

# Joint genotypic and phenotypic outcome modeling improves base editing variant effect quantification

Jayoung Ryu<sup>1-3</sup>, Sam Barkal<sup>4</sup>, Tian Yu<sup>4</sup>, Martin Jankowiak<sup>3</sup>, Yunzhuo Zhou<sup>5,6</sup>, Matthew Francoeur<sup>4</sup>, Quang Vinh Phan<sup>4</sup>, Zhijian Li<sup>1,3</sup>, Manuel Tognon<sup>1,3,7</sup>, Lara Brown<sup>4</sup>, Michael I. Love<sup>8</sup>, Guillaume Lettre<sup>9,10</sup>, David B. Ascher<sup>5,6</sup>, Christopher A. Cassa<sup>4†</sup>, Richard I. Sherwood<sup>4†</sup>, Luca Pinello<sup>1,3,11†</sup>

## Affiliations

<sup>1</sup>Molecular Pathology Unit, Center for Cancer Research, Massachusetts General Hospital, Boston, MA, USA

<sup>2</sup>Department of Biomedical Informatics, Harvard Medical School, Boston, MA, USA

<sup>3</sup>Broad Institute of Harvard and MIT, Cambridge, MA, USA

<sup>4</sup>Division of Genetics, Department of Medicine, Brigham and Women's Hospital and Harvard Medical School, Boston, MA, USA

<sup>5</sup>School of Chemistry and Molecular Biosciences, The University of Queensland, Brisbane, Australia

<sup>6</sup>Computational Biology and Clinical Informatics, Baker Heart and Diabetes Institute, Melbourne, Victoria, Australia

<sup>7</sup>Computer Science Department, University of Verona, Verona, Italy

<sup>8</sup>Department of Genetics, Department of Biostatistics, University of North Carolina at Chapel Hill, Chapel Hill, NC

<sup>9</sup>Montreal Heart Institute, Montréal, QC H1T 1C8, Canada

<sup>10</sup>Faculté de Médecine, Université de Montréal, Montréal, QC H3T 1J4, Canada

<sup>11</sup>Department of Pathology, Harvard Medical School, Boston, MA, USA

†Corresponding authors: [ccassa@bwh.harvard.edu](mailto:ccassa@bwh.harvard.edu) (C.A.C.), [rsherwood@bwh.harvard.edu](mailto:rsherwood@bwh.harvard.edu) (R.I.S.), [lpinello@mgh.harvard.edu](mailto:lpinello@mgh.harvard.edu) (L.P.)

## Abstract

CRISPR base editing screens are powerful tools for studying disease-associated variants at scale.

However, the efficiency and precision of base editing perturbations vary, confounding the assessment of variant-induced phenotypic effects. Here, we provide an integrated pipeline that improves the estimation of variant impact in base editing screens. We perform high-throughput ABE8e-SpRY base editing screens with an integrated reporter construct to measure the editing efficiency and outcomes of each gRNA alongside their phenotypic consequences. We introduce BEAN, a Bayesian network that accounts for per-guide editing outcomes and target site chromatin accessibility to estimate variant impacts. We show this pipeline attains superior performance compared to existing tools in variant classification and effect size quantification. We use BEAN to pinpoint common variants that alter LDL uptake, implicating novel genes. Additionally, through saturation base editing of *LDLR*, we enable accurate quantitative prediction of the effects of missense variants on LDL-C levels, which aligns with measurements in UK Biobank individuals, and identify structural mechanisms underlying variant pathogenicity. This work provides a widely applicable approach to improve the power of base editor screens for disease-associated variant characterization.

## 1 Introduction

2 Genetic variation contributes substantially to complex disease risk. While well-powered genome-wide  
3 association studies (GWAS)<sup>1</sup> and rare variant analyses from cohort studies such as the UK Biobank (UKB)<sup>2</sup>  
4 have associated thousands of loci and genes with clinical phenotypes, these observational approaches  
5 are often insufficient to identify causal variants. Perturbation-based methods enable evaluation of the  
6 impact of an individual variant in a common genetic background, isolated from genetically linked  
7 variants, and such testing can be performed in high throughput through multiplex assays of variant  
8 effect (MAVEs)<sup>3</sup>. Numerous types of MAVEs have been developed, including deep mutational scanning  
9 (DMS)<sup>4</sup>, saturation mutagenesis<sup>5</sup>, massively parallel reporter assays (MPRA)<sup>6</sup>, and CRISPR-based  
10 screens<sup>7-9</sup>.

11 CRISPR base editing screens have emerged as a uniquely powerful method to study variants in their  
12 endogenous genomic context. Base editors, fusions of Cas9-nickase and single-stranded cytosine or  
13 adenine deaminase enzymes<sup>10,11</sup>, enable site-specific installation of transition variants. As the majority  
14 of disease-associated variants are single-nucleotide transitions<sup>12</sup>, base editors enable the installation of  
15 functionally relevant variants in a precise and scalable way. Base editing screens have been employed to  
16 dissect coding variant effects as well as to evaluate GWAS-associated variant functions<sup>13-28</sup>.

17 However, base editing efficiency varies substantially depending on the local sequence context  
18 surrounding the target base, the specific Cas9 variant and deaminase used, and the cellular context<sup>29</sup>.

19 Moreover, base edits can occur at multiple positions within the single-stranded DNA bubble created by  
20 the guide RNA (gRNA)-DNA binding on the opposite strand, therefore a single gRNA can install a variety  
21 of alleles, each with distinct efficiencies. While there have been efforts to predict editing outcomes  
22 using massively parallel base editor reporter assay data<sup>29</sup>, these predictions do not generalize well to  
23 unprofiled base editors and cellular contexts<sup>20</sup>.

24 In previous base editing screens, analysis of phenotypic outcomes is confounded by variable editing  
25 efficiencies and outcomes. Phenotypic effects of gRNAs with robust editing are exaggerated, and effects  
26 of variants that are not installed as efficiently are underestimated. Such confounding is especially  
27 pernicious when the target elements are coding variants, as a single gRNA may install distinct coding  
28 variants with different frequencies, and current analysis methods are unable to deconvolve such data.

29 Existing base editing screens have dealt with the heterogeneity in gRNA efficiency and genotypic  
30 outcomes in several ways. One approach that has been employed is to assume all editable nucleotides  
31 within the editing window are edited with uniform efficiency<sup>13</sup>. Two recent studies have profiled the  
32 gRNAs used in phenotypic base editing screening using a base editor reporter (or sensor) assay<sup>20,21</sup> to  
33 filter gRNAs with low editing efficiency when analyzing their phenotypic data. Despite these initial  
34 efforts, the computational analyses of these screens have not yet been formalized, often relying on  
35 existing tools that were not designed specifically for base editor data with or without the target site  
36 reporter.

37 Here, we design an experimental-computational pipeline to improve the accuracy of variant effect  
38 estimation in base editing screens. By incorporating a target site reporter sequence into the gRNA  
39 construct, we simultaneously measure the editing efficiency of a gRNA and its phenotypic impact. We  
40 develop a computational pipeline, BEAN, that normalizes the phenotypic scores of target variants using  
41 genotypic outcome information collected from the target site reporter. Moreover, we extend BEAN to  
42 analyze densely tiled coding sequence base editing screen data, sharing information among neighboring  
43 gRNAs to obtain accurate phenotypic scores for each coding variant. BEAN provides a first-in-class  
44 integrated solution to experimental assessment of variant effects through base editing screens. We  
45 systematically benchmark BEAN against current state-of-the-art methods for the analyses of pooled  
46 CRISPR screens and show substantially improved performance of BEAN.

47 To leverage activity-normalized base editing screening, we have conducted screens assessing the impact  
48 of low-density lipoprotein cholesterol (LDL-C)-associated GWAS variants and low-density lipoprotein  
49 receptor (*LDLR*) coding variants on LDL-C uptake in HepG2 hepatocellular carcinoma cells. Genetic  
50 differences in LDL-C levels contribute substantially to coronary artery disease risk. Serum LDL-C  
51 measurements are quantitative and nearly uniformly measured in most biobanks, and thus they provide  
52 among the highest quality human phenotypic data for any trait. A trans-ancestry GWAS meta-analysis  
53 from the Global Lipids Genetics Consortium (GLGC) has identified >900 genome-wide significant loci  
54 associated with blood lipid levels, including >400 loci associated with LDL-C<sup>30</sup>. LDL-C GWAS loci overlap  
55 strongly with liver-enriched gene expression, nominating liver as the primary tissue driving LDL-C variant  
56 effects<sup>31,32</sup>. Yet, the causal variants and mechanisms by which many of these loci modulate LDL-C levels  
57 remain unknown.

58 LDL-C levels are also impacted by rare coding variants. In the most severe instances, inherited  
59 monogenic variants in several genes cause Familial Hypercholesterolemia (FH), a disease associated with  
60 extremely elevated LDL-C levels and premature cardiovascular disease<sup>33</sup>. The majority of genetic  
61 mutations known to cause FH occur in *LDLR*, a cell surface receptor that uptakes LDL, thus removing it  
62 from circulation<sup>34</sup>. Despite the effectiveness of lipid lowering therapies, FH patients are still 2-4-fold  
63 more likely to have coronary events than the general population<sup>35</sup>. Elevated LDL-C levels increase  
64 cardiovascular disease risk throughout life, so the early identification of at-risk individuals would have  
65 immense clinical utility<sup>33</sup>. However, many *LDLR* variants currently lack clinical interpretation. Of the  
66 1,427 *LDLR* missense variants in the ClinVar database<sup>36</sup>, 50% are classified as variants of unknown  
67 significance (VUS) or to have conflicting interpretations of pathogenicity (“conflicting”), thus impeding  
68 FH diagnosis. Likewise, of the 758 unique *LDLR* missense variants carried by sequenced individuals in the  
69 UKB cohort, 69% are either unreported or have an uncertain annotation in ClinVar. Altogether,

70 improved understanding of *LDLR* variant impacts would enable earlier diagnosis and treatment for a  
71 large number of at-risk individuals.

72 We have modeled the impacts of both common GWAS-associated and rare *LDLR* coding variants through  
73 base editing installation followed by cellular uptake of fluorescent LDL-C in HepG2 cells, which provides  
74 a scalable flow cytometric assay to measure a key contributing factor of serum LDL-C levels<sup>37</sup> given the  
75 majority of serum LDL-C is cleared in liver<sup>38</sup>. By applying our experimental-computational pipeline to this  
76 screen model, we identify LDL uptake-altering GWAS-associated variants and characterize their  
77 downstream impact on chromatin accessibility, transcription factor binding, and gene expression that  
78 leads to differential LDL uptake. We nominate causal variants that alter LDL-C uptake through impacting  
79 the genes *OPRL1*, *VTN*, and *ZNF329*, which have not previously been connected with LDL-C levels.

80 Through saturation tiled base editing of *LDLR*, not only do we accurately distinguish known pathogenic  
81 vs. benign variants, we find strong correlation between missense variant functional scores and the LDL-C  
82 levels of patients in the UKB who carry these variants. We combine functional scores with structural  
83 modeling to mechanistically classify deleterious variant impacts, revealing a key, conserved tyrosine  
84 residue in each *LDLR* class B repeat that interacts with the neighboring repeat to maintain structural  
85 integrity. Altogether, BEAN provides a widely applicable tool to characterize single-nucleotide variant  
86 functions.

## 87 **Results**

### 88 **A base editing reporter profiles endogenous editing outcomes**

89 To enable accurate interrogation of variant effects at scale, we built a platform to perform dense, high-  
90 coverage base editing screens that accounts for variable editing efficiency and genotypic outcomes. To  
91 maximize coverage of variants in base editing screens, we built lentiviral adenine (ABE8e)<sup>11,39</sup> and  
92 cytosine (AID-BE5)<sup>29</sup> deaminase base editor (BE) constructs using the near-PAM-less SpCas9 variant,

93 SpRY<sup>40</sup>. Both BEs showed native genomic editing activity, as measured in HepG2 cells by ASGR1 splice  
94 site editing followed by flow cytometric anti-ASGR1 antibody staining, with ABE8e-SpRY showing  
95 considerably more robust maximal activity (**Supplementary Fig. 1a**). Editing efficiency was increased by  
96 5-10% by prior lentiviral integration of constitutively expressed BEs and by transient dosing of cells with  
97 the histone deacetylase valproic acid immediately after BE and gRNA transduction (**Supplementary Fig.**  
98 **1b-c**), and thus these treatments were implemented in all screens.

99 Base editing efficiency is known to vary depending on Cas9 binding efficiency as well as the local sequence  
100 and chromatin context surrounding the target base<sup>29,41,42</sup>, and thus we expected gRNAs to vary  
101 substantially in editing efficiency across target sites. To account for this variability, we synthesized and  
102 cloned each gRNA paired with a 32-nt reporter sequence comprising the genomic target sequence of that  
103 gRNA into lentiviral base editor vectors (**Fig. 1a, Methods**), akin to previously published CRISPR mutational  
104 outcome reporter constructs<sup>20,21</sup>. When introduced into cells, the gRNA can edit both its native genomic  
105 target site and the adjacent target site (reporter) in the lentiviral vector, which can be read out using next-  
106 generation sequencing (NGS).

107 We designed two gRNA libraries using this approach to improve understanding of the genetics of LDL-C  
108 levels. The first library (LDL-C GWAS library) targets 583 variants associated with LDL-C levels from the UK  
109 Biobank GWAS cohort (**Methods, Supplementary Table 1**). We included fine-mapped variants with  
110 posterior inclusion probability (PIP) > 0.25 from either the SUSIE or Polyfun fine-mapping pipelines<sup>43,44</sup>,  
111 and also variants with PIP > 0.1 within 250 kb of any of 490 genes found to significantly alter LDL-C uptake  
112 from recent CRISPR-Cas9 knockout screens<sup>37,45</sup>. We designed five tiled gRNAs for each variant allele that  
113 place the variant in positions shown to induce most efficient editing with ABE8e (**Fig. 1e**)<sup>46</sup>. Positive control  
114 gRNAs which ablate splice donor and acceptor consensus sites in six genes found to have significantly  
115 altered LDL-C uptake upon knockout<sup>37</sup>, and 100 non-targeting negative control gRNAs that tile 20 synthetic  
116 variants were included, for a total of 3,455 gRNAs.

117 The second library (*LDLR* tiling library) targeted the *LDLR* gene (**Supplementary Table 2**). Taking  
118 advantage of the flexible PAM recognition of SpRY, every possible gRNA targeting the *LDLR* coding  
119 sequence on both strands was included. Lower density gRNAs, tiled every 2-3-nt, targeted the the 50-nt  
120 flanking each *LDLR* exon, the *LDLR* 5' and 3' UTR, promoter, and two intronic enhancers (**Fig. 1f**). This  
121 library also contained 150 non-targeting negative control gRNAs, for a total of 7,500 gRNAs.

122 We first assessed editing outcomes through lentiviral transduction of each library in HepG2 cells  
123 followed by NGS of gRNA-reporter pairs 10-14 days afterwards. We developed an end-to-end  
124 computational toolkit for base-editing screens, BEAN, which includes the ability to perform quality  
125 control and quantify editing outcomes from raw reads among other functionalities. Importantly, the  
126 quantification step is designed to account for self-editing of the spacer sequence, which we found to  
127 occur at appreciable frequency and with modest correlation with reporter editing frequency (LDL-C  
128 GWAS library median 31%, Pearson  $r=0.36$ , *LDLR* tiling library median 18%,  $r=0.31$ , **Supplementary Fig.**  
129 **2**). We used BEAN to profile the previously uncharacterized PAM-less base editors ABE8e-SpRY and AID-  
130 BE5-SpRY on reporter data from the >10,000 gRNAs in both libraries (**Fig. 1b-c, Supplementary Fig. 3**).

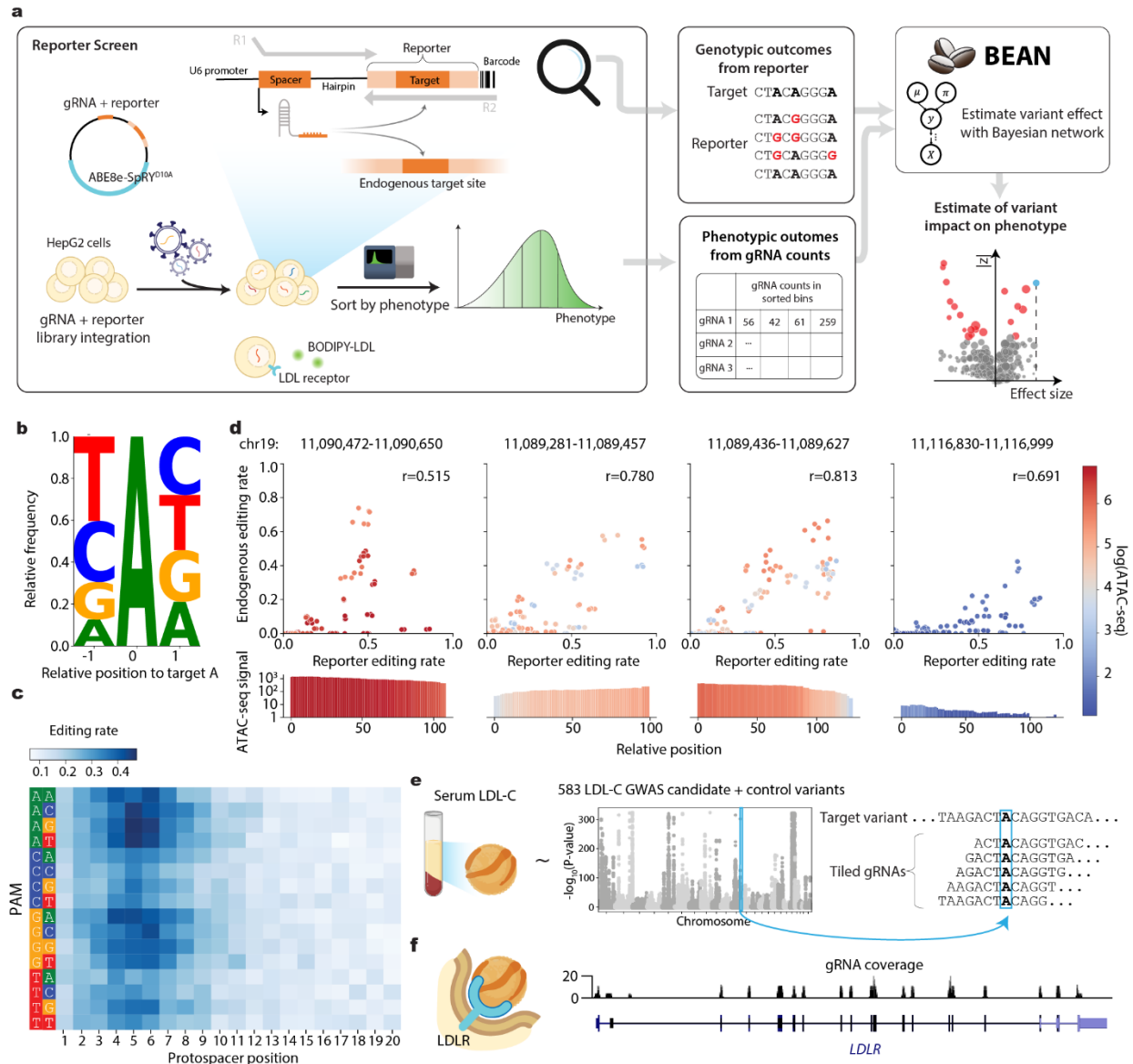
131 The result clearly recapitulated the hallmark positional preferences of these base editors<sup>5,9</sup> the NRY PAM  
132 preference of the SpRY enzyme<sup>10,11</sup>, and the relative depletion of editing at AA dinucleotides by ABE8e.  
133 Notably, the average maximal positional ABE8e-SpRY editing frequency at protospacer positions 3-8  
134 across dinucleotide PAM sequences ranges from 32% to 46%, indicating the ability of this enzyme to  
135 install variants efficiently across a wide variety of genomic locations.

136 To validate that editing of the reporter provides an accurate surrogate for endogenous editing, we  
137 generated a library where both the reporter and endogenous target site are sequenced following the  
138 editing by 49 gRNAs across four loci surrounding *LDLR* with varying levels of HepG2 chromatin  
139 accessibility (**Supplementary Table 3**). We demonstrate that nucleotide-level and allele-level reporter  
140 editing fractions correlate well with endogenous target site editing fractions (**Fig. 1d, Supplementary**



141 **Fig. 4**, average Pearson correlation across 4 loci is  $r=0.70$  for per-nucleotide editing rate  $r=0.70$ , per-  
142 allele editing rate  $r=0.69$ ), and the reporter shows higher correspondence than BE-Hive predictions<sup>29</sup>  
143 (Nucleotide  $r=0.44$ , allele  $r=0.64$ ) (**Supplementary Fig. 5**). Notably, while reporter editing correlates with  
144 endogenous editing at all four loci, we found that endogenous editing frequency also depends on the  
145 accessibility of the target region, as has been previously reported for Cas9-nuclease<sup>47-49</sup> and base  
146 editors<sup>41,42</sup>. Yet, current computational analyses do not model these dependencies, motivating the  
147 development of a tailored modeling framework.

148 We then performed fluorescent LDL uptake screens with each library in  $\geq 5$  biological replicates, ensuring  
149  $>500$  cells per gRNA at all stages. We used simulation to determine the optimal flow cytometric sorting  
150 scheme, accounting for variability in gRNA editing rate, gRNA coverage, gDNA sampling and PCR  
151 amplification (<https://github.com/pinellolab/screen-simulation/>). Based on our simulation result that  
152 finer bin widths improves sensitivity (**Supplementary Fig. 6, Supplementary Note 1**), we flow  
153 cytometrically isolated four populations per replicate with the very low (0-20% percentile), low (20-  
154 40%), high (60-80%), and very high (80-100%) LDL uptake (**Fig. 1a**), performing NGS on gRNA and  
155 reporter pairs in each sorted population. We observed robust replicability (median Spearman  $\rho=0.84$  for  
156 LDL-C GWAS library, 0.88 for *LDLR* tiling library) in gRNA counts across replicates (**Supplementary Fig. 7**),  
157 indicating technical reproducibility.



