of PRIME to functionally characterize some nonsense mutations.



469

470    **Fig 4. Functional clinical variants identified using PRIME.** (a) Functional clinical variants (red)

471    with either a positive or a negative impact on cell growth were determined by relative effects on

472    cell fitness between Alt and Ref alleles. Blue dots represent significant iSTOPs, and black dots

473    represent negative controls. The red dashed line indicates 5% FDR. (b) Effect sizes of identified

474    functional iSTOPs and clinical variants are larger than that of negative controls ($P$ values were

475    calculated by two-tailed two-sample t-test). Box plots indicate the median, IQR, Q1 − 1.5 × IQR,

476    and Q3 + 1.5 × IQR. Red dots indicate the mean. (c) CADD scores for iSTOPs and clinical

477 variants. (d) Number of identified functional VUS causing each amino acid group transition. (N, 478 Nonpolar; P, Polar; Pc, Positively charged; Nc, Negatively charged). (e,f) Lollipop plots of 479 functional VUS in *RAD51C* and *BARD1* mapped to their canonical isoforms. The identified 480 significant VUSs are labeled in red. Their effects on cell growth are indicated by fold changes. 481 (g,h) Lollipop plots of the nonsense variants in *BRCA1* and *BRCA2* mapped to their canonical 482 isoforms. The identified significant hits are labeled in blue. Their effects on cell growth are 483 indicated by fold changes.



484

485 **Supplementary Figure 5. Examples of functional VUS with their potential consequences.**

486 (a) Sequence conservation of RAD51 family proteins. Alignment of RAD51 family proteins using

487 MUSCLE. Functional VUS identified by PRIME in RAD51C are labeled. (b) Graphic showing the

488 binding regions between BARD1 and BRCA1. (c) The Alphafold predicted protein structure of the

489 BARD1 and BRCA1 complex. Two hydrogen bonds were identified between wild type His36 in

490 BARD1 and Asp96 in BRCA1, but lost following the BARD1 His36Pro mutation.

491

492 **Discussion**

493 In this study, we describe a new genomic screening method, PRIME, to interrogate DNA

494 function at base-pair resolution by adopting and optimizing 'search-and-replace' prime editing[5, 9].

495 We demonstrate the success of pooled prime-editing screens to identify essential nucleotides in

496 a *MYC* enhancer via saturation mutagenesis screen, the functional characterization of 1,304

497 breast cancer-associated risk SNPs, and provide accurate annotation for 3,699 clinical variants.

498 Our study offers a novel strategy to elucidate genome function at an unprecedented precision and

499 scale. The broad applications demonstrated in this work suggest that PRIME can significantly

500 augment the functional characterization toolbox and advance our ability to elucidate the roles of

501 disease-associated variants in the human genome.

502 Our analyses show that lentiviral installation of PE yields long-lasting expression of nCas9,

503 pegRNA, and ngRNAs, but can result in unwanted sequence-specific repression similar to

504 CRISPRi. This bias must be corrected to produce accurate base-pair resolution annotations.

505 When assessing the functional impact of a variant, pegRNA controls should be included to

506 introduce other alleles at the same locus. Our study normalized sequence-specific repression

507 bias by comparing the differential effects on cell survival of all base pair substitutions at each

508 locus in the *MYC* enhancer, and between Alt and Ref alleles for disease variants. Additional

509 improvement could be achieved through controlled nCas9 expression duration. For example, a

510 doxycycline-inducible nCas9 could be selectively expressed when editing is needed and

511 reversibly turned off afterwards. In addition to establishing and optimizing PRIME, we defined

512 sensitive base pairs (SBPs) and core sequences for a *MYC* enhancer's function. We generated

513 a functional PWM for this enhancer by leveraging effect sizes for all possible substitutions at each

514 base from the screens. The functional PWM enabled us to accurately predict TF binding sites

515 within the enhancer, providing critical annotations for delineating *MYC* activation in MCF7 cells.

516 Interpreting the effect of inherited genetic variations will dramatically advance our ability

517 to predict an individual's disease risk. However, utilizing GWAS data for risk prediction is still

518 limited without substantial functional annotation. In this study, 7.9% of the 1,304 tested GWAS

519    breast cancer variants, and 6.2% of the 2,532 tested VUS were identified as significant hits with
520    functions linked to MCF7 growth phenotypes. Our results demonstrate the feasibility of PRIME
521    for functionally characterizing individual variants. The impact of variants was context-specific and
522    our findings were limited to assessing variants with growth phenotype related functions in MCF7
523    cells. Other ClinVar did not show changes in our functional assay likely have functional
524    consequences for breast cancer susceptibility genes in a different cell type or other biological
525    processes.

526    Future work employing different phenotypic screening readouts across multiple cell lines
527    will provide new insights into variant function. For example, screens that identify variants
528    associated with differential drug treatment responses will help construct better predictive models
529    for an individual's unique benefits and risks from therapeutics. Screens of variants with readouts
530    directly linked to physiological functions e.g. endolysosomal activities in microglia or synaptic
531    activities in neurons using iPSC models will uncover functional variants associated with
532    neuropsychiatric diseases. In summary, our study provides a roadmap to advance functional
533    genomics toward the actionable disease prediction, prevention and treatment necessary to realize
534    personalized medicine.

535

543

544    **Author contributions**
545    X.R. H.Y., and Y.S. conceived the study. Y.S. and E.Z. supervised the study. X.R. and H.Y.
546    designed PRIME screens. X.R. H.Y. C.B. Y.S. M.N. M.A.T. and V.N. performed experiments
547    under the supervision of Y.S. X.R. HY, J.L.N, Y.S. and J.C. performed computational analysis
548    under the supervision of Y.S. Y.L. and E.Z.  Y.S. X.R. and H.Y. prepared the manuscript with
549    input from all other authors.

550

551    **Competing interests statement**
552    X.R., H.Y., and Y.S. have filed a patent application related to pooled prime editing screens.

553

**Code availability statement**

A copy of the custom code used for data analysis and figure generation in this study is available upon request.

