

1 **Title:** Decoding Heterogenous Single-cell Perturbation Responses

2
3 **One sentence summary:** We present a method to quantify diverse perturbation responses and discover
4 novel biological insights in single-cell perturbation datasets.

5
6 Bicna Song^{1,2}, Dingyu Liu^{3,4}, Weiwei Dai^{5,6}, Natalie McMyn⁵, Qingyang Wang⁷, Dapeng Yang³, Adam
7 Krejci⁸, Anatoly Vasilyev⁸, Nicole Untermoser⁸, Anke Loregger⁸, Dongyuan Song⁹, Breanna Williams³,
8 Bess Rosen^{3,10}, Xiaolong Cheng^{1,2}, Lumen Chao^{1,2}, Hanuman T. Kale³, Hao Zhang⁵, Yarui Diao¹¹,
9 Tilmann Bürckstümmer⁸, Jenet M. Siliciano⁵, Jingyi Jessica Li^{7,9,12-14}, Robert Siliciano^{5,6}, Danwei
10 Huangfu³, Wei Li^{1,2,#}

11
12 1 Center for Genetic Medicine Research, Children’s National Hospital, Washington DC, USA

13 2 Department of Genomics and Precision Medicine, George Washington University, Washington DC,
14 USA

15 3 Developmental Biology Program, Sloan Kettering Institute, New York City, NY, USA

16 4 Louis V. Gerstner Jr. Graduate School of Biomedical Sciences, Memorial Sloan Kettering Cancer
17 Center, New York City, NY, USA.

18 5 Department of Medicine, Johns Hopkins University School of Medicine, Baltimore, MD, USA

19 6 Howard Hughes Medical Institute, Johns Hopkins University School of Medicine, Baltimore, MD,
20 USA.

21 7 Department of Statistics and Data Science, University of California, Los Angeles, CA, USA

22 8 Myllia Biotechnology GmbH. Am Kanal 27, 1110 Vienna Austria.

23 9 Bioinformatics Interdepartmental Ph.D. Program, University of California, Los Angeles, CA, USA

24 10 Weill Cornell Graduate School of Medical Sciences, Weill Cornell Medicine, New York, NY, USA.

25 11 Department of Cell Biology, Duke University Medical Center, Durham, NC, USA

26 12 Department of Human Genetics, University of California, Los Angeles, CA, USA

27 13 Department of Biostatistics, University of California, Los Angeles, CA, USA

28 14 Department of Computational Medicine, University of California, Los Angeles, CA, USA

29
30 # Correspondences should be addressed. wli2@childrensnational.org;

31
32 **Keywords:** Perturb-seq, CRISPR-based genetic perturbations, single-cell RNA-seq, computational
33 model

34

35

36 Abstract

37
38 Understanding diverse responses of individual cells to the same perturbation is central to many
39 biological and biomedical problems. Current methods, however, do not precisely quantify the strength of
40 perturbation responses and, more importantly, reveal new biological insights from heterogeneity in
41 responses. Here we introduce the perturbation-response score (PS), based on constrained quadratic
42 optimization, to quantify diverse perturbation responses at a single-cell level. Applied to single-cell
43 transcriptomes of large-scale genetic perturbation datasets (e.g., Perturb-seq), PS outperforms existing
44 methods for quantifying partial gene perturbation responses. In addition, PS presents two major
45 advances. First, PS enables large-scale, single-cell-resolution dosage analysis of perturbation, without
46 the need to titrate perturbation strength. By analyzing the dose-response patterns of over 2,000 essential
47 genes in Perturb-seq, we identify two distinct patterns, depending on whether a moderate reduction in
48 their expression induces strong downstream expression alterations. Second, PS identifies intrinsic and
49 extrinsic biological determinants of perturbation responses. We demonstrate the application of PS in
50 contexts such as T cell stimulation, latent HIV-1 expression, and pancreatic cell differentiation. Notably,
51 PS unveiled a previously unrecognized, cell-type-specific role of coiled-coil domain containing 6
52 (CCDC6) in guiding liver and pancreatic lineage decisions, where CCDC6 knockouts drive the
53 endoderm cell differentiation towards liver lineage, rather than pancreatic lineage. The PS approach
54 provides an innovative method for dose-to-function analysis and will enable new biological discoveries
55 from single-cell perturbation datasets.

56

57

58 Introduction

59 Perturbation is essential for understanding the functions of the mammalian genome that encodes protein-
60 coding genes and non-coding elements (e.g., enhancers). Single-cell profiling of cells undergoing
61 genetic, chemical, environmental or mechanical perturbations is commonly used to examine
62 perturbation responses at the single-cell level. Recently, high-throughput approaches of perturbation
63 have been developed using single-cell RNA-seq (scRNA-seq) readout, including multiplexing of
64 perturbations and single-cell CRISPR screen (e.g., Perturb-seq, CROP-seq)¹⁻⁷. This concept has been
65 extended to study changes in single-cell chromatin accessibility^{8,9}, spatial transcriptomics¹⁰ upon
66 perturbations or perturbation combinations¹¹⁻¹³, and other phenomena.

67

68 Understanding how perturbations lead to different responses within cells is critical to understanding the
69 fundamental biology behind perturbation. Technical factors, including single-cell assays used to profile
70 the response, and the on-target/off-target effects of perturbations, are known drivers that lead to
71 differences of single-cell profiles in the data¹⁴⁻¹⁶. In Perturb-seq experiments that use CRISPR/Cas9 for
72 knockouts, both in-frame deletions¹⁶ and chromosomal losses¹⁷ contribute to different expression
73 profiles and clustering patterns of single cells.

74

75 Perhaps more interestingly, the heterogeneity perturbation responses may be driven by underlying
76 biological factors (**Fig. 1a**). These factors may be either cell-intrinsic (e.g., the activities of other coding-
77 and non-coding genomic elements) or cell-extrinsic (e.g., cell states or types, environment factors), all of
78 which define the context of perturbation response. For example, combined expressions of transcription
79 factors (TFs) are critical for many cellular state conversions. Therefore, to properly decode the functions
80 of these TFs via perturbation, one must consider the effect of the cell state and the activities of other
81 companion TFs. For this reason, defining the heterogeneity of perturbation response and identifying

82 factors that contribute to these outcomes is important for understanding how cells respond to
83 perturbation.

84
85 Unfortunately, computational frameworks are currently lacking to decode the diverse outcomes of
86 perturbations. For technical factors, mixscape is the only method to detect and mitigate confounding
87 variations (e.g., incomplete knockouts from CRISPR/Cas9)¹⁶. However, its performance has not been
88 rigorously benchmarked, especially when partial gene functions are perturbed using techniques like
89 CRISPR interference (CRISPRi). More importantly, no methods have been developed to reveal new
90 biological insights from the heterogenous perturbation outcomes, including studying how partial gene
91 perturbations affect a phenotype of interest (i.e., “dosage” analysis), and discovering biological
92 determinants that govern differential perturbation responses.

93
94 Here we present a computational framework, the perturbation-response score (PS), to quantify
95 heterogenous perturbation outcomes in single-cell transcriptomics datasets. The PS, estimated from
96 constrained quadratic optimization, quantifies the strength of the perturbation outcome for a single cell.
97 We performed comprehensive benchmark studies that demonstrated the outstanding performance of PS
98 over existing methods, including simulated datasets, genome-scale Perturb-seq, and published Perturb-
99 seq datasets that cover various CRISPR-based technologies. More importantly, PS analysis presents two
100 major conceptual advances for analyzing single-cell perturbation data: it enables analysis of the dose of
101 perturbation, and identification of novel biological determinants that govern the heterogeneity of
102 perturbation responses. First, we analyzed essential gene Perturb-seq and found two patterns of dose
103 response, based on whether moderate perturbation leads to strong expression changes of downstream
104 genes. Second, we identified intrinsic and extrinsic biological factors governing critical gene functions
105 in latent HIV-1 expression and pancreatic/liver development. Based on PS analysis results, we identified
106 and confirmed a novel function of CCDC6, wherein perturbation drives duodenum cell differentiation
107 towards liver commitment. Collectively, PS analysis provides a powerful tool to decode heterogenous
108 perturbation outcomes from single-cell assays.

109 110 **Results**

111 **Using PS to detect heterogenous perturbation outcomes within and across datasets.**

112 Perturbing the same gene (or non-coding elements) may result in different phenotypic changes or
113 transcriptional outcomes (**Fig. 1a**), depending on technical factors (e.g., perturbation efficiency) and
114 biological factors (e.g., cell type, cell state, activities of cofactors). Unfortunately, existing methods can
115 detect only technical factors¹⁶, while biological factors remain unexplored. To bridge this gap, we built a
116 computational framework to quantify perturbation outcomes in single-cell datasets using scRNA-seq as
117 readout. Corresponding assays include single-cell CRISPR screens (e.g., Perturb-seq), or simply
118 multiplex scRNA-seq profiling of various perturbations (e.g., sci-Plex; **Fig. 1b, c**). We define the
119 perturbation-response score (PS) to quantify the strength of perturbation, where PS=0 indicates no
120 perturbation effect (e.g., effects corresponding to unperturbed, wild-type gene functions) and PS=1
121 indicates the maximum perturbation effect observed within a dataset; for example, effects that
122 correspond to homozygous knockouts on both alleles of a gene. We utilize the expressions of multiple
123 downstream targets of a perturbed gene to infer the (unknown) values of PS (**Fig. 1b**). For example, if
124 one cell has dramatic expression changes on the known downstream target genes, then its value of PS
125 should be higher than cells with weak expression changes of these genes.

126

127 We built a computational model, based on a constrained quadratic optimization, to automatically
128 identify the downstream targets of perturbed genes and calculate PS (**Fig. 1c**). This model, named
129 “scMAGeCK-PS”, is based on our previously published scMAGeCK algorithm¹⁵ and consists of three
130 steps. First, scMAGeCK-PS identifies differentially expressed genes (DEGs) upon perturbation (e.g.,
131 perturbing the function of gene X), by comparing the transcriptome profiles between perturbed cells and
132 unperturbed cells. These DEGs are served as “signature” target genes of X . Second, scMAGeCK-PS
133 used a previously developed scMAGeCK model to estimate the average effect of perturbation on these
134 target genes, which can be estimated from the first step. Third, scMAGeCK-PS uses a constraint
135 optimization procedure to find the value of PS that minimizes the sum of mean squared errors between
136 predicted and measured expression changes of all downstream targets (see Methods). The constraints are
137 established such that any PS is non-negative for cells with X perturbed, and is exactly zero in cells
138 without perturbation. Such constraints can be established based on the prior information of
139 perturbations; for example, the expression matrix of single-guide RNAs (sgRNAs).

140
141 **PS outperforms mixscape in quantifying partial perturbations.**

142 mixscape¹⁶ is currently the only method to detect and remove technical factors that affect perturbation
143 outcomes, especially incomplete gene knockouts that are generated from CRISPR/Cas9. However, the
144 performance of mixscape on partial gene perturbations has not been fully evaluated. Here we compare
145 PS with mixscape using multiple benchmark datasets. We first used synthetic datasets to evaluate the
146 performances of different methods, because finding a real scRNA-seq dataset that contains ground truth
147 (*i.e.*, accurate measurements of loss-of-function upon perturbation) is challenging. For synthetic data
148 generation, we used scDesign3¹⁸ to simulate the single-cell transcriptomic responses upon perturbing the
149 50% and 100% functions of Nelfb, based on a real scRNA-seq dataset that deletes Nelfb in mouse T
150 cells¹⁹ (**Supplementary Fig. S1a**; see Methods). We specified different numbers of DEGs (from 10 to
151 500) and simulated their expression changes upon 50% or 100% perturbations of Nelfb functions. In all
152 the cases, PS correctly estimated partial perturbation, where the median PSs range from 0.32-0.34 for
153 50% perturbation, and greater than 0.8 for 100% perturbation, respectively (**Fig. 1d-e**; **Supplementary**
154 **Fig. S1b-e**). In contrast, mixscape uniformly assigned the posterior probability of perturbation to 1 in all
155 cases, an indication that mixscape is not suited to analyze the outcome of partial gene perturbations (**Fig.**
156 **1d-e**; **Supplementary Fig. S1b-e**), possibly due to the bimodal statistic model it uses, which only
157 considers 100% knockout effects¹⁶.