**Figure 1. Activity-normalized base editing screening pipeline.** **a)** Schematic of activity-normalized base editing screening process and analysis by BEAN. A library of gRNAs, each paired with a reporter sequence encompassing its genomic target sequence, is cloned into a lentiviral base editor expression vector. Lentiviral transduction is performed in HepG2, followed by flow cytometric sorting of four populations based on fluorescent LDL-cholesterol (BODIPY-LDL) uptake. The gRNA and reporter sequences are read out by paired-end NGS to obtain gRNA counts and reporter editing outcomes in each flow cytometric bin. BEAN models the reporter editing frequency and allelic outcomes and gRNA enrichments among flow cytometric bins using BEAN to estimate variant phenotypic effect

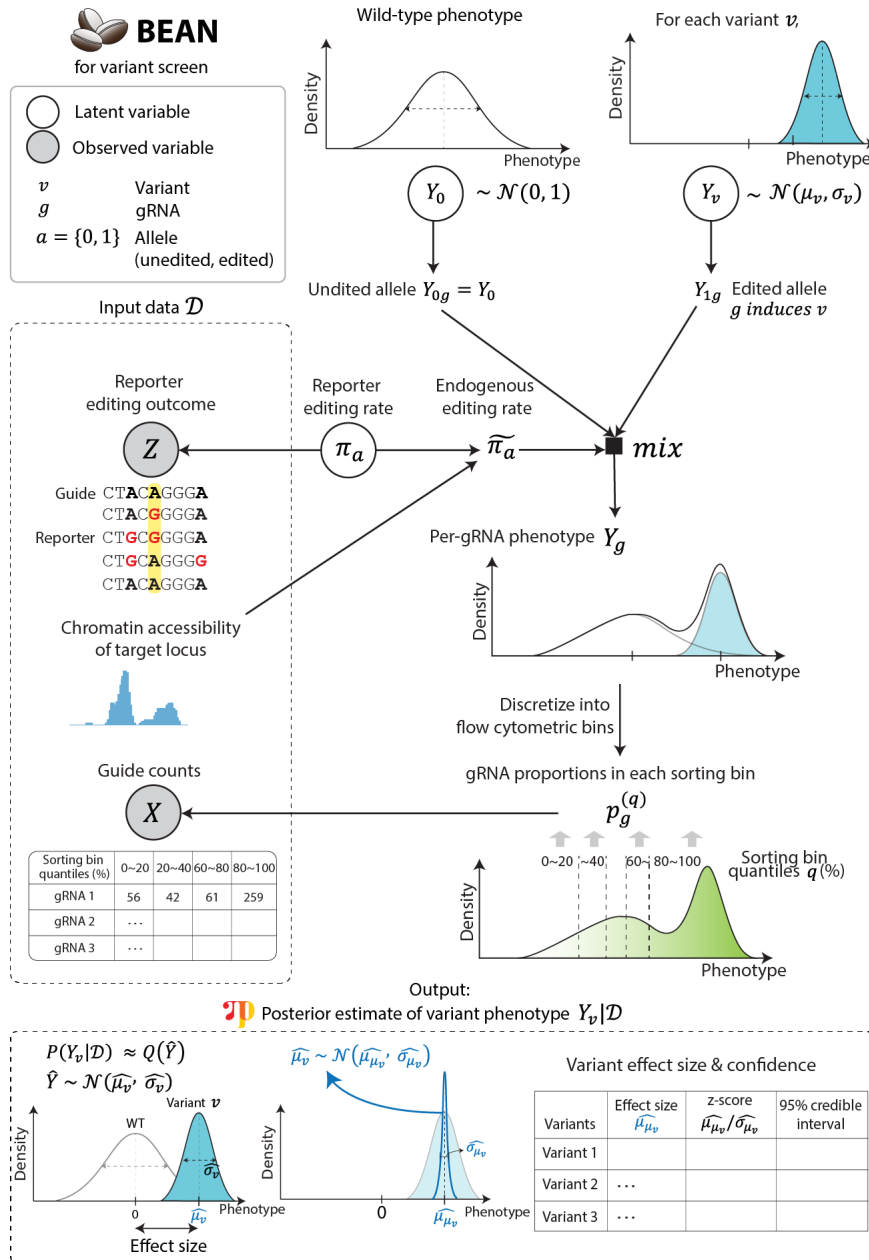
sizes. **b)** Adjacent nucleotide specificity of ABE8e-SpRY editing represented as a sequence logo from 7,320 gRNAs; the height of each base represents the relative frequency of observing each base given an edit at position 0. **c)** Average editing efficiency of ABE8e-SpRY by protospacer position and PAM sequence **d)** Scatterplots comparing nucleotide-level editing efficiency between the reporter and endogenous target sites for a total of 49 gRNAs across four loci across 3 experimental replicates. The accessibility of the four loci as measured by ATAC-seq signal in HepG2 is shown in the top panel, and the scatterplot markers are colored by the accessibility of each nucleotide. Pearson correlation coefficients are shown as  $r$ . **e)** Schematic of the LDL-C variant library gRNA design for selected GWAS candidate variants with a Manhattan plot showing variant P-values from a recent GWAS study<sup>50</sup>. gRNAs tile the variant at five positions with maximal editing efficiency (protospacer positions 4-8). **f)** gRNA coverage of the LDLR tiling library across LDLR coding sequence along with 5' and 3' UTRs and several regulatory regions.

## 158 **Activity-normalized base editing screen analysis with BEAN**

159 We postulated that the gRNA editing outcomes provided by the reporter together with the accessibility  
160 of the target region could improve the quantification of variant phenotypic effects in our pooled BE  
161 screens. To do so, we developed a novel analysis method, BEAN (Base Editor screen analysis with  
162 Activity Normalization), to quantify the effect of each variant from gRNA abundance in sorted  
163 populations along with genotypic outcome information provided by reporter editing. BEAN assumes that  
164 the observed phenotypic distribution in a population of cells for each gRNA derives from a mixture of  
165 cells with unedited and edited alleles (**Fig. 2**). The proportion of cells carrying a given gRNA that possess  
166 a particular genotype is inferred based on the editing outcome observed in reporter as well as  
167 chromatin accessibility of the target locus using a Bayesian network. The distribution of cells with each  
168 gRNA prior to sorting is modeled as a Gaussian mixture for each underlying genotype produced by that  
169 gRNA. Because multiple gRNAs may induce the same genotypic outcome at different frequencies, BEAN  
170 uses this redundancy to build confidence in the predicted phenotypic impacts of a given genotype. As  
171 the output for each variant, BEAN provides its effect size i.e. the posterior mean phenotypic shift along  
172 with the corresponding z-score, and 95% credible interval (CI). We also note that BEAN can be adapted

173 to an arbitrary number and arrangement of sorting bins and other base editing enzymes including those  
174 with uncharacterized editing preferences, and can accommodate screens without reporter or  
175 accessibility information (**Methods**).

176 BEAN only assumes population-level consistency between editing of the reporter and endogenous  
177 target site. We hypothesized that variation in editor expression or cellular state may lead certain cells to  
178 be more amenable to editing than others. In this scenario, “jackpot” cells would be more likely to have  
179 editing at both endogenous and reporter loci. To assess this possibility, we compared the enrichment of  
180 a gRNA in the highest vs. lowest sorted LDL uptake quantile bin with the difference in reporter editing  
181 observed in cells sorted into these bins, reasoning that endogenous editing should be highest in the cells  
182 sorted into the enriched bin. We indeed observed such correlation for *LDLR* and *MYLIP* splice-ablating  
183 gRNAs (Spearman  $\rho=0.32$ , **Supplementary Fig. 8**), suggesting the existence of cell-level factors leading to  
184 “jackpot” cells with higher editing at both endogenous and reporter loci. However, the correlation  
185 between phenotypic and reporter editing enrichment was weaker when considering all positive control  
186 gRNAs (Spearman  $\rho=0.13$ ). We thus concluded that incorporating the jackpot effect into BEAN would be  
187 unlikely to improve model performance.



**Figure 2. BEAN models variant effects from activity-normalized base editing screens.** Simplified schematic of BEAN Bayesian network that models input reporter editing outcomes and gRNA counts. The Bayesian network model recapitulates the data generation process starting from a variant-level phenotype  $Y_v$  and models per-gRNA phenotypes as a Gaussian mixture distribution of edited and unedited (wild-type) allele phenotypes. The weights of the mixture components are modeled to generate reporter editing outcomes. gRNA abundance in each sorting bin is then calculated by discretizing the gRNA phenotype based on the experimental design into the phenotypic

quantiles, and is modeled to generate the observed gRNA counts using an overdispersed multivariate count distribution (see Methods). BEAN outputs the parameters of the posterior distribution of mean phenotypic shift as Gaussian distribution with mean  $\widehat{\mu}_\mu$  (effect size), along with negative-control adjusted z-score and credible interval (CI), where  $\mathcal{D}$  is the input data.

## 188 **BEAN identifies LDL uptake altering GWAS variants**

189 We applied BEAN to the LDL-C GWAS library screen. From the reporter data, variant editing efficiency  
190 per gRNA is highly variable with average edit fraction of 34.0%. Encouragingly, most target variants are  
191 edited at high efficiency by at least one of the five targeting gRNAs (median maximal editing of 60.4%,  
192 **Supplementary Fig. 9**).

193 First, we compared the performance of BEAN and five published CRISPR screen analysis methods at  
194 distinguishing the effects of positive control splice-altering variants versus negative control non-  
195 targeting gRNAs<sup>23,51–54</sup> (**Methods, Fig. 3a**). To dissect the contributions of individual features to BEAN  
196 performance, we included two reduced versions of BEAN: one that considers reporter editing but not  
197 chromatin accessibility (BEAN-Reporter), and another that ignores the reporter, assuming uniform gRNA  
198 editing efficiency (BEAN-Uniform) (**Methods, Supplementary Fig. 10**). BEAN outperforms other  
199 evaluated methods at this classification task (**Fig. 3b, Supplementary Fig. 11**), and this improved  
200 performance is accentuated when the data is subsampled for fewer replicates, demonstrating its ability  
201 to maintain robustness even with less data. Importantly, BEAN shows improved performance (mean  
202 AUPRC=0.90 across 15 2-replicate subsamples) over BEAN-Reporter (mean AUPRC=0.87), which in turn  
203 outperforms BEAN-Uniform (mean AUPRC=0.85), supporting the value of accurately modeling target site  
204 editing. Intriguingly, even BEAN-Uniform outperforms alternative approaches, likely due to more  
205 accurate modeling of sorting bins, suggesting the utility of BEAN in sorting screens without reporter.

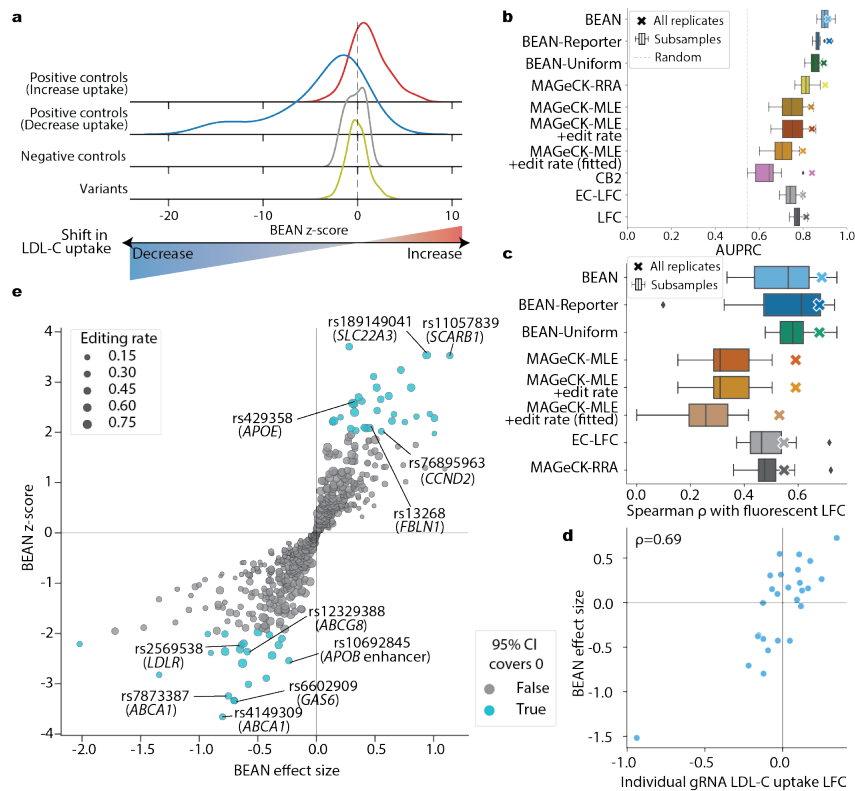
206 Having demonstrated robust performance of BEAN, we evaluated our ability to characterize common  
207 variants that alter LDL-C uptake. We identified 54 variants that significantly alter LDL-C uptake (95% CI  
208 does not contain 0, **Supplementary Table 4**). These variants include intronic variants in well-known LDL-  
209 C uptake mediators whose knockout altered LDL-C uptake in a recent genome-scale CRISPR screen<sup>37</sup>  
210 such as *ABCA1*, *LDLR*, and *SCARB1* (**Fig. 3e**). We additionally identified coding/intronic variants in *APOE*,  
211 *CCND2*, *GAS6*, and *FBLN1* with strong genetic likelihood of causality (UKBB SUSIE fine-mapping PIP > 0.99  
212 and/or the only variant in a fine-mapped credible set<sup>30</sup>), indicating that the effect of these variants on  
213 serum LDL-C is at least partially mediated by LDL-C uptake.

214 To validate the inferred effect sizes, we performed individual lentiviral ABE8e-SpRY transduction of  
215 HepG2 cells with gRNAs targeting 22 variants and 4 positive controls (**Supplementary Table 5**). We  
216 performed fluorescent LDL-C uptake profiling of each edited cell line mixed with an in-well control cell  
217 line in 6 biological replicates (**Methods**), allowing us to compare changes in LDL-C uptake with matched  
218 data from the screen. The individual LDL-C uptake log-fold-change (LFC) values showed strong  
219 correlation to the BEAN effect sizes ( $\mu$ , Spearman R=0.69, **Fig. 3c-d**, **Supplementary Fig. 12**), showing  
220 more robust correlation than BEAN-Uniform (R=0.68), log fold change based on MAGeCK-RRA (R=0.51),  
221 and regression coefficients  $\beta$  of MAGeCK-MLE (R=0.61, see **Methods**). These data demonstrate that  
222 BEAN enables accurate inference of variant effects on LDL-C uptake over a wide dynamic range.

223

224

It is made available under a [CC-BY-NC-ND 4.0 International license](https://creativecommons.org/licenses/by-nc-nd/4.0/).



**Figure 3. BEAN improves variant impact estimation from the LDL-C GWAS library screen. a)** Ridge plot of BEAN z-score distributions of positive controls, negative controls, and test variants. **b)** AUPRC of classifying *LDLR* and *MYLIP* splicing variants vs. negative controls. Metrics for all 5 replicates are shown as markers and metrics of 15 two-replicate subsamples among the 5 replicates are shown as box plots. **c)** Spearman correlation coefficient between BEAN effect size and log fold change of LDL-C uptake following individual testing of 26 gRNAs. Metrics for all 5 replicates are shown as marker and metrics of 15 two-replicate subsamples among the 5 replicates are shown as box plots. **d)** Scatterplot of BEAN effect size and log fold change of LDL-C uptake following individual testing of 26 gRNAs. Spearman correlation coefficient is denoted as  $\rho$ . **e)** Scatterplot of variant effect size and significance estimated by BEAN. Labels show rsIDs of selected variants and dbSNP gene annotations and a manual annotation for *APOB* enhancer in the parenthesis.

225 To gain insight into a set of 20 variants for which the mechanism of LDL-C uptake alteration is less clear,  
 226 we developed a pipeline to assess the cellular effects of variant installation (**Fig. 4a**). First, we asked  
 227 which of these variants impact chromatin accessibility. We established an approach to perform pooled



228 variant editing followed by ATAC-seq. High multiplicity of infection (MOI) lentiviral delivery of a pool of  
229 20 ABE8e-SpRY gRNAs to HepG2 cells was followed by ATAC-seq and paired genomic DNA collection in  
230 three biological replicates in standard and serum-starved conditions. We performed multiplexed PCR  
231 enrichment of the regions surrounding each of the 20 edited variants followed by targeted amplicon  
232 sequencing by NGS. Differential representation of an alternate allele in ATAC-seq relative to gDNA  
233 sequencing implies differential accessibility of the alternate allele than the reference (**Fig. 4b**).

234 Eight of the 20 variants are heterozygous in HepG2, and thus we could assess whether these variants  
235 reside in chromatin accessibility quantitative trait loci (caQTL)<sup>55</sup>, showing differential relative  
236 accessibility of the two haplotypes irrespective of base editing. We found five of these eight variants to  
237 be caQTLs (**Fig. 4c**). Two of these loci (rs35081008 and rs2618566) were also identified as caQTLs in a  
238 recent analysis of 20 human liver tissue samples<sup>56</sup>. Importantly, caQTL analysis cannot address the  
239 causality of the evaluated variant due to the presence of linked variants which could contribute to the  
240 differential ATAC-seq signal.

241 To assess whether individual variants alter chromatin accessibility, we evaluated whether base editing  
242 induces differential accessibility for any of the 20 tested variants. Technical issues including insufficient  
243 representation of the region, insufficient editing, and inability to phase heterozygous loci prevented  
244 assessment of five of the variants (see **Methods**). Of the 15 remaining variants, eight significantly altered  
245 chromatin accessibility when edited (family-wise error rate 0.1, **Fig. 4c**). Four such variants (rs11149612,  
246 rs35081008, rs8126001, and rs2618566) were in loci identified as liver tissue caQTLs. Because base  
247 editing only alters a single variant in a locus, this analysis establishes at least partial causality to the  
248 tested variant.

249 We performed deeper characterization of three variants whose editing alters both LDL-C uptake and  
250 chromatin accessibility. Rs704 is a missense coding variant in *VTV* and is the only variant in a fine-

251 mapped credible set from LDL-C GWAS<sup>30</sup>, suggesting high likelihood of causality. The other two variants  
252 are in gene promoters—rs35081008 is in the *ZNF329* promoter, and rs8126001 is in the shared  
253 *OPRL1/RGS19* promoter (**Fig. 4d**). Both variants have moderate probability of causality from GWAS  
254 evidence (SUSIE PIP=0.49 for rs35081008, PIP=0.25 for rs8126001), with the remaining probability in the  
255 rs35081008 locus deriving from a linked variant (rs34003091) 19-nt upstream in the *ZNF329* promoter.  
256 None of the putative target genes has been previously found to alter LDL-C uptake, nor do they show  
257 significance in LDL-C burden analyses.

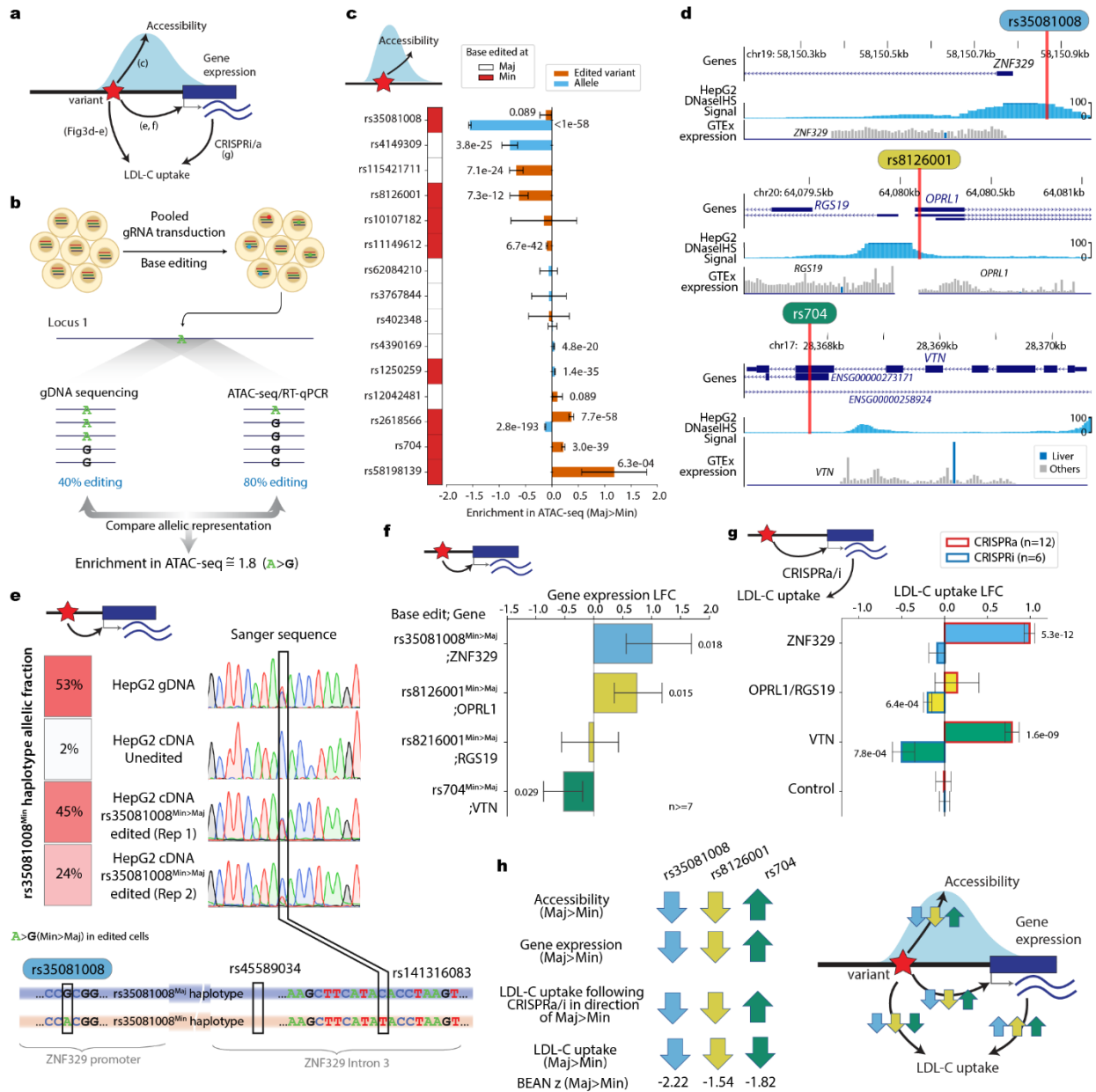
258 To investigate how the prioritized variants might affect transcription factors (TF) binding sites and  
259 thereby regulate proximal genes involved in LDL-C uptake, we adapted the MotifRaptor pipeline<sup>57</sup>.  
260 Briefly, MotifRaptor predicts both the binding strength and disruption scores of TF motifs at specified  
261 SNP loci (see **Methods**). Using the human TFs from the CIS-BP motif database<sup>58</sup>, for each variant, we  
262 ranked the TFs with high binding and potential disruption by the variant (**Supplementary Table 6**). For  
263 rs8126001, our approach prioritized two zinc finger TFs, ZNF333 and ZNF770 with enhanced binding site  
264 sequences due to the heterozygous minor allele in HepG2 cells (**Supplementary Fig. 13**). HepG2 ChIP-  
265 seq data<sup>59</sup> further support the binding of these TFs at this locus, although the variant lies at the edge of  
266 the peaks (**Supplementary Fig. 14**). While definitive conclusions about these factors will require further  
267 experimental validation, our observations align with previous research<sup>60</sup> suggesting that only a minority  
268 of causal variants directly alter canonical TF binding sequences and instead affect non-canonical  
269 sequences.