557

**Supplementary Tables**

**Supplementary Table 1. pegRNA and ngRNA oligo sequences and their fold changes in *MYC* enhancer.**

**Supplementary Table 2. TF motif analysis for alleles based on functional data from PRIME.**

**Supplementary Table 3. pegRNA and ngRNA oligo sequences for SNP and ClinVar.**

**Supplementary Table 4. PRIME results for breast cancer-associated variants.**

**Supplementary Table 5. PRIME results for clinical variants.**

**Supplementary Table 6. TF motif analysis for alleles with functional SNPs for breast cancer.**

567

**Methods**

**Cell culture**

MCF7 cells were cultured in Dulbecco's Modified Eagle Medium (DMEM) (Gibco, 10569010) supplemented with 10% fetal bovine serum (FBS) (HyClone, SH30396.03), and were passaged with trypsin-EDTA (Gibco, 25200072). All cells were cultured with 5% $CO_2$ at 37°C and verified to be free of mycoplasma using the MycoAlert Mycoplasma Detection Kit (Lonza, LT07-218). Wild type MCF7 cells were a gift from Howard Y. Chang's lab. The MCF7-nCas9/RT cell line was generated by lentiviral transduction of cells with a cassette expressing the nickase Cas9 (nCas9) Moloney murine leukemia virus reverse transcriptase (M-MLV RT) fusion protein. The infected MCF7 cell pool was treated with puromycin (2.5 μg/ml) for two weeks. Then, single cells were sorted into 96-well plates with one cell per well by fluorescence-activated cell sorting (FACS) to generate a clonal MCF7-nCas9/RT cell line. nCas9/RT expression levels were quantified in each clone via RT-qPCR, and normalized to the dCas9 expression level in a WTC11 doxycycline-inducible dCas9-KRAB iPSC line[45, 46].

582

**Functional characterization of a *MYC* enhancer by CRISPR deletion**

Two sgRNAs were designed to knock out a MCF7 enhancer (chr8:128,141,747-128,142,627, hg38) (sg1: GAAGTTGTAAGTATAGCGAG, sg2: AGTGCCTGGCACAAGGCAGA). sgRNAs were synthesized *in vitro* using the Precision gRNA Synthesis Kit (Invitrogen, A29377) according

587 to the manufacturer protocol and concentrations were quantified with Nanodrop. To deliver
588 genome editing machinery, 100 pmol of Cas9-NLS protein (QB3 MacroLab in University of
589 California, Berkeley) and 120 pmol of *in vitro* synthesized gRNA were electroporated into 250,000
590 MCF7 cells with the P3 primary nucleofection solution (Lonza, V4XP-3024), using the DN-100
591 Lonza 4D-Nucleofector program. Cells were then plated into 6-well plates and cultured for 2 days,
592 followed by plating into 96-well plates to pick single clones. Successful knockout clones were
593 identified by genomic PCR with the primers forward: CACCAGGACTTGAAGGCAGC and
594 reverse: CACTTCCCAACCTCAGTTTCC. RT-qPCR was used to quantify *MYC* expression (*MYC*
595 forward primer: GTCCTCGGATTCTCTGCTCT, reverse primer
596 ATCTTCTTGTTCCTCCTCAGAGTC) and normalized to the *GAPDH* expression level (*GAPDH*
597 forward primer: ATTCCATGGCACCGTCAAGG, reverse primer
598 TTCTCCATGGTGGTGAAGACG).

599

**Cloning of prime editing plasmids**

601 To construct the lentiV2-EF1a-nCas9/RT plasmid, we first excised the U6-sgRNA cassette from
602 the lentiCRISPR v2 plasmid (Addgene, 52961) by dual KpnI and EcoRI digestion followed by
603 blunt end ligation. We further replaced the Cas9 cassette with an nCas9/M-MLV-RT cassette from
604 the pCMV-PE2 plasmid (Addgene, 132775). The lentiV2-pegRNA and lentiV2-ngRNA plasmids
605 were constructed by replacing the Cas9 and Puromycin sequences in the lentiCRISPR v2 plasmid
606 (Addgene, 52961), with hygromycin B and EGFP sequences. RNA motifs and sgRNA scaffolds
607 were further integrated by Gibson assembly.

608

**Testing prime editing efficiency**

610 To assess prime editing efficiencies at the *EMX1* and *FANCF* loci, we cloned paired
611 pegRNAs/ngRNAs into individual vectors. For lentivirus co-infection testing, we first infected
612 MCF7 cells with EF1a-nCas9/RT lentivirus followed by treatment with puromycin (2.5 µg/ml;
613 Sigma-Aldrich, P8833) for 2 weeks to eliminate uninfected cells. Then, EF1a-nCas9/RT-infected
614 cells were seeded in 24-well plates at 12,500 cells per well for pegRNA and ngRNA co-infection.
615 The infected cells were treated with hygromycin B (200 µg/ml; Gibco, 10687010) 48 hours after
616 infection, and were collected one week after infection for editing efficiency assessment. For
617 testing in the MCF7-nCas9/RT clonal line, we seeded cells in 24-well plates at 12,500 cells per
618 well, followed by lentiviral infection (pegRNA-mCherry and ngRNA-EGFP). Two days after
619 infection, mCherry and EGFP double-positive cells were isolated by FACS and cultured. Cultured
620 cells were then collected at 2 weeks and 4 weeks post-infection for editing efficiency assessment.

621   Genomic DNA was then extracted from each sample using the Wizard genomic DNA purification
622   kit (Promega, A1120). Genomic sites of interest were amplified from purified genomic DNA and
623   amplicons were sequenced on the Illumina NovaSeq 6000 platform. Briefly, sequencing libraries
624   were prepared using DNA primers amplifying target genomic loci of interest for the first round of
625   PCR (PCR1). Then, DNA primers containing index adapters were used for the second round of
626   PCR (PCR2) to add these adapters to PCR1 amplicons. Finally, dual indexing primers were used
627   for the third round PCR (PCR3) to add Illumina indexes to each PCR2 amplicon. Alignment of
628   amplicons to reference sequences was performed using CRISPResso2[47]. For all prime editing
629   efficiency quantification, wild-type and edited amplicon frequencies were quantified using a 21 bp
630   window centered on either the 1 bp wild-type or edited sequence. The remaining amplicons were
631   classified as indels.