158
159 We next evaluated different methods using real single-cell perturbation datasets. We chose CRISPRi-
160 based Perturb-seq datasets (**Fig. 1f-i**) because the CRISPRi system directly modulates the expression
161 levels of perturbed genes, and the perturbation efficiency can be accessed using single-cell
162 transcriptomic data. We use two published K562 CROP-seq datasets:²⁰ in the first, only 1 gRNA is
163 expressed within each cell (*i.e.*, low multiplicity of infection or MOI), and in the second multiple
164 gRNAs are expressed (*i.e.*, high MOI). We examine cells where the transcription starting sites (TSS) of
165 highly expressed protein-coding genes are targeted (23 and 342 genes, respectively; **Fig. 1g-i**,
166 **Supplementary Fig. S1f-g**). If the TSS of gene X is perturbed, we first removed the expression of X
167 from expression matrix, and used the rest of gene expressions to measure the perturbation efficiency of
168 X . The scores of different methods were then compared with the expression of X , producing a direct
169 measurement of perturbation efficiency (**Fig. 1f**). In over 40% of these genes (10 out of 23 for low MOI,
170 139 out of 342 for high MOI), PS has a strong negative correlation with the expression of X (**Fig. 1g**),
171 defined as Pearson correlation coefficient < -0.1 and p value < 0.01 . In contrast, mixscape scores
172 correlate with X expression in none of the genes (for low MOI dataset; **Fig. 1g**), or in less than 5% of all

173 the genes (for high MOI dataset; **Supplementary Fig. S1f-g**). PS detects a much greater number of cells
174 that have a strong perturbation effect (PS or mixscape score >0.5 ; **Fig. 1h**), whose scores are strongly
175 negatively correlated with gene expression (**Fig. 1i**). We also tested both methods in another CRISPRi-
176 based Perturb-seq dataset, where sgRNAs with mismatches reduce efficiency, leading to partial
177 perturbation effects²¹ (**Supplementary Fig. S1h-i**). PS has a high sensitivity and a good balance
178 between sensitivity and specificity, evidenced by the higher Pearson correlation coefficients
179 (**Supplementary Fig. S1h**) and areas under the receiver-operating characteristic (ROC) curve (AUC)
180 values (**Supplementary Fig. S1i**).

181
182 To further benchmark methods in terms of a phenotype of interest, we designed and performed a
183 genome-scale CRISPRi Perturb-seq on both unstimulated and stimulated Jurkat, a T lymphocyte cell
184 model (**Fig. 2a**), and evaluated the performances of different methods in identifying known regulators of
185 T cell activation. We designed Perturb-seq library that contains sgRNAs targeting the TSS of 18,595
186 genes (4-6 guides per gene) and used a TAP-seq-based²² multiplex primer panel to detect the
187 expressions of 374 genes with high sensitivity (see **Supplementary Table S1** and Methods). We
188 obtained high-quality scRNA-seq data on over 586,000 single cells after quality control, and the UMAP
189 clustering of Perturb-seq datasets clearly demonstrated the differences between stimulated and non-
190 stimulated cells (**Fig. 2b**). Next, we ran PS or mixscape to calculate the scores of all perturbations at a
191 single-cell level; and for each perturbed gene, we calculated its overall perturbation-response score, by
192 adding the scores of all cells that express a corresponding sgRNA targeting that gene. Because our
193 system focuses on T cell stimulation, perturbing a gene that reaches a highest (and lowest) cumulative
194 score should have the strongest (and no) effect on T cell stimulation, respectively. For an independent
195 evaluation, we extracted 385 (and 1297) positive (and negative) hits whose perturbation impairs (or does
196 not impair) the stimulation of T cells from a published genome-scale CRISPR screen²³. Both Perturb-seq
197 and pooled CRISPR screen identified many known positive regulators of T cell activation, such as
198 components of the T cell receptor complex (e.g., CD3D) and proximal signaling components (e.g., LCK;
199 **Fig. 2c**). For many positive genes, cells with higher values of PS or mixscape score are skewed towards
200 non-stimulating state, consistent with their negative selections in pooled CRISPR screens using T cell
201 stimulation as readout (**Fig. 2c; Supplementary Fig. S2**). However, when comparing the ROC score,
202 PS reaches a higher AUC score than mixscape (**Fig. 2d**), indicating its better performance in accurately
203 separating positive from negative hits.

204
205 Finally, we tested different methods on a published ECCITE-seq, which simultaneously measures
206 single-cell transcriptomes, surface proteins, and perturbations¹⁶. PDL1 protein expression was used as an
207 independent metric of evaluation (**Fig. 2e**), because PDL1 is a well-studied gene whose protein
208 expression is well understood. Among 25 perturbed genes in the ECCITE-seq perturbation library, 17
209 are known to regulate PDL1 expression (**Fig. 2f**). We compared PS with mixscape in terms of predicting
210 changes in PDL1 expression (**Fig. 2f; Supplementary Fig. S3**). In addition, the expression of the
211 perturbed gene is included in the comparison as a naïve method. In 19 out of 25 genes (76%), PS
212 outperformed mixscape and perturbed gene expression in predicting PDL1 expression (**Fig. 2e**),
213 including 12 out of 17 (71%) known PDL1 regulators. Notably, for genes whose perturbations led to
214 strong transcriptomic changes (e.g., IFNGR1, IFNGR2, JAK2, STAT1), both PS and mixscape work
215 well, reaching $AUC > 0.8$ (**Fig. 2f**). For other genes whose perturbation only leads to moderate or weak
216 expression changes, as described previously¹⁶, PS outperforms mixscape, including those that are
217 confirmed to be PDL1 regulators (*i.e.*, genes marked in red in **Fig. 2f**).

218

219 **Analyzing dose-dependent effects of perturbation.**

220 Traditionally, dosage analysis requires a careful, time-consuming adjustment of perturbation strength,
221 including changing drug concentrations or designing sgRNA sequences to achieve various editing
222 efficiencies^{21,24}. Since the quantifying partial gene perturbation by PS is highly accurate (**Fig. 1-2**), we
223 can use PS to perform dose-response analysis of perturbation, without the need to titrate the strength of
224 perturbation. By examining ECCITE-seq data in which PDL1 expression was measured directly (**Fig.**
225 **2e**), we found correlations between PDL1 expression and the PS of known PDL1 regulators (**Fig. 3a**).
226 The PSs of positive PDL1 regulators (e.g., IFNGR1/2, STAT1; **Fig. 3b**) are negatively correlated with
227 PDL1 expression, while the scores of negative regulators (e.g., CUL3, BRD4) are negatively correlated
228 (**Supplementary Fig. S3; Fig. 3c**). One example is CUL3, which is known to destabilize and degrade
229 PDL1 protein expression²⁵. Consequently, higher CUL3 PSs, indicating higher CUL3 functional
230 perturbation, correspond to higher PDL1 protein expressions (**Fig. 3a**). Compared with mixscape, PS
231 more accurately predicts the quantitative changes in PDL1 expression, evidenced by stronger Pearson
232 correlations between the two (**Supplementary Fig. S3**).

233
234 We further investigated the relationships between perturbation efficiency and the strength of
235 perturbation responses, which is measured by PS (**Fig. 3d**). In particular, we are interested in genes that
236 show one of two different patterns of PSs upon perturbation: “buffered” distribution, where genes have
237 high PSs only when stronger perturbation efficiency is achieved; and “sensitive” distribution, where the
238 PSs are high, even with moderate or weak perturbation efficiency. Both “buffered” or “sensitive” terms
239 have been coined previously to describe the effects of transcription factor dosages to chromatin
240 accessibility²⁶. CRISPRi-based Perturb-seq datasets are used, as the efficiencies of CRISPR inhibition
241 can be directly evaluated by examining perturbed gene expressions (**Fig. 2e**).

242
243 We calculated PS for every gene in a published essential-wide Perturb-seq²⁷, which uses CRISPRi to
244 inhibit the expressions of 2,285 common essential genes. We classified genes based on their PS
245 quantiles that correspond to around 50% perturbation efficiency (**Fig. 3e**): a gene is classified as
246 “buffered” if its median PS is smaller than 0.1; or “sensitive”, if its 95% quantile is greater than 0.75.
247 Among over 2,000 essential genes, we classified 613 genes as either buffered or sensitive. The majority
248 are buffered (529 out of 613), indicating high robustness to perturbation, possibly due to their essential
249 roles in cellular functions that require compensations on expression reductions. Many buffered genes
250 belong to essential protein complexes, including proteasomes (e.g., PSMA3; **Fig. 3f**) and ribosomal
251 subunits (e.g., RPL4; **Supplementary Fig. S4a**). 30% of the genes (185 out of 613) belong to
252 “sensitive” category, showing strong transcriptome responses even with moderate or weak efficiencies
253 upon perturbing gene expression (**Supplementary Fig. S4b-c**). Many of the sensitive genes are also
254 displaying buffering effect, a demonstration of complex, heterogenous responses of cells undergoing the
255 same perturbation of essential genes. Notably, 50% reduction of HSPA5 and GATA1 expression
256 achieved near-maximal transcriptional response (and the associated growth defect), as in previous
257 studies²².

258
259 We further examined possible mechanisms by which buffered genes resist perturbation, especially those
260 that belong to the same functional protein complex. Interestingly, perturbing one member of the protein
261 complex usually leads to the expression up-regulation of other members of the complex, indicating a
262 possible mechanism for compensation. For example, perturbing proteasome subunits led to a strong
263 expression reduction of the perturbed gene (e.g., PSMA5; blue squares in **Fig. 3g**) and concurrent up-
264 regulation of other members of the proteasomes (e.g., PSMB7, PSMD2). Perturbing many other protein

265 complexes, including ribosomal subunit, mediator, and RNA polymerases, also leads to similar up-
266 regulation of some members of the same functional unit (**Supplementary Fig. S5a-c**), indicating that
267 compensation occurs by up-regulation of other subunits of the same molecular machine. To confirm our
268 findings on a different cellular system, we examined the effects of perturbing proteasomes in our
269 genome-scale Perturb-seq dataset (**Fig. 2a**). The TAP-seq approach used in this dataset provides a
270 sensitive and accurate measurement of gene expression changes upon perturbation²⁷. Indeed, perturbing
271 members of the proteasome subunits leads to the up-regulation of other proteasomes (**Supplementary**
272 **Fig. S5d**), consistent with the known transcriptional feedback loop that is observed between proteasome
273 genes²⁸. Overall, the widespread existence of such compensatory effect may explain the perturbation-
274 expression phenotype of buffered genes, where a strong perturbation efficiency is needed to achieve
275 strong expression changes.

276 277 **PS reveals intrinsic and extrinsic biological factors that regulate gene functions in latent HIV** 278 **expression.**

279 We next perform Perturb-seq experiment and use PS to investigate the functions of key genes regulating
280 latent HIV-1 expression. We used a Jurkat HIV cell model that we previously established for pooled
281 CRISPR screening²⁹, where cells stably express Cas9 and are latently infected with HIV-GFP viral
282 vector. We designed a Perturb-seq library that targets 10 protein-coding genes (**Supplementary Table**
283 **S2**), which are either (1) known factors in HIV-1 virus expression and T cell activation (e.g., BIRC2), or
284 (2) top hits from genome-scale CRISPR screens that we previously performed (e.g., BRD4)²⁹. We
285 performed Perturb-seq experiments in three different conditions, including stimulated Jurkat (by PMA/I)
286 followed by GFP expression sorting (GFP+ or GFP-), and unstimulated cells (**Fig. 4a**). The single-cell
287 transcriptomes were profiled via the 10X Genomics Chromium platform, and expressed guide RNAs can
288 be captured directly. After quality controls, we received 7,063-8,811 single cells per sample, where the
289 mean reads per cell (and median genes expressed per cell) in each sample is at least 69,888 (and 4,744),
290 respectively (**Supplementary Fig. S6a**). Guide RNAs were detected in over 96% of the cells, and over
291 85% of these cells could be assigned a unique guide RNA (**Supplementary Table S3**). The
292 transcriptome profiles of cells are primarily clustered by cell states (stimulated vs. unstimulated),
293 indicating that the primary sources of expression variation are coming from cell states (**Fig. 4b**).

294
295 We investigated gene functions using our PS framework. Among all perturbed genes, the PS of BRD4
296 (bromodomain containing 4) demonstrates a strong cell state-specific pattern, where a subset of cells
297 with BRD4 perturbation has strong BRD4 PSs (named “BRD4-PS+ cells”) than other BRD4-perturbed
298 cells (or “BRD4-PS- cells”; **Fig. 4c**). BRD4-PS+ cells overexpress genes that are involved in known
299 functions of BRD4^{30,31} including NF-kB/TNF-alpha signaling, hypoxia and apoptosis (**Fig. 4c-d**,
300 **Supplementary Fig. S6b-d**). We examined whether the differences in BRD4 PS reflects the degree of
301 BRD4 functional perturbation. We first checked the expressions of BRD4 “signature” genes from
302 another published study³². Compared with BRD4-PS- cells, BRD4-PS+ cells have a much lower
303 expressions of these signature genes (**Supplementary Fig. S6e**), indicating a stronger functional BRD4
304 perturbation. In addition, BRD4 has been shown to inhibit HIV transcription and activation in many
305 studies, including our previous CRISPR screens^{29,33}, consistent with the fact that HIV-GFP is one of the
306 strongest up-regulated genes in BRD4-PS+ cells (**Supplementary Fig. S6f**). Furthermore, BRD4-PS+
307 cells have a stronger GFP expression (**Fig. 4d**) than other cells, confirming a stronger BRD4 functional
308 perturbation in these cells.