270 We confirmed through RT-qPCR analysis that editing the minor to major alleles of rs35081008 and  
271 rs8126001 leads to increased expression of *ZNF329* and *OPRL1* respectively (**Fig. 4f**), which is consistent  
272 with the increased chromatin accessibility induced by these edits. rs35081008 is heterozygous in HepG2,  
273 and we used two linked *ZNF329* intronic variants to assess allele-specific expression. In wild-type HepG2,  
274 only 2% of *ZNF329* transcripts derive from the minor allele haplotype (**Fig. 4e**), consistent with the

275 diminished chromatin accessibility of this allele (**Fig. 4c**) and the status of rs35081008 as a liver eQTL.  
276 Editing rs35081008 from minor to major allele restores expression of this haplotype to 35% of total  
277 transcripts (**Fig. 4e**), providing further evidence that rs35081008<sup>Maj</sup> results in increased expression of  
278 *ZNF329*.

279 We then performed CRISPRa and CRISPRi targeting to assess whether altered expression of the four  
280 candidate target genes alters LDL-C uptake. CRISPRa induction of *VTN* and *ZNF329* significantly increased  
281 LDL-C uptake, and CRISPRi repression of *VTN* and *OPRL1/RGS19* reduced LDL-C uptake (**Fig. 4g**). In our  
282 base editing experiments, rs704<sup>Min</sup> shows decreased LDL-C uptake, so we surmise that this allele must  
283 have decreased expression or function, given that decreased *VTN* expression decreases LDL-C uptake  
284 (**Fig. 4h**). Prior biochemical characterization has shown decreased cellular binding capacity of rs704<sup>Min</sup>,  
285 suggesting a possible mechanistic explanation. Our data are consistent with rs35081008<sup>Min</sup> decreasing  
286 *ZNF329* expression, which in turn decreases LDL-C uptake. Finally, our data are most consistent with  
287 rs8126001<sup>Min</sup> decreasing *OPRL1* expression, which leads to decreased LDL-C uptake. This observation  
288 aligns with the higher predictive binding affinity of ZNF333, a transcriptional repressor<sup>61,62</sup>, with  
289 rs8126001<sup>Min</sup> and the potential disruption of its binding with rs8126001<sup>Maj</sup>. In summary, through  
290 accurately quantifying impacts of disease-associated variants on LDL-C uptake, BEAN reveals genetic  
291 mechanisms underlying control of LDL-C levels.

It is made available under a [CC-BY-NC-ND 4.0 International license](https://creativecommons.org/licenses/by-nc-nd/4.0/).



**Figure 4.** Functional characterization of LDL-C GWAS variants. **a)** Schematic of potential variant mechanisms and the figure panels showing data from each mechanistic experiment. **b)** Schematic of pooled ATAC-seq analysis to identify variants impacting accessibility. Differential representation of allele in gDNA and ATAC-seq reflects differential accessibility induced by the base edit or heterozygous reference allele. **c)** Change in ATAC-seq enrichment from the pooled ATAC-seq experiment. 95% confidence intervals are shown as the error bars. “Edited variant” denotes the enrichment by base edit and “Allele” denotes the enrichment by either of the heterozygous alleles, when available, translated uniformly to major (Maj) to minor (Min) allele direction. Variants where the

base edit is conducted from minor allele are denoted as red in the color bar. Family-wise error rate (FWER) with Benjamini-Hochberg multiple correction is shown for each enrichment value where  $FWER < 0.1$ . **d)** Genomic tracks for three selected variants. DNaseIHS; DNase 1 Hypersensitivity. Multiple transcript variants of RGS19 and OPRL1 are shown in the middle panel. **e)** Fraction of *ZNF329* minor (Min) allele haplotype reads in gDNA and cDNA from untreated HepG2 and HepG2 with rs35081008<sup>Min>Maj</sup> base editing. **f)** Change in gene expression following base editing of three selected variants from minor (Min) to major (Maj) allele. P-values of the one sample Student's t-test of LFC vs. mean of 0 that are smaller than 0.05 are shown above each bar. **g)** Change in cellular LDL-C uptake following CRISPRa/i of proximal genes for three selected variants. P-values of the one sample Student's t-test of LFC vs. mean of 0 that are smaller than 0.05 are shown above each bar. **h)** Summary schematic of characterization results.

292 **Saturation *LDLR* coding sequence tiling screening enables quantitative assessment of rare variant**  
293 **deleteriousness**

294 We next adapted BEAN to the *LDLR* tiling library, enhancing the model to specifically assess the  
295 contributions of individual amino acid mutations rather than SNVs, by enabling a more comprehensive  
296 understanding of coding region alterations. Previous coding sequence base editing analyses have assumed  
297 that all editable bases within a window are edited, which leads to erroneous amino acid mutation  
298 assignments, or have analyzed gRNA-level signal only<sup>15</sup>. We aimed to exploit the combination of dense  
299 tiling afforded by ABE8e-SpRY and reporter editing outcomes to model the effects of coding variants more  
300 accurately.

301 The *LDLR* tiling screen showed high coverage of edited nucleotides and amino acids (92% of targetable  
302 nucleotides and 74% of the 860 *LDLR* amino acids in the *LDLR* coding sequence were edited at >10%  
303 frequency by at least one gRNA in the reporter, **Supplementary Fig. 15**). A total of 2,182 distinct variants  
304 were assessed, of which 874 are missense coding variants. Each gRNA produced an average of 2.6  
305 distinct alleles, and each variant was covered by 5.8 gRNAs on average (**Supplementary Fig. 16**). Thus,

306 ABE8e-SpRY tiling of *LDLR* resulted in a rich dataset of coding variants for the evaluation of their  
307 phenotypic impacts.

308 As opposed to the LDL-C GWAS analysis in which each gRNA was evaluated based on its editing  
309 frequency at a single target position, we adapted BEAN to account for multi-allelic outcomes. First, BEAN  
310 translates the edited alleles, i.e., aggregates nucleotide-level allele counts that leads to the identical  
311 amino acid transition into a single amino-acid level allele counts, while preserving nucleotide transition  
312 in non-coding regions. BEAN then filters for the translated alleles that are robustly observed (see  
313 **Methods**) for each gRNA (**Fig. 5a**). BEAN uses a Bayesian network to combine phenotypic information  
314 from all the gRNAs that produce a given allele. Importantly, the phenotype attributed to each gRNA is  
315 modeled as a mixture distribution of the alleles it generates, with the contribution of each allele  
316 weighted by its corresponding editing frequency.

317 BEAN assigned significant z-scores ( $<-1.96$ , equivalent to 95% credible interval not covering 0) to 145  
318 among 2,182 variants assessed from the *LDLR* tiling library (**Supplementary Table 7**), 131 of which  
319 decrease LDL-C uptake. 47 variants that significantly decrease LDL-C uptake are annotated in ClinVar as  
320 pathogenic/likely pathogenic, while 17 are ClinVar VUS/conflicting variants and none are ClinVar  
321 benign/likely benign (**Fig. 5g**), indicating that BEAN can reliably predict the pathogenicity for variants  
322 without a pathogenic or benign classification (**Fig. 5b-c**).

323 We compared the performance of BEAN at distinguishing ClinVar-annotated pathogenic from  
324 benign/likely benign *LDLR* variants to other available screen analysis methods<sup>51-54</sup>. To allow comparison  
325 of methods that do not account for editing outcomes, we assigned outcomes to each gRNA either by  
326 assuming all editable bases within the maximal editing window are perfectly edited<sup>13</sup> ("All") or by using  
327 the most frequent predicted outcome from BE-Hive<sup>29</sup>("BE-Hive"). As in the LDL-C GWAS screen, BEAN  
328 showed better performance than any other method (**Fig. 5d, Supplementary Fig. 17**), and BEAN also

329 outperformed the model variants that do not account for accessibility (BEAN-Reporter) or reporter  
330 editing outcomes (BEAN-Uniform), further justifying the modeling decision to explicitly leverage editing  
331 outcomes and accessibility. BEAN achieves an AUPRC of 0.88 at this task, indicating highly effective  
332 distinction of pathogenic and benign *LDLR* variants through scoring HepG2 LDL-C uptake proficiency.

333 To gain insight into mechanisms by which these variants disrupt LDL-C uptake, we examined BEAN z-  
334 scores for variants that reside within conserved functional domains (**Fig. 5e-f, Supplementary Fig. 18**).  
335 *LDLR* contains seven highly conserved *LDLR* class A repeats that bind to LDL. The *LDLR* class A repeat is  
336 structurally anchored by six highly conserved cysteines that form three disulfide bonds<sup>63</sup>. As expected,  
337 many of the missense variants with the strongest effects on *LDLR* function disrupt these cysteines (**Fig.**  
338 **5f**). We find that cysteine mutating edits in each of the seven *LDLR* class A repeats disrupt *LDLR* activity  
339 (**Supplementary Fig. 19**), suggesting that structural integrity of all repeats is required for efficient LDL  
340 binding, although disruption is most impactful in repeats 3-7. Truncation experiments have reported  
341 that repeats 1 and 2 are dispensable for LDL binding<sup>64</sup>, in partial accord with our results<sup>65</sup>. To examine  
342 the relationship between conservation and function more comprehensively in these repeats, we  
343 compared the BEAN z-score of every installed variant with its change in amino acid conservation score  
344 from the Pfam profile HMM<sup>66</sup> (see **Methods**). We observed strong concordance (Pearson  $r = 0.57$  **Fig.**  
345 **5h**), with N-terminal hydrophobic residues and C-terminal calcium-coordinating acidic residues within  
346 the repeats also showing particular functional importance, as expected from the known function of  
347 these domains.

348 Encouraged by the concordance between our screen and conservation scores within the *LDLR* class A  
349 repeats, we asked whether BEAN scores could predict functional impairment across the entire *LDLR*  
350 gene. We examined statin-adjusted LDL-C levels<sup>67</sup> in the UK Biobank (UKB) for individuals with paired  
351 exome sequencing and lipid level data. To control for the contribution of other variants in genes known  
352 to impact serum LDL-C level, we filtered out individuals who harbor nonsynonymous *APOB* or *PCSK9*

353 variants or multiple *LDLR* missense variants, leading to 9,819 individuals harboring 358 distinct *LDLR*  
354 missense variants. There are 76 distinct *LDLR* missense variants observed in our base editing data with  
355 UKB carriers. We observe robust concordance between the average carrier LDL-C and BEAN scores for  
356 these variants (Spearman  $\rho = 0.40$ , Pearson  $r = 0.45$ , **Fig. 5i**), suggesting that BEAN provides accurate  
357 quantitative prediction of the impact of *LDLR* missense variants on control of serum LDL-C levels in the  
358 human population.

359 As our base editing screen does not exhaust possible mutation types per position, we used the FUSE<sup>68</sup>  
360 pipeline to impute the impact of unobserved variants at positions at which a different missense variant  
361 is scored. FUSE uses an amino acid substitution matrix derived from 24 deep mutational scanning  
362 datasets to impute functional scores for all possible missense variants at positions observed in base  
363 editing data (BEAN+FUSE score, see **Methods**). Applying FUSE to the 76 UKB variants with observed base  
364 editing data, BEAN-FUSE shows improved correlation with UKB carrier LDL-C (Spearman  $\rho = 0.50$ ,  
365 Pearson  $r = 0.51$ , **Fig. 5j**). BEAN-FUSE correlation with UKB carrier LDL-C was robust but lower at all 358  
366 missense *LDLR* variants with lipid measurements (Spearman  $\rho = 0.37$ , Pearson  $r = 0.35$ , **Supplementary**  
367 **Fig. 20a**). Altogether, BEAN-FUSE provides a pipeline to extend base editing screening to predict  
368 functional impairment for unobserved missense variants, although our data suggest that accuracy does  
369 decrease for unobserved variants.

370 As base editing provides orthogonal functional assessment to conservation, we asked whether the LDL-C  
371 levels of UKB variant carriers could be predicted with BEAN-FUSE scores and PhyloP 100way vertebrate  
372 conservation scores. Using XGBoost regression<sup>69</sup>, we achieved more robust correlation with UKB carrier  
373 LDL-C than either BEAN-FUSE or PhyloP alone at the 76 variants observed in the base editing screen  
374 (Spearman  $\rho = 0.48$ , Pearson  $r = 0.61$ , RMSE=39.0, **Fig. 5k**, **Supplementary Fig. 20c**) and at 358 variants  
375 with BEAN-FUSE score (Spearman  $\rho = 0.37$ , Pearson  $r = 0.31$ , RMSE=51.1, **Supplementary Fig. 20b,d**).

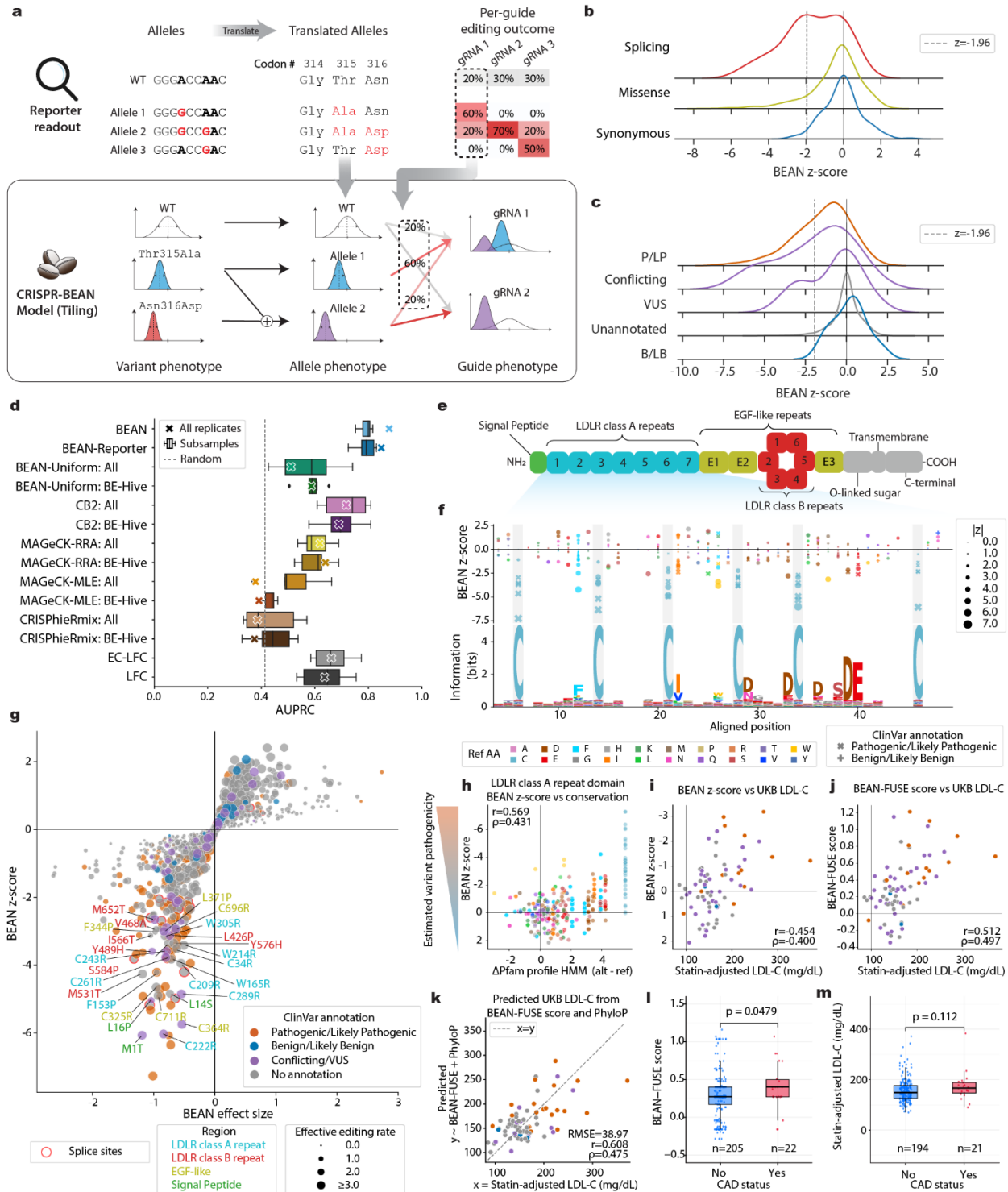


376 This result demonstrates the potential utility of base editing data to improve quantitative phenotype  
377 prediction combined with computational prediction methods.

378 Individuals with pathogenic FH variants are at higher risk of coronary artery disease (CAD), even after  
379 controlling for LDL-C levels<sup>70</sup>. However, the vast majority of rare *LDLR* missense variants lack ClinVar  
380 pathogenic/likely pathogenic designations, preventing information about these potentially disease-  
381 causing variants from being shared with patients. Therefore, we asked whether CAD incidence within  
382 *LDLR* variant carriers could be stratified by functional scores. We found that for individuals with rare  
383 *LDLR* variants, functional scores processed by BEAN-FUSE were significantly higher for patients with  
384 prevalent or incident CAD (Wilcoxon rank-sum test,  $p = 0.0479$ , **Fig. 5l**). BEAN-FUSE scores provided  
385 more robust stratification of individuals with CAD than statin-adjusted LDL-C values for individuals with  
386 variants covered in the screen (Wilcoxon rank-sum test,  $p = 0.112$ , **Fig. 5m**). This demonstrates the  
387 advantage of quantifying genetic risk, which has a lifelong impact on LDL-C levels, over the snapshot  
388 provided by a single LDL-C measurement. Overall, we show that activity-normalized base editing  
389 screening can yield accurate quantitative estimation of *LDLR* variant pathogenicity in a large human  
390 cohort.

391

It is made available under a [CC-BY-NC-ND 4.0 International license](https://creativecommons.org/licenses/by-nc-nd/4.0/).