632

633   **SNP prioritization**
634   We selected 14 MCF7 growth-related genes overlapping with GWAS identified breast cancer
635   susceptibility genes[26]. For each gene, we selected SNPs using the GWAS results from the Breast
636   Cancer Association Consortium[25]. We identified genome-wide significant SNPs with GWAS $P <$
637   $1 \times 10^{-5}$, minor allele frequency < 0.02, and odds ratios < 0.9 or > 1.2 (representing approximately
638   the top and bottom quartiles of the odds ratio distribution for SNPs meeting the location, $P$ value,
639   and MAF thresholds) for association with breast cancer within the locus +/- 500 kb of each
640   transcription start site. We also separately selected SNPs with GWAS $P < 1 \times 10^{-5}$ in the *ESR1*
641   locus using GWAS results from a Latina population[48]. We determined linkage disequilibrium (LD)
642   clumps among the selected SNPs using the LD Link R package[49] with an LD threshold of $R^2 >$
643   0.1. We then prioritized the most likely causal variants using CAVIAR[50], as those with a causal
644   posterior probability (> 0.1), the highest posterior probability (≤ 0.1), or most extreme odds ratio
645   in each haplotype block. We ran CAVIAR twice for each locus, once assuming only one causal
646   variant per LD clump, and again allowing for more than one causal variant in each LD clump.

647

648   **Clinical variant prioritization**
649   We retrieved clinical variants from the ClinVar database (accessed 2021-12-25), and all single
650   nucleotide variants (SNVs) were kept for the PRIME design (**Supplementary Fig. 3c**). We first
651   selected only the SNVs whose genes overlapped with breast cancer risk and MCF7 growth-
652   related genes. Next, we only retained SNVs in the benign, pathogenic and uncertain significance
653   categories. Further, for SNVs associated with *BARD1*, *BRCA1*, *BRCA2*, *RAD51C*, *RAD51D*, and
654   *PTEN*, we only retained the SNVs with more than three submitters, as there are thousands of

655    identified variants for these genes. Finally, our selection criteria yielded 5310 SNVs, of which we

656    successfully designed pegRNA/ngRNA pairs for 3699 SNVs.

657

658    **Design and construction of prime-editing libraries**

659    For nucleotide-resolution analyses of *MYC* enhancer function**,** paired pegRNAs/ngRNAs targeting

660    a 716 bp enhancer region were first designed using PrimeDesign's PooledDesign-Saturation

661    mutagenesis tool[51]. We optimized pegRNAs/ngRNAs pairs based on ngRNA pegRNA proximity

662    (more than 50 bp) and primer binding site (PBS) length (near 14 nt), redesigning the sequence

663    containing the BsmBI cutting sites (GAGACG, CGTCTC) or TTTTT. Next, we used GuideScan2

664    to assess the specificity and efficiency of each pegRNA and ngRNA spacer sequence. Spacer

665    sequences with low specificity were redesigned to improve the specificity. Finally, three different

666    pegRNA/ngRNA pairs were designed to target the same base pair for 93.0% (666/716) of the

667    substitutions. Each replicate pegRNA/ngRNA pair shared the same pegRNA and sgRNA spacer

668    sequences, and only the substitution alleles differed in the pegRNA extension sequence. To

669    design  positive control guides, we used pegIT[52] to generate pegRNA/ngRNA pairs which alter a

670    single base pair to introduce a stop codon within the *MYC* coding region. We selected the best

671    pegRNA/ngRNA pair for each position suggested by pegIT[52]. The *AAVS1* locus was selected as

672    the targeting pegRNA/ngRNA pair negative control region based on previous work[53], and guides

673    were designed as described above using PrimeDesign[51]. For non-targeting pegRNA/ngRNA

674    pairs, pegRNA and ngRNA spacer sequences and pegRNA extension sequences were selected

675    from the ENCODE non-targeting sgRNA reference data set

676    (https://www.encodeproject.org/files/ENCFF058BPG/). A guanine nucleotide was added to the 5'

677    end of all pegRNAs/ngRNAs with leading nucleotides other than G, to increase transcription

678    efficiency from the U6 promoter. We used the following template to link these component

679    sequences: 5'- CTTGGAGAAAAGCCTTGTTT[ngRNA-spacer]GTTTAGAGACG[5nt-random-

680    sequence]CGTCTCACACC[pegRNA-

681    spacer]GTTTTAGAGCTAGAAATAGCAAGTTAAAATAAGGCTAGTCCGTTATCAACTTGAAAA

682    AGTGGCACCGAGTCGGTGC[pegRNA extension]CCTAACACCGCGGTTC-3'.

683

684    Library oligos for the *MYC* enhancer screen were synthesized by Twist Bioscience and amplified

685    using the NEBNext High-Fidelity 2× PCR Master Mix (NEB, M0541L), forward primer:

686    GTGTTTTGAGACTATAAATATCCCTTGGAGAAAAGCCTTGTTT and reverse primer

687    CTAGTTGGTTTAACGCGTAACTAGATAGAACCGCGGTGTTAGG. To amplify paired

688    PegRNA/ngRNA library oligos for enhancer saturation mutagenesis, we employed emulsion PCR

689   (ePCR) to reduce recombination of similar amplicons during PCR. Briefly, ninety-six 20 µl ePCR

690   reactions were performed using 0.01 fmol of pooled oligos with NEBNext High-Fidelity 2× PCR

691   Master Mix (NEB, M0541S). Each 20 µl PCR mix was combined with 40 µl of oil-surfactant mixture

692   (containing 4.5 % Span 80 (v/v), 0.4 % Tween 80 (v/v) and 0.05 % Triton X-100 (v/v) in mineral

693   oil)[54]. This mixture was vortexed at maximum speed for 5 min, briefly centrifuged, and placed into

694   the PCR machine for amplification. Thermocycler settings were: 98 °C for 30 s, then 26 cycles

695   (98 °C 10 s, 60 °C 20 s, 72 °C 30 s), then 72 °C for 5 min, and finally a 4 °C hold. The ramp rate

696   for each step was 2°C/s. After PCR, individual reactions were combined and purified using the

697   QIAQuick PCR Purification Kit (Qiagen, 28104) following previously established guidelines[55].

698   Purified PCR products were then treated with Exonuclease I (NEB, M0568L) and purified using

699   1× AMPure XP beads (Beckman Coulter, A63881). The isolated ePCR products were then

700   inserted into a BsmBI–digested lentiV2-mU6-evopreQ1 vector via Gibson assembly (NEB,

701   E2621L). The assembled products were electroporated into Endura electrocompetent

702   Escherichia coli cells (Biosearch Technologies, 60242) and approximately 4,000 independent

703   bacterial colonies were cultured for each library. The resulting plasmid DNA was linearized by