309

310 To build a quantitative perturbation-expression relationship, we recalculated BRD4 PS without using
311 HIV-GFP expression and examined how the scores are associated with a phenotype of interest (i.e.,
312 latent HIV-GFP expression) in different conditions (**Fig. 4e**). BRD4 PS correlation with HIV-GFP
313 expression is cell-state dependent: in stimulated T cells (PMA/I treatment), a linear, positive correlation
314 is observed regardless of the GFP expression. In contrast, a nonlinear relationship exists in unstimulated
315 T cells (DMSO), where stronger BRD4 PS (>0.5) leads to a sharp increase in HIV-GFP expression (**Fig.**
316 **4e**).

317
318 Another gene, cyclin T1 (CCNT1), also displays heterogeneity in PS distribution: cells with CCNT1
319 perturbation have a high PS distribution only in stimulated cells (**Fig. 4f**). This is different from CCNT1
320 gene expression or guide distribution, which do not show such pattern differences between cell states
321 (**Supplementary Fig. S7a**). Confirming our findings, the number of DEGs (cells with CCNT1
322 perturbation vs. cells expressing non-targeting guides) is over 100 in stimulated cells, but only 1 in non-
323 stimulated cells (adjusted p value <0.001; **Supplementary Fig. S7b**). In particular, HIV-GFP is the
324 strongest DEG in cells with CCNT1 perturbation, consistent with the known role of CCNT1 in
325 activating HIV transcription.

326
327 CCNT1 is a key subunit of P-TEFb (positive transcription elongation factor b)/CDK9 complex that
328 drives RNA transcription, including the transcription of HIV. The transcription elongation control of P-
329 TEFb/CDK9 is a complicated process that is regulated by multiple mechanisms, including various T cell
330 signaling pathways (e.g., NF- κ B signaling), translation control, and epigenetic modification (reviewed in
331 ³⁴). The activities of these factors are different in different states of T cells (e.g., NF- κ B;
332 **Supplementary Fig. S7c**), which may explain the differences of CCNT1 PSs. Despite the strong cell
333 state dependency of CCNT1 PS, PS shows weak correlation with HIV-GFP within one cell state
334 (**Supplementary Fig. S7d**), which is different from BRD4 PS (**Fig. 4e**).

335
336 To further confirm our finding that different cellular states affect the transcriptomic responses of
337 CCNT1 perturbation, we stimulated Jurkat cells using a different agonist (TNF-alpha). To measure the
338 downstream effect of CCNT1 perturbation, we sorted cells by expression of HIV-GFP, which is the
339 strongest down-regulated gene upon CCNT1 knockout (**Supplementary Fig. S7b**), and whose
340 expression is known to be regulated by CCNT1^{35,36}. Indeed, with the presence of TNF-alpha, CCNT1
341 knockout leads to a strong reduction in HIV-GFP expression (over 50% reduction), while such reduction
342 is much smaller (<5% reduction) in cellular states without TNF-alpha stimulation (**Fig. 4g**).

343 Collectively, these results demonstrated that PS is a powerful computational framework for investigating
344 cofactors (cell states, other genes) that drive transcriptomic responses upon gene perturbation.

345 346 **PS enables identification of novel cell-type dependent gene functions in regulating pancreatic cell** 347 **differentiation from multiplex single-cell transcriptomics.**

348 Besides Perturb-seq, multiplexing cells with different perturbations are also used to measure single-cell
349 responses to perturbation^{2,19}. A mixture of cells from different perturbations can be sequenced at the
350 same time, and the identity of cells can be established using various methods including cell hashing³⁷,
351 the expressions of pre-defined barcodes³⁸, or a combination of random barcodes³⁹. We therefore tested
352 our PS framework on pooled single-cell transcriptomics of different perturbations to study the functions
353 of lineage regulators during human pancreatic differentiation. By using an established in vitro human
354 embryonic stem cell (hESC) pancreatic differentiation system, we generated cells corresponding to early
355 stage (definitive endoderm, DE) and middle stage (pancreatic progenitor, PP) pancreas development. To

356 test the performance of PS framework and uncover the functions of unknown regulators, we picked ten
357 clonal hESC lines with the homozygous knockout of four genes (**Supplementary Table S4**), including
358 two known pancreatic lineage regulators (HHEX, FOXA1) and two uncharacterized candidate regulators
359 from previous genetic screens (OTUD5, CCDC6)^{40,41}. These clones are then labelled with different
360 LARRY (Lineage and RNA recovery) DNA barcodes³⁸, pooled together and differentiated into DE and
361 PP stages using established protocols⁴⁰. Finally, the single-cell expressions of these cells were profiled
362 via 10X genomics Chromium platform (**Fig. 5a**). The clone information of each cell was identified from
363 LARRY barcodes. Among 26,286 single cells that passed the quality control measurements, over 97%
364 (25,694/26,286) of the cells had at least one barcode detected, and over 80% (20,678/25,694) were
365 identified as singlets and retained for downstream analysis. UMAP clustering revealed different known
366 cell types during pancreatic differentiation, based on the expression markers of known cell types (**Fig.**
367 **5b; Supplementary Fig. S8**), including DE, PP, liver/duodenum progenitor (LV/DUO), endocrine
368 precursor (EP), and cells in transition stages (e.g., DE in transition, PP in transition).

369
370 We next applied the PS framework to the pooled single-cell RNA-seq datasets containing different
371 knockout clones. Among all knockout genes, HHEX PS is high in cells whose type is between two
372 different differentiated cell types (PP and LP/DP; **Fig. 5c; Supplementary Fig. S8**), consistent with the
373 known function of HHEX as a key determinant of cell fate decision, whose deletion drives DE cell
374 differentiation towards LP/DP, rather than PP⁴⁰. Indeed, *HHEX* knockout led to a much fewer
375 percentage of cells that are annotated as PP (**Fig. 5d**). The PS of FOXA1, another key transcription
376 factor during PP differentiation, is strong in DE and PP cell types, consistent with the specific
377 expression pattern of FOXA1 in DE/PP cell types (**Supplementary Fig. S9a-c**).

378
379 As in our previous genome-wide CRISPR screens, CCDC6 is one of the top hits whose perturbation
380 hinders PP differentiation^{40,42}. However, the exact function of CCDC6 during pancreatic differentiation
381 is largely unknown. CCDC6 may have different functions at different cell types, evidenced by the few
382 overlaps of DEGs between different cell types (**Supplementary Fig. S9d-f**). To investigate these
383 different functions, we calculated PSs from the DEGs from four major cell types in the dataset (DE in
384 transition, DE, PP/PP in transition, and LV/DUO). An unbiased clustering on these CCDC6 PSs
385 demonstrated two distinct distributions across cell types (**Fig. 5e**), where scores calculated from late-
386 stage cell types including PP/PP in transition/LV/DUO (“pattern 1”) are distributed differently from
387 scores calculated from early-stage cell types including DE in transition/DE (“pattern 2”; **Fig. 5f;**
388 **Supplementary Fig. S10a-b**), implying different behaviors of CCDC6 perturbation at different cell
389 types. Indeed, functional analysis on DEGs leading to both patterns have distinct enrichment terms. In
390 early-stage cell types, DEG genes are enriched in the targets of stem cell transcription factors (e.g.,
391 SOX2, POU5F1, NANOG) and cell cycle regulation (**Supplementary Fig. S10c-e**), consistent with the
392 known function of *CCDC6* as a cell cycle regulator^{43,44}. In contrast, DEGs in late-stage cell types are
393 primarily the targets of HNF4A, a key transcription factor that drives LP/DP differentiation (**Fig. 5g;**
394 **Supplementary Fig. S10f**). The expressions of these transcription factors (SOX2, HNF4A) are among
395 the up-regulated genes in both programs, respectively (**Supplementary Fig. S9d-e**). Furthermore,
396 compared with wild-type cells, *CCDC6* knockout cells have a much lower percentage of PP cells and a
397 higher percentage of LP/DP cells (**Fig. 5h**). Collectively, these results imply that *CCDC6* has different
398 functions for early vs. late-stage cell types. Especially in late-stage cell types, *CCDC6* knockout drives
399 cell differentiation towards LV/DUO cell types rather than PP cell types.

400

401 To further validate the prediction results of *CCDC6*, we performed flow cytometry analysis to evaluate
402 the effects of *CCDC6* knockout on the composition of late-stage cell types (PP/LV/DUO). We examined
403 the percentage of HNF4A⁺ cells, a marker for LV population, and PDX1⁺ cells, a marker for PP
404 population. Indeed, both clones of *CCDC6* knockout greatly reduced PDX1⁺ population and increased
405 HNF4A⁺ population in three biological replicates (**Fig. 5i; Supplementary Fig. S11**), confirming our
406 finding on the enrichment of *CCDC6* PS in LP/DP populations (**Fig. 5f-g**).

407

408 Discussion

409 Understanding cellular responses to perturbations is a central task in modern biology, from studying
410 tumor heterogeneity to developing personalized medicine. These perturbations may be genetic (e.g.,
411 knocking out genes or non-coding elements), chemical (e.g., drug treatments), mechanical (e.g.,
412 pressure) or environmental (e.g., temperature changes). Single-cell genomics profiles of perturbations
413 are commonly used to investigate the mechanisms of perturbations. Many technologies, including
414 Perturb-seq and sci-Plex, provide a high-content readout of the results of systematically perturbing many
415 genes or non-coding elements. Despite rapid technological advancements, a major bottleneck is the lack
416 of a computational model to fully unlock the potential of high-content perturbation, especially for
417 discovering novel biological insights from the data. Here we introduce the PS framework to model the
418 heterogenous transcriptomic responses of perturbations and to enable novel biological discovery from
419 modeling perturbation heterogeneity.

420

421 Partial gene perturbation is common in perturbation experiments. Partial perturbations may come from
422 dose-controlled drug treatment, gene editing technology that does not fully knockout gene function (e.g.,
423 RNA- or CRISPR-interference, epigenome editors), or from CRISPR/Cas9 that generates random DNA
424 editing outcomes. We demonstrated the outstanding performance of our PS method over existing
425 methods in quantifying partial gene perturbation. Specifically, partial perturbation identification enables
426 the analysis of dose-dependent effect, which is demonstrated in this study using various datasets.

427

428 More importantly, PS enables novel biological investigations, including analysis of perturbation dosage
429 without the need to titrate perturbation strength and identification of cell-intrinsic and extrinsic
430 biological factors that regulate perturbation responses. In the latter case, the PS, ranging between 0 and
431 1, no longer represents the quantity of partial perturbation, but instead represents the strength of the
432 perturbation outcome. Therefore, PS becomes a convenient tool to identify cell context that determines
433 perturbation outcome. We demonstrated the application of PS in various biological problems, including
434 T cell activation, essential gene function, latent HIV-1 virus expression, and pancreatic cell
435 differentiation. Importantly, our PS model leads the discovery of novel *CCDC6* functions that are cell
436 type dependent, whose role as a regulator during pancreatic and liver cell fate decision is experimentally
437 validated.

438

439 Partial perturbations of gene functions contribute to the complexity of many biological processes. For
440 example, “haploinsufficient” genes are able to cause disease phenotypes when 50% of their functions
441 are disrupted, while “haplosufficient” genes will require a nearly complete gene knockout. However, we
442 currently lack a method to investigate the phenotypes of partial gene perturbations or to efficiently
443 perform dosage analysis at a large scale. Current approaches, such as introducing mismatches to guide
444 RNAs to modulate the effects of CRISPRi²¹ or Cas13²⁸, require a complex design of a specific CRISPR
445 system. Here we demonstrated that both CRISPR knockout (e.g., **Fig. 2f, Fig. 4e**) and CRISPRi
446 naturally introduce partial perturbation effects, which can be used to study the dose effect of partial gene

447 perturbations on downstream gene expressions or a phenotype of interest. Our PS framework is
448 versatile, enabling the dosage analysis using various perturbation methods (e.g., CRISPRi or CRISPR
449 knockout) and assays (e.g., Perturb-seq or multiplex scRNA-seq).

450
451 Results from genetic perturbations (e.g., via CRISPR/Cas9) are informative for drug development, and
452 confirmations from genetic perturbation experiments are usually required to demonstrate the feasibility
453 of candidate drug targets. However, titrating pharmaceutical interventions are easy (e.g., by using
454 different doses of drugs), while it is much more difficult to precisely control the degree of genetic
455 perturbations. Our PS framework provides a convenient alternative to dose-dependent perturbations,
456 especially genetic perturbations, and their associations with phenotypic changes, which will be
457 informative in designing drugs. For example, BRD4 is the primary target of bromodomain inhibitors
458 (BETi), many of which have been proposed as candidates of latency reversing agents (LRAs) to
459 reactivate latent HIV-1 expression. The distribution of BRD4 PSs (**Fig. 3**) reveals that stronger
460 perturbation effects are needed to induce the desired phenotype, in this case, the expression of HIV-GFP
461 (**Fig. 3**). Since BRD4 is an essential gene, a strong BRD4 perturbation may lead to unexpected toxicity,
462 thereby limiting the efficacy of BETi. Indeed, our previous study²⁹ demonstrated that 10-1000x higher
463 doses of JQ1, a commonly used BETi, are needed to induce latent HIV-1 expression at a similar level
464 with other potent LRAs. Our results further warrant the development of synergistic drug combinations to
465 mitigate the narrow therapeutic window of BETi, which is currently tested in many studies.