**Figure 5. Dissection of LDLR variant effects through BEAN modeling of a saturation tiled base editing screen. a)** BEAN model for coding sequence tiling screens. Reporter editing efficiencies are calculated at the amino acid-level when the edited nucleotides are in coding region. Phenotypes of gRNAs with multi-allelic outcomes are modeled

as the Gaussian mixture of allelic phenotypes. If an allele consists of more than one variant, the phenotype of the allele is modeled as the sum of the component variants. A Bayesian network is used to model variant-level phenotypes, sharing phenotypic information from all available gRNAs. **b)** Ridge plot of BEAN z-score distributions of positive controls, negative controls, and variants. **c)** Ridge plot of BEAN z-score distributions of Clinvar variants annotated as pathogenic/likely pathogenic (P/LP), benign/likely benign (B/LB), conflicting interpretation of pathogenicity (conflicting), and Uncertain significance (VUS), and unannotated variants. **d)** AUPRC of classifying ClinVar pathogenic/likely pathogenic vs. benign/likely benign variants. The marker shows the metrics of each method run on 4 replicates with no failing samples. Boxplot shows the metrics of 6 2-replicate combinations of the 4 replicates. **e)** *LDLR* domain structure adopted from Oommen et al<sup>71</sup>. **f)** BEAN z-scores for variants in the 7 *LDLR* class A repeat domains aligned with the Pfam profile HMM logo. Highly conserved cysteines are highlighted in grey. **g)** Scatterplot of estimated variant effect sizes and z-scores. Labels of selected deleterious variants without ClinVar pathogenic/likely pathogenic annotations are shown. **h)** Scatterplot of *LDLR* class A repeat missense variant BEAN z-scores and  $\Delta$ Pfam profile HMM scores. Higher  $\Delta$ Pfam scores correspond to substitution from highly conserved to rarely observed amino acids. **i)** Comparison of mean statin-adjusted LDL-C level and BEAN z-score for variants observed in UKB and base editing. **j)** Comparison of mean statin-adjusted LDL-C level and BEAN-FUSE scores for variants observed in UKB and base editing. **k)** LDL-C levels of observed missense variants predicted by a regression model using BEAN-FUSE and PhyloP scores with 10-fold cross validation, compared with mean statin-adjusted LDL-C level in UKB. **l)** Boxplots of BEAN-FUSE functional scores for UKB individuals with variants observed in our base editing screen with or without CAD. **m)** Boxplots of statin-adjusted LDL-C levels of UKB individuals with variants observed in our base editing screen with or without CAD. P-value of two-sided Wilcoxon rank-sum test is denoted. *r*; Pearson correlation coefficient,  $\rho$ ; Spearman correlation coefficient

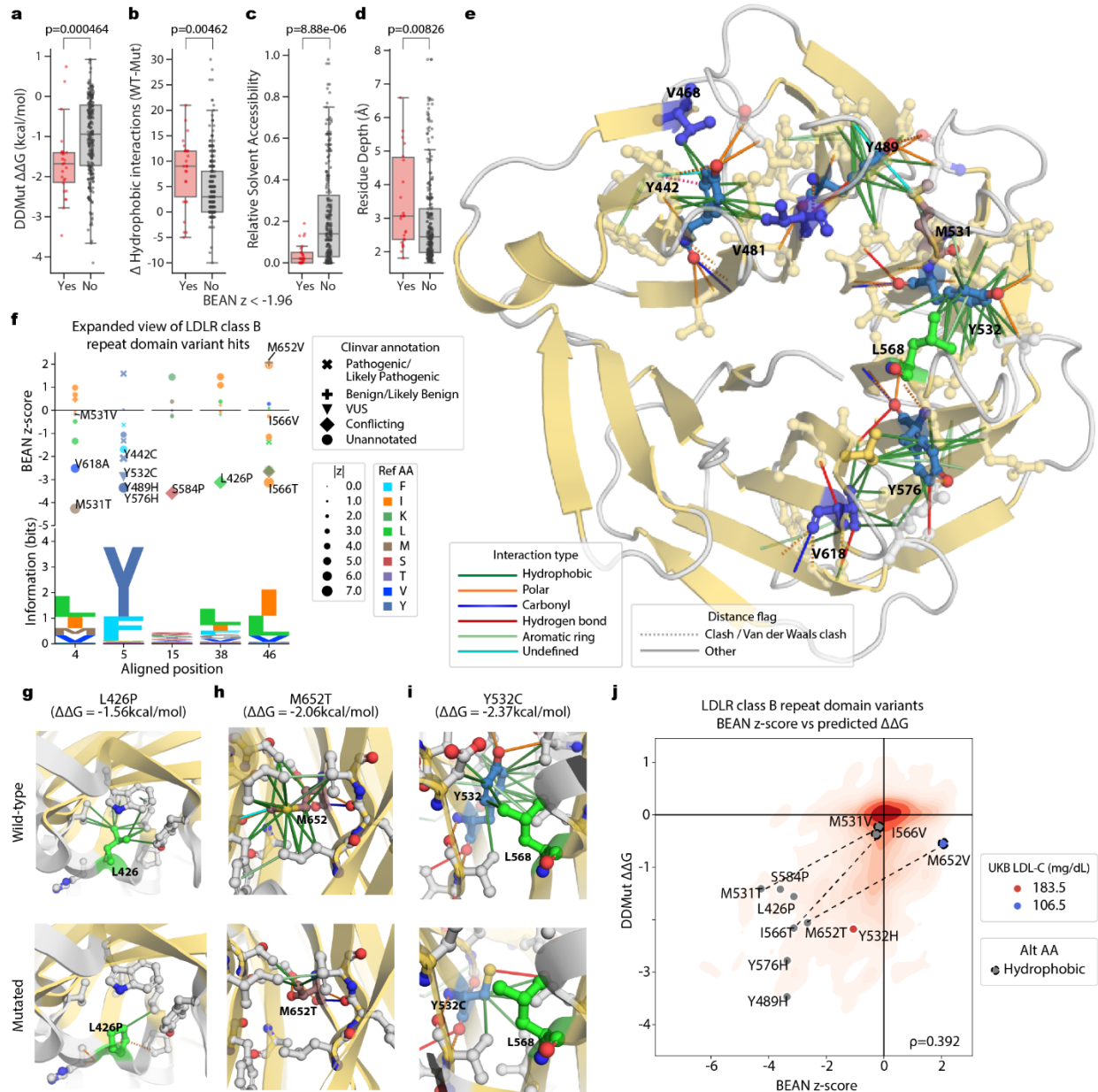
## 392 **Structural basis of *LDLR* missense variants**

393 We further analyzed *LDLR* missense variants identified to significantly impair LDL-C uptake by BEAN to  
394 gain insight into mechanisms of their pathogenicity. We first examined variants with top z-scores that  
395 are unannotated or annotated as conflicting, or VUS in ClinVar. The top ranked variant, which shows  
396 even more significant loss-of-function than splice-ablating variants, alters the start codon, preventing

397 full-length LDLR translation. Other top variants such as C222R, C261R, C289R, and C364R disrupt  
398 conserved disulfide bond-forming cysteines in LDLR class A repeats and EGF-like domains. Top-ranked  
399 variant in the of the signal peptide L16P substitute hydrophobic leucines with prolines in the  
400 transmembrane alpha helix, which is likely to distort the alpha helix<sup>72</sup> and the neighboring L15P has  
401 been shown to reduce LDLR transport to the plasma membrane<sup>73</sup>. Neighboring L14S that also ranks high  
402 substitutes hydrophobic leucine with serine in the hydrophobic h-region central to the signal peptide<sup>74</sup>.  
403 Additionally, multiple variants disrupt calcium ion binding, which is key to LDLR class A repeat folding<sup>75</sup>  
404 through the conversion of negatively charged amino acids (D/E) to glycine (G), thereby disrupting ionic  
405 interactions with side-chain carboxylates and calcium ions (D94G, E101G, E179G, D307G) in LDLR class A  
406 repeats (**Supplementary Fig. 22**). We also found that L371P, a VUS, disrupts a calcium ion interaction in  
407 the EGF-like domain by breaking the coordinate covalent bond between the calcium ion and the  
408 carbonyl group within the L371 main chain due to backbone distortion. Finally, we found that F153P  
409 significantly interfered with hydrophobic interactions between the aromatic ring and the attached  
410 saccharide on Q182.

411 We noticed that an appreciable number of deleterious variants that lack ClinVar pathogenic designation  
412 reside in the six LDLR class B repeats. The LDLR class B repeats, also known as YWTD repeats, form a  
413 propeller-like structure involved in the release of LDL following its endocytosis. To gain insights into  
414 unannotated variant impact, focusing on the LDLR class B repeats, we used the full wild-type LDLR  
415 structure from the AlphaFold Protein Structure Database<sup>76,77</sup> and the MODELLER<sup>78</sup>-generated mutant  
416 structures to calculate changes in interatomic interactions using Arpeggio<sup>79</sup>. Additionally, we predicted  
417 the effects of variants on protein stability ( $\Delta\Delta G$ , negative value indicates destabilization) with DDMut<sup>80</sup>  
418 (**Supplementary Table 8**). We found that the 26 significant LDLR class B variants induce more  
419 destabilizing effects, disrupt more hydrophobic interactions, have lower relative solvent accessibility<sup>81</sup>  
420 (0.041 of maximum residue solvent accessibility), and have higher wild-type residue depth as compared

421 to the other observed variants in this region (**Fig. 6a-d, Supplementary Fig. 23**). Collectively, these  
422 observations strongly indicate that these significant LDLR class B repeat variants are predominantly  
423 buried within the protein core where they engage in extensive hydrophobic interactions essential for  
424 protein folding. Moreover, we found a conserved interaction across repeat domains in which a tyrosine  
425 (aligned position 5 in **Supplementary Fig. 18b**) holds neighboring propeller blades together through  
426 interactions with a hydrophobic residue of the neighboring repeat (**Fig. 6e-f**). We identified five of these  
427 variant pairs (Y442C with V481A, Y442C with V468A, Y489H with M531T, Y532C with L568P, and Y576H  
428 with V618A), where all nine positions have at least one variant that weakens their hydrophobic  
429 interaction and has a significant BEAN z-score. Among the top-ranked unannotated or ClinVar VUS and  
430 conflicting variants within LDLR class B repeats, the six most significant variants (L426P, Y489H, M531T,  
431 I566T, S584P, and Y576H) all disrupt residues that hold the propeller blades together through  
432 hydrophobic interactions (**Fig. 6g-i, Supplementary Fig. 24**). Further supporting the importance of  
433 hydrophobic interactions, the base editing screen installed additional missense variants at positions 531,  
434 566, and 652 that conserve hydrophobicity. In all cases, mutation into hydrophobic residues has less  
435 severe impact from the base editing screen and DDMut-predicted destabilization than mutation into  
436 non-hydrophobic residues (**Fig. 6j**). For example, while we find M652T to be highly deleterious (BEAN z=-  
437 2.65), we find no functional disruption from the hydrophobicity-conserving M652V variant (BEAN  
438 z=+2.06). This analysis is supported clinically, as M652V is designated in ClinVar as “Likely Benign,” and  
439 the average UKB carrier LDL-C is below average (106mg/dL). In summary, structural analysis of rare *LDLR*  
440 variants identified by BEAN provides a basis for the missense variant impact through affecting structural  
441 integrity of LDLR, highlighting a central role for hydrophobic interactions that hold together adjacent  
442 beta blades of the LDLR class B repeat domain.



**Figure 6. Deleterious variants in LDLR class B repeats weaken hydrophobic interactions. a-d)** Boxplots of 26 significant ( $z < -1.96$ ) and the rest of 259 variants observed in LDLR class B repeats. P-values of two-sided Wilcoxon rank-sum test are denoted. WT; wild-type, Mut; mutated. **e)** Conserved interactions involving tyrosine of which mutation showed significant BEAN scores. Simplified interaction types and distance flags as annotated by Arpeggio are shown in the legend. **f)** BEAN z-scores of positions with conserved hydrophobic residues are shown along with the LDLR class B repeat PFAM HMM logo. **g-i)** Local atomic interaction in wild-type and mutated structure for

ClinVar conflicting variants or VUS L426P, M652T, and Y532C. Residues in the variant positions are colored by the reference amino acids. Residues that interact with the variant position are shown. Variant position and interacting residues are colored by elements (O: red, N: blue, S: yellow). **j**) Contour plot of BEAN z-score against  $\Delta\Delta G$  predicted by DDMut for 872 missense variants. Positions with distinct observed missense variants that disrupt and conserve hydrophobic sidechains are connected by dashed line.

## 443 **Discussion**

444 In this work, we develop a framework to improve variant effect quantification in base editing screens  
445 through accounting for variable genotypic outcomes estimated by gRNA efficiencies and chromatin  
446 accessibility. The activity-normalized approach is straightforward to apply, simply requiring synthesis of  
447 gRNA-reporter pairs, which can be cloned and screened using standard experimental procedures. This  
448 approach should prove particularly useful when employing new CRISPR enzymes and deaminases, as it  
449 allows for simultaneous characterization of editing preferences and phenotypic screening. We note that,  
450 while the screens described herein utilized flow cytometric phenotypic readouts, BEAN should also be  
451 suitable in dropout and enrichment paradigms by performing reporter analysis at an early timepoint  
452 prior to extensive phenotypic selection. We provide an open source implementation of BEAN in the  
453 comprehensive Python package *bean* with end-to-end functionality from base editing screen sequencing  
454 data to variant effect quantification at <https://pypi.org/project/crispr-bean/>.

455 Our results show that careful Bayesian modeling of the data generation process can substantially  
456 improve analytical power. We show that incorporation of reporter editing outcomes (BEAN-Reporter)  
457 and accessibility (BEAN) improves classification over the minimal model (BEAN-Uniform). In our work,  
458 we take into account the dependence of editing on loci accessibility from measurements taken from 49  
459 gRNAs at four loci and use the relationship to improve our model (we provide a guide on fitting this  
460 relationship in **Supplementary Note 2**). Higher-throughput measurements and incorporation of  
461 additional epigenetic features influencing base editing could refine the inference of endogenous editing

462 efficiency. It's also worth noting that BEAN currently does not consider off-target editing, an omission  
463 that may affect the evaluation of phenotypic impacts for certain gRNAs.

464 BEAN also makes certain assumptions regarding how data is distributed. It assumes that phenotypic  
465 readout follows a Gaussian distribution and that alleles with multiple variants show additive effects. We  
466 also present an approach to modeling multivariate gRNA and allele count data by employing a Dirichlet-  
467 Multinomial distribution, building on the Negative Binomial modeling of counts used in prior  
468 methods<sup>82,83</sup> (**Supplementary Note 3**).

469 We show that BEAN outperforms existing analysis methods at classifying GWAS variants with phenotypic  
470 impacts. Applying BEAN to the LDL-C GWAS library screen, we found it superior to existing methods at  
471 distinguishing positive control splice-altering variants. Furthermore, BEAN enabled accurate inference of  
472 variant effects on LDL-C uptake that were recapitulated using individual lentiviral gRNA transduction.

473 We used BEAN to uncover common variants that modulate LDL-C levels through altering liver cell  
474 expression/function of three previously unappreciated genes, *OPRL1*, *VTN*, and *ZNF329*. It is unclear why  
475 individuals with rare deleterious variants in these genes do not show altered LDL-C levels, but given that  
476 these genes have evidence of selective constraint<sup>84</sup>, we speculate that variants that are tolerated in the  
477 population may have weak functional and phenotypic effects<sup>85</sup>. The GWAS effect sizes for these variants  
478 are small, so they are unlikely to represent therapeutic targets; nonetheless, their elucidation  
479 contributes to understanding of the complex genetic underpinnings of lipid homeostasis.

480 We additionally demonstrate that dense activity-normalized base editing screens can improve  
481 characterization of coding variants in *LDLR*. Prior coding sequence-targeted base editing screens have  
482 used more restrictive PAMs<sup>13,15–18,20–23,25–27</sup>, thus limiting their breadth. By combining ABE8e-SpRY, which  
483 we show to have robust activity across the vast majority of PAMs, with BEAN, which enables sharing of  
484 information across adjacent gRNAs to boost power, we obtain accurate phenotypic measurements for



485 an average of one variant per amino acid. This resolution is less than that of DMS, which can evaluate all  
486 possible missense variants at each position<sup>4,86</sup>. However, base editing is considerably less work-intensive  
487 and less expensive to perform and is far more scalable, allowing assessment of sets of genes in a single  
488 experiment. We provide pathogenicity assessment of 874 *LDLR* missense variants, most of which do not  
489 have prior clinical designation. Structural characterization of identified *LDLR* variants reveals distinct  
490 domain-specific characteristics of the most deleterious variants, including a central role of hydrophobic  
491 interactions gluing adjacent LDLR class B repeat domain's beta blades.

492 Past coding variant editing screens have focused on binary classification of pathogenic and benign  
493 variants, and BEAN effectively distinguishes these variant classes in *LDLR*, while also making predictions  
494 about dozens of variants of unknown significance and conflicting annotation. Notably, we also show that  
495 BEAN scores associate with quantitative serum LDL-C levels measured in patients with rare *LDLR* variants  
496 in the UKB. Functional assays are accepted as evidence of pathogenicity or benignity in clinical variant  
497 interpretation guidelines such as those published by the American College of Medical Genetics and the  
498 Association for Molecular Pathology<sup>87,88</sup>. However, these frameworks have focused on classifying larger  
499 effect pathogenic and benign variants in a binary fashion, and these approaches have not been designed  
500 to offer risk predictions for quantitative traits. Given that CAD risk depends proportionally on the level  
501 of lifelong serum LDL-C exposure<sup>89,90</sup>, our ability to assign quantitative estimates of clinical risk to *LDLR*  
502 variants has strong clinical significance that, if replicable in other disease-associated genes, may merit  
503 adopting an assessment paradigm for clinical risk based on quantitative traits. We additionally show,  
504 albeit on a small cohort, that *LDLR* functional scores associate more closely with CAD incidence than  
505 LDL-C levels. This result is consistent with the residual CAD risk of FH patients even after LDL-reductive  
506 therapy<sup>91</sup> and reinforces the value of *LDLR* pathogenicity analysis above and beyond LDL-C monitoring.  
507 We show that base editing-derived functional scores correlate moderately with evolutionary  
508 conservation (Pfam, PhyloP) as well as structural disruption (DDMut) analyses. We show preliminary

509 evidence that integrating BEAN and PhyloP scores improves prediction of LDL-C levels. Thus, functional  
510 screening offers at least partially orthogonal information to other forms of evidence that have been  
511 used to build computational pathogenicity predictors<sup>92-94</sup>. We anticipate that principled integration of  
512 distinct lines of evidence including MAVE data will improve pathogenicity prediction, since each  
513 evidence type is an imperfect surrogate for the effects of a variant over the lifespan of an individual. For  
514 example, our LDL-C uptake screening fails to measure how variants impact uptake of other lipoproteins  
515 by LDLR<sup>95</sup>, interaction with PCSK9<sup>96</sup>, or expression or function across environmental conditions or cell  
516 types, and thus its accuracy may be improved by integrating orthogonal methods. In conclusion, activity-  
517 normalized base editor reporter screening with BEAN markedly improves variant impact quantification  
518 from CRISPR base editing screens. Given the importance of variants of weak effect in complex disease,  
519 this approach promises to accelerate the characterization of human disease-associated variants in their  
520 native genomic context.

## References

1. Tam, V. *et al.* Benefits and limitations of genome-wide association studies. *Nat. Rev. Genet.* **20**, 467–484 (2019).
2. Bycroft, C. *et al.* The UK Biobank resource with deep phenotyping and genomic data. *Nature* **562**, 203–209 (2018).
3. Gasperini, M., Starita, L. & Shendure, J. The power of multiplexed functional analysis of genetic variants. *Nat. Protoc.* **11**, 1782–1787 (2016).
4. Araya, C. L. & Fowler, D. M. Deep mutational scanning: assessing protein function on a massive scale. *Trends Biotechnol.* **29**, 435–442 (2011).
5. Myers, R. M., Tilly, K. & Maniatis, T. Fine structure genetic analysis of a beta-globin promoter. *Science* **232**, 613–618 (1986).
6. Inoue, F. & Ahituv, N. Decoding enhancers using massively parallel reporter assays. *Genomics* **106**, 159–164 (2015).
7. Bock, C. *et al.* High-content CRISPR screening. *Nature Reviews Methods Primers* **2**, 1–23 (2022).
8. Shalem, O. *et al.* Genome-scale CRISPR-Cas9 knockout screening in human cells. *Science* **343**, 84–87 (2014).
9. Wang, T., Wei, J. J., Sabatini, D. M. & Lander, E. S. Genetic screens in human cells using the CRISPR-Cas9 system. *Science* **343**, 80–84 (2014).
10. Komor, A. C., Kim, Y. B., Packer, M. S., Zuris, J. A. & Liu, D. R. Programmable editing of a target base in genomic DNA without double-stranded DNA cleavage. *Nature* **533**, 420–424 (2016).
11. Gaudelli, N. M. *et al.* Programmable base editing of A•T to G•C in genomic DNA without DNA cleavage. *Nature* **551**, 464–471 (2017).
12. Rees, H. A. & Liu, D. R. Base editing: precision chemistry on the genome and transcriptome of living cells. *Nat. Rev. Genet.* **19**, 770–788 (2018).

13. Hanna, R. E. *et al.* Massively parallel assessment of human variants with base editor screens. *Cell* **184**, 1064-1080.e20 (2021).
14. Morris, J. A. *et al.* Discovery of target genes and pathways at GWAS loci by pooled single-cell CRISPR screens. *Science* **380**, eadh7699 (2023).
15. Martin-Rufino, J. D. *et al.* Massively parallel base editing to map variant effects in human hematopoiesis. *Cell* **186**, 2456-2474.e24 (2023).
16. Cuella-Martin, R. *et al.* Functional interrogation of DNA damage response variants with base editing screens. *Cell* **184**, 1081-1097.e19 (2021).
17. Pablo, J. L. B. *et al.* Scanning mutagenesis of the voltage-gated sodium channel Nav1.2 using base editing. *Cell Rep.* **42**, 112563 (2023).
18. Coelho, M. A. *et al.* Base editing screens map mutations affecting interferon- $\gamma$  signaling in cancer. *Cancer Cell* **41**, 288-303.e6 (2023).
19. Cheng, L. *et al.* Single-nucleotide-level mapping of DNA regulatory elements that control fetal hemoglobin expression. *Nat. Genet.* **53**, 869–880 (2021).
20. Sánchez-Rivera, F. J. *et al.* Base editing sensor libraries for high-throughput engineering and functional analysis of cancer-associated single nucleotide variants. *Nat. Biotechnol.* 1–12 (2022).
21. Kim, Y. *et al.* High-throughput functional evaluation of human cancer-associated mutations using base editors. *Nat. Biotechnol.* (2022) doi:10.1038/s41587-022-01276-4.
22. Kweon, J. *et al.* A CRISPR-based base-editing screen for the functional assessment of BRCA1 variants. *Oncogene* **39**, 30–35 (2020).
23. Huang, C., Li, G., Wu, J., Liang, J. & Wang, X. Identification of pathogenic variants in cancer genes using base editing screens with editing efficiency correction. *Genome Biol.* **22**, 80 (2021).
24. Sangree, A. K. *et al.* Benchmarking of SpCas9 variants enables deeper base editor screens of BRCA1 and BCL2. *Nat. Commun.* **13**, 1318 (2022).

25. Lue, N. Z. *et al.* Base editor scanning charts the DNMT3A activity landscape. *Nat. Chem. Biol.* **19**, 176–186 (2023).
26. Després, P. C., Dubé, A. K., Seki, M., Yachie, N. & Landry, C. R. Perturbing proteomes at single residue resolution using base editing. *Nat. Commun.* **11**, 1871 (2020).
27. Garcia, E. M. *et al.* Base Editor Scanning Reveals Activating Mutations of DNMT3A. *bioRxiv* 2023.04.12.536656 (2023) doi:10.1101/2023.04.12.536656.
28. Lue, N. Z. & Liau, B. B. Base editor screens for in situ mutational scanning at scale. *Mol. Cell* **83**, 2167–2187 (2023).
29. Arbab, M. *et al.* Determinants of Base Editing Outcomes from Target Library Analysis and Machine Learning. *Cell* **182**, 463-480.e30 (07/2020).
30. Graham, S. E. *et al.* The power of genetic diversity in genome-wide association studies of lipids. *Nature* **600**, 675–679 (2021).
31. Wang, R., Lin, D.-Y. & Jiang, Y. EPIC: Inferring relevant cell types for complex traits by integrating genome-wide association studies and single-cell RNA sequencing. *PLoS Genet.* **18**, e1010251 (2022).
32. Finucane, H. K. *et al.* Heritability enrichment of specifically expressed genes identifies disease-relevant tissues and cell types. *Nat. Genet.* **50**, 621–629 (2018).
33. Bouhairie, V. E. & Goldberg, A. C. Familial hypercholesterolemia. *Cardiol. Clin.* **33**, 169–179 (2015).
34. Brown, M. S. & Goldstein, J. L. How LDL receptors influence cholesterol and atherosclerosis. *Sci. Am.* **251**, 58–66 (1984).
35. Mundal, L. J. *et al.* Impact of age on excess risk of coronary heart disease in patients with familial hypercholesterolaemia. *Heart* **104**, 1600–1607 (2018).
36. Landrum, M. J. *et al.* ClinVar: public archive of relationships among sequence variation and human phenotype. *Nucleic Acids Res.* **42**, D980-5 (2014).