704   BsmbI digestion, gel-purified, and ligated using T4 ligase (NEB, M0202M) to a DNA fragment

705   containing an sgRNA scaffold and the human U6 promoter. The resulting library was

706   electroporated into Endura electrocompetent Escherichia coli cells (Biosearch Technologies,

707   60242) and cultured as described above. The final plasmid library was extracted using the Qiagen

708   EndoFree Plasmid Mega Kit (Qiagen, 12381).

709

710   For the SNP and clinical variant screen Alt library, pegRNA/ngRNA pairs were designed using

711   PrimeDesign[51]. The sequences 200 bp upstream and downstream of each variant or iSTOP were

712   used as inputs for PrimeDesign. We generated initial pegRNA/ngRNA pairs using the following

713   parameters: number of pegRNAs per edit: 10, length of homology downstream: 10 nt, PBS length:

714   13 nt, maximum reverse transcription template (RTT) length: 50 nt, number of ngRNAs per

715   pegRNA: 10, ngRNA to pegRNA nicking distance: 50 and 75 bp. Next, a guanine nucleotide was

716   added to the 5' end of all pegRNAs/ngRNAs with leading nucleotides other than G to increase

717   transcription efficiency from the U6 promoter. pegRNA/ngRNA pairs containing BsmBI sites

718   (GAGACG, CGTCTC) or a TTTTT sequence in the pegRNA spacer, ngRNA spacer or pegRNA

719   extension were eliminated. pegRNA/ngRNA pairs were further selected to maximize specificity,

720   efficiency, and ngRNA to pegRNA distance while minimizing pegRNA to edit distance when

721   multiple pairs were available for the same locus. For non-targeting pegRNA/ngRNA pairs,

722   pegRNA spacer, ngRNA spacer and pegRNA extension sequences were selected from the

723 ENCODE non-targeting sgRNA reference data set

724 (https://www.encodeproject.org/files/ENCFF058BPG/). To design the Ref library, we used the

725 same pegRNA/ngRNA pairs as the Alt library, but replaced the alternative alleles in the pegRNA

726 extension sequences with the reference allele sequences. The final oligos adhered to the

727 following template architecture: 5′-CTTGTGGAAAGGACGAAACACC[ngRNA-

728 spacer]GTTTCGAGACG[6nt-random-sequence]CGTCTCTTGTTT[pegRNA-

729 spacer]gttttagagctagaaatagcaagttaaaataaggctagtccgttatcaacttgaaaaagtggcaccgagtcggtgc[pegR

730 NA extension]TTGACGCGGTTCTATCTAGTTAC-3′.

731

732 The Alt and Ref library oligos were synthesized by Twist Bioscience. The Alt and Ref plasmid

733 libraries were cloned separately using two-step cloning. First, the oligo pool for each library was

734 amplified with NEBNext High-Fidelity 2× PCR Master Mix (NEB, M0541L) and the following

735 primers: Forward primer: TCGATTTCTTGGCTTTATATATCTTGTGGAAAGGACGAAACAC,

736 Reverse primer: ATTTCTAGTTGGTTTAACGCGTAACTAGATAGAACCGCGTCAA. PCR

737 products were purified via gel excision and column purification (Promega, A9282), followed by

738 insertion into the BsmBI–digested lentiV2-hU6-evopreQ1 vector by Gibson assembly (NEB,

739 E2621L). The assembled products were electroporated into Endura electrocompetent

740 Escherichia coli cells (Biosearch Technologies, 60242). About 25 million bacterial colonies were

741 cultured for each library, followed by purification with the QIAGEN Plasmid Maxi Kit (GIAGEN,

742 12163). For the second step, the resulting plasmid libraries from the first cloning step were

743 linearized by BsmbI digestion, gel-purified, and ligated using T4 ligase (NEB, M0202M) to a DNA

744 fragment containing an sgRNA scaffold and the mouse U6 promoter. The ligated products were

745 electroporated into Endura electrocompetent Escherichia coli cells (Biosearch Technologies,

746 60242), and about 40 million bacterial colonies were cultured for each library. The final plasmid

747 libraries were extracted with the Qiagen EndoFree Plasmid Mega Kit (Qiagen, 12381).

748

749 **Lentivirus production and titration**

750 To produce the lentiviral library, we used our previously described method[46]. Briefly, 5 µg of

751 plasmid library, with 3 µg of psPAX (Addgene, 12260) and 1 µg of pMD2.G (Addgene, 12259)

752 packaging plasmids were cotransfected into 8 million HEK293T cells in a 10-cm dish

753 supplemented with 36 µl PolyJet (SignaGen Laboratories, SL100688). The medium was replaced

754 12 hours after transfection and harvested every 24 hours thereafter for a total of three harvests.

755 Harvested viral media was filtered through a Millex-HV 0.45-µm polyvinylidene difluoride filter

756 (Millipore, SLHV033RS) and further concentrated via centrifugation using 100,000 NMWL

757 (nominal molecular weight limit) Ultra-15 centrifugal filter units (Amicon, UFC910008).

758

759 The lentiviral titer was determined by transducing 400,000 cells with increasing volumes (0, 1, 2,

760 5, 10, 20, and 40 µl) of concentrated virus and polybrene (6 µg/ml; Millipore, TR-1003-G). 48

761 hours after the transduction, cells were dissociated with Trypsin-EDTA (0.25%; Gibco, 25200056)

762 and seeded as two separate replicates; one treated with hygromycin B (200 µg/ml; Gibco,

763 10687010) for four days, and another that was not. Finally, hygromycin-resistant and control cells

764 were counted to calculate the infected cell ratios and viral titers.

765

766 **Prime-editing screens**

767 We performed *MYC* enhancer screens in triplicate. We transfected MCF7-dCas9/RT cells with

768 lentivirus libraries at a multiplicity of infection (MOI) of 0.3 with a coverage of 1,000 transduced

769 cells per paired pegRNA/ngRNA. 48 hours later, approximately 10 million cells were harvested

770 as controls and the remaining cells were treated with hygromycin B (200 µg/ml; Gibco, 10687010)

771 for 7 days. After antibiotic selection, the cells were maintained in DMEM supplemented with 10%

772 FBS for 30 days post infection, and 10 million cells were collected from the final cell population.

773

774 We performed Alt and Ref library screens in quadruplicate. We separately infected about 24

775 million MCF7-nCas9/RT cells with the lentivirus library for each replicate of the Alt and Ref

776 screens at an MOI of 0.5, with a cell coverage of 2,000 infected cells per pegRNA/ngRNA pair.