466
467 Our PS analysis provides a general framework to analyze several major sources that contribute to the
468 heterogeneity of perturbation responses: the strength of perturbation *per se* (e.g., **Fig. 1i, 3d**; BRD4 in
469 **Fig. 4c**), compensations to perturbation especially on essential genes (e.g., proteasomes; **Fig. 3g**), and
470 cell type/state specificity (e.g., T cell states in **Fig. 4**; differentiation cell types in **Fig. 5**). Importantly,
471 cell type/state is linked to perturbation responses in three distinct ways: cell type/state may change as a
472 result of perturbation (e.g., CCDC6 and HHEX in Fig. 5); cell type/state serving a critical context to
473 define perturbation responses (e.g., T cell states in response to CCNT1 perturbation in **Fig. 4f-g**); and
474 cell type/state as a confounding factor that drives perturbation responses (e.g., BRD4 perturbation
475 heterogeneity in unstimulated T cells in **Fig. 4c**). Compared with other methods, PS is currently the only
476 method to analyze heterogeneity of perturbation responses from all these aspects.

477
478 Confounding factors are the major sources of variation when analyzing single-cell perturbation effects.
479 These confounding factors can be modeled explicitly (e.g., using generalized linear models) if
480 confounding source is known; or be detected and corrected using mathematical or statistical approaches
481 including matrix factorization (e.g., using GSFA⁴⁵) or independent component analysis (e.g., using
482 CINEMA-OT⁴⁶). In contrast, PS does not explicitly model confounding factors. Instead, PS scores can
483 be used in combination with methods that remove confounding sources of variation, or to detect these
484 confounding factors that contribute to the heterogeneity in perturbation responses (e.g., **Fig. 4c**).
485 Importantly, many confounding factors defined in previous methods^{16,46} are not always confounding;
486 instead, they can be used to discover novel biological insights, as are shown in this study (e.g.,
487 perturbation efficiency, cell type/state). The orthogonal algorithmic design of PS compared with existing
488 methods also allows the combination of PS with these methods to simultaneously remove confounding
489 factors and measure the strength of perturbation responses.

490
491 One limitation of PS is its power in detecting drastic changes in cell types or states. For example, even
492 moderate perturbations on essential gene functions affect cellular viability^{47,48}. In this case, single-cell

493 profiling only captures surviving cells that are resistant to essential gene perturbations in various
494 mechanisms (e.g., expression compensation in **Fig. 3**), and largely misses dead cells due to essential
495 gene dysfunction. Consequently, due to this “survival bias”, PS probably only reflects the perturbation
496 responses in a fraction of cells, rather than the full spectrum of perturbations. To overcome this
497 limitation, PS can combine with recently developed prediction methods that predict the responses of
498 perturbations, even if cells between perturbed/non-perturbed states are unevenly distributed⁴⁹.

499

500 **Methods**

501

502 **The Perturbation-response Score (PS) framework**

503 Estimating PS proceeds in three steps, as illustrated in Figure 1c: target gene identification (Step 1),
504 average perturbation effect estimation using a previously published scMAGeCK (Step 2), and PS
505 estimation using constrained optimization (Step 3).

506

507 Step 1: target gene identification. We first performed differential expression analysis between cells with
508 certain perturbation (e.g., knocking out gene X) and negative control cells. In most cases, negative
509 control cells are cells that express non-targeting guide RNAs (in Perturb-seq), or wild-type cells (in
510 pooled scRNA-seq). In Perturb-seq with high MOI condition, these cells may come from cells that do
511 not have a particular perturbation. We used Wilcoxon rank sum test (implemented in Seurat) to identify
512 and rank differentially expressed genes. Top genes were then selected as potential target genes of the
513 specific perturbation. The maximum and minimum numbers of top genes can be specified by the user.
514 Alternatively, users can provide the list of target genes for each perturbation, based on prior knowledge,
515 therefore skipping the differential expression analysis in this step.

516

517 Step 2: average perturbation effect estimation. We used the linear regression module in scMAGeCK
518 (scMAGeCK-LR) to estimate the average perturbation effect. scMAGeCK-LR takes the expressions of
519 all target genes (identified in Step 1) in all cells as input and outputs a β score, which is conceptually
520 similar to log fold change. There are two advantages of using β score, instead of simply using the log
521 fold changes in Step 1. First, scMAGeCK-LR naturally supports datasets from high MOI Perturb-seq,
522 where one cell may express multiple guides targeting different genes. Second, scMAGeCK-LR is able to
523 estimate average perturbation effects of multiple perturbations (e.g., genome-scale perturbations) in one
524 step, while a naïve DEG analysis can only calculate LFC for each perturbation.

525

526 The mathematical model of scMAGeCK-LR is described as follows. Let Y be the log-transformed, $M \times N$
527 expression matrix of M single cells and N target genes. These genes are the union of all target genes for
528 all K perturbations, extracted from Step 1. Let D be the $M \times K$ binary cell identity matrix of M single cells
529 and K perturbations, where $d_{jX} = 1$ if single cell j contains sgRNAs targeting gene X ($j =$
530 $1, 2, \dots, M; X = 1, 2, \dots, K$), and $d_{jX} = 0$ otherwise. D can be obtained from the detected guide RNA
531 expression matrix from Perturb-seq or from the prior sample information from pooled scRNA-seq. The
532 effect of target gene knockout on all expressed genes is indicated as a β score in a matrix B with size
533 $K \times N$, where $\beta_{XA} > 0$ (< 0) indicates gene X is positively (or negatively) selected on gene A expression,
534 respectively. In other words, gene X knockout increases (or decreases) gene A expression if $\beta_{XA} > 0$ ($<$
535 0), respectively.

536

537 The log-transformed expression matrix Y is modeled as follows:

538

539
$$Y = Y_0 + D \times B + \epsilon, \quad \text{Eq (1)}$$

540

541 where Y_0 is the basal expression level of all genes in an unperturbed state, and ϵ is a noise term
 542 following a Gaussian distribution with zero means. Y_0 can be estimated from negative control cells (e.g.,
 543 wild-type cells or cells expressing non-targeting guides), or be modeled using the expressions of
 544 neighboring negative control cells (e.g., the approach used by mixscape¹⁶). The value of B can be
 545 estimated using ridge regression:

546

547
$$B = (D^T D + \lambda I)^{-1} D^T Y, \quad \text{Eq (2)}$$

548

549 where I is the identity matrix, and λ is a small positive value (default 0.01).

550

551 Step 3: PS estimation using constrained optimization. We revise Eq (1) to incorporate PS. Here, the log-
 552 transformed expression matrix Y is modelled as follows:

553
$$Y = Y_0 + \Psi \times B + \epsilon, \quad \text{Eq (3)}$$

554

555 Where Ψ is the non-negative, raw PS matrix with the same size as D in Step 2 ($M \times K$). Each element ψ_{jx}
 556 in Ψ indicates the raw PS of cell j of perturbing gene X . Here, B is the β score matrix which is estimated
 557 in Step 2. We find the value of Ψ to minimize the squared error of predicted and observed expressions
 558 of all genes within all cells, subject to constraints and regularization terms:

559

560
$$\min \sum_{ji} (y_{ji} - y_{ji}^0 - \sum_k \psi_{jk} \beta_{ki})^2 + \lambda \sum_{jk} |\psi_{jk}|, \quad \text{Eq (4)}$$

561

562 subject to the following constraints:

563

564
$$\begin{cases} 0 \leq \psi_{jk} \leq U, & \text{if } d_{jk} = 1 \\ \psi_{jk} = 0 & \text{if } d_{jk} = 0 \end{cases}.$$

565

566 Here, U is a positive value indicating the upper bound of raw Ψ values, and d_{ik} is the value of the binary
 567 cell identity matrix in Step 2. $1 \leq j \leq M$ is the index of single cells, $1 \leq i \leq N$ is the index of target
 568 genes, and $1 \leq k \leq K$ is the index of perturbations.

569

570 Because we are imposing non-negative constraints to Ψ , the absolute operator can be removed from the
 571 objective function in Eq (4) and can be rewritten as

572
$$\min \sum_{ji} (y_{ji} - y_{ji}^0 - \sum_k \psi_{jk} \beta_{ki})^2 + \lambda \sum_{jk} \psi_{jk}. \quad \text{Eq (4)}$$

573

574 This becomes a constrained quadratic optimization problem where the best solution can be easily
 575 achieved using methods like Newton's method. The final, normalized PS is to scale values of ψ_{ik} to
 576 $[0,1]$:

577

578
$$PS_{ik} = \psi_{ik}/U.$$

579

580 We implemented this framework as part of the scMAGeCK pipeline¹⁸. The PS source code,
 581 documentation and tutorials can be found on Github: <https://github.com/davidliwei/PS>

582

583 **Simulated datasets**

582 The eight simulated datasets are generated by the simulator scDesign3⁴⁶ with modifications for Perturb-
583 seq. The simulation utilizes scDesign3's parametric model to capture the characteristics of the user-
584 inputted reference data, specify the desired ground truth, and simulate synthetic cells via sampling from
585 the model (to be detailed in **Steps 1-4** below). The reference data is the real scRNA-seq dataset with the
586 gene *Nelfb* perturbed in some mouse T cells⁴⁷; the cells with *Nelfb* perturbed are referred to as *knockout*
587 *cells*, and the cells with *Nelfb* unperturbed serve as the negative control and are referred to as *wild-type*
588 *cells*. Based on the same reference data, the eight simulated datasets are generated under eight different
589 settings. Each setting corresponds to a combination of two simulation parameters' values: the number of
590 *Nelfb*'s downstream genes (i.e., the genes whose expression levels are affected by *Nelfb*'s knockout;
591 with candidate values 0, 10, 200, and 500) and the perturbation efficiency (with candidate values 50%
592 and 100%). The candidate downstream genes of *Nelfb* are the top differentially expressed (DE) genes
593 identified from the bulk RNA-seq data of the same biological sample (from the second sheet in the
594 Excel file from Wu et al.'s Supplementary Data 1⁴⁸). Thus, we have $4 \times 2 = 8$ simulated datasets in total.
595

596 Before running the simulation, we pre-process the scRNA-seq dataset and the bulk DE gene rank list.

- 597 1. First, we perform the same quality control as in the dataset's original publication⁴⁹. Specifically,
598 cells are retained only if their numbers of detected genes are between 1,000 and 5,000, and their
599 UMI counts have less than 12% mitochondrial counts.
- 600 2. Second, we impute and amplify the gene-by-cell count matrix of the wild-type mouse cells to
601 enhance the perturbation effects in the simulated data. Specifically, we first impute the wild-type
602 count matrix using scImpute⁵⁰ (default version 0.0.9) to reduce the sparsity. Then we multiply
603 the imputed count matrix by an amplification factor of 10 to increase the range of gene
604 expression levels.
- 605 3. Third, we construct a gene-by-cell count matrix by combining the wild-type cells in the post-
606 imputation-and-amplification wild-type count matrix and the knockout cells in the knockout
607 count matrix. By the end of this step, the dimension of this combined matrix is $(P+1) \times N$, with
608 rows corresponding to $P+1$ genes (*Nelfb* and P other genes) and columns corresponding to N
609 cells, which consist of N^{wt} wild-type cells and N^{ko} knockout cells.
- 610 4. Fourth, we extract the row corresponding to *Nelfb* as a vector, which contains *Nelfb*'s counts in
611 all cells (an N -dimensional vector denoted as C , where C_j is *Nelfb*'s count in cell j), and we
612 denote the remaining P rows as a $P \times N$ matrix \mathbf{Y} , where Y_{ij} is gene i 's count in cell j .
- 613 5. Fifth, using \mathbf{Y} , we refine the list of bulk DE genes by excluding the DE genes that correspond to
614 zero rows in \mathbf{Y} or do not correspond to any rows in \mathbf{Y} .
- 615 6. Lastly, to reduce the computation time for data simulation, we use the scran package⁵¹ to select
616 3,000 highly variable genes in \mathbf{Y} . We only keep the union of these 3,000 highly variable genes
617 and the refined bulk DE genes as the rows in \mathbf{Y} . The number of the kept genes is 3,390, so the
618 dimension of \mathbf{Y} is $3,390 \times N$.