37. Hamilton, M. C. *et al.* Systematic elucidation of genetic mechanisms underlying cholesterol uptake. *Cell Genomics* **3**, 100304 (2023).
38. Spady, D. K. Hepatic clearance of plasma low density lipoproteins. *Semin. Liver Dis.* **12**, 373–385 (1992).
39. Richter, M. F. *et al.* Phage-assisted evolution of an adenine base editor with improved Cas domain compatibility and activity. *Nat. Biotechnol.* **38**, 883–891 (2020).
40. Walton, R. T., Christie, K. A., Whittaker, M. N. & Kleinstiver, B. P. Unconstrained genome targeting with near-PAMless engineered CRISPR-Cas9 variants. *Science* **368**, 290–296 (2020).
41. Shin, H. R. *et al.* Small-molecule inhibitors of histone deacetylase improve CRISPR-based adenine base editing. *Nucleic Acids Res.* **49**, 2390–2399 (2021).
42. Yang, C. *et al.* HMGN1 enhances CRISPR-directed dual-function A-to-G and C-to-G base editing. *Nat. Commun.* **14**, 2430 (2023).
43. Wang, G., Sarkar, A., Carbonetto, P. & Stephens, M. A simple new approach to variable selection in regression, with application to genetic fine mapping. *J. R. Stat. Soc. Series B Stat. Methodol.* **82**, 1273–1300 (2020).
44. Weissbrod, O. *et al.* Functionally informed fine-mapping and polygenic localization of complex trait heritability. *Nat. Genet.* **52**, 1355–1363 (2020).
45. Emmer, B. T. *et al.* Genome-scale CRISPR screening for modifiers of cellular LDL uptake. *PLoS Genet.* **17**, e1009285 (2021).
46. Arbab, M. *et al.* Base editing rescue of spinal muscular atrophy in cells and in mice. *Science* **380**, eadg6518 (2023).
47. Schep, R. *et al.* Impact of chromatin context on Cas9-induced DNA double-strand break repair pathway balance. *Mol. Cell* **81**, 2216–2230.e10 (2021).

48. Ding, X. *et al.* Improving CRISPR-Cas9 Genome Editing Efficiency by Fusion with Chromatin-Modulating Peptides. *CRISPR J* **2**, 51–63 (2019).
49. Liu, G., Yin, K., Zhang, Q., Gao, C. & Qiu, J.-L. Modulating chromatin accessibility by transactivation and targeting proximal dsRNAs enhances Cas9 editing efficiency in vivo. *Genome Biol.* **20**, 145 (2019).
50. Klimentidis, Y. C. *et al.* Phenotypic and genetic characterization of lower LDL cholesterol and increased type 2 diabetes risk in the UK Biobank. *Diabetes* **69**, 2194–2205 (2020).
51. Li, W. *et al.* MAGeCK enables robust identification of essential genes from genome-scale CRISPR/Cas9 knockout screens. *Genome Biol.* **15**, 554 (2014).
52. Li, W. *et al.* Quality control, modeling, and visualization of CRISPR screens with MAGeCK-VISPR. *Genome Biol.* **16**, 281 (2015).
53. Jeong, H.-H., Kim, S. Y., Rousseaux, M. W. C., Zoghbi, H. Y. & Liu, Z. Beta-binomial modeling of CRISPR pooled screen data identifies target genes with greater sensitivity and fewer false negatives. *Genome Res.* **29**, 999–1008 (2019).
54. Daley, T. P. *et al.* CRISPhierMix: a hierarchical mixture model for CRISPR pooled screens. *Genome Biol.* **19**, 159 (2018).
55. Degner, J. F. *et al.* DNase I sensitivity QTLs are a major determinant of human expression variation. *Nature* **482**, 390–394 (2012).
56. Currin, K. W. *et al.* Genetic effects on liver chromatin accessibility identify disease regulatory variants. *Am. J. Hum. Genet.* **108**, 1169–1189 (2021).
57. Yao, Q. *et al.* Motif-Raptor: a cell type-specific and transcription factor centric approach for post-GWAS prioritization of causal regulators. *Bioinformatics* **37**, 2103–2111 (2021).
58. Weirauch, M. T. *et al.* Determination and inference of eukaryotic transcription factor sequence specificity. *Cell* **158**, 1431–1443 (2014).

59. Partridge, E. C. *et al.* Occupancy maps of 208 chromatin-associated proteins in one human cell type. *Nature* **583**, 720–728 (2020).
60. Farh, K. K.-H. *et al.* Genetic and epigenetic fine mapping of causal autoimmune disease variants. *Nature* **518**, 337–343 (2015).
61. Jing, Z., Liu, Y., Dong, M., Hu, S. & Huang, S. Identification of the DNA binding element of the human ZNF333 protein. *J. Biochem. Mol. Biol.* **37**, 663–670 (2004).
62. Witzgall, R., O’Leary, E., Leaf, A., Onaldi, D. & Bonventre, J. V. The Krüppel-associated box-A (KRAB-A) domain of zinc finger proteins mediates transcriptional repression. *Proc. Natl. Acad. Sci. U. S. A.* **91**, 4514–4518 (1994).
63. Fass, D., Blacklow, S., Kim, P. S. & Berger, J. M. Molecular basis of familial hypercholesterolaemia from structure of LDL receptor module. *Nature* **388**, 691–693 (1997).
64. Russell, D. W., Brown, M. S. & Goldstein, J. L. Different combinations of cysteine-rich repeats mediate binding of low density lipoprotein receptor to two different proteins. *J. Biol. Chem.* **264**, 21682–21688 (1989).
65. Jeon, H. & Blacklow, S. C. Structure and physiologic function of the low-density lipoprotein receptor. *Annu. Rev. Biochem.* **74**, 535–562 (2005).
66. Bateman, A. *et al.* Pfam 3.1: 1313 multiple alignments and profile HMMs match the majority of proteins. *Nucleic Acids Res.* **27**, 260–262 (1999).
67. Discovery and refinement of loci associated with lipid levels. *Nat. Genet.* **45**, 1274–1283 (2013).
68. Yu, T., Fife, J. D., Adzhubey, I., Sherwood, R. & Cassa, C. A. Joint estimation and imputation of variant functional effects using high throughput assay data. *medRxiv* (2023)  
doi:10.1101/2023.01.06.23284280.
69. Chen, T. & Guestrin, C. XGBoost: A Scalable Tree Boosting System. *arXiv [cs.LG]* (2016).



70. Clarke, S. L. *et al.* Coronary artery disease risk of familial hypercholesterolemia genetic variants independent of clinically observed longitudinal cholesterol exposure. *Circ. Genom. Precis. Med.* **15**, e003501 (2022).
71. Oommen, D., Kizhakkedath, P., Jawabri, A. A., Varghese, D. S. & Ali, B. R. Proteostasis Regulation in the Endoplasmic Reticulum: An Emerging Theme in the Molecular Pathology and Therapeutic Management of Familial Hypercholesterolemia. *Front. Genet.* **11**, 570355 (2020).
72. Kim, M. K. & Kang, Y. K. Positional preference of proline in alpha-helices. *Protein Sci.* **8**, 1492–1499 (1999).
73. Pavloušková, J., Réblová, K., Tichý, L., Freiberger, T. & Fajkusová, L. Functional analysis of the p.(Leu15Pro) and p.(Gly20Arg) sequence changes in the signal sequence of LDL receptor. *Atherosclerosis* **250**, 9–14 (2016).
74. von Heijne, G. Signal sequences. The limits of variation. *J. Mol. Biol.* **184**, 99–105 (1985).
75. Pena, F., Jansens, A., van Zadelhoff, G. & Braakman, I. Calcium as a crucial cofactor for low density lipoprotein receptor folding in the endoplasmic reticulum. *J. Biol. Chem.* **285**, 8656–8664 (2010).
76. Jumper, J. *et al.* Highly accurate protein structure prediction with AlphaFold. *Nature* **596**, 583–589 (2021).
77. Varadi, M. *et al.* AlphaFold Protein Structure Database: massively expanding the structural coverage of protein-sequence space with high-accuracy models. *Nucleic Acids Res.* **50**, D439–D444 (2022).
78. Webb, B. & Sali, A. Comparative protein structure modeling using MODELLER. *Curr. Protoc. Protein Sci.* **86**, 2.9.1-2.9.37 (2016).
79. Jubb, H. C. *et al.* Arpeggio: A web server for calculating and visualising interatomic interactions in protein structures. *J. Mol. Biol.* **429**, 365–371 (2017).

80. Zhou, Y., Pan, Q., Pires, D. E. V., Rodrigues, C. H. M. & Ascher, D. B. DDMut: predicting effects of mutations on protein stability using deep learning. *Nucleic Acids Res.* **51**, W122–W128 (2023).
81. Rose, G. D., Geselowitz, A. R., Lesser, G. J., Lee, R. H. & Zehfus, M. H. Hydrophobicity of amino acid residues in globular proteins. *Science* **229**, 834–838 (1985).
82. Anders, S. & Huber, W. Differential expression analysis for sequence count data. *Genome Biol.* **11**, R106 (2010).
83. 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).
84. Karczewski, K. J. *et al.* The mutational constraint spectrum quantified from variation in 141,456 humans. *Nature* **581**, 434–443 (2020).
85. Cassa, C. A. *et al.* Estimating the selective effects of heterozygous protein-truncating variants from human exome data. *Nat. Genet.* **49**, 806–810 (2017).
86. Fowler, D. M. & Fields, S. Deep mutational scanning: a new style of protein science. *Nat. Methods* **11**, 801–807 (2014).
87. Richards, S. *et al.* Standards and guidelines for the interpretation of sequence variants: a joint consensus recommendation of the American College of Medical Genetics and Genomics and the Association for Molecular Pathology. *Genet. Med.* **17**, 405–424 (2015).
88. Brnich, S. E. *et al.* Recommendations for application of the functional evidence PS3/BS3 criterion using the ACMG/AMP sequence variant interpretation framework. *Genome Med.* **12**, 3 (2019).
89. Domanski, M. J. *et al.* Time Course of LDL Cholesterol Exposure and Cardiovascular Disease Event Risk. *J. Am. Coll. Cardiol.* **76**, 1507–1516 (2020).
90. Duncan, M. S., Vasan, R. S. & Xanthakis, V. Trajectories of Blood Lipid Concentrations Over the Adult Life Course and Risk of Cardiovascular Disease and All-Cause Mortality: Observations From the Framingham Study Over 35 Years. *J. Am. Heart Assoc.* **8**, e011433 (2019).

91. Mundal, L. & Retterstøl, K. A systematic review of current studies in patients with familial hypercholesterolemia by use of national familial hypercholesterolemia registries. *Curr. Opin. Lipidol.* **27**, 388–397 (2016).
92. Gao, H. *et al.* The landscape of tolerated genetic variation in humans and primates. *Science* **380**, eabn8153 (2023).
93. Frazer, J. *et al.* Disease variant prediction with deep generative models of evolutionary data. *Nature* 1–5 (2021).
94. Ioannidis, N. M. *et al.* REVEL: An Ensemble Method for Predicting the Pathogenicity of Rare Missense Variants. *Am. J. Hum. Genet.* **99**, 877–885 (2016).
95. Go, G.-W. & Mani, A. Low-density lipoprotein receptor (LDLR) family orchestrates cholesterol homeostasis. *Yale J. Biol. Med.* **85**, 19–28 (2012).
96. Lagace, T. A. PCSK9 and LDLR degradation: regulatory mechanisms in circulation and in cells. *Curr. Opin. Lipidol.* **25**, 387–393 (2014).

521 **Figure 1. Activity-normalized base editing screening pipeline. a)** Schematic of activity-normalized base  
522 editing screening process and analysis by BEAN. A library of gRNAs, each paired with a reporter  
523 sequence encompassing its genomic target sequence, is cloned into a lentiviral base editor expression  
524 vector. Lentiviral transduction is performed in HepG2, followed by flow cytometric sorting of four  
525 populations based on fluorescent LDL-cholesterol (BODIPY-LDL) uptake. The gRNA and reporter  
526 sequences are read out by paired-end NGS to obtain gRNA counts and reporter editing outcomes in  
527 each flow cytometric bin. BEAN models the reporter editing frequency and allelic outcomes and gRNA  
528 enrichments among flow cytometric bins using BEAN to estimate variant phenotypic effect sizes. **b)**  
529 Adjacent nucleotide specificity of ABE8e-SpRY editing represented as a sequence logo from 7,320  
530 gRNAs; the height of each base represents the relative frequency of observing each base given an edit at  
531 position 0. **c)** Average editing efficiency of ABE8e-SpRY by protospacer position and PAM sequence **d)**  
532 Scatterplots comparing nucleotide-level editing efficiency between the reporter and endogenous target  
533 sites for a total of 49 gRNAs across four loci across 3 experimental replicates. The accessibility of the four  
534 loci as measured by ATAC-seq signal in HepG2 is shown in the top panel, and the scatterplot markers are  
535 colored by the accessibility of each nucleotide. Pearson correlation coefficients are shown as  $r$ . **e)**  
536 Schematic of the LDL-C variant library gRNA design for selected GWAS candidate variants with a  
537 Manhattan plot showing variant P-values from a recent GWAS study<sup>50</sup>. gRNAs tile the variant at five  
538 positions with maximal editing efficiency (protospacer positions 4-8). **f)** gRNA coverage of the LDLR tiling  
539 library across LDLR coding sequence along with 5' and 3' UTRs and several regulatory regions.

540 **Figure 2. BEAN models variant effects from activity-normalized base editing screens.** Simplified  
541 schematic of BEAN Bayesian network that models input reporter editing outcomes and gRNA counts.  
542 The Bayesian network model recapitulates the data generation process starting from a variant-level  
543 phenotype  $Y_v$  and models per-gRNA phenotypes as a Gaussian mixture distribution of edited and  
544 unedited (wild-type) allele phenotypes. The weights of the mixture components are modeled to

545 generate reporter editing outcomes. gRNA abundance in each sorting bin is then calculated by  
546 discretizing the gRNA phenotype based on the experimental design into the phenotypic quantiles, and is  
547 modeled to generate the observed gRNA counts using an overdispersed multivariate count distribution  
548 (see Methods). BEAN outputs the parameters of the posterior distribution of mean phenotypic shift as  
549 Gaussian distribution with mean  $\widehat{\mu}_{\mu}$  (effect size), along with negative-control adjusted z-score and  
550 credible interval (CI), where  $\mathcal{D}$  is the input data.

551 **Figure 3. BEAN improves variant impact estimation from the LDL-C GWAS library screen. a)** Ridge plot  
552 of BEAN z-score distributions of positive controls, negative controls, and test variants. **b)** AUPRC of  
553 classifying *LDLR* and *MYLIP* splicing variants vs. negative controls. Metrics for all 5 replicates are shown  
554 as markers and metrics of 15 two-replicate subsamples among the 5 replicates are shown as box plots.  
555 **c)** Spearman correlation coefficient between BEAN effect size and log fold change of LDL-C uptake  
556 following individual testing of 26 gRNAs. Metrics for all 5 replicates are shown as marker and metrics of  
557 15 two-replicate subsamples among the 5 replicates are shown as box plots. **d)** Scatterplot of BEAN  
558 effect size and log fold change of LDL-C uptake following individual testing of 26 gRNAs. Spearman  
559 correlation coefficient is denoted as  $\rho$ . **e)** Scatterplot of variant effect size and significance estimated by  
560 BEAN. Labels show rsIDs of selected variants and dbSNP gene annotations and a manual annotation for  
561 APOB enhancer in the parenthesis.

562 **Figure 4. Functional characterization of LDL-C GWAS variants. a)** Schematic of potential variant  
563 mechanisms and the figure panels showing data from each mechanistic experiment. **b)** Schematic of  
564 pooled ATAC-seq analysis to identify variants impacting accessibility. Differential representation of allele  
565 in gDNA and ATAC-seq reflects differential accessibility induced by the base edit or heterozygous  
566 reference allele. **c)** Change in ATAC-seq enrichment from the pooled ATAC-seq experiment. 95%  
567 confidence intervals are shown as the error bars. “Edited variant” denotes the enrichment by base edit

568 and “Allele” denotes the enrichment by either of the heterozygous alleles, when available, translated  
569 uniformly to major (Maj) to minor (Min) allele direction. Variants where the base edit is conducted from  
570 minor allele are denoted as red in the color bar. Family-wise error rate (FWER) with Benjamini-Hochberg  
571 multiple correction is shown for each enrichment value where  $FWER < 0.1$ . **d)** Genomic tracks for three  
572 selected variants. DNaseIHS; DNase 1 Hypersensitivity. Multiple transcript variants of RGS19 and OPRL1  
573 are shown in the middle panel. **e)** Fraction of *ZNF329* minor (Min) allele haplotype reads in gDNA and  
574 cDNA from untreated HepG2 and HepG2 with rs35081008<sup>Min>Maj</sup> base editing. **f)** Change in gene  
575 expression following base editing of three selected variants from minor (Min) to major (Maj) allele. P-  
576 values of the one sample Student’s t-test of LFC vs. mean of 0 that are smaller than 0.05 are shown  
577 above each bar. **g)** Change in cellular LDL-C uptake following CRISPRa/i of proximal genes for three  
578 selected variants. P-values of the one sample Student’s t-test of LFC vs. mean of 0 that are smaller than  
579 0.05 are shown above each bar. **h)** Summary schematic of characterization results.

580 **Figure 5. Dissection of *LDLR* variant effects through BEAN modeling of a saturation tiled base editing**  
581 **screen. a)** BEAN model for coding sequence tiling screens. Reporter editing efficiencies are calculated at  
582 the amino acid-level when the edited nucleotides are in coding region. Phenotypes of gRNAs with multi-  
583 allelic outcomes are modeled as the Gaussian mixture of allelic phenotypes. If an allele consists of more  
584 than one variant, the phenotype of the allele is modeled as the sum of the component variants. A  
585 Bayesian network is used to model variant-level phenotypes, sharing phenotypic information from all  
586 available gRNAs. **b)** Ridge plot of BEAN z-score distributions of positive controls, negative controls, and  
587 variants. **c)** Ridge plot of BEAN z-score distributions of Clinvar variants annotated as pathogenic/likely  
588 pathogenic (P/LP), benign/likely benign (B/LB), conflicting interpretation of pathogenicity (conflicting),  
589 and Uncertain significance (VUS), and unannotated variants. **d)** AUPRC of classifying ClinVar  
590 pathogenic/likely pathogenic vs. benign/likely benign variants. The marker shows the metrics of each  
591 method run on 4 replicates with no failing samples. Boxplot shows the metrics of 6 2-replicate

592 combinations of the 4 replicates. **e)** *LDLR* domain structure adopted from Oommen et al<sup>71</sup>. **f)** BEAN z-  
593 scores for variants in the 7 *LDLR* class A repeat domains aligned with the Pfam profile HMM logo. Highly  
594 conserved cysteines are highlighted in grey. **g)** Scatterplot of estimated variant effect sizes and z-scores.  
595 Labels of selected deleterious variants without ClinVar pathogenic/likely pathogenic annotations are  
596 shown. **h)** Scatterplot of *LDLR* class A repeat missense variant BEAN z-scores and  $\Delta$ Pfam profile HMM  
597 scores. Higher  $\Delta$ Pfam scores correspond to substitution from highly conserved to rarely observed amino  
598 acids. **i)** Comparison of mean statin-adjusted LDL-C level and BEAN z-score for variants observed in UKB  
599 and base editing. **j)** Comparison of mean statin-adjusted LDL-C level and BEAN-FUSE scores for variants  
600 observed in UKB and base editing. **k)** LDL-C levels of observed missense variants predicted by a  
601 regression model using BEAN-FUSE and PhyloP scores with 10-fold cross validation, compared with  
602 mean statin-adjusted LDL-C level in UKB. **l)** Boxplots of BEAN-FUSE functional scores for UKB individuals  
603 with variants observed in our base editing screen with or without CAD. **m)** Boxplots of statin-adjusted  
604 LDL-C levels of UKB individuals with variants observed in our base editing screen with or without CAD. P-  
605 value of two-sided Wilcoxon rank-sum test is denoted. *r*; Pearson correlation coefficient,  $\rho$ ; Spearman  
606 correlation coefficient

607 **Figure 6. Deleterious variants in *LDLR* class B repeats weaken hydrophobic interactions. a-d)** Boxplots  
608 of 26 significant ( $z < -1.96$ ) and the rest of 259 variants observed in *LDLR* class B repeats. P-values of  
609 two-sided Wilcoxon rank-sum test are denoted. WT; wild-type, Mut; mutated. **e)** Conserved interactions  
610 involving tyrosine of which mutation showed significant BEAN scores. Simplified interaction types and  
611 distance flags as annotated by Arpeggio are shown in the legend. **f)** BEAN z-scores of positions with  
612 conserved hydrophobic residues are shown along with the *LDLR* class B repeat PFAM HMM logo. **g-i)**  
613 Local atomic interaction in wild-type and mutated structure for ClinVar conflicting variants or VUS  
614 L426P, M652T, and Y532C. Residues in the variant positions are colored by the reference amino acids.  
615 Residues that interact with the variant position are shown. Variant position and interacting residues are

616 colored by elements (O: red, N: blue, S: yellow). **j**) Contour plot of BEAN z-score against  $\Delta\Delta G$  predicted  
617 by DDMut for 872 missense variants. Positions with distinct observed missense variants that disrupt and  
618 conserve hydrophobic sidechains are connected by dashed line.