777 48 hours post infection, one-third of the infected cells were collected from each cell pool as control

778 samples (Day 2). The remaining cells were treated with hygromycin B (200 µg/ml; Gibco,

779 10687010) for 7 days and cultured until 32 days post infection (Day 32).

780

781 **Generation of Illumina sequencing libraries**

782 Genomic DNA was extracted from each sample via cell lysis and digestion [100 mM tris-HCl (pH

783 8.5), 5 mM EDTA, 200 mM NaCl, 0.2% SDS, and proteinase K (100 µg/ml)], phenol:chloroform

784 (Thermo Fisher Scientific, 17908) extraction, and isopropanol (Thermo Fisher Scientific,

785 BP2618500) precipitation. For the *MYC* enhancer screen, we applied ePCR during library

786 preparation to amplify the paired pegRNA/ngRNA sequences from each sample and reduce

787 recombination between similar sequences. Briefly, thirty 20 µl ePCRs were performed using 400

788 ng of DNA for each reaction and NEBNext High-Fidelity 2× PCR Master Mix (NEB, M0541S) with

789 the                          following                          primers:                          Enh-lib-Forward:

790    TCCCTACACGACGCTCTTCCGATCTNNNNNCCTTGGAGAAAAGCCTTGTTT, Enh-lib-

791    Reverse: GGAGTTCAGACGTGTGCTCTTCCGATCTNNNNNGAACCGCGGTGTTAGG. ePCR

792    was performed as described previously to amplify pegRNA/ngRNA pairs from genomic DNA.

793    Thermocycler settings were 98 °C for 30 s, then 25 cycles (98 °C 10 s, 60 °C 20 s, 72 °C 1 min),

794    then 72 °C 5 min, and finally a 4 °C hold. The ramp rate for each step was 2°C/s. After PCR,

795    individual reactions were combined and purified using the QIAQuick PCR Purification Kit (Qiagen

796    28104) following previously established guidelines[55]. Purified PCR products were then treated

797    with Exonuclease I (NEB, M0568L) and purified using 1× AMPure XP beads (Beckman Coulter,

798    A63881). Round one PCR amplicons were used in the 2nd round of PCR to add Illumina adapter

799    and index sequences. For the 2nd round PCR, we performed 6 ePCR reactions containing 0.023

800    ng of purified DNA each, using NEBNext High-Fidelity 2× PCR Master Mix (NEB, M0541S). The

801    2nd round PCR mixture was prepared and purified similarly to the 1st. Thermocycler settings were

802    98 °C for 30 s, then 12 cycles (98 °C 10 s, 60 °C 20 s, 72 °C 1 min), then 72 °C 5 min, and finally

803    a 4 °C hold. The ramp rate for each step was 2°C/s. For Alt and Ref screens, we amplified

804    pegRNA/ngRNA pair sequences from each sample using NEBNext High-Fidelity 2× PCR Master

805    Mix (NEB, M0541L) and the following primers: Alt-Ref-lib-Forward:

806    TCCCTACACGACGCTCTTCCGATCTNNNNNCTTGTGGAAAGGACGAAACACC, Alt-Ref-lib-

807    Reverse:

808    GGAGTTCAGACGTGTGCTCTTCCGATCTNNNNNCGTAACTAGATAGAACCGCGTCAA.

809    Twenty-four 50 µl PCR reactions, each containing 600 ng genomic DNA, were performed for each

810    sample. Individual reactions were combined for each sample and column purified (Promega,

811    A9282). The purified products were then amplified by indexing PCR to add Illumina TruSeq

812    adaptors and sample index sequences with the following primers: Index-Forward:

813    aatgatacggcgaccaccgagatctacac[8 bp index]acactctttccctacacgacgctcttccgatct, Index-Reverse:

814    caagcagaagacggcatacgagat[8 bp index]gtgactggagttcagacgtgtgctcttccgatct. The final libraries

815    were gel purified and sequenced with 150 bp paired-ends on the Illumina NovaSeq 6000 platform.

816

817    **Data processing and analysis of prime-editing data**

818    Sequencing libraries were first trimmed with 5 bp random sequences from read1 and read2, and

819    low quality reads were filtered out with the fastp tool before formal mapping. To calculate the read

820    counts, each pegRNA/ngRNA pair was included if it met the following criteria: (1) Read 1 exactly

821    matched the sequence containing a 20-21 nt ngRNA spacer and 5 bp flanking sequences; (2)

822    Read 2 exactly matched the reverse complementary sequence containing the full pegRNA

823    extension and 5 bp flanking sequences.

824

825    For PRIME of *MYC* enhancer, the MAGeCK (0.5.9) pipeline[13] was used to estimate the statistical
826    significance and fold change for each pegRNA/ngRNA pair at the sgRNA level, and for each
827    substitution at the gene level in the cell population relative to controls. The non-targeting and
828    AAVS1 targeting pegRNAs were used as negative controls for normalization. To identify the core
829    enhancer region for the *MYC* enhancer based on the screening results, we first identified base
830    pairs with three significant substitutions (FDR < 0.05), and calculated the slopes for each
831    continuous bin (moving step = 1 bp, bin size = 30 bp, x axis: the position of each base pair, y axis:
832    the accumulation number of  SBPs with three significant substitutions) (**Supplementary Fig. 2e**).
833    The slopes were then transformed into Z score-derived *P* values accordingly. The core enhancer
834    region was identified by merging overlapping significant bins (*P* value < 0.05).

835

836    For Alt and Ref library screens, oligos with zero reads for any sample were removed before the
837    following analysis. Oligo counts from all samples were passed into DESeq2 (1.38.0)[31] and a
838    median-of-ratios method was used to normalize samples for varying sequencing depths.
839    Normalized read counts for each oligo were then modeled by DESeq2 as a negative binomial
840    distribution. We then used DESeq2 to check the fold changes for each oligo in Alt and Ref libraries
841    by comparing Day 32 to Day 2 data (design= ~ Replicate + Condition). We further estimated
842    relative effects between the reference and alternate alleles by adding an interaction term (design=
843    ~ Replicate + Condition + Allele + Condition:Allele). Condition refers to the collection timepoint
844    (i.e. Day 32 or Day 2), and Allele refers to the allele category (i.e. Alt or Ref). Finally, a Wald test
845    was performed via DESeq2 to calculate the *P* value. To minimize false positive hits and achieve
846    an empirical FDR less than 5%, we then selected a *P* value cutoff corresponding to the fifth
847    percentile of *P* values from non-targeting control oligos.