620 Additionally, we know which cells have *Nelfb* perturbed; thus, we have another N -dimensional binary
621 vector denoted as K , where K_j indicates whether the j -th cell has *Nelfb* perturbed or not; that is, $K_j = 0$
622 means the j -th cell is a wild-type cell, and $K_j = 1$ means the j -th cell is a knockout cell. K and C are used
623 as two covariate vectors, and \mathbf{Y} is used as the reference count matrix for scDesign3. Finally, we modify
624 scDesign3 by using \mathbf{Y} , C , K , the refined DE genes, the number of *Nelfb*'s downstream genes, and the
625 perturbation efficiency to simulate data in the following four steps:
626

627 **Step 1: modeling each gene's marginal distribution independently.** For each gene i , if it is a downstream
 628 gene of Nelfb, we assume that Y_{ij} , conditional on C_j , follows a zero-inflated negative binomial (ZINB)
 629 distribution with the mean parameter μ_{ij} , the dispersion parameter ϕ_i , and the zero-inflation probability
 630 parameter ν_{ij} . Otherwise, if gene i is not a downstream gene of Nelfb, we assume that Y_{ij} follows a
 631 ZINB distribution with the mean parameter μ_i , the dispersion parameter ϕ_i , and the zero-inflation
 632 probability parameter ν_i . This marginal distribution for each gene is specified by a generalized additive
 633 model for location, scale, and shape (GAMLSS). Without loss of generality, we define the first D genes
 634 in \mathbf{Y} to be the top D DE genes in the refined DE gene list ($D \in \{0, 10, 200, 500\}$); we treat these top D
 635 DE genes as the D downstream genes of Nelfb. Then we modify scDesign3's original code
 636 implementation so Nelfb's downstream genes and non-downstream genes have different marginal
 637 distributions: a downstream gene's marginal distribution in each cell j depends on C_j , Nelfb's count in
 638 cell j ; a non-downstream gene's marginal distribution in each cell j is irrelevant to C_j .

639 For Nelfb's downstream gene $i = 1, \dots, D$:

$$640 \begin{cases} Y_{ij} | C_j \sim \text{ZINB}(\mu_{ij}, \phi_i, \nu_{ij}) \\ \log(\mu_{ij}) = \alpha_i + \beta_i \times C_j \\ \log(\phi_i) = \omega_i \\ \text{logit}(\nu_{ij}) = \gamma_i + \eta_i \times C_j \end{cases} .$$

641

642 For Nelfb's non-downstream gene $i = D + 1, \dots, P$:

$$643 \begin{cases} Y_{ij} \sim \text{ZINB}(\mu_i, \phi_i, \nu_i) \\ \log(\mu_i) = \alpha_i \\ \log(\phi_i) = \omega_i \\ \text{logit}(\nu_i) = \gamma_i \end{cases} .$$

644

645 After parameter estimation by the R package gamlss, the fitted distribution of $Y_{ij} | C_j$, for $i = 1, \dots, D$, is
 646 denoted as $\text{ZINB}(\hat{\mu}_{ij}, \hat{\phi}_i, \hat{\nu}_{ij})$ with the CDF \hat{F}_{ij} ; the fitted distribution of Y_{ij} , for $i = D + 1, \dots, P$, is
 647 denoted as $\text{ZINB}(\hat{\mu}_i, \hat{\phi}_i, \hat{\nu}_i)$ with the CDF \hat{F}_i . The other parameters including α_i , β_i , γ_i , and η_i are
 648 estimated as $\hat{\alpha}_i$, $\hat{\beta}_i$, $\hat{\gamma}_i$, and $\hat{\eta}_i$ for each i respectively.

649
 650
 651 **Step 2: modeling genes' joint distribution using the Gaussian copula.** To approximate the pairwise gene-
 652 gene correlations in the reference dataset, scDesign3 utilizes a multivariate statistical technique, the
 653 Gaussian copula. Given each gene's marginal distribution fitted in Step 1, scDesign3 approximates the
 654 multivariate joint distribution of the P genes in cell j as

$$655 (\Phi^{-1}(\hat{F}_{1j}(Y_{1j})), \dots, \Phi^{-1}(\hat{F}_{Dj}(Y_{Dj})), \Phi^{-1}(\hat{F}_{D+1}(Y_{(D+1)j})), \dots, \Phi^{-1}(\hat{F}_P(Y_{Pj}))) \sim N(\mathbf{0}, \hat{\mathbf{R}}(K_j)) ,$$

656 where $\Phi^{-1}(\cdot)$ denotes the inverse of the cumulative distribution function (CDF) of the standard
 657 Gaussian distribution, $\mathbf{0}$ is the P -dimensional zero vector, and $\hat{\mathbf{R}}(K_j)$ is the estimated $P \times P$ gene-gene
 658 correlation matrix of the Gaussian copula conditional on the value of K_j . Specifically, since K_j is binary,
 659 we have two estimated gene-gene correlation matrices, one for the wild-type cells ($K_j = 0$) and the other
 660 for the knockout cells ($K_j = 1$). For $\hat{F}_{1j}(Y_{1j}), \dots, \hat{F}_{Dj}(Y_{Dj}), \hat{F}_{D+1}(Y_{(D+1)j}), \dots, \hat{F}_P(Y_{Pj})$, a technique called
 661 "distributional transform" is used to make the CDFs continuous; see Sun et al.⁵² for a detailed
 662 explanation.

663

664 Step 3: modifying the fitted parameters. Since we want to generate synthetic datasets with two
 665 perturbation efficiencies, for each downstream gene $i = 1, \dots, D$, we modify the mean parameters for all
 666 downstream genes in the knockout cells to reflect the user-specified perturbation efficiency. Without
 667 loss of generality, we assume the first $N^{\text{ko}} = \sum_{j=1}^N I(K_j = 1)$ of the N cells as the knockout cells. Then,
 668 we update the mean parameters for Nelfb's D downstream genes in the N^{ko} knockout cells (i.e., $\hat{\mu}_{ij}$ for
 669 $i \in \{1, \dots, D\}, j \in \{1, \dots, N^{\text{ko}}\}$) based on the user-specified perturbation efficiency as follows.

670
 671 For the 50% perturbation efficiency: We randomly sample N^{ko} values from $\{C_j, j \in \{N^{\text{ko}} + 1, \dots, N\}\}$
 672 (i.e., Nelfb's counts in the wild-type cells) and multiply the sampled C_j values by 0.5 to represent the
 673 50% perturbation efficiency. We store these sampled and scaled values by $C^* = (C_1^*, \dots, C_{N^{\text{ko}}}^*)^T$ as
 674 Nelfb's counts in the N^{ko} synthetic knockout cells to be simulated. Then, we modify the mean
 675 parameters for the D downstream genes in the N^{ko} synthetic knockout cells (for $i \in \{1, \dots, D\}, j \in$
 676 $\{1, \dots, N^{\text{ko}}\}$) as

$$\hat{\mu}_{ij} = \hat{\alpha}_i + \hat{\beta}_i \cdot C_j^* .$$

677
 678
 679 For the 100% perturbation efficiency: C^* becomes a zero vector with length N^{ko} , and we modify $\hat{\mu}_{ij}$ for
 680 $i \in \{1, \dots, D\}, j \in \{1, \dots, N^{\text{ko}}\}$ in the same way as above.

681
 682 We do not change any estimated mean parameters for the D downstream genes in the N^{wt} wild-type
 683 cells, any estimated mean parameters for the non-downstream ($P - D$) genes in all N cells, any
 684 estimated dispersion parameters, or any estimated zero-inflation probability parameters.

685
 686 Moreover, we use S to denote an N -dimensional vector representing Nelfb's counts in the N synthetic
 687 cells, with the counts in the first N^{ko} synthetic knockout cells set above based on the perturbation
 688 efficiency, and the counts in the last N^{wt} synthetic wild-type cells same as those in the real N^{wt} wild-
 689 type cells. That is, $S_j = C_j^*$ for $j \in \{1, \dots, N^{\text{ko}}\}$, and $S_j = C_j$ for $j \in \{N^{\text{ko}} + 1, \dots, N\}$.

690
 691 Step 4: generating synthetic data with the fitted model and modified parameters. First, we independently
 692 sample N^{wt} Gaussian vectors of length P from the estimated P -dimensional multivariate Gaussian
 693 distribution $N(0, \hat{\mathbf{R}}(K_j = 0))$ and N^{ko} Gaussian vectors of length P from the estimated P -dimensional
 694 multivariate Gaussian distribution $N(0, \hat{\mathbf{R}}(K_j = 1))$. Together, we stack these $N = N^{\text{wt}} + N^{\text{ko}}$ vectors
 695 $(\tilde{Z}_{11}, \dots, \tilde{Z}_{P1})^T, \dots, (\tilde{Z}_{1N}, \dots, \tilde{Z}_{PN})^T$ by row into a $P \times N$ Gaussian matrix $\tilde{\mathbf{Z}}$.

696
 697 Given the parameter estimates (modified or not) from Step 3, we convert the $P \times N$ Gaussian
 698 matrix $\tilde{\mathbf{Z}}$ into a $P \times N$ ZINB count matrix $\tilde{\mathbf{Y}}$ as

$$\tilde{\mathbf{Y}} = \begin{bmatrix} \tilde{Y}_1 = (\hat{F}_{11}^{-1}(\tilde{Z}_{11} | \hat{\mu}_{11}, \hat{\phi}_1, \hat{v}_{11}), \dots, \hat{F}_{1N}^{-1}(\tilde{Z}_{1N} | \hat{\mu}_{1N}, \hat{\phi}_1, \hat{v}_{1N})) \\ \vdots \\ \tilde{Y}_D = (\hat{F}_D^{-1}(\tilde{Z}_{D1} | \hat{\mu}_{D1}, \hat{\phi}_D, \hat{v}_{D1}), \dots, \hat{F}_{DN}^{-1}(\tilde{Z}_{DN} | \hat{\mu}_{DN}, \hat{\phi}_D, \hat{v}_{DN})) \\ \tilde{Y}_{D+1} = (\hat{F}_{D+1}^{-1}(\tilde{Z}_{(D+1)1} | \hat{\mu}_{D+1}, \hat{\phi}_{D+1}, \hat{v}_{D+1}), \dots, \hat{F}_{D+1}^{-1}(\tilde{Z}_{(D+1)N} | \hat{\mu}_{D+1}, \hat{\phi}_{D+1}, \hat{v}_{D+1})) \\ \vdots \\ \tilde{Y}_P = (\hat{F}_P^{-1}(\tilde{Z}_{P1} | \hat{\mu}_P, \hat{\phi}_P, \hat{v}_P), \dots, \hat{F}_{PN}^{-1}(\tilde{Z}_{PN} | \hat{\mu}_P, \hat{\phi}_P, \hat{v}_P)) \end{bmatrix} .$$

700
 701 Lastly, we combine \tilde{Y} with S by row into a $(P+1) \times N$ matrix, obtaining the final $(P+1) \times N$ synthetic
 702 count matrix $\begin{pmatrix} \tilde{Y} \\ S \end{pmatrix}$.

703
 704 **Genome-scale Perturb-seq on Jurkat cells**

705 Perturb-seq. We performed genome-scale Perturb-seq on Jurkat E6 cell line expressing dCas9-KRAB as
 706 our model to study. We transduced them with a genome-wide CRISPRi CROP-seq library at a high
 707 MOI. After infection, we split the cells into two populations, including untreated cells and activated cells
 708 (cells treated with anti-TCR and anti-CD28 antibodies for approximately 24 hours to stimulate TCR
 709 signaling). Cells were then labelled with cell hashing antibodies. Multiple labels were used for the
 710 activated population to help with cell multiplet detection. Cells were loaded on 16 channels of a 10x
 711 Chromium X instrument. We loaded 115 000 cells per channel, and the expected recovery rate was 60
 712 000 cells per channel, including 24% multiples. Samples were pooled unequally before they were loaded
 713 on the ChromiumX: 10% untreated cells, 90% treated cells. A sequencing library was prepared using
 714 3'Chemistry with a targeted primer panel: custom multiplex PCR step to enrich for specific transcripts.
 715 Libraries were sequenced on NovaSeq S4 PE100 in asymmetric read mode (R1: 28 cycles; R2: 172
 716 cycles), with PhiX concentration of 1%. The expected coverage is around 9 000 ~ 10 000 input reads per
 717 cell.

718
 719 Hash oligos.

oligo	condition
CGGCTCGTGCTGCGTCGTCTCAAGTCCAGAACTCCGTGTATCCT	untreated
CTCCCTGGTGTCAATACCCGATGTGGTGGGCAGAATGTGGCTGG	activated
TTACCCGCAGGAAGACGTATACCCCTCGTGCCAGGCGACCAATGC	activated
TGTCTACGTCGGACCGCAAGAAGTGAGTCAGAGGCTGCACGCTGT	activated
CCCCACCAGGTTGCTTTGTCGGACGAGCCCGCACAGCGCTAGGAT	activated
GTGATCCGCGCAGGCACACATACCGACTCAGATGGGTTGTCCAGG	activated
GCAGCCGGCGTTCGTACGAGGCACAGCGGAGACTAGATGAGGCCCC	activated

720
 721 sgRNA library design. The genome-wide CRISPRi sgRNA library was designed to target the
 722 transcription start site (TSS) coordinates, calculated from publicly available FANTOM CAGE peaks
 723 data. In total, 18 595 genes were targeted, with 4 sgRNAs per gene. On top of that, we designed another
 724 CRISPRi library targeting 3220 genes with 4 sgRNAs per gene. This library was designed using Jurkat-
 725 specific TSS, which were calculated from public Jurkat CAGE-seq datasets. Both libraries were
 726 combined into a final library targeting 3220 genes with 8 sgRNA/gene and 15 375 genes with 4
 727 sgRNA/gene.