## 619 **Methods**

### 620 **Establishing Cell Lines**

621 HepG2 cells were obtained from American Type Culture Collection (ATCC). HepG2 cells were infected  
622 with lentiviral constitutive base editor vectors pLenti\_ABE8e-SpRY-P2A-BFP\_HygroR and pLenti\_AID-  
623 BE5-SpRY-P2A-BFP\_HygroR. After Hygromycin selection, cells were sorted twice to enrich for BFP  
624 expression.

### 625 **Screen library design**

626 The LDL-C GWAS library was constructed to include gRNAs targeting variants associated with LDL-C  
627 levels from the UK Biobank GWAS cohort. Fine-mapped variants with posterior inclusion probability  
628 (PIP) > 0.25 from either the SUSIE or Polyfun fine-mapping pipelines (updated in December 2019,  
629 downloaded from <https://www.finucanelab.org/data>) were included, as well as variants with PIP > 0.1  
630 within 250 kb of any of 490 genes found to significantly alter LDL-C uptake from recent CRISPR-Cas9  
631 knockout screens<sup>1,2</sup>. All HepG2 haplotypes derived from the phased HepG2 genome sequence were  
632 targeted, and thus multiple allelic variants were targeted at certain genomic locations. Variants were  
633 assigned to ABE or CBE sub-libraries according to the identity of the affected nucleotide, and variants  
634 were included even in cases such as transversions and variable-length variants in which the edited  
635 variant would not exactly match the alternate allele identity. Five gRNAs that position the variant in  
636 protospacer positions 4-8 were included. gRNAs that contain TTTT stretches were removed from the  
637 library. Positive control gRNAs were designed to ablate all feasible splice donor and acceptor sites for  
638 ABE and to install 16 stop/gain variants for CBE, using five gRNAs per target site using the same logic as  
639 for variants. A total of 100 negative control non-targeting gRNAs were included in each sub-library,  
640 designed as 20 sets of five tiled gRNAs using the same logic as for variants. Paired reporter sequences  
641 were designed as the 32-nt genomic target sequence centered on the 20-nt gRNA.

642 The LDLR tiling library was constructed to include all gRNAs targeting coding regions (tiling density of 1)  
643 on both strands. The 26-nt intronic region surrounding each LDLR exon was tiled at 1/2 density (on both  
644 strands as in all subsequent cases), and the 24-nt distal to this region was tiled at 1/3 density. 5' and 3'  
645 UTR regions were tiled at 1/3 density. The two strongest LDLR intronic enhancers, both in intron 1, were  
646 also tiled at 1/3 density. In all cases, gRNAs were designed to match the HepG2 genomic sequence, and  
647 gRNAs were designed to target all HepG2 haplotypes when HepG2 is heterozygous at that sequence.  
648 gRNAs that contain TTTT stretches were removed from the library. 150 negative control non-targeting  
649 gRNAs were included. Paired reporter sequences were designed as the 32-nt genomic target sequence  
650 centered on the 20-nt gRNA.

651 Libraries were designed in the following template:

652 TGGAAAGGACGAAACACCG[19-20-nt gRNA]

653 GTTTAAGAGCTATGCTGGAAACAGCATAGCAAGTTTAAATAAGGCTAGTCCGTTATCAACTTGAAAAAGTGGCACC

654 GAGTCGGTGCTTTTTTT[32-nt target (6-nt upstream, 20-nt gRNA, 6-nt PAM)][4-nt barcode]

655 AGATCGGAAGAGCACACGNNNNNNNNNNNNNNNN

656 Where the final 14 N's are a variable primer sequence enabling pooling of sublibraries into a single  
657 synthesis order. The gRNA libraries were ordered from Agilent.

### 658 **Base editing screening**

659 The gRNA libraries were cloned into either CRISPRv2FE-ABE8e-SpRY-BsrGI or CRISPRv2FE-AIDBE5-SpRY-  
660 BsrGI. Libraries were amplified using NEBNext Ultra II Q5 mastermix, cloned using NEBuilder HiFi DNA  
661 Assembly mastermix, and electroporated into Endura competent cells (Biosearch Technologies) for  
662 propagation. Lentivirus was produced in HEK293T cells for each library using TransIT-Lenti transfection  
663 reagent, titered, and incubated with  $6.25 \times 10^6$  ABE8e-SpRY stable HepG2 and AID-BE5-SpRY stable  
664 HepG2 cells respectively per replicate at an MOI of 0.3-0.5. Five or more biological replicate screens  
665 were performed for each library. After 24 hours, media containing the lentivirus was removed and fresh

666 DMEM media + 2mM VPA was added to allow for construct integration and promote base editing. After  
667 another 48 hours of media + VPA treatment, media with 1:20,000 Puromycin was added and cells  
668 underwent selection for the next 5-7 days and were split as needed.  
669 Following complete selection as defined by complete death of a concurrent untreated control, cells  
670 were started on the 2-day LDL uptake protocol: on day 0, library cells were split, counted, and replated  
671 onto 15-cm plates at  $1 \times 10^5$  cells/cm<sup>2</sup>. On the evening of day 1, DMEM media was removed and  
672 replaced with Optimem to induce overnight serum starvation. On the morning of day 2, 1:400 BODIPY  
673 FL-LDL (Thermo Fisher Scientific) + Optimem was added and incubated with cells for 4-6 hours. After this  
674 incubation period, plates were trypsinized, stained with 50 ng/mL DAPI, and sorted based on BODIPY-  
675 LDL fluorescence levels into 4 bins (top 20%, top 20-40%, bottom 20%, and bottom 20-40%). gDNA was  
676 then collected from each sorted population as well as an unsorted bulk population and prepared for  
677 next generation sequencing.

#### 678 **Library preparation and next generation sequencing**

679 To prepare for next generation sequencing of samples, genomic DNA collected from each sorted  
680 population and an unsorted bulk population was used as the input for a PCR1 reaction aimed at  
681 amplifying the integrated construct spanning the gRNA, as well as adding different inline PE1 barcodes  
682 to specific samples for downstream analysis. The ideal total genomic DNA input per sample for PCR1  
683 was 20ug of DNA, but if less than 20ug of genomic DNA was collected for a specific sorting population,  
684 then the total genomic DNA yield was input into PCR1. NEBNext Ultra II Q5 mastermix was used. Please  
685 refer to **Supplementary Table 9** for specific PCR1 primers used depending on the specific reporter and  
686 corresponding sub-experiment described above in the paper. Following PCR1, reactions were  
687 individually PCR purified using a standard QIAquick PCR Purification Kit. Next, a qPCR2 was performed  
688 from 0.25ul of each purified sample in a 15ul qPCR reaction to determine the optimum number of cycles  
689 for PCR2, and the products of qPCR2 were run on a gel to confirm lack of primer dimer and the

690 confidence of CT cycles from qPCR2. PCR2 cycle counts were chosen to be 2-3 cycles less than the qPCR2  
691 CT for the corresponding sample, with a minimum number of PCR2 cycles being 7. In order to set up  
692 PCR2, half of each sample's purified PCR1 product was then used with NEBNext Ultra II Q5 mastermix  
693 (please refer to the **Supplementary Table 9** for specific PCR2 primers used. It is important to note that,  
694 despite differences in primers between sub-experiments, NEBNext i7 primers were always used in PCR1,  
695 and NEBNext i5 primers were always used in PCR2. The unique combination of these two barcodes for  
696 each specific sample is what allows for downstream identification of reads post sequencing. Following  
697 PCR2, samples were again PCR purified with a QIAquick PCR Purification Kit, and then purified samples  
698 were run on a 2200 Agilent Tapestation to identify the quantity of product as well as any unwanted  
699 byproducts and primer dimers that might have occurred throughout NGS preparation. Samples were  
700 then pooled based on their molarity, and the pool was purified to remove unwanted products using  
701 SPRI-Select beads from Beckman Coulter. The SPRI bead: sample ratio for pool purification was chosen  
702 based on the quantity and size of unwanted byproducts relative to the desired product, and varied  
703 among 0.8-0.9 for all experiments. Following bead purification, the pooled sample was sequenced using  
704 Illumina Nextseq used paired-end sequencing with >50-nt read 1 and >36-nt read 2.

#### 705 **Comparison of endogenous and reporter editing**

706 Endogenous sublibraries A-D targeting four distinct loci across *LDLR* (**Supplementary Table 3**) were  
707 cloned and run as the screens (see Screen Procedure above), but were run in six-well format with  
708 400,000 cells per sample six well being infected. Following the standard lentiviral infection protocol and  
709 puromycin selection process, samples were harvested at the end of selection as only editing data was  
710 desired. After gDNA isolation, a reporter-based library preparation for next generation sequencing was  
711 performed (see Library preparation for next generation sequencing protocol above), together with an  
712 additional library preparation for amplifying endogenous editing sites: PCR1 samples to determine  
713 endogenous editing were set up with Ultra II Q5 mastermix and 2.5ug of genomic DNA input and a

714 unique primer mix for each of the four LDLR endogenous libraries. These primer mixes contained three  
715 unique R1 forward primers and one common R2 reverse primer that would allow amplification segments  
716 from the endogenous LDLR containing all desired editing sites targeted within that library. Upon  
717 completion of PCR1, samples were PCR purified using a QIAquick PCR Purification Kit and prepared  
718 following the standard library preparation protocol above. For specific primers used, please see  
719 **Supplementary Table 9**). Following the preparation of these endogenous samples, both endogenous  
720 and reporter samples were run on an 2200 Agilent Tapestation and pooled + purified accordingly to  
721 prepare for next generation sequencing (see Library preparation and next generation sequencing  
722 section).

723 Endogenous target site reads are mapped to the reference amplicon sequences using CRISPResso2 with  
724 base editing mode and custom mismatch score matrix that tolerates A to G mutation generated by  
725 ``bean-count`` command of *bean* software that implements all the computational functions of the  
726 proposed BEAN workflow as a Python package. Bean is available at [https://pypi.org/project/crispr-](https://pypi.org/project/crispr-bean/)  
727 *bean/*. The paired gRNA and reporter library is mapped to the expected gRNA and reporter sequences  
728 using ``bean-count`` function of the *bean* package. Position-wise reporter base edits are tested for  
729 significance against control data without editing using the function  
730 ``bean.annotate.filter_alleles.filter_alleles`` from the *bean* package. This function conducts Fisher's exact  
731 test and was used to filter for edits with Bonferroni-corrected P-value < 0.05 and odds ratio > 5.

### 732 **Read mapping**

733 FASTQ files are first demultiplexed by matching the corresponding 8nt pair-end index sequences (12  
734 sequences "NNNNNAG" are partially degenerative). Then, the FASTQ files are further demultiplexed by  
735 exact matching of the 3-6nt barcodes and the U6 hairpin stub ("GGAAAGGACGAAACACCG"). Paired end  
736 reads in the demultiplexed FASTQ files are mapped to gRNA sequences and read for the reporter editing  
737 outcome using the ``bean-map-samples`` command from *bean*. There, read pairs where either R1 or R2

738 with average quality Phred score lower than score 30 are discarded. Read pairs with good quality are  
739 mapped to the spacer sequence where the spacer positions in R1 perfectly match with any of the  
740 provided gRNA sequences while being masked for the editable positions to account for self-editing. As  
741 26%-54% of the reads are shown to have undergone R1-R2 recombination (**Supplementary Fig. 25**), we  
742 assign the read to two categories, where both R1 spacer sequence and R2 barcode sequence uniquely  
743 matches to a gRNA in the library and where on the R1 spacer sequence has the unique match to the  
744 gRNA sequence. For the reads that have the correct R1-R2 matching (barcode matched reads), we count  
745 the allele-level editing outcome by comparing the reference reporter sequence to the R2 reporter  
746 sequence by globally aligning two sequences using the modified `CRISPResso2Align.global_align``  
747 function of CRISPResso2<sup>3</sup> for base editing. The matched gRNA count, all gRNA count regardless of R1-R2  
748 matching, and per-guide reporter allele counts for the matched gRNAs are stored as the output per each  
749 sequencing sample.

#### 750 **Quality control of reporter screen data**

751 To exclude failing samples and outlier guides, `bean-qc`` from *bean* is run on the mapped gRNAs and  
752 edited allele counts. Specifically, samples with median Spearman correlation of gRNA count smaller than  
753 0.8, median log fold change of positive control guides (gRNAs targeting splicing sites for both LDL-C  
754 GWAS and LDLR tiling library) in top 20% and bottom 20% quantile smaller than -0.1, or median gRNA  
755 editing rate in reporter smaller than 0.1 are labeled as low-quality. The gRNA editing rate is calculated as  
756 the target variant cognate (A to G in ABE) editing rate in LDL-C GWAS library and mean cognate (A to G  
757 in ABE) editing rate in editable base (A for ABE) protospacer position 3-8. Outlier gRNAs and replicate  
758 pairs are defined as the gRNAs with a median absolute deviation larger than 5 and RPM (reads per  
759 million)  $\geq 10000$  among replicates.

760

761 **Profiling and visualization of base editing preference**

762 LDLR tiling library editing outcome in bulk samples were analyzed to profile base editing preferences,  
763 and avoid bias in sequence context that comes from having A in designated position in LDL variant  
764 library. The heatmap of editing preference across spacer position and PAM was generated by the  
765 function ``bean.pl.editing_patterns.plot_by_pos_pam`` from *bean*. This considered the mean editing  
766 efficiency of A to G transitions within protospacer positions 1-20 for 7320 targeting gRNAs with more  
767 than 10 reads in any bulk sample. The average editing range per dinucleotide in the PAM was calculated  
768 as the maximal editing rate in protospacer positions 3-8, where the the relative editing rate is the  
769 highest. To quantify the context preference, the mean A to G editing efficiency of the -1 and +1 bases of  
770 the intended target base in protospacer position 3-8 was calculated. The context preference logo was  
771 generated from the normalized mean editing efficiency across replicates by the nucleotide in -1 and +1  
772 positions with the Logomaker package.

773 **Prediction of editing outcomes with BE-Hive**

774 We used the Python implementation of BE-Hive<sup>4</sup> ([https://github.com/maxwshen/be\\_predict\\_bystander](https://github.com/maxwshen/be_predict_bystander))  
775 to predict editing outcomes of the *LDLR* tiling library. We initialized the model with “mES” as cell type  
776 and “ABE8” as the editor due to the lack of HepG2 cell type model. For each spacer, we extracted a 50nt  
777 long sequence around the starting position of the spacer in the hg38 genome, with 20nt before the  
778 spacer start and 30nt after, on the same strand. These 50nt sequences were used as the input to the BE-  
779 Hive to predict likely editing outcomes. To calculate allele-level edit rates for each spacer, we summed  
780 up the probability of any editing outcomes with the same editing patterns in position 0-18 (0-based)  
781 relative to the start of the spacer. Similarly, to calculate base-level edit rate, we summed up the  
782 probability of any editing outcomes with identical base edits in position 3-8, relative to the start of the  
783 spacer.

784

785 **BEAN models**

786 *Variant and gRNA-level phenotypic modeling with reporter*

787 BEAN models the phenotype of cells with variant  $v$ ,  $Y_v$ , as Normal distribution, where the wild-type cells  
788 have standard normal phenotypic distribution  $Y_0$  and the variant effects are quantified in a relative  
789 scale, using  $Y_0$  as reference. For the cells with a gRNA, their phenotype is modeled as a mixture of allelic  
790 distributions produced by the gRNA, reflecting the heterogeneous outcome from a gRNA.  
791 For variant screens (LDL-C GWAS library), we aggregate alleles into two categories: alleles with or  
792 without observed edits at the target variant. The non-edited component in these models is fixed to have  
793 a wild-type phenotypic distribution. That is, the phenotype  $Y_g$  of cells with gRNA  $g$  that induces variant  
794  $v$  with editing rate  $\pi$  is modeled as follows:

795 
$$f_{Y_g}(y) = (1 - \pi_g)f_{Y_0}(y) + \pi_g f_{Y_v}(y)$$

796 
$$Y_0 \sim \mathcal{N}(0, 1)$$

797 
$$\mu_v \sim \text{Laplace}(0, 1)$$

798 
$$\sigma_v \sim \text{LogNormal}(0, 0.01)$$

799 
$$Y_v \sim \mathcal{N}(\mu_v, \sigma_v)$$

800 
$$\mu_v \sim \text{Laplace}(0, 1)$$

801 
$$\sigma_v \sim \text{LogNormal}(0, 0.01)$$

802 , where  $f_Y$  indicates the probability density function of  $Y$ . The prior for  $\mu_v$  and  $\sigma_v$  are set to be narrow  
803 based on the assumption that most variant would have close to wild-type effect size of mean 0 and  
804 standard deviation 1.

805 For saturation tiling screen, as bystander edits are more likely to have phenotypic effect, BEAN accounts  
806 for more than one non-wild-type allele where each allele may include one or more variants. Here, we  
807 use the term “allele” to refer to the multiple editing outcome produced by base editing, and we  
808 aggregate multiple nucleotide-level variants that lead to the same coding sequence amino acid



809 mutations together. To account for splicing and noncoding region variants, variants that fall outside  
810 coding regions are not aggregated.

811 We denote with  $A(g) = \{a | \text{Allele } a \text{ is produced by } g\}$  the set of alleles produced by gRNA  $g$  that is  
812 robustly observed, here in at least 10% of the gRNA read counts across 30% of the samples after above-  
813 mentioned aggregation. However, we note that users have the flexibility to set their own robustness  
814 thresholds in `'bean-filter'` of `bean` package. The phenotype of a given allele  $a$  is defined as the sum of  
815 phenotypic effect of non-wild-type nucleotide and amino acid level variants. Finally, the phenotype of  
816 cells with gRNA  $g$  is modeled again as the mixture distribution of allelic phenotype for the alleles it  
817 induces ( $a \in A(g)$ ) as follows:

$$818 \quad f_{Y_g}(y) = \sum_{a \in A(g)} \tilde{\pi}_a f_{Y_a}(y)$$

$$819 \quad Y_a = \sum_{v \in a} Y_v, \quad Y_a = Y_0 \text{ if } |a| = 0$$

820 , where  $\tilde{\pi}_a$  is the endogenous editing rate, estimated from  $\pi_a$ , the reporter editing rate, of allele  $a$ . The  
821 non-edited allele phenotype and the priors for  $\mu_v$  and  $\sigma_v$  are identical to the variant screen modeling.  
822 The identity of the alleles and their frequency in reporters are learned from per-gRNA reporter allele  
823 counts in pre-sort (bulk) sample. Within this modeling framework, the allelic editing frequency in the  
824 reporters is proportionately adjusted based on chromatin accessibility of the intended gRNA target locus  
825 to better estimate the endogenous allele frequency, while allowing for deviation from the scaled values.  
826 For each gRNA, allele editing rate  $\boldsymbol{\pi}_g = (\pi_{g0}, \dots, \pi_{g|A(g)|})$  and per-gRNA allele count  $\mathbf{Z}_g =$   
827  $(A_{g0}, \dots, A_{g|A(g)|})$  are modeled as the Dirichlet and Multinomial distributions :

$$828 \quad \tilde{\boldsymbol{\alpha}}_{g\pi} = \frac{\boldsymbol{\alpha}_{g\pi} + \epsilon}{\sum \boldsymbol{\alpha}_{g\pi} + \epsilon} \boldsymbol{\alpha}_{g\pi}^\circ$$

$$829 \quad \boldsymbol{\pi}_g \sim \text{Dir}(\tilde{\boldsymbol{\alpha}}_{g\pi})$$

$$830 \quad \mathbf{Z}_g \sim \text{Multinomial}(\boldsymbol{\pi}_g)$$

831 Where  $\alpha_\pi$  is initialized as  $\vec{1}$ ,  $\epsilon = 1e^{-5}$  and  $\alpha_{g\pi}^o$  is the precision parameter that is fitted from the data  
832 (**Supplementary Note 3**). This approach partially follows DESeq2<sup>5,6</sup> approach of dispersion parameter  
833 estimation for the Negative Binomial distribution. The reporter editing rate  $\pi_g$  is further scaled by  
834 accessibility to be used as the endogenous editing rate  $\widetilde{\pi}_g$  through a function  $f$ . This function  $f$  is  
835 learned a priori from the paired reporter and endogenous editing rate data while the deviation of  $\widetilde{\pi}_g$   
836 from  $f(\pi_g)$  is fitted per gRNA. The deviation  $\epsilon_{g\pi}$  below accounts for the incomplete correlation  
837 between endogenous and reporter editing rates.

$$\begin{aligned} 838 \quad \widetilde{\pi}_{gj} &= \frac{f(\pi_{gj})}{\sum_{j \in \{1, \dots, |A(g)|\}} f(\pi_{gj})} + \epsilon_\pi, & f(\pi) &= \pi e^{bw^a} \\ 839 \quad \epsilon_{g\pi} &= \text{logit}^{-1}(l_{g\pi}), & l_{g\pi} &\sim \mathcal{N}(0, \sigma_\pi) \\ 840 \quad \widetilde{\pi}_g &= \left( 1 - \sum_{j \in \{1, \dots, |A(g)|\}} \widetilde{\pi}_{gj}, \quad \widetilde{\pi}_{g1}, \quad \dots, \quad \widetilde{\pi}_{gn} \right) \end{aligned}$$

841  $f(\pi)$  is fitted from the data generated for comparison of endogenous and reporter editing based on the  
842 regression  $E \left[ \log \left( \frac{\pi_{endo}}{\pi_{reporter}} \right) \right] = aw + b$  where  $w$  is  $\log(\text{accessibility signal} + 1)$  and the resulting  
843 coefficients  $a = 0.2513$  and  $b = -1.9458$  are used for the analyses presented in this paper. The  
844 residual of the regression is fitted as the Normal distribution, which is used as the prior for the logit-  
845 scale deviation  $l_\pi$  (see full detail in **Supplementary Note 2**).

846 In modeling base editing data screens with reporter data, we have built and evaluated two variants of  
847 the BEAN model that utilize less information than the original model. The first variant, BEAN-Uniform  
848 assumes a single component Normal distribution of cellular phenotype, reflecting the assumption that  
849 all gRNAs would have the same editing efficiency.

$$850 \quad \text{BEAN - Uniform: } Y_g \sim \mathcal{N}(\mu_v, \sigma_v), g \text{ induces } v$$

851 While BEAN utilizes both reporter data and accessibility to estimate endogenous editing efficiency from  
 852 the reporter data, the second variant, BEAN-Reporter focuses only on the incorporation of the reporter  
 853 data without accessibility information. That is, for BEAN-Reporter,  $\tilde{\pi} = \pi$ .