848

849    **Motif matrix comparison analysis**
850    To identify potential transcription factor (TF) binding sites within the target *MYC* enhancer, we
851    established a new method based on motif comparison[56] to directly compare known TF motifs with
852    our base-pair resolution functional data. We first calculated the $\log_2$(fold change) for each
853    substitution at each base pair with MAGeCK  (0.5.9)[13]. The $\log_2$(fold changes) of the wild type
854    alleles were set to 0. We then transformed the $\log_2$(fold change) of each substitution into the
855    corresponding fold change value. We further constructed the position weight matrix by
856    normalizing the fold change of each allele per base pair to the sum of all unique alleles' fold
857    change per base pair. We further partitioned the enhancer sequence into multiple bins with

858    lengths of 5 and 10 base pairs. We only retained bins with an information content (IC) over 3 and

859    an 'N' content less than 10%. We then collected all TF motifs from JASPAR, HOCOMOCO, and

860    SwissRegulon databases with high expression in MCF7 cells (TPM > 10, GSE175204). Next, we

861    compared the filtered TF motif matrices with the enhancer bin matrix using Tomtom (*P* value <

862    0.05) to identify the potential TF binding sites at the enhancer. Finally, we only retained positive

863    TF motif hits overlapping at least 95% of the input sequences' essential base pairs (positions with

864    maximum probabilities > 0.5). Details about the best matching motifs are summarized in

865    **Supplementary Table 2**.

866

867    **Predicting base pair contribution to enhancer activity with BPNet**

868    We trained a convolutional neural network using BPNet consistent with the published approach[24]

869    to explain the GATA3, ELF1, FOXM1, MTA3, and RCOR1 ChIP-seq data from ENCODE projects.

870    Briefly, the model inputs were 1kb sequences across each ChIP-seq peak locus, and

871    corresponding ChIP-seq control peaks were used as the bias track for training. The region from

872    chromosome 2 was used as the tuning set, and chromosomes 5, 6, 7, 10, and 14 were used as

873    the test set. The X and Y chromosomes were excluded. The remaining regions from other

874    chromosomes were used to train the model with default parameters. Once models were acquired

875    for each TF's ChIP-seq data, DeepLIFT was used to calculate each input sequence base pair's

876    contribution to enhancer activity. TF-MoDISco contribution scores were finally used to cluster and

877    determine consolidated TF motifs and map these to input peak regions.

878

879    **MCF7 genotyping analysis**

880    Sequence Read Archive (SRA) files for SRR7707725 and SRR7707726 (paired-end, two reads

881    per loci) were retrieved from BioProject PRJNA486532. We used bwa-mem v.0.7.17 to align

882    sequenced reads to the human reference genome hg38 for each run separately. The Picard tools,

883    SortSam, MarkDuplicates, AddOrReplaceReadGroups were then used to process the BAM files.

884    Finally, GATK v.4.2.5.0 was used to call SNPs and indels via local haplotype re-assembly

885    (HaplotypeCaller) followed by joint genotyping on a single-sample GVCF from HaplotypeCaller

886    (GenotypeGVCFs). Finally, CalcMatch v.1.1.2 was used to verify genotype consistency between

887    two runs.

888

889    **Motif scan and TF identification for alleles with functional breast cancer SNPs**

890    The sequences 20 bp upstream and downstream of each SNP (Alt and Ref alleles) were used as

891    input sequences for TF motif analysis. FIMO software (version 5.5.0)[57] was used to identify

892    matching motifs centered on the SNP regions against the human TF motif database HOCOMOCO

893    (v11 FULL)[19]. All FIMO motif scans were performed using default settings. Finally, TFs (FPKM

894    >1) with binding motifs overlapping target SNP loci were selected (FDR < 0.05, *P* value < 0.0001).

895

896    **Functional validation of rs10956415 using prime editing and RT-qPCR**

897    To validate the function of rs10956415 in MCF7 cells, we converted the alternative allele (A) to

898    the reference allele (C) at this locus using PE. To clone the ngRNA/pegRNA expression plasmid,

899    we amplified the fragment containing the ngRNA-mU6-pegRNA for the rs10956415 reference

900    allele (C) from the screening plasmid library, and inserted this fragment into the BsmBI–digested

901    lentiV2-hU6-evopreQ1 vector using Gibson assembly (NEB, E2621L). We verified the cloned

902    ngRNA/pegRNA plasmid sequence using Primordium whole-plasmid sequencing.

903

904    To perform PE, we transfected two million MCF7-dCas9/RT cells with 2000 ng of ngRNA/pegRNA

905    plasmid containing an EGFP marker using PolyJet (SignaGen Laboratories, SL100688). Five

906    days after transfection, we sorted the cells with the highest EGFP expression level (top 2%) into

907    96-well plates with 100 cells per well using FACS. Approximately two weeks later, we extracted

908    genomic DNA from half of the cells in each well and maintained the other half by seeding them in

909    a 24-well plate. We estimated the PE efficiency for each well by performing genotyping PCR

910    followed by Sanger sequencing. We then expanded the cells in the wells with the highest editing

911    efficiency to isolate clonal PE edited cell lines. We sorted the cell pool into 96-well plates with one

912    cell per well using FACS. Approximately two weeks later, we performed genotyping PCR followed

913    by Sanger sequencing to identify successfully edited clones. Deep sequencing was then

914    performed to quantify the copy number of edited alleles.