728
 729 Targeted primer panel. The primer panel for targeted transcriptomic readout consisted of 374 target
 730 genes from several categories:

Source	Number of targets (genes)
bulk RNA-seq DEG, 50 top up/downregulated	100
T-cell related genes	51
KEGG TCR signaling pathway	35
RNA-binding proteins	27

Replogle*: unfolded protein response	51
Replogle: proteasome	41
Replogle: NFkB	20
Replogle: cell cycle	29
Replogle: targets of nonsense-mediated decay	5
Controls (cell cycle, mitochondrial, Cas9)	15
Total number of targets	374

732

733 * “Replogle” prefix means that this category of targets was derived from the published genome-scale
734 Perturb-seq dataset⁵⁰.

735

736 **Data preprocessing.** For each of the 16 channels, all 3 kinds of sequencing libraries (mRNA, sgRNA,
737 cell hashing) were indexed using the same Illumina index sequence. We obtained high-quality scRNA-
738 seq data of over 586,000 single cells after quality control, with a median 13 guides detected per cell. We
739 obtained an average of 400 cells per gene perturbation. STAR and STAR solo 2.7.10a were used to map
740 transcriptomic reads against a custom gtf annotation, which was based on gencode.v34.annotation.gtf
741 (hg38). Reads that did not map to the transcriptome reference were then mapped & counted using STAR
742 solo against a custom fasta reference with guide sequences and a fasta reference with hash label
743 sequences. STAR Solo output transcriptome matrices were first filtered using an approach similar to 10x
744 cellranger EmptyDrops filtering, which retained cells with at least 10 % UMI count of the 99th percentile
745 UMI counts of the top expected cells number. Then, an initial Seurat object was created from those
746 filtered transcriptome matrices using the CreateSeuratObject function with following parameters:
747 min.cells = 5; min.features = 10; all other parameters at default values. Outlier cells were filtered out by
748 mitochondrial and mRNA content (percent.mt, nCount_RNA). In order to detect cell multiplets and
749 determine cell population (untreated or activated), cell labels (also known as hashes) were called using
750 the MULTI-seq approach (deMULTIplex::classifyCells in R). Only cells with exactly 1 known label
751 were kept. Then, sgRNA calling was conducted using a binomial test, with total sgRNA UMI counts
752 used to derive background frequencies. A threshold of 0.05 on Benjamini-Hochberg corrected p-values
753 (per channel) was used to generate the final calls. The sgRNA assays are sparse matrices containing 1,
754 where the respective cell is considered to be carrying the respective sgRNA and 0 elsewhere. Following
755 that, all results from steps above from 16 channels were merged together, and merged counts were
756 normalized using NormalizeData and scaled using ScaleData. Cell cycle scoring was performed via
757 CellCycleScoring, PCA was calculated using RunPCA, and UMAP was calculated on first 30 principal
758 components using RunUMAP in Seurat.

759

760 **HIV latency Perturb-seq**

761 We used a previously established cell line model of HIV latency⁴⁰. In this model, Jurkat cells were
762 infected with an HIV vector with GFP tied to the LTR promoter, resulting in a positive GFP signal as a
763 measurement of viral transcription reactivation and HIV latency reversal. These cells, which already
764 express Cas9, were transduced with a lenti-sgRNA library. The lenti-sgRNA library (MilliporeSigma;
765 LV14, U6-gRNA-10x:EF1a-Puro-2a-BFP) was designed to target 10 genes, with 3 gRNAs per gene. In
766 addition to non-targeting controls, the library contained five positive regulators (NFKB1, CCNT1,
767 PRKCA, TLR1, MAP3K14) and five negative regulators (NFKBIA, NELFE, HDAC2, BRD4, BIRC2)
768 of HIV transcription. Transduction was carried out on 850,000 cells at an MOI of 0.3 using 8ug/ml
769 polybrene in 2 ml of RPMI containing 10% FBS and 1% penicillin–streptomycin. The media was
770 replaced 24 hours later with fresh media without polybrene. Two days after transduction, the cells were

771 selected for using 1.5 ug/ml puromycin for 5 days. After selection, the cells were split evenly into three
772 groups. One-third of the cells were kept in culture with no drug added, and two-thirds of the cells were
773 stimulated with PMA/I (50ng/ml PMA in combination with 1 μ M Ionomycin). After 16 hours, the
774 stimulated cells were sorted into GFP+ and GFP- populations. All three samples were then analyzed
775 following the 10x Genomics single-cell sequencing protocol. Sequencing data, encompassing gene
776 expression and CRISPR guide capture libraries, underwent demultiplexing and processing using Cell
777 Ranger (version 6.1.2). The resulting feature-barcode matrices from three samples were then merged,
778 and subsequent analysis was carried out utilizing the Seurat R package (version 4.3.1). To ensure data
779 quality, cells were excluded if the number of expressed genes was greater than 7,500 or fewer than 200.
780 Additionally, cells were removed if the percentage of mitochondrial reads exceeded 15%. Single cells
781 harboring more than one detected sgRNA sequence, attributable to either multiple sgRNA transductions
782 or the presence of multiple cells in a single-cell droplet, were also excluded from the analysis. Following
783 quality control measures, merged counts underwent normalization and scaling. PCA was computed
784 based on the top 2,000 highly variable genes. Subsequently, clustering and UMAP embeddings were
785 performed using default parameters. To gain further insights into the biological significance of the
786 obtained clusters, enrichment analysis was conducted utilizing Enrichr (PMID: 27141961).

787 788 **Pancreatic differentiation clones and pooled single-cell RNA-seq**

789 Culture of hESC. Generation of KO hESCs was described in published studies, including HHEX KO H1
790 and HUES8 cell lines⁵³, FOXA1 KO HUES8 cell lines⁵⁴, OTUD5 KO HUES8 cell lines, and CCDC6
791 KO H1 cell lines⁴¹. Cells were regularly confirmed to be mycoplasma-free by the Memorial Sloan
792 Kettering Cancer Center (MSKCC) Antibody & Bioresource Core Facility. KO and WT hESCs were
793 maintained in Essential 8 (E8) medium (Thermo Fisher Scientific, A1517001) on vitronectin (Thermo
794 Fisher Scientific, A14700) pre-coated plates at 37 °C with 5% CO₂. The Rho-associated protein kinase
795 (ROCK) inhibitor Y-276325 (5 μ M; Selleck Chemicals, S1049) was added to the E8 medium the first
796 day after passaging or thawing of hESCs.

797
798 hESC-directed pancreatic differentiation. hESCs were seeded at a density of 2.3×10^5 cells/cm² on
799 vitronectin-coated plates in E8 medium with 10 μ M Y-27632. After 24 hours, cells were washed with
800 PBS and differentiated to DE (stage 1), primitive gut tube (stage 2), PP1 (stage 3) and PP2 (stage 4)
801 stages following previously described 4-stage protocol⁴⁰. In brief, stage 1 (3 d): S1/2 medium
802 supplemented with 100 ng ml⁻¹ Activin A (Bon Opus Biosciences) and 5 μ M CHIR99021 (04-0004-10,
803 Stemgent) for 1 d. S1/2 medium supplemented with 100 ng ml⁻¹ Activin A for the next 2 d. Stage 2 (2 d):
804 S1/2 medium supplemented with 50 ng ml⁻¹ KGF (AF-100-19, PeproTech) and 0.25 mM vitamin C
805 (VitC) (Sigma-Aldrich, A4544). Stage 3 (2 d): S3/4 medium supplemented with 50 ng ml⁻¹ KGF,
806 0.25 mM VitC and 1 μ M retinoic acid (R2625, MilliporeSigma). Stage 4 (4 d): S3/4 medium
807 supplemented with 50 ng ml⁻¹ KGF, 0.1 μ M retinoic acid, 200 nM LDN (Stemgent, 04-0019), 0.25 μ M
808 SANT-1 (Sigma, S4572), 0.25 mM VitC and 200 nM TPB (EMD Millipore, 565740). The base
809 differentiation medium formulations used in each stage were as follows. S1/2 medium: 500 ml MCDB
810 131 (15-100-CV, Cellgro) supplemented with 2 ml 45% glucose (G7528, MilliporeSigma), 0.75 g
811 sodium bicarbonate (S5761, MilliporeSigma), 2.5 g BSA (68700, Proliant), 5 ml GlutaMAX (35050079,
812 Invitrogen). S3/4 medium: 500 ml MCDB 131 supplemented with 0.52 ml 45% glucose, 0.875 g sodium
813 bicarbonate, 10 g BSA, 2.5 ml ITS-X, 5 ml GlutaMAX.

814
815 Cell infection with LARRY barcode virus. Individual LARRY barcode constructs were cloned from the
816 LARRY barcode library (Addgene:140024) and transfected to 293T cells to generate lentivirus. Next,

817 each KO and WT hESC clone was infected with a unique LARRY barcode at low MOI. One week after
818 lentiviral infection, the barcoded cells, which expressed GFP, were sorted out and cultured in E8
819 medium as described in previous section.

820
821 Pooled single-cell RNA-seq. One day before differentiation, each of 10 hESC barcoded clones were
822 counted, mixed at the same cell number ratio, and then seeded at a density of 2.3×10^5 cells/cm² onto a
823 12-well cell culture plate. At DE and PP2 stages, pooled differentiating cells were dissociated into single
824 cell suspension by TrypLE Select for 5 min at 37 °C. Cells were then stored in BAMBANKER™
825 freezing medium for future experiments. For scRNA-seq, frozen cells were thawed and sorted to collect
826 live GFP+ cells. Cellular suspensions were then loaded on a Chromium Controller following the
827 manufacturer's instructions (10x Genomics Chromium Single Cell 3' Reagent Kit v3.1 User Guide).
828 cDNA libraries and targeted LARRY barcode libraries were generated separately using 10ul cDNA
829 each. cDNA libraries were made under manufacturer's instructions and targeted LARRY barcode
830 libraries were amplified using specific primers (F: CTACACGACGCTCTTCCGATCT; R:
831 GTGACTGGAGTTCAGACGTGTGCTCTTCCGATCTtaaccgttgctaggagacataT).

832
833 Data analysis. The sequencing data which included transcriptome and LARRY barcode libraries,
834 underwent demultiplexing and processing via Cell Ranger (version 6.1.2). Subsequent analysis was
835 conducted using the Seurat R package (version 4.3.1). Quality control measures were implemented to
836 ensure robust data analysis. Cells were excluded if the number of expressed genes exceeded 7,000 or fell
837 below 200. Additionally, cells were removed if the percentage of mitochondrial read exceeded 20%.
838 Singlet cells were defined by considering the highest feature barcode count, ensuring it was at least
839 twice as large as the second highest feature barcode count. Single cells containing more than one
840 detected barcode sequence were excluded from the dataset. This process resulted in a final set of 20,678
841 cells for downstream analysis. After quality control measures, the count matrix underwent normalization
842 and scaling. PCA was performed using the top 2,000 highly variable genes. Subsequently, clustering and
843 UMAP embeddings were generated using default parameters to elucidate the underlying structure and
844 relationships within the dataset.

845
846 Flow cytometry. Cells were dissociated using TrypLE Select and resuspended in FACS buffer (5% FBS
847 in PBS). Live/Dead Fixable Violet cell stain (Invitrogen, L34955) was used to discriminate dead cells
848 from live cells. Permeabilization/fixation was performed at room temperature for 1 h. Antibody staining
849 was performed in permeabilization buffer. Antibodies for this study include HNF4A, Novus Biologicals,
850 NBP2-67679, 1:200; PDX1, R&D Systems, AF2419, 1:500, Donkey anti-Rabbit IgG (H+L) Highly
851 Cross-Adsorbed Secondary Antibody, Thermo Fisher Scientific, 1:500; Donkey anti-goat IgG (H+L)
852 Highly Cross-Adsorbed Secondary Antibody, Thermo Fisher Scientific, 1:500. Cells were then analysed
853 using BD LSRFortessa. Flow cytometry analysis and figures were generated using FlowJo v.10.

854 855 **Acknowledgements**

856 The authors thank all members from the Li and Huangfu laboratory for comments and discussions. The
857 authors thank Jake P. Taylor-King for discussions. This study is supported by NIH R01 HG010753,
858 HL168174 (to W.L., B.S., L.C.), District of Columbia Center for AIDS (DC-CFAR) Research
859 Transitioning Investigator Award (AI117970, to W.L.), and startup support from the Center for Genetic
860 Medicine Research at Children's National Hospital. D.H is supported by NIH UM1 HG012654, U01
861 HG012051. J.J. Li is supported by National Science Foundation DBI-1846216 and DMS-2113754, NIH
862 R35 GM140888, Johnson and Johnson WiSTEM2D Award, Sloan Research Fellowship, UCLA David

863 Geffen School of Medicine W.M. Keck Foundation Junior Faculty Award, and Chan-Zuckerberg
864 Initiative Single-Cell Biology Data Insights [Silicon Valley Community Foundation Grant Number:
865 2022-249355]. W.D. and R.F.S. is supported by the Howard Hughes Medical Institute.