#### 854 *Sorting screen and gRNA count data modeling*

855 Sorting screens sorts the pool of cells with different gRNA and editing outcomes into distinct bins based  
 856 on the phenotype they're sorted on prior to sequencing. To model the sorting procedure, the proportion  
 857 of cells that falls within sorting quantile bins for each gRNA is calculated analytically. This process allows  
 858 for the determination of the relative fraction of cells with the gRNA that falls into each sorting bin, which  
 859 is then used as the concentration parameter of Dirichlet-Multinomial distribution. Dirichlet-Multinomial  
 860 distribution is chosen to model the gRNA read count across sorting bins that is over-dispersed  
 861 multinomial count distribution, which we confirm from our data (see **Supplementary Note 3**). The gRNA  
 862 read counts across sorting bins  $\mathbf{X}_{gr} = (X_{gr}^{(0.0,0.2)}, X_{gr}^{(0.2,0.4)}, X_{gr}^{(0.6,0.8)}, X_{gr}^{(0.8,1.0)})$  and the barcode-matched  
 863 gRNA read count  $\mathbf{X}_{gr}^b$  for gRNA  $g$  and replicate  $r$  are modeled as following:

$$864 \quad p_g^{(q_l, q_h)} = P(q_l \leq Y_g \leq q_h) = \sum_a \Phi\left(\frac{q_h - \mu_a}{\sigma_a}\right) - \Phi\left(\frac{q_l - \mu_a}{\sigma_a}\right)$$

$$865 \quad \mathbf{p}_g = (p_g^{(0.0,0.2)}, p_g^{(0.2,0.4)}, p_g^{(0.6,0.8)}, p_g^{(0.8,1.0)})$$

$$866 \quad \tilde{\mathbf{p}}_g = \frac{\mathbf{p}_g}{\sum \mathbf{p}_g} p_g^\circ$$

$$867 \quad \mathbf{X}_{gr} \sim \text{DirichletMultinomial}(\tilde{\mathbf{p}}_g \odot \mathbf{s}_r)$$

$$868 \quad \mathbf{X}_{gr}^b \sim \text{DirichletMultinomial}(\tilde{\mathbf{p}}_g^b \odot \mathbf{s}_r^b)$$

869 ,where  $\odot$  denotes element-wise multiplication. Here,  $\mathbf{p}_g$  is scaled as  $\alpha_\pi$  by the data-fitted precision  
 870 parameter  $p_g^\circ$  (**Supplementary Note 3**) then scaled by the sample-specific size factor  $\mathbf{s}_r =$

871  $(s_r^{(0.0,0.2)}, s_r^{(0.2,0.4)}, s_r^{(0.6,0.8)}, s_r^{(0.8,1.0)})$  , where the sample size factor is calculated as in DESeq2<sup>5,6</sup>. For

872 sample  $j$ ,  $s_j = \text{median}_g \frac{X_{gj}}{(\prod_{v=1}^m X_{gv})^{1/m}}$  and the same function is used to calculate size factor for barcode-

873 matched read counts for sample  $j$ ,  $s_j^b$  with  $X_{gj}^b$ . We note that  $X_{gr}^b$  is not used in the inference when  
874 benchmarking against other methods to make sure we provide the same input. Samples marked as low-  
875 quality and gRNA and replicate pair with  $\leq 10$  total reads, or identified as outliers are excluded from  
876 inference (see **Quality control of reporter screen data**).

877 The parameters  $\mu_g, \sigma_g, \alpha_\pi, l_\pi$  of posterior distributions are fitted using stochastic variational inference  
878 (SVI) of Pyro using 2000 steps with decaying learning rate starting from 0.01 with variational distribution  
879 that mirrors the model. Specifically, the posterior phenotypic distribution of each variant is fitted as a  
880 Normal distribution with a posterior standard deviation parameter and mean parameter which has  
881 Normal posterior distribution:

$$\begin{aligned} 882 \quad P(Y_v | \mathcal{D}) &\approx Q(\hat{Y}) \\ 883 \quad \hat{Y} &\sim \mathcal{N}(\widehat{\mu}_v, \widehat{\sigma}_v) \\ 884 \quad \widehat{\mu}_v &\sim \mathcal{N}(\widehat{\mu}_{\mu_v}, \widehat{\sigma}_{\mu_v}) \end{aligned}$$

885 Where  $\mathcal{D}$  is observed data for the model and  $Q$  is the variational distribution. Negative control variants  
886 are used to control the significance of variant effect, by fitting the shared phenotypic distribution of  
887 negative controls as a single normal distribution. Subsequently the results are scaled so that the fitted  
888 negative control distribution is transformed to a standard normal.

$$\begin{aligned} 889 \quad \widehat{Y}_{ctrl} &\sim \mathcal{N}(\widehat{\mu}_{ctrl}, \widehat{\sigma}_{ctrl}) \\ 890 \quad Y_v^{(scaled)} &= \frac{\widehat{Y}_v - \widehat{\mu}_{ctrl}}{\widehat{\sigma}_{ctrl} | \mathcal{D}} \sim \mathcal{N}(\mu_v^{(scaled)}, \sigma_v^{(scaled)}) \\ 891 \quad \mu_v^{(scaled)} &= \frac{\widehat{\mu}_{\mu_v} - \widehat{\mu}_{ctrl}}{\widehat{\sigma}_{ctrl}} \sim \mathcal{N}(\mu_{\mu_v}^{(scaled)}, \sigma_{\mu_v}^{(scaled)}) \end{aligned}$$

892 In order to control for false discovery with negative control variants, the standard deviations of variants  
893  $\sigma_{\mu_v}^{(scaled)}$  are scaled so that the standard deviation of  $\mu_n$ , where  $n$  are the negative control variants, is  
894 equal to 1. In the LDL-C variant screen, 20 negative control variants (each tiled by 5 gRNAs) are used and  
895 for LDLR tiling screen 175 synonymous variants are used as the negative control variants:

896

$$\sigma_{\mu_v}^{(adj)} = \sigma_{\mu_v}^{(scaled)} * \widehat{\sigma}_n$$

897

Where  $\widehat{\sigma}_n$  is fitted as the standard deviation estimate of  $z_{\mu_n}^{(scaled)} = \mu_{\mu_n}^{(scaled)} / \sigma_{\mu_n}^{(scaled)}$  with

898

``stats.norm.fit`` of Python's SciPy package with setting location parameter to 0.

899

The model's output includes various parameters relating to the phenotype of the variant, such as the

900

mean and standard deviation of variant phenotype  $\mu_v^{(scaled)}$ ,  $\sigma_v^{(scaled)}$  and scaled and significance-

901

adjusted phenotypic mean distribution parameters  $\mu_{\mu_v}^{(scaled)}$ ,  $\sigma_{\mu_v}^{(scaled)}$ ,  $\sigma_v^{(adj)}$ ,  $z_{\mu_v}^{(adj)}$  where  $z_{\mu_v}^{(adj)} =$

902

$\mu_{\mu_v}^{(scaled)} / \sigma_{\mu_v}^{(adj)}$  are reported together with estimated endogenous editing efficiency for each variant.

903

For ``variant`` mode, the mean targeting gRNA editing rate is reported and for ``tiling`` mode, effective

904

editing efficiency is reported and calculated as  $\sum_{g \in \{g | g \text{ induces } v\}} \sum_{a \in \{g \text{ induces } a, a \text{ has } v\}} \frac{\pi_{ga}}{|a|}$ . The model,

905

variational distribution and inference procedure are available as the default options of ``bean-run``

906

command of *bean* software. Specifically, BEAN-Uniform is run with ``--uniform-edit`` and full BEAN model

907

is run by specifying ``--scale-by-acc`` argument.

908

### **Benchmarking of CRISPR Pooled Screen Analysis Methods**

909

We reviewed and selected several CRISPR pooled screen analysis methods for benchmarking against

910

BEAN, based on their availability and applicability to our experimental design and sorting screens.

911

BAGEL<sup>7</sup> was not applicable as it required positive and negative control target genes as the input. ACE<sup>8</sup>

912

was designed for gene essentiality screens. Gscreend<sup>9</sup> required a single unsorted population to be

913

compared against the multiple treatment samples. Consequently, we chose MAGeCK-RRA<sup>10</sup>, three

914

running modes of MAGeCK-MLE<sup>11</sup>, CRISPRBetaBinomial (CB2)<sup>12</sup>, and CRISPhieRmix<sup>13</sup>. CRISPhieRmix is

915

only used for benchmarking the *LDLR* tiling screen as it requires the negative control gRNAs as LDL-C

916

GWAS library benchmarking uses negative control variant label to evaluate classification

917

performance. Huang et al., (2021)<sup>14</sup> didn't offer their method available as software, but we incorporated

918

their efficiency correction concept as efficiency-corrected log fold change (EC-LFC). Log fold change of

919 variants calculated from MAGeCK-RRA was included as the baseline. We believe that these methods  
920 represent the state-of-the-art based on multiple recent benchmark studies.

921 We used two modes of MAGeCK v0.5.9.4, MAGeCK-RRA<sup>10</sup> and MAGeCK-MLE<sup>11</sup>. MAGeCK-RRA takes  
922 treatment and control samples and evaluates if the rank of log fold change of gRNA abundance is not  
923 uniformly distributed. We used the paired mode with bottom 20% and top 20% quantile bin samples of  
924 each replicate as the paired treatment and control.

925 MAGeCK-MLE<sup>11</sup> uses Negative Binomial generalized linear model with log link to output the coefficients,  
926 that can be interpreted as the log fold change of the gRNA abundance following the unit increase in the  
927 covariate that is provided in the input design matrix. MAGeCK-MLE is the only method that we  
928 benchmarked against that can use all 4 quantile bins of our sorting screens. We assigned 0, 1, 3, 4 to 0-  
929 20%, 20-40%, 60-80%, 80-100% quantile bin samples as the input covariate values. We further  
930 benchmarked MAGeCK-MLE where it uses the gRNA activity (`--guide_efficiency_file`) or fits the gRNA  
931 activity (`--guide_efficiency_file --update-efficiency`). As MAGeCK-MLE assumes gRNA efficiency in  
932 `guide_efficiency_file` to scale from -1 to 0.25, the editing efficiency is normalized to the range. In case  
933 the gRNA has not enough reads and is not assigned of the editing rate, it is assigned to the editing rate  
934 of 0.5 before the scaling. All runs are ran both with (`--genes-varmodeling 1000`) or without (default)  
935 dispersion fitting.

936 CB2<sup>12</sup> models the gRNA counts using the beta-binomial distribution in which the variance can be either  
937 large or smaller than the mean to quantify gRNA abundance for CRISPR pooled screen data analysis. It  
938 uses Fisher's combined probability test to estimate the gene-level significance. We installed the CB2  
939 package (v1.3.4) and benchmarked the performance by comparing the bottom 20% and top 20%  
940 quantile bins.

941 CRISPhieRmix<sup>13</sup> is a hierarchical mixture model for analyzing CRISPR pooled screen data by assuming  
942 that the majority of genes does not impact phenotype. It builds a two-group mixture model to identify

943 the impactful target genes. Specifically, log 2 fold change between bottom and top 20% quantile bins is  
944 calculated from DESeq2<sup>6</sup> and used as the input to compare the bottom and top samples.  
945 EC-LFC was calculated by dividing the variant log fold change calculated by MAGeCK-RRA<sup>10</sup> by the variant  
946 editing efficiency. For the LDL-C GWAS library, mean editing efficiency of the targeting gRNAs are used  
947 as the variant editing efficiency. For *LDLR* tiling library, effective editing efficiency was used as the  
948 variant editing efficiency.

#### 949 **Benchmark on LDL-C GWAS and *LDLR* tiling library**

950 For both LDL-C GWAS library and *LDLR* tiling library, the classification performance AUPRC of  
951 distinguishing positive controls against negative controls are evaluated. For the benchmark, 6 biological  
952 replicates of LDL-C GWAS library and 4 *LDLR* tiling library with no failing samples are used, and barcode-  
953 matched reads are ignored during inference in BEAN runs. For replicate subsample analysis, all possible  
954 2-replicate combinations are subset to be analyzed by each method. For LDL-C GWAS library, its positive  
955 control variants, which are the splice sites of the genes that changes the LDL-C uptake is used as the  
956 positive control variants and the 20 non-targeting negative control variants are used as the negative  
957 control variants. For *LDLR* tiling library, ClinVar “pathogenic” or “pathogenic/likely pathogenic”  
958 annotated variants are classified against ClinVar “benign” or “benign/likely benign” annotated variants.  
959 As each method has different strategy to assign gRNA to variant thus scores different set of variants, we  
960 evaluate the recall as how much of the all ABE-discoverable Pathogenic/Likely Pathogenic variants are  
961 identified as Pathogenic.

#### 962 **Cloning and testing of individual gRNAs**

##### 963 **Base Edit**

964 Oligonucleotides including protospacer sequences were ordered in the following format:  
965 GGAAAGGACGAAACACCG [19-20-bp protospacer —remove initial G for any 20-bp protospacer with one  
966 natively] GTTTAAGAGCTATGCTGGAAAC (see Supplementary Table 9). Using NEBuilder HiFi DNA

967 assembly, ABE8e-Cas9NG designated oligonucleotides were cloned into CRISPRv2FE-ABE8e-Cas9NG,  
968 while ABE8e-SPRY designated oligonucleotides were cloned into CRISPRv2FE-ABE8e-SpRY-BsrGI. To  
969 make base edited cell lines, the gRNA constructs were packaged into lentivirus and transduced into  
970 HepG2 ABE8e-SPRY-BFP cells seeded at  $4 \times 10^4$  cells/cm<sup>2</sup> on 6-well plates in two replicates with 8 µg/ml  
971 polybrene. Two days post-transduction, cells were treated with 500 ng/ul puromycin and selected for  
972 approximately one week. HepG2 Base Edited cells were seeded 1:1 with HepG2-mcherry cells to achieve  
973 a total density of  $1.08 \times 10^5$  cells/cm<sup>2</sup> on a 96 well plate in at least two technical replicates of two  
974 biological replicates and incubated overnight. The next day, the media was replaced with optiMEM and  
975 cells were incubated overnight. Approximately 4-6 hours prior to flow cytometric analysis, cells were  
976 treated with 2.5 mg/mL BODIPY-LDL in optiMEM. Cells were trypsinized and analyzed for presence of  
977 mCherry and LDL uptake using a Beckman CytoFLEX flow cytometer. LDL uptake of each base edited cell  
978 line was normalized to the LDL uptake of the mCherry cells within the same well. Differential LDL uptake  
979 between base edited and control cells was further normalized using data from the ABE8e and SPRY  
980 sgCTRL lines.

#### 981 **CRISPRi**

982 Oligonucleotides including protospacer sequences (**Supplementary Table 9**) were ordered in the  
983 following format: GGAAAGGACGAAACACCG [19-20-bp protospacer—remove initial G for any 20-bp  
984 protospacer with one natively] GTTTAAGAGCTATGCTGGAAAC were cloned into a pHR-U6-gRNAFE-Zim3-  
985 dCas9-P2A-Hygro backbone through NEBuilder HiFi DNA assembly. To make CRISPRi cell lines, the gRNA  
986 constructs were packaged into lentivirus and transduced into HepG2 cells seeded at  $4 \times 10^4$  cells/cm<sup>2</sup> on  
987 48-well plates in two replicates with 8 µg/mL of polybrene. Two days post-transduction, cells were  
988 treated with 125 µg/ml Hygromycin B and were selected for approximately one week. LDL uptake  
989 experiments were performed as described above, seeding CRISPRi cell lines 1:1 with HepG2-tTA-BFP  
990 cells as the internal control.



991 **CRISPRa**

992 Oligonucleotides including protospacer sequences (**Supplementary Table 9**) were ordered in the  
993 following format: GGAAAGGACGAAACACCG [19-20-bp protospacer —remove initial G for any 20-bp  
994 protospacer with one natively] GTTTAAGAGCTAGGCCAACATG. Using NEBuilder HiFi DNA assembly,  
995 oligonucleotides were cloned into a pLenti U6-2xMS2gRNA MCPp65 PuroR backbone. To make CRISPRa  
996 cell lines, the gRNA constructs were packaged into lentivirus and transduced into HepG2 dCas9-  
997 10xGcn4-mChe + scFv-Sbno1-Nfe2l1-Krt40-BFP cells and seeded at  $4 \times 10^4$  cells/cm<sup>2</sup> on 6-well plates in  
998 two replicates with 8 µg/ml polybrene. Two days post-transduction, cells were treated with 500 ng/ml  
999 puromycin and selected for approximately one week. LDL uptake experiments were performed as  
1000 described above, seeding CRISPRa cell lines 1:1 with HepG2 wt cells as the internal control.

1001 **Pooled ATAC-seq**

1002 A pool of 20 gRNAs was cloned into CRISPRv2FE-ABE8e-Cas9NG or CRISPRv2FE-ABE8e-SpRY-BsrGI (see  
1003 Cloning and testing of individual gRNAs), packaged into lentivirus and transduced into HepG2 ABE8e-  
1004 SpRY-BFP cells seeded at  $4 \times 10^4$  cells/cm<sup>2</sup> on 6-well plates in three replicates with 8 µg/mL of polybrene.  
1005 Cells were treated with VPA and selected with Puromycin as in screens. Once selected, two sets of  $1 \times 10^6$   
1006 cells for each of the three replicates as well as an unedited control replicate were seeded to 6-well  
1007 plates. The next day, one well per replicate was fed with DMEM + FBS and the other with optiMEM  
1008 (serum-starved). 24 hours later, the wells were trypsinized, and  $1 \times 10^5$  cells were used for ATAC-seq  
1009 while the remaining cells were used for bulk genomic DNA isolation using the Purelink Genomic DNA  
1010 mini kit (Life Technologies). ATAC-seq was performed using the Active Motif ATAC-Seq kit according to  
1011 manufacturer's instructions.

1012 To obtain valid primers to amplify the loci surrounding the 20 target variants, Primer3<sup>15</sup> was used to  
1013 generate 5 candidate primer sets within +-150-nt from each variant. Primer-Dimer.com was used to  
1014 calculate a deltaG (dG) interaction matrix for all candidate primers. Primers with average dG  $\leq -7$  were

1015 removed. Then, recursive pairwise filtering was performed to iteratively remove the primer with the  
1016 worst dG interaction until no pairwise  $dG \leq -7$  remains. This recursive filtering was performed 300  
1017 times, and the run with the most primers remaining was used. The primer set for each variant with  
1018 highest minimum dG was selected. Primers were all ordered from IDT preceded by NNN to randomize  
1019 initial nucleotides in NGS. We provide the amplicon sequence of 20 loci in **Supplementary Table 9**.  
1020 Genomic DNA and ATAC-Seq products from 8 total samples (3 experimental and 1 control, both in serum  
1021 and starved conditions) were amplified using two primer pools, each composed of 10 primer sets to  
1022 synchronize annealing temperature. 2.5 ug of gDNA/half of ATAC-Seq product was used in 100 uL  
1023 reactions for 32 cycles (gDNA) or 35 cycles (ATAC-Seq). Tapestation was used to pool the two PCR  
1024 products for each sample, and these 16 pools were used as input to the NEBNext UltraII DNA Library  
1025 Prep to prepare NGS libraries. Libraries were sequenced using 150-nt single-end sequencing using  
1026 Illumina Nextseq.

### 1027 **Pooled ATAC-seq analysis**

1028 For each sample, the ATAC-seq reads were mapped to amplicon sequences (**Supplementary Table 9**)  
1029 from 20 loci via Bowtie2<sup>16</sup> (v2.5.1). `bowtie2-build` was used to build indices for the amplicon sequences  
1030 for each of 20 loci and reads are mapped onto the indices with default parameters. In-house Perl script  
1031 was used to parse the SAM output from Bowtie2 and to demultiplex the reads by the locus they mapped  
1032 to with default options. Demultiplexed reads are then profiled for the target base editing rate using  
1033 CRISPResso2<sup>3</sup> (v.2.2.9) using average read quality cutoff of Phred score 30 and assigned of base `N` if  
1034 per-base quality is lower than Phred score 20. For each variant, reads are assigned to the reference  
1035 allele or alternate allele based on the base identity at the target SNP position. In case there exists a  
1036 neighboring variant that allows phasing as HepG2 is heterozygous for the variant, the reads are counted  
1037 per phase based on the identity of the neighboring variant. We note that for two of the variants  
1038 examined (rs3767844 and rs4390169), whether the base is the result of editing or is the reference allele

1039 was ambiguous (that is, variants are heterozygous in HepG2 and two reference alleles of A and G, then  
1040 we cannot assign reads with G in the variant position to the edited reference A or unedited reference G).  
1041 For the variants, we simply compare two observed bases and treat the effect as caQTL. For  
1042 rs771555783, rs76895963, and rs116734477, edited reads are not detected due to insufficient  
1043 representation of the loci, and thus excluded from the enrichment analysis.  
1044 We first identified the variants with significant editing observed in treatment samples compared to the  
1045 control samples where base editors are not treated. This is done by assessing the significance of  
1046 coefficient for *is\_treatment* in the following Binomial regression with `GLM` module of Python  
1047 `statsmodels` package<sup>17</sup>, where  $Edited_j$  and  $Unedited_j$  is the read counts of edited and unedited  
1048 variants in sample  $j$ , and  $is\_treatment_j$  is the indicator variable for the sample  $j$  being treatment  
1049 sample.

$$1050 \quad Edited_j \sim \text{Binomial}(p_j, Edited_j + Unedited_j)$$

$$1051 \quad \text{logit}(p_j) \sim 1 + is\_treatment_j$$

1052 Significantly edited variants should show higher proportion of edited reads ( $p_j$ ) in treatment samples  
1053 compared to the control samples. For all significance testing, Benjamini-Hochberg family-wise error rate  
1054 (FWER) value of 0.1 is used as the threshold, where multiple testing correction is performed with  
1055 `stats.multitest.multipletests` function of Python `statsmodels` package<sup>17</sup>.