915

916    To assess the effect of rs10956415 on *MYC* expression, we used seven PE edited clones with

917    four copies of the C allele and one copy of the A allele. About two million cells from each sample

918    were used to extract total RNA with the RNeasy Plus Mini Kit (Qiagen, #74134), and 1 μg of RNA

919    was used to generate cDNA with the iScript cDNA Synthesis Kit (Bio-Rad, #1708890). We used

920    RT-qPCR to quantify *MYC* expression (forward primer: GTCCTCGGATTCTCTGCTCT, reverse

921    primer: ATCTTCTTGTTCCTCCTCAGAGTC), which was normalized to the *GAPDH* expression

922    level        (forward        primer:        CCACTCCTCCACCTTTGACG,        reverse        primer:

923    ATGAGGTCCACCACCCTGTT).

924

925    **Protein structure prediction with AlphaFold**

926    To explore the impact of the BARD1 His36Pro mutation on BARD1/BRCA1 complex structure,

927    we predicted the wild type BRAD1/BRCA1 and BARD1(His36Pro)/BRCA1 complex structures

928    with AlphaFold. We used the same amino acid chain which is used in the BARD1/BRCA1 complex

929    structure determined by NMR spectroscopy[42] (BARD1, residues 26-122; BRCA1, residues 1-103)

930    as input for complex structure predictions. The amino acid chains of BARD1 and BRCA1 were

931    imported into the Google Colab Version of AlphaFold V2.2.4[58, 59], powered by Python 3 Google

932    Compute Engine. AlphaFold applied a multimer model in response to the duo-sequence

933    imputation, then searched the genetic database to determine the best suited multiple sequence

934    alignment (MSA) for the imported sequence and initiated structural prediction. To avoid

935    stereochemical violations, all structures are relaxed with AMBER model (Assisted Model Building

936    with Energy Refinement) using GPU acceleration. The resulting PDB files were imported into

937    UCSF Chimera X[60, 61] for structure visualization. Protein chains were assigned different colors to

938    distinguish individual chains, and selected amino acid atomic structures and hydrogen bonds were

939    illustrated for interaction analysis. Finally, the real-time rendered complex structures were

940    exported using the snapshot function in Chimera X at the optimal visualization angle.

941

## References

1. Taliun, D. et al. Sequencing of 53,831 diverse genomes from the NHLBI TOPMed Program. *Nature* **590**, 290-299 (2021).

2. Shalem, O., Sanjana, N.E. & Zhang, F. High-throughput functional genomics using CRISPR-Cas9. *Nat Rev Genet* **16**, 299-311 (2015).

3. Anzalone, A.V., Koblan, L.W. & Liu, D.R. Genome editing with CRISPR-Cas nucleases, base editors, transposases and prime editors. *Nat Biotechnol* **38**, 824-844 (2020).

4. Chen, P.J. & Liu, D.R. Prime editing for precise and highly versatile genome manipulation. *Nat Rev Genet* (2022).

5. Anzalone, A.V. et al. Search-and-replace genome editing without double-strand breaks or donor DNA. *Nature* **576**, 149-157 (2019).

6. Erwood, S. et al. Saturation variant interpretation using CRISPR prime editing. *Nat Biotechnol* **40**, 885-895 (2022).

7. Anzalone, A.V., Lin, A.J., Zairis, S., Rabadan, R. & Cornish, V.W. Reprogramming eukaryotic translation with ligand-responsive synthetic RNA switches. *Nat Methods* **13**, 453-458 (2016).

8. Houck-Loomis, B. et al. An equilibrium-dependent retroviral mRNA switch regulates translational recoding. *Nature* **480**, 561-564 (2011).

9. Nelson, J.W. et al. Engineered pegRNAs improve prime editing efficiency. *Nat Biotechnol* **40**, 402-410 (2022).

10. Dang, Y. et al. Optimizing sgRNA structure to improve CRISPR-Cas9 knockout efficiency. *Genome Biol* **16**, 280 (2015).

11. Chen, P.B. et al. Systematic discovery and functional dissection of enhancers needed for cancer cell fitness and proliferation. *Cell Rep* **41**, 111630 (2022).

12. Cho, S.W. et al. Promoter of lncRNA Gene PVT1 Is a Tumor-Suppressor DNA Boundary Element. *Cell* **173**, 1398-1412 e1322 (2018).

13. Li, W. et al. MAGeCK enables robust identification of essential genes from genome-scale CRISPR/Cas9 knockout screens. *Genome Biol* **15**, 554 (2014).

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

15. Baluapuri, A., Wolf, E. & Eilers, M. Target gene-independent functions of MYC oncoproteins. *Nat Rev Mol Cell Biol* **21**, 255-267 (2020).

16. Vitsios, D., Dhindsa, R.S., Middleton, L., Gussow, A.B. & Petrovski, S. Prioritizing non-coding regions based on human genomic constraint and sequence context with deep learning. *Nat Commun* **12**, 1504 (2021).

17. Villar, D. et al. Enhancer evolution across 20 mammalian species. *Cell* **160**, 554-566 (2015).

18. Fornes, O. et al. JASPAR 2020: update of the open-access database of transcription factor binding profiles. *Nucleic Acids Res* **48**, D87-D92 (2020).

19. Kulakovskiy, I.V. et al. HOCOMOCO: towards a complete collection of transcription factor binding models for human and mouse via large-scale ChIP-Seq analysis. *Nucleic Acids Res* **46**, D252-D259 (2018).

20. Pachkov, M., Balwierz, P.J., Arnold, P., Ozonov, E. & van Nimwegen, E. SwissRegulon, a database of genome-wide annotations of regulatory sites: recent updates. *Nucleic Acids Res* **41**, D214-220 (2013).

21. Consortium, E.P. An integrated encyclopedia of DNA elements in the human genome. *Nature* **489**, 57-74 (2012).

22. Schreiber, J., Durham, T., Bilmes, J. & Noble, W.S. Avocado: a multi-scale deep tensor factorization method learns a latent representation of the human epigenome. *Genome Biol* **21**, 81 (2020).

23. Behan, F.M. et al. Prioritization of cancer therapeutic targets using CRISPR-Cas9 screens. *Nature* **568**, 511-516 (2019).

24. Avsec, Z. et al. Base-resolution models of transcription-factor binding reveal soft motif syntax. *Nat Genet* **53**, 354-366 (2021).

25. Michailidou, K. et al. Association analysis identifies 65 new breast cancer risk loci. *Nature* **551**, 92-94 (2017).

26. Fachal, L. et al. Fine-mapping of 150 breast cancer risk regions identifies 191 likely target genes. *Nat Genet* **52**, 56-73 (2020).

27. Hanna, R.E. et al. Massively parallel assessment of human variants with base editor screens. *Cell* **184**, 1064-1080 e1020 (2021).