866

867 **Author contributions**

868 W.L. conceived the project. W.L. and B.S. developed the method. W.L., B.S., and D.L. designed and
869 performed the experiments and analyzed the data. W.D., N.M. and B.S. performed and analyzed HIV
870 Perturb-seq under the supervision of J.M.S., R.S. and W.L. B.S., Q.W. and D.S. performed synthetic
871 experiments under the supervision of W.L. and J.J.L. D.L., D.Y., B.W. and B.R. generated pancreatic
872 differentiation dataset and performed validations under the supervision of D.H. A.K., A.V., N.U., and
873 A.L. generated and analyzed genome-scale Perturb-seq under the supervision of T.B. X.C., L.C. and
874 Y.D. performed the analysis and interpretation of the results. W.L., and B.S. wrote the manuscript with
875 input from all the authors. W.L. T.B., J.J.L., R.S. and D.H. supervised the study.

876

877 **Competing interests**

878 T.B. is a co-founder and Managing Director of Myllia Biotechnology. A.K., A.V., N.U. and A.L. are
879 employees of Myllia Biotechnology. Other authors declare that they have no competing interest.

880

881 **Data and materials availability**

882 The Perturb-seq scRNA-seq data have been deposited to Gene Expression Omnibus (GEO) under the
883 accession number GSE247601. The source code of the PS method, and the documentation and demos
884 are available on GitHub: (<https://github.com/davidliwei/PS>).

885

886

887

888

889

890

891 **References**

892

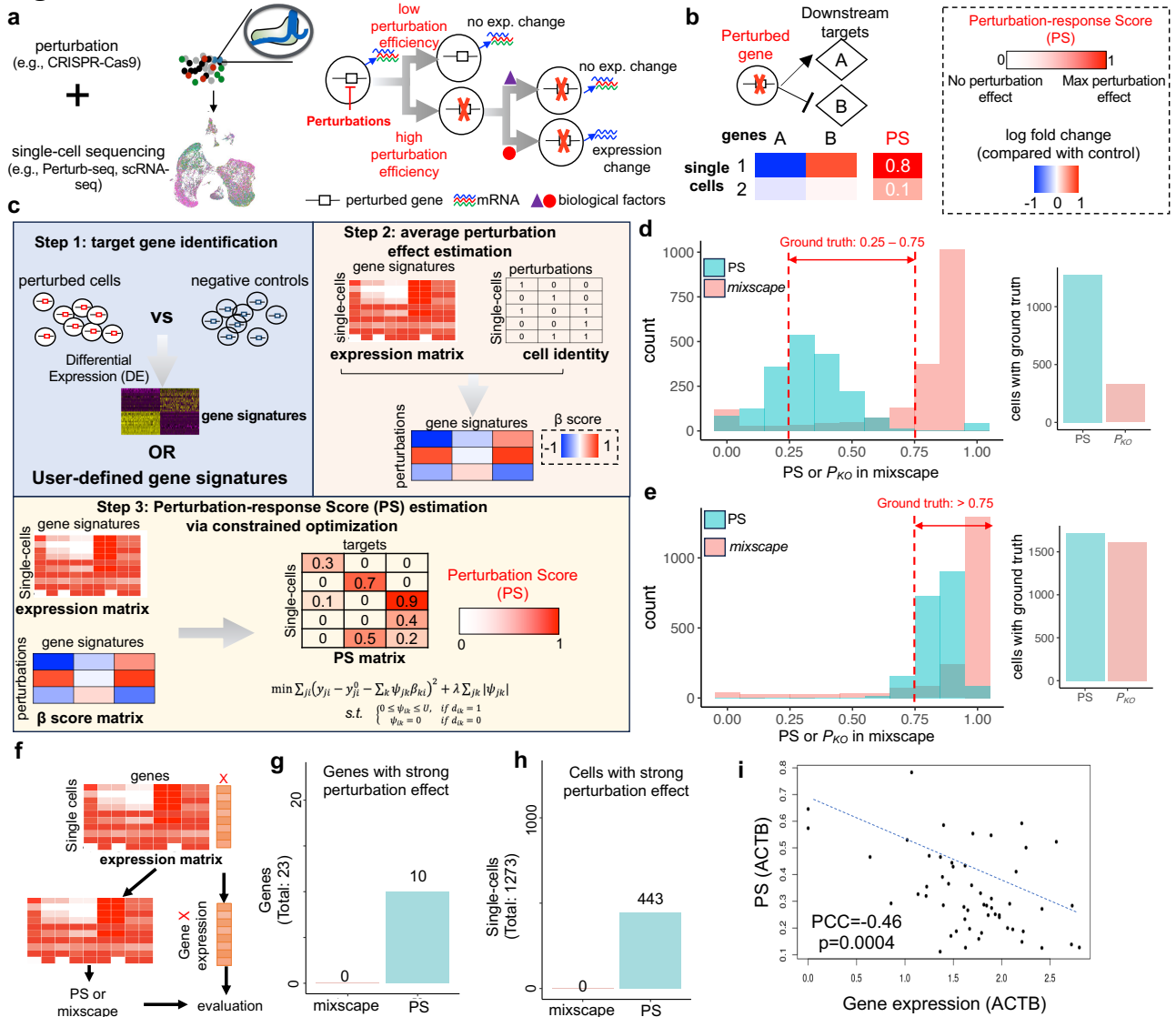
- 893 1 High-content CRISPR screening. *Nat Rev Methods Primers* 2022;**2**:.
894 <https://doi.org/10.1038/s43586-022-00098-7>.
- 895 2 Srivatsan SR, McFaline-Figueroa JL, Ramani V, Saunders L, Cao J, Packer J, *et al*. Massively
896 multiplex chemical transcriptomics at single-cell resolution. *Science* 2020;**367**:45–51.
- 897 3 Dixit A, Parnas O, Li B, Chen J, Fulco CP, Jerby-Arnon L, *et al*. Perturb-Seq: Dissecting
898 Molecular Circuits with Scalable Single-Cell RNA Profiling of Pooled Genetic Screens. *Cell*
899 2016;**167**:1853-1866.e17.
- 900 4 Adamson B, Norman TM, Jost M, Cho MY, Nuñez JK, Chen Y, *et al*. A Multiplexed Single-Cell
901 CRISPR Screening Platform Enables Systematic Dissection of the Unfolded Protein Response. *Cell*
902 2016;**167**:1867-1882.e21.
- 903 5 Jaitin DA, Weiner A, Yofe I, Lara-Astiaso D, Keren-Shaul H, David E, *et al*. Dissecting Immune
904 Circuits by Linking CRISPR-Pooled Screens with Single-Cell RNA-Seq. *Cell* 2016;**167**:1883-
905 1896.e15.
- 906 6 Datlinger P, Rendeiro AF, Schmidl C, Krausgruber T, Traxler P, Klughammer J, *et al*. Pooled
907 CRISPR screening with single-cell transcriptome readout. *Nat Methods* 2017;**14**:297–301.

- 908 7 Xie S, Duan J, Li B, Zhou P, Hon GC. Multiplexed Engineering and Analysis of Combinatorial
909 Enhancer Activity in Single Cells. *Mol Cell* 2017;**66**:285-299.e5.
- 910 8 Rubin AJ, Parker KR, Satpathy AT, Qi Y, Wu B, Ong AJ, *et al.* Coupled Single-Cell CRISPR
911 Screening and Epigenomic Profiling Reveals Causal Gene Regulatory Networks. *Cell*
912 2019;**176**:361-376.e17.
- 913 9 Liscovitch-Brauer N, Montalbano A, Deng J, Méndez-Mancilla A, Wessels H-H, Moss NG, *et al.*
914 Profiling the genetic determinants of chromatin accessibility with scalable single-cell CRISPR
915 screens. *Nat Biotechnol* 2021;**39**:1270–7.
- 916 10 Dhainaut M, Rose SA, Akturk G, Wroblewska A, Nielsen SR, Park ES, *et al.* Spatial CRISPR
917 genomics identifies regulators of the tumor microenvironment. *Cell* 2022;**185**:1223-1239.e20.
- 918 11 Norman TM, Horlbeck MA, Replogle JM, Ge AY, Xu A, Jost M, *et al.* Exploring genetic
919 interaction manifolds constructed from rich single-cell phenotypes. *Science* 2019;**365**:786–93.
- 920 12 Replogle JM, Norman TM, Xu A, Hussmann JA, Chen J, Cogan JZ, *et al.* Combinatorial single-cell
921 CRISPR screens by direct guide RNA capture and targeted sequencing. *Nat Biotechnol*
922 2020;**38**:954–61.
- 923 13 Wessels H-H, Méndez-Mancilla A, Hao Y, Papalexi E, Mauck WM 3rd, Lu L, *et al.* Efficient
924 combinatorial targeting of RNA transcripts in single cells with Cas13 RNA Perturb-seq. *Nat*
925 *Methods* 2023;**20**:86–94.
- 926 14 Hill AJ, McFaline-Figueroa JL, Starita LM, Gasperini MJ, Matreyek KA, Packer J, *et al.* On the
927 design of CRISPR-based single-cell molecular screens. *Nat Methods* 2018;**15**:271–4.
- 928 15 Yang L, Zhu Y, Yu H, Cheng X, Chen S, Chu Y, *et al.* scMAGeCK links genotypes with multiple
929 phenotypes in single-cell CRISPR screens. *Genome Biol* 2020;**21**:19.
- 930 16 Papalexi E, Mimitou EP, Butler AW, Foster S, Bracken B, Mauck WM 3rd, *et al.* Characterizing
931 the molecular regulation of inhibitory immune checkpoints with multimodal single-cell screens.
932 *Nat Genet* 2021;**53**:322–31.
- 933 17 Tsuchida CA, Brandes N, Bueno R, Trinidad M, Mazumder T, Yu B, *et al.* Mitigation of
934 chromosome loss in clinical CRISPR-Cas9-engineered T cells. *Cell* 2023;**186**:4567-4582.e20.
- 935 18 Song D, Wang Q, Yan G, Liu T, Sun T, Li JJ. scDesign3 generates realistic in silico data for
936 multimodal single-cell and spatial omics. *Nat Biotechnol* 2023. [https://doi.org/10.1038/s41587-023-](https://doi.org/10.1038/s41587-023-01772-1)
937 [01772-1](https://doi.org/10.1038/s41587-023-01772-1).
- 938 19 Wu B, Zhang X, Chiang H-C, Pan H, Yuan B, Mitra P, *et al.* RNA polymerase II pausing factor
939 NELF in CD8+ T cells promotes antitumor immunity. *Nat Commun* 2022;**13**:2155.
- 940 20 Gasperini M, Hill AJ, McFaline-Figueroa JL, Martin B, Kim S, Zhang MD, *et al.* A Genome-wide
941 Framework for Mapping Gene Regulation via Cellular Genetic Screens. *Cell* 2019;**176**:1516.
- 942 21 Jost M, Santos DA, Saunders RA, Horlbeck MA, Hawkins JS, Scaria SM, *et al.* Titrating gene
943 expression using libraries of systematically attenuated CRISPR guide RNAs. *Nat Biotechnol*
944 2020;**38**:355–64.
- 945 22 Schraivogel D, Gschwind AR, Milbank JH, Leonce DR, Jakob P, Mathur L, *et al.* Targeted
946 Perturb-seq enables genome-scale genetic screens in single cells. *Nat Methods* 2020;**17**:629–35.
- 947 23 Shifrut E, Carnevale J, Tobin V, Roth TL, Woo JM, Bui CT, *et al.* Genome-wide CRISPR screens
948 in primary human T cells reveal key regulators of immune function. *Cell* 2018;**175**:1958-1971.e15.
- 949 24 Wessels H-H, Stirn A, Méndez-Mancilla A, Kim EJ, Hart SK, Knowles DA, *et al.* Prediction of on-
950 target and off-target activity of CRISPR-Cas13d guide RNAs using deep learning. *Nat Biotechnol*
951 2023. <https://doi.org/10.1038/s41587-023-01830-8>.
- 952 25 Zhang J, Bu X, Wang H, Zhu Y, Geng Y, Nihira NT, *et al.* Cyclin D–CDK4 kinase destabilizes
953 PD-L1 via cullin 3–SPOP to control cancer immune surveillance. *Nature* 2018;**553**:91–5.