1056 For the variants with significant observed editing, we calculated the enrichment of the editing in ATAC-  
1057 seq compared to the gDNA sample, which indicates the editing opened the chromatin at the variant loci  
1058 and increased its capture rate for ATAC-seq. The enrichment of edited allele is calculated as the Binomial  
1059 regression coefficient of edited and unedited read counts for each variant. The proportion of edited  
1060 read count is regressed on whether the sequencing sample  $j$  is from ATAC-seq ( $is\_ATAC_j = 1$ ) or gDNA  
1061 ( $is\_ATAC_j = 0$ ), and the regression coefficient of  $is\_ATAC_j$  is used as the accessibility enrichment of the  
1062 variant editing. We condition for replicate and condition specific effect, along with the interaction effect

1063 between condition and ATAC-seq sample to examine if the variant only alters accessibility under either  
1064 one of two conditions.

$$1065 \quad Edited_j \sim \text{Binomial}(p_j, Edited_j + Unedited_j)$$

$$1066 \quad \text{logit}(p_j) \sim 1 + is\_ATAC_j + condition_j * is\_ATAC_j + condition_j + replicate_j$$

1067 We also calculated the caQTL effect of the variants that are heterozygous in HepG2. Here, whether one  
1068 allele has higher enrichment in ATAC-seq sample is examined as the regression coefficient for  $is\_ATAC_j$   
1069 of following regression, again conditioned on experimental condition and replicate.

$$1070 \quad Allele0_j \sim \text{Binomial}(q_j, Allele0_j + Allele1_j)$$

$$1071 \quad \text{logit}(q_j) \sim 1 + is\_ATAC_j + condition_j * is\_ATAC_j + condition_j + replicate_j$$

1072 Here,  $Allele0_j$  and  $Allele1_j$  are the read counts of alleles 0 and 1. The regression coefficient for  
1073  $is\_ATAC_j$  is used as the accessibility enrichment of allele 0. When the enrichment is shown uniformly in  
1074 major to minor allele, enrichment values and confidence intervals calculated for the opposite direction is  
1075 inverted of their sign.

## 1076 **MotifRaptor**

1077 We adapted the MotifRaptor<sup>18</sup> pipeline to investigate how prioritized genetic variants may influence  
1078 nearby genes involved in LDL-C uptake. For each variant, we retrieved genomic sequences spanning 61  
1079 bp centered around the SNP location, using the hg38 genome assembly as a reference. Each sequence  
1080 was mutated by substituting the major allele with the minor allele at the SNP position, yielding both a  
1081 reference and an alternative sequence for each variant.

1082 Subsequently, to evaluate the potential for transcription factor (TF) binding, we employed all the human  
1083 TF position weight matrices (PWMs) from the CIS-BP database<sup>19</sup> to scan each pair of reference and  
1084 alternative sequences. This motif scanning generated binding scores at each sequence position, serving  
1085 as predictive indicators of TF binding potential.

1086 We then compared these scores for each TF across the reference and alternative alleles within every  
1087 sequence pair. This comparative step is crucial for determining a variant's impact on TF binding.  
1088 Specifically, higher binding scores for the alternative sequence indicate an increase in TF binding  
1089 potential, while lower scores suggest a decrease. To quantify these changes, we calculated a 'disruption  
1090 score' as follows:

$$1091 \text{ disruption} = \text{score}(s_{alt}) - \text{score}(s_{ref}).$$

1092 This score helps capture the directional change each variant induces, where a negative value signifies  
1093 reduced TF binding potential and a positive value indicates an increase.

#### 1094 **Pfam profile HMM scores**

1095 Pfam profile HMM files of PF00057, PF00058 and PF00008 for LDLR class A repeat, LDLR class B repeat,  
1096 and EGF-like domain, respectively, are downloaded from Pfam<sup>20</sup> to generate sequence logo through  
1097 Skyline<sup>21</sup>. Match emission score from the profile HMMs is used to calculate  $\Delta Pfam$  score. Match  
1098 emission score is the negative log probability to observe the amino acid from multiple sequence  
1099 alignment for a given position, thus lower score corresponds to high conservation and lower  
1100  $\Delta Pfam(ref - alt) = -(Pfam_{alt} - Pfam_{ref})$  corresponds to higher reference amino acid  
1101 conservation and lower chance to observe alt amino acid.

#### 1102 **LDLR repeat domain alignment**

1103 LDLR class A repeats is aligned as shown in a previous study to align for all Cysteine residues. Alignments  
1104 for LDLR class B repeats and EGF-like domain were obtained with Clustal Omega<sup>22-24</sup> by aligning domain  
1105 sequences with seed alignments from Pfam PF00058 and PF00008.

#### 1106 **UK Biobank data processing**

##### 1107 *Study participants*

1108 The UK Biobank<sup>25</sup> is a prospective cohort of over 500,000 individuals recruited between 2006 and 2010  
1109 of ages 40-69. Drawing from 469,803 participants with whole exome sequencing (WES) data, we

1110 included 443,353 participants with available LDL cholesterol measurements in this study. Patients with  
1111 homozygous variants and participants with more than one rare variant across *LDLR* and with any rare  
1112 variant in *APOB* and *PCSK9* were not considered for these analyses.

1113 *Variant inclusion and quality control*

1114 Exon coordinates were determined for *LDLR*, *APOB*, and *PCSK9* using MANE transcripts<sup>26</sup>, with an  
1115 additional 5nt retained upstream and downstream of each coding region to capture splice-site variants.  
1116 Exome sequencing was performed for UKB participants as previously described. Analysis was conducted  
1117 on the Research Analysis Platform ([ukbiobank.dnanexus.com](https://ukbiobank.dnanexus.com)). We extracted gene-level VCF files from  
1118 the WES joint-called pVCFs using *bcftools*<sup>27</sup> (v1.15.1) using the Swiss Army Knife app, then normalized to  
1119 flatten multiallelic sites and align variants to the GRCh38 reference genome.

1120 Variants in low complexity regions, segmental duplications, or other regions known to be challenging for  
1121 next generation sequencing alignment or calling were removed from analysis (National Institute of  
1122 Standards and Technology Genome in a Bottle Consortium<sup>28</sup> difficult regions), as were variants with an  
1123 alternate allele frequency greater than 0.1% in the UK Biobank cohort. Further filtering removed  
1124 variants in which more than 10% of samples were missing genotype calls and variants that did not  
1125 appear in the UK Biobank cohort. To mitigate differences in sequencing coverage between individuals  
1126 who were sampled at different phases of the UK Biobank project, variants were only retained in the final  
1127 set if at least 90% of their called genotypes had a read depth of at least 10.

1128 The canonical functional consequence of each variant was calculated using Variant Effect Predictor (VEP,  
1129 v99)<sup>29</sup>. Non-coding variants outside of essential splice sites were not considered in the analysis.

1130 Computational scores are provided by VEP, including the PhastCons conservation score<sup>30</sup>  
1131 (PhastCons100way Vertebrate). When multiple PhastCons conservation scores are available for a coding  
1132 variant, the mean of the available scores was used.

1133

1134 *Clinical endpoints and endophenotypic data*

1135 Coronary artery disease and myocardial infarction cases were aggregated from hospital records (primary  
1136 or secondary diagnosis), death registries (primary or secondary cause of death), and self-reported data.

1137 Age of onset was estimated based on date of onset and birth date when not directly provided, and cases  
1138 with uncertain or unavailable onset data were excluded.

1139 Patient-level LDL-C values were ascertained from the UK Biobank data files. Estimated untreated LDL-C  
1140 levels were obtained using adjustments for lipid-lowering therapies were used in analyses, as described  
1141 in the supplement of this manuscript.

1142 *ClinVar assertions*

1143 ClinVar clinical assessments were identified from the tab delimited version of ClinVar released on  
1144 04/04/2023. In this analysis, we use ‘pathogenic’, ‘likely pathogenic’, and ‘pathogenic/likely pathogenic’  
1145 classifications as ‘P/LP’ collectively, and ‘benign’, ‘likely benign’, and ‘benign/likely benign’ classifications  
1146 as B/LB.

1147 **BEAN-FUSE scores**

1148 We make use of the FUSE (Functional Substitution Estimation) pipeline<sup>31</sup> to improve the estimation of  
1149 variant functional effects and to impute effects of variants which have not been screened. FUSE makes  
1150 use of related measurements within and across experimental assays to jointly estimate variant impacts.  
1151 After functional scores have been estimated by the BEAN pipeline, the full set of scores are processed by  
1152 FUSE, which first collectively estimates the mean functional effect per amino acid residue position  
1153 within the assay, using shrinkage estimation. FUSE then makes estimates for individual allelic variants  
1154 within the amino acid residue position, based on a functional substitution matrix derived from deep  
1155 mutational scanning data across many genes. The result is a full set of estimated variant functional  
1156 effects for both: 1) the original variants screened in the assay, and 2) other possible variants which were  
1157 not screened but fall within amino acid residues which had variants covered in the screen.

1158 **Prediction of UKB LDL-C level**

1159 UKB LDL-C level of variants observed in the base editing and had high confidence ( $\sigma_{\mu} < 0.5$ ) was  
1160 predicted using XGBoost<sup>32</sup> Python package and with default option with 10-fold cross validation  
1161 implemented in scikit-learn<sup>33</sup> `model\_selection.cross\_val\_predict`. The LDL-C levels of UKB variants that  
1162 are unobserved or not observed with enough confidence ( $\sigma_{\mu} > 0.5$ ) was predicted by the XGBoost  
1163 model that is trained on the variants observed with  $\sigma_{\mu} < 0.5$ .

1164 **Structural analysis**

1165 The protein structures were visualized and the screenshots were generated using PyMOL<sup>34</sup> (v2.5.2).  
1166 Relative solvent accessibility (RSA)<sup>35</sup> and residue depth were calculated using the DSSP module in  
1167 BioPython (v1.79)<sup>36</sup> to capture the local 3D accessibility of residues. The wild-type atomic interactions  
1168 between residues were calculated by Arpeggio using the LDLR AlphaFold2 structure (position 1-860)  
1169 from the AlphaFold Protein Structure Database<sup>37</sup>. Additionally, interactions with calcium ions and  
1170 saccharides were calculated using PDB structure 1N7D (position 65-714 after renumbering according to  
1171 Uniprot P01130). These interactions were also computed for mutant structures generated using  
1172 MODELLER<sup>38</sup> 10.3. Subsequently, the change in interactions was determined by subtracting the  
1173 interactions in the mutant from those in the wild-type. Within each of the LDLR class B, LDLR class A, and  
1174 EGF-like domains, two-sided Wilcoxon rank-sum tests were conducted to compare the features  
1175 calculated for deleterious variants (identified by BEAN z-scores below -1.96) against those of other  
1176 variants. DDMut<sup>39</sup> is a deep learning model that predicts protein stability change induced by mutation,  
1177  $\Delta\Delta G$ , based on the local atomic environment and interactions in wild type and mutated residue. LDLR  
1178 AlphaFold2 structure is used as the input to predict the  $\Delta\Delta G$  of variants observed in our LDLR tiling  
1179 screen.  
1180 Molecular interactions were visually represented using a color-coded scheme to differentiate between  
1181 interaction types. Hydrophobic interactions were depicted in 'Forest'; polar interactions were depicted



1182 in 'Orange'; Carbonyl interactions were depicted in 'Blue'; hydrogen bonds were depicted in 'Red';  
1183 Aromatic ring interactions including Methionine Sulfur- $\pi$ , Donor- $\pi$ , Cation-  $\pi$ , and Amide-Ring  
1184 interactions, were depicted in 'Pale Green'; Undefined interactions were depicted in 'Cyan'; Coordinate  
1185 covalent bonds were depicted in 'Purple'; and Ionic interactions were depicted in 'Yellow'. Moreover,  
1186 line type specifies distance flag from Arpeggio output, where thin dashed dashed lines represent Van  
1187 der Waals (vdW) Clashes, where the vdW radii between two atoms cause steric clashes. When such  
1188 clashes co-occurred with other interactions mentioned earlier, they were portrayed with dashed lines,  
1189 using the color code corresponding to the additional interaction except for undefined interaction type.  
1190 Additionally, all ionic interactions and coordinate covalent bonds with ions were consistently  
1191 represented by yellow and purple dashed thick lines. Other interactions were all represented by solid  
1192 lines.

## Methods References

1. Hamilton, M. C. *et al.* Systematic elucidation of genetic mechanisms underlying cholesterol uptake. *Cell Genomics* **3**, 100304 (2023).
2. Emmer, B. T. *et al.* Genome-scale CRISPR screening for modifiers of cellular LDL uptake. *PLoS Genet.* **17**, e1009285 (2021).
3. Clement, K. *et al.* CRISPResso2 provides accurate and rapid genome editing sequence analysis. *Nat. Biotechnol.* **37**, 224–226 (2019).
4. Arbab, M. *et al.* Determinants of Base Editing Outcomes from Target Library Analysis and Machine Learning. *Cell* **182**, 463-480.e30 (07/2020).
5. Anders, S. & Huber, W. Differential expression analysis for sequence count data. *Genome Biol.* **11**, R106 (2010).
6. 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).
7. Hart, T. & Moffat, J. BAGEL: a computational framework for identifying essential genes from pooled library screens. *BMC Bioinformatics* **17**, 164 (12/2016).
8. Hutton, E. R., Vakoc, C. R. & Siepel, A. ACE: a probabilistic model for characterizing gene-level essentiality in CRISPR screens. *Genome Biol.* **22**, 278 (2021).
9. Imkeller, K., Ambrosi, G., Boutros, M. & Huber, W. gscreeend: modelling asymmetric count ratios in CRISPR screens to decrease experiment size and improve phenotype detection. *Genome Biol.* **21**, 53 (12/2020).
10. Li, W. *et al.* MAGeCK enables robust identification of essential genes from genome-scale CRISPR/Cas9 knockout screens. *Genome Biol.* **15**, 554 (2014).
11. Li, W. *et al.* Quality control, modeling, and visualization of CRISPR screens with MAGeCK-VISPR. *Genome Biol.* **16**, 281 (2015).

12. Jeong, H.-H., Kim, S. Y., Rousseaux, M. W. C., Zoghbi, H. Y. & Liu, Z. Beta-binomial modeling of CRISPR pooled screen data identifies target genes with greater sensitivity and fewer false negatives. *Genome Res.* **29**, 999–1008 (2019).
13. Daley, T. P. *et al.* CRISPhieRmix: a hierarchical mixture model for CRISPR pooled screens. *Genome Biol.* **19**, 159 (2018).
14. Huang, C., Li, G., Wu, J., Liang, J. & Wang, X. Identification of pathogenic variants in cancer genes using base editing screens with editing efficiency correction. *Genome Biol.* **22**, 80 (2021).
15. Untergasser, A. *et al.* Primer3--new capabilities and interfaces. *Nucleic Acids Res.* **40**, e115 (2012).
16. Langmead, B. & Salzberg, S. L. Fast gapped-read alignment with Bowtie 2. *Nat. Methods* **9**, 357–359 (2012).
17. Seabold, S. & Perktold, J. Statsmodels: Econometric and statistical modeling with python. in *Proceedings of the 9th Python in Science Conference (SciPy, 2010)*. doi:10.25080/majora-92bf1922-011.
18. Yao, Q. *et al.* Motif-Raptor: a cell type-specific and transcription factor centric approach for post-GWAS prioritization of causal regulators. *Bioinformatics* **37**, 2103–2111 (2021).
19. Weirauch, M. T. *et al.* Determination and inference of eukaryotic transcription factor sequence specificity. *Cell* **158**, 1431–1443 (2014).
20. Mistry, J. *et al.* Pfam: The protein families database in 2021. *Nucleic Acids Res.* **49**, D412–D419 (2021).
21. Wheeler, T. J., Clements, J. & Finn, R. D. Skylign: a tool for creating informative, interactive logos representing sequence alignments and profile hidden Markov models. *BMC Bioinformatics* **15**, 7 (2014).
22. McWilliam, H. *et al.* Analysis tool Web Services from the EMBL-EBI. *Nucleic Acids Res.* **41**, W597-600 (2013).

23. Sievers, F. *et al.* Fast, scalable generation of high-quality protein multiple sequence alignments using Clustal Omega. *Mol. Syst. Biol.* **7**, 539 (2011).
24. Goujon, M. *et al.* A new bioinformatics analysis tools framework at EMBL-EBI. *Nucleic Acids Res.* **38**, W695-9 (2010).
25. Sudlow, C. *et al.* UK biobank: an open access resource for identifying the causes of a wide range of complex diseases of middle and old age. *PLoS Med.* **12**, e1001779 (2015).
26. Morales, J. *et al.* A joint NCBI and EMBL-EBI transcript set for clinical genomics and research. *Nature* **604**, 310–315 (2022).
27. Danecek, P. *et al.* Twelve years of SAMtools and BCFtools. *Gigascience* **10**, (2021).
28. Krusche, P. *et al.* Best practices for benchmarking germline small-variant calls in human genomes. *Nat. Biotechnol.* **37**, 555–560 (2019).
29. McLaren, W. *et al.* The ensembl variant effect predictor. *Genome Biol.* **17**, (2016).
30. Siepel, A. *et al.* Evolutionarily conserved elements in vertebrate, insect, worm, and yeast genomes. *Genome Res.* **15**, 1034–1050 (2005).
31. Yu, T., Fife, J. D., Adzhubey, I., Sherwood, R. & Cassa, C. A. Joint estimation and imputation of variant functional effects using high throughput assay data. *medRxiv* (2023)  
doi:10.1101/2023.01.06.23284280.
32. Chen, T. & Guestrin, C. XGBoost: A Scalable Tree Boosting System. *arXiv [cs.LG]* (2016).
33. Pedregosa, F. *et al.* Scikit-learn: Machine Learning in Python. *J. Mach. Learn. Res.* **12**, 2825–2830 (2011).
34. Schrödinger, LLC. The PyMOL Molecular Graphics System, Version 1.8. Preprint at (2015).
35. Rose, G. D., Geselowitz, A. R., Lesser, G. J., Lee, R. H. & Zehfus, M. H. Hydrophobicity of amino acid residues in globular proteins. *Science* **229**, 834–838 (1985).

36. Cock, P. J. A. *et al.* Biopython: freely available Python tools for computational molecular biology and bioinformatics. *Bioinformatics* **25**, 1422–1423 (2009).
37. Varadi, M. *et al.* AlphaFold Protein Structure Database: massively expanding the structural coverage of protein-sequence space with high-accuracy models. *Nucleic Acids Res.* **50**, D439–D444 (2022).
38. Webb, B. & Sali, A. Comparative protein structure modeling using MODELLER. *Curr. Protoc. Protein Sci.* **86**, 2.9.1-2.9.37 (2016).
39. Zhou, Y., Pan, Q., Pires, D. E. V., Rodrigues, C. H. M. & Ascher, D. B. DDMut: predicting effects of mutations on protein stability using deep learning. *Nucleic Acids Res.* **51**, W122–W128 (2023).

## **Acknowledgments**

The authors thank Grigoriy Losyev and Allison James for technical assistance and funding from UM1HG012010 (R.I.S., L.P.), 1R01HL164409 (C.A.C., R.I.S., L.P.), 1R01GM143249 (R.I.S.), R01HG010372 (C.A.C., T.Y.), American Cancer Society (R.I.S.), American Heart Association (R.I.S.), National Organization for Rare Diseases (R.I.S.), 1R35HG010717-01 (L.P.), National Health and Medical Research Council of Australia (GNT1174405; D.B.A. and Y.Z.), and the Victorian Government's Operational Infrastructure Support Program (Y.Z. and D.B.A.). We are indebted to the UK Biobank and its participants (UK Biobank application #41250 and IRB protocol 2020P002093). We thank Qian Qin, Cameron Smith, and Logan Blaine for help on BEAN model implementation and representation, Kendell Clement for advice on CRISPResso2 usage, Zain Patel for insights on transcription factor binding analysis, Soojung Yang for advice in structural analysis, and Hanna Boen for helping endogenous editing comparison.

## **Author contributions**

R.I.S. conceived the experimental design and J.R. and L.P. conceptualized BEAN. S.B. collected screen data. J.R. developed BEAN and M.J., M.I.L. and L.P. advised on design and implementation of BEAN. J.R. and T.Y. processed and analyzed data. T.Y. performed BE-Hive and FUSE analysis. M.F., Q.V.P., and R.I.S. performed downstream characterization of LDL-C GWAS variants. T.Y., L.B., and C.A.C. obtained and analyzed UKB data. Y.Z. led structural analysis of LDLR variants with J.R. and D.B.A. J.R. and Z.L. benchmarked classification performance. M.T. and L.P. performed analysis on variant impact on transcription factor binding. G.L. advised on library design. J.R. and R.I.S. drafted the manuscript. R.I.S., L.P., and C.A.C. provided guidance and supervised this project. All the authors wrote and approved the final manuscript.

## **Competing interests**

L.P. has financial interests in Edilytics, Inc., Excelsior Genomics, and SeQure Dx, Inc. L.P.'s interests were reviewed and are managed by Massachusetts General Hospital and Partners HealthCare in accordance with their conflict of interest policies. The remaining authors declare no competing interests.

### **Code availability**

Bean source code is available at <https://github.com/pinelloab/crispr-bean>. The scripts used to generate the figures and analyses presented in the manuscript have been deposited here: [https://github.com/pinelloab/bean\\_manuscript](https://github.com/pinelloab/bean_manuscript).

### **Data Availability**

The data used in this manuscript have been provisionally deposited on Zenodo (doi: 10.5281/zenodo.8270605). Controlled access, patient-level data from the UKB may be requested at <https://ams.ukbiobank.ac.uk/ams/>

### **Additional information**

Supplementary Information is available for this paper. Correspondence and requests for materials should be addressed to [ccassa@bwh.harvard.edu](mailto:ccassa@bwh.harvard.edu) (C.A.C.), [rsherwood@bwh.harvard.edu](mailto:rsherwood@bwh.harvard.edu) (R.I.S.), [lpinello@mgh.harvard.edu](mailto:lpinello@mgh.harvard.edu) (L.P.).