28. Landrum, M.J. et al. ClinVar: improvements to accessing data. *Nucleic Acids Res* **48**, D835-D844 (2020).

29. Cuella-Martin, R. et al. Functional interrogation of DNA damage response variants with base editing screens. *Cell* **184**, 1081-1097 e1019 (2021).

30. Qi, L.S. et al. Repurposing CRISPR as an RNA-guided platform for sequence-specific control of gene expression. *Cell* **152**, 1173-1183 (2013).

31. Love, M.I., Huber, W. & Anders, S. Moderated estimation of fold change and dispersion for RNA-seq data with DESeq2. *Genome Biol* **15**, 550 (2014).

32. Bruna, A. et al. TGFbeta induces the formation of tumour-initiating cells in claudinlow breast cancer. *Nat Commun* **3**, 1055 (2012).

33. Bossone, S.A., Asselin, C., Patel, A.J. & Marcu, K.B. MAZ, a zinc finger protein, binds to c-MYC and C2 gene sequences regulating transcriptional initiation and termination. *Proc Natl Acad Sci U S A* **89**, 7452-7456 (1992).

34. Wang, X. et al. MAZ drives tumor-specific expression of PPAR gamma 1 in breast cancer cells. *Breast Cancer Res Treat* **111**, 103-111 (2008).

35. Ghandi, M. et al. Next-generation characterization of the Cancer Cell Line Encyclopedia. *Nature* **569**, 503-508 (2019).

36. Wang, Y.H. et al. Knockdown of c-Myc expression by RNAi inhibits MCF-7 breast tumor cells growth in vitro and in vivo. *Breast Cancer Res* **7**, R220-228 (2005).

37. Kircher, M. et al. A general framework for estimating the relative pathogenicity of human genetic variants. *Nat Genet* **46**, 310-315 (2014).

38. Pollard, K.S., Hubisz, M.J., Rosenbloom, K.R. & Siepel, A. Detection of nonneutral substitution rates on mammalian phylogenies. *Genome Res* **20**, 110-121 (2010).

39. Li, W. et al. A synergetic effect of BARD1 mutations on tumorigenesis. *Nat Commun* **12**, 1243 (2021).

40. UniProt, C. UniProt: the universal protein knowledgebase in 2021. *Nucleic Acids Res* **49**, D480-D489 (2021).

41. Prakash, R. et al. Homologous recombination-deficient mutation cluster in tumor suppressor RAD51C identified by comprehensive analysis of cancer variants. *Proc Natl Acad Sci U S A* **119**, e2202727119 (2022).

42. Brzovic, P.S., Rajagopal, P., Hoyt, D.W., King, M.C. & Klevit, R.E. Structure of a BRCA1-BARD1 heterodimeric RING-RING complex. *Nat Struct Biol* **8**, 833-837 (2001).

43. Densham, R.M. et al. Human BRCA1-BARD1 ubiquitin ligase activity counteracts chromatin barriers to DNA resection. *Nat Struct Mol Biol* **23**, 647-655 (2016).

44. Spain, B.H., Larson, C.J., Shihabuddin, L.S., Gage, F.H. & Verma, I.M. Truncated BRCA2 is cytoplasmic: implications for cancer-linked mutations. *Proc Natl Acad Sci U S A* **96**, 13920-13925 (1999).

45. Mandegar, M.A. et al. CRISPR Interference Efficiently Induces Specific and Reversible Gene Silencing in Human iPSCs. *Cell Stem Cell* **18**, 541-553 (2016).

46. Ren, X. et al. Parallel characterization of cis-regulatory elements for multiple genes using CRISPRpath. *Sci Adv* **7**, eabi4360 (2021).

1043 47.    Clement, K. et al. CRISPResso2 provides accurate and rapid genome editing sequence
1044        analysis. *Nat Biotechnol* **37**, 224-226 (2019).
1045 48.    Fejerman, L. et al. Genome-wide association study of breast cancer in Latinas identifies
1046        novel protective variants on 6q25. *Nat Commun* **5**, 5260 (2014).
1047 49.    Machiela, M.J. & Chanock, S.J. LDlink: a web-based application for exploring population-
1048        specific haplotype structure and linking correlated alleles of possible functional variants.
1049        *Bioinformatics* **31**, 3555-3557 (2015).
1050 50.    Hormozdiari, F., Kostem, E., Kang, E.Y., Pasaniuc, B. & Eskin, E. Identifying causal
1051        variants at loci with multiple signals of association. *Genetics* **198**, 497-508 (2014).
1052 51.    Hsu, J.Y. et al. PrimeDesign software for rapid and simplified design of prime editing guide
1053        RNAs. *Nat Commun* **12**, 1034 (2021).
1054 52.    Anderson, M.V., Haldrup, J., Thomsen, E.A., Wolff, J.H. & Mikkelsen, J.G. pegIT - a web-
1055        based design tool for prime editing. *Nucleic Acids Res* **49**, W505-W509 (2021).
1056 53.    Chen, C.H. et al. Improved design and analysis of CRISPR knockout screens.
1057        *Bioinformatics* **34**, 4095-4101 (2018).
1058 54.    Williams, R. et al. Amplification of complex gene libraries by emulsion PCR. *Nat Methods*
1059        **3**, 545-550 (2006).
1060 55.    Verma, V., Gupta, A. & Chaudhary, V.K. Emulsion PCR made easy. *Biotechniques* **69**,
1061        421-426 (2020).
1062 56.    Gupta, S., Stamatoyannopoulos, J.A., Bailey, T.L. & Noble, W.S. Quantifying similarity
1063        between motifs. *Genome Biol* **8**, R24 (2007).
1064 57.    Grant, C.E., Bailey, T.L. & Noble, W.S. FIMO: scanning for occurrences of a given motif.
1065        *Bioinformatics* **27**, 1017-1018 (2011).
1066 58.    Jumper, J. et al. Highly accurate protein structure prediction with AlphaFold. *Nature* **596**,
1067        583-589 (2021).
1068 59.    Mirdita, M. et al. ColabFold: making protein folding accessible to all. *Nat Methods* **19**, 679-
1069        682 (2022).
1070 60.    Goddard, T.D. et al. UCSF ChimeraX: Meeting modern challenges in visualization and
1071        analysis. *Protein Sci* **27**, 14-25 (2018).
1072 61.    Pettersen, E.F. et al. UCSF ChimeraX: Structure visualization for researchers, educators,
1073        and developers. *Protein Sci* **30**, 70-82 (2021).
1074