- 954 26 Naqvi S, Kim S, Hoskens H, Matthews HS, Spritz RA, Klein OD, *et al.* Precise modulation of
955 transcription factor levels identifies features underlying dosage sensitivity. *Nat Genet* 2023;**55**:841–
956 51.
- 957 27 Replogle JM, Saunders RA, Pogson AN, Hussmann JA, Lenail A, Guna A, *et al.* Mapping
958 information-rich genotype-phenotype landscapes with genome-scale Perturb-seq. *Cell*
959 2022;**185**:2559-2575.e28.
- 960 28 Radhakrishnan SK, Lee CS, Young P, Beskow A, Chan JY, Deshaies RJ. Transcription factor Nrfl
961 mediates the proteasome recovery pathway after proteasome inhibition in mammalian cells. *Mol*
962 *Cell* 2010;**38**:17–28.
- 963 29 Dai W, Wu F, McMyn N, Song B, Walker-Sperling VE, Varriale J, *et al.* Genome-wide CRISPR
964 screens identify combinations of candidate latency reversing agents for targeting the latent HIV-1
965 reservoir. *Sci Transl Med* 2022;**14**:eabh3351.
- 966 30 Yin M, Guo Y, Hu R, Cai WL, Li Y, Pei S, *et al.* Potent BRD4 inhibitor suppresses cancer cell-
967 macrophage interaction. *Nat Commun* 2020;**11**:1833.
- 968 31 Tan Y-F, Wang M, Chen Z-Y, Wang L, Liu X-H. Inhibition of BRD4 prevents proliferation and
969 epithelial-mesenchymal transition in renal cell carcinoma via NLRP3 inflammasome-induced
970 pyroptosis. *Cell Death Dis* 2020;**11**:239.
- 971 32 Shu S, Wu H-J, Ge JY, Zeid R, Harris IS, Jovanović B, *et al.* Synthetic Lethal and Resistance
972 Interactions with BET Bromodomain Inhibitors in Triple-Negative Breast Cancer. *Mol Cell*
973 2020;**78**:1096-1113.e8.
- 974 33 Li Z, Guo J, Wu Y, Zhou Q. The BET bromodomain inhibitor JQ1 activates HIV latency through
975 antagonizing Brd4 inhibition of Tat-transactivation. *Nucleic Acids Res* 2013;**41**:277–87.
- 976 34 Mbonye U, Kizito F, Karn J. New insights into transcription elongation control of HIV-1 latency
977 and rebound. *Trends Immunol* 2023;**44**:60–71.
- 978 35 Wei P, Garber ME, Fang S-M, Fischer WH, Jones KA. A novel CDK9-associated C-type cyclin
979 interacts directly with HIV-1 tat and mediates its high-affinity, loop-specific binding to TAR RNA.
980 *Cell* 1998;**92**:451–62.
- 981 36 Peng J, Zhu Y, Milton JT, Price DH. Identification of multiple cyclin subunits of human P-TEFb.
982 *Genes Dev* 1998;**12**:755–62.
- 983 37 Stoeckius M, Zheng S, Houck-Loomis B, Hao S, Yeung BZ, Mauck WM 3rd, *et al.* Cell Hashing
984 with barcoded antibodies enables multiplexing and doublet detection for single cell genomics.
985 *Genome Biol* 2018;**19**:224.
- 986 38 Weinreb C, Rodriguez-Fraticelli A, Camargo FD, Klein AM. Lineage tracing on transcriptional
987 landscapes links state to fate during differentiation. *Science* 2020;**367**:eaaw3381.
- 988 39 Vitak SA, Torkenczy KA, Rosenkrantz JL, Fields AJ, Christiansen L, Wong MH, *et al.* Sequencing
989 thousands of single-cell genomes with combinatorial indexing. *Nat Methods* 2017;**14**:302–8.
- 990 40 Yang D, Cho H, Tayyebi Z, Shukla A, Luo R, Dixon G, *et al.* CRISPR screening uncovers a central
991 requirement for HHEX in pancreatic lineage commitment and plasticity restriction. *Nat Cell Biol*
992 2022;**24**:1064–76.
- 993 41 Rosen BP, Li QV, Cho H, Liu D, Yang D, Graff S, *et al.* Parallel genome-scale CRISPR screens
994 distinguish pluripotency and self-renewal. *BioRxivorg* 2023.
995 <https://doi.org/10.1101/2023.05.03.539283>.
- 996 42 Li QV, Dixon G, Verma N, Rosen BP, Gordillo M, Luo R, *et al.* Genome-scale screens identify
997 JNK-JUN signaling as a barrier for pluripotency exit and endoderm differentiation. *Nat Genet*
998 2019;**51**:999–1010.

- 999 43 Thanasopoulou A, Stravopodis DJ, Dimas KS, Schwaller J, Anastasiadou E. Loss of CCDC6
1000 affects cell cycle through impaired intra-S-phase checkpoint control. *PLoS One* 2012;**7**:e31007.
- 1001 44 Morra F, Luise C, Merolla F, Poser I, Visconti R, Ilardi G, *et al.* FBXW7 and USP7 regulate
1002 CCDC6 turnover during the cell cycle and affect cancer drugs susceptibility in NSCLC. *Oncotarget*
1003 2015;**6**:12697–709.
- 1004 45 Zhou Y, Luo K, Liang L, Chen M, He X. A new Bayesian factor analysis method improves
1005 detection of genes and biological processes affected by perturbations in single-cell CRISPR
1006 screening. *Nat Methods* 2023. <https://doi.org/10.1038/s41592-023-02017-4>.
- 1007 46 Dong M, Wang B, Wei J, de O Fonseca AH, Perry CJ, Frey A, *et al.* Causal identification of single-
1008 cell experimental perturbation effects with CINEMA-OT. *Nat Methods* 2023;**20**:1769–79.
- 1009 47 Hart T, Brown KR, Sircoulomb F, Rottapel R, Moffat J. Measuring error rates in genomic
1010 perturbation screens: gold standards for human functional genomics. *Mol Syst Biol* 2014;**10**:733.
- 1011 48 Morgens DW, Deans RM, Li A, Bassik MC. Systematic comparison of CRISPR/Cas9 and RNAi
1012 screens for essential genes. *Nat Biotechnol* 2016;**34**:634–6.
- 1013 49 Bunne C, Stark SG, Gut G, Del Castillo JS, Levesque M, Lehmann K-V, *et al.* Learning single-cell
1014 perturbation responses using neural optimal transport. *Nat Methods* 2023;**20**:1759–68.
- 1015 50 Li WV, Li JJ. An accurate and robust imputation method scImpute for single-cell RNA-seq data.
1016 *Nat Commun* 2018;**9**:997.
- 1017 51 Lun ATL, Bach K, Marioni JC. Pooling across cells to normalize single-cell RNA sequencing data
1018 with many zero counts. *Genome Biol* 2016;**17**:75.
- 1019 52 Sun T, Song D, Li WV, Li JJ. scDesign2: a transparent simulator that generates high-fidelity
1020 single-cell gene expression count data with gene correlations captured. *Genome Biol* 2021;**22**:163.
- 1021 53 Hao Y, Hao S, Andersen-Nissen E, Mauck WM 3rd, Zheng S, Butler A, *et al.* Integrated analysis of
1022 multimodal single-cell data. *Cell* 2021;**184**:3573-3587.e29.
- 1023 54 Lee K, Cho H, Rickert RW, Li QV, Pulecio J, Leslie CS, *et al.* FOXA2 is required for enhancer
1024 priming during pancreatic differentiation. *Cell Rep* 2019;**28**:382-393.e7.
- 1025
1026
1027
1028
1029

Figures



1031

1032

1033

1034

1035

1036

1037

1038

1039

1040

1041

1042

1043

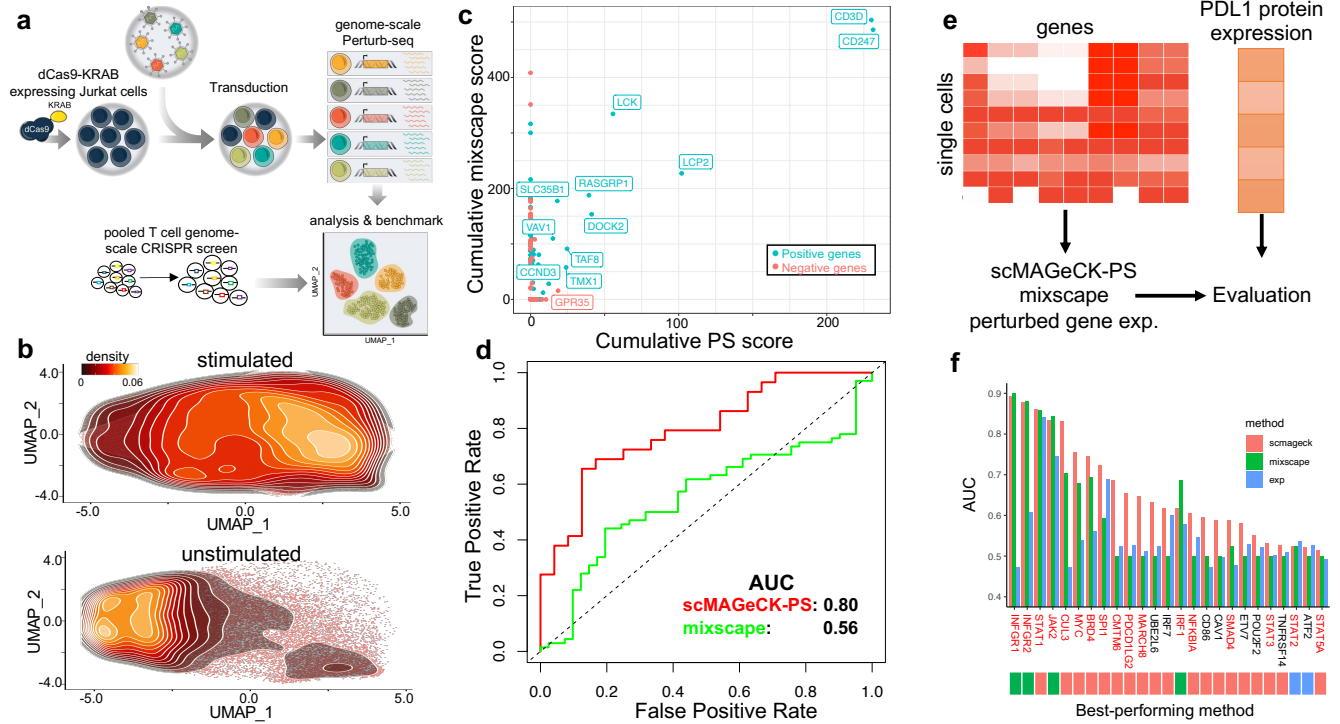
1044

1045

Figure 1. The Perturbation-response Score (PS) framework and benchmark. **a**, Overview of different technical and biological factors that contribute to heterogeneous perturbation outcomes from single-cell perturbation datasets. **b**, Using downstream gene expressions to infer the value of PSs. **c**, Overview of the scMAGeCK-PS that estimates PS value. **d-e**, Benchmark results of both PS and mixscape using simulated datasets, where 50% (**d**) and 100% (**e**) gene perturbation effects are simulated using scDesign3. Here, the expressions of 200 differentially expressed genes (DEGs) from bulk RNA-seq (Nelf knockout vs. wild-type) are simulated, and ground truth efficiency value is indicated in red color. **f**, Benchmark pipeline using real CRISPRi-based Perturb-seq datasets, where the perturbation efficiency can be evaluated directly via gene expression. **g-h**, Benchmark results of mixscape and scMAGeCK-PS using a published Perturb-seq dataset, by counting the numbers of cells or genes with strong perturbation effects. A gene is considered to have strong perturbation effect, if a strong negative correlation (Pearson correlation coefficient < -0.1) is observed between PS and the expression of that gene across all perturbed cells. A cell is considered to be strongly perturbed, if its predicted efficiency score (by scMAGeCK-PS or mixscape) within one cell is greater than 0.5. The Perturb-seq experiment

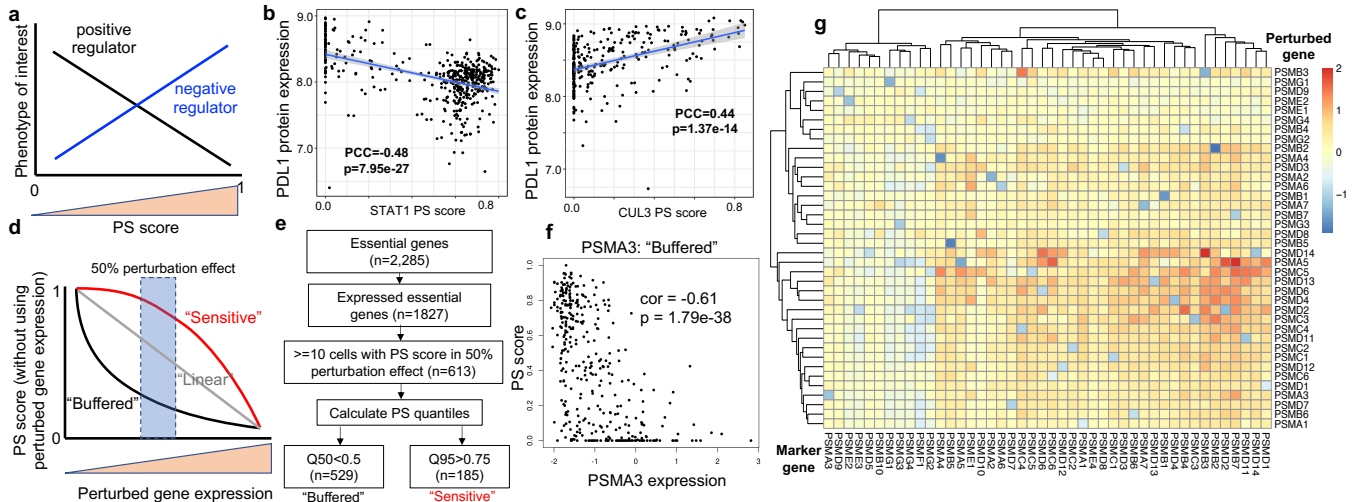
1046
1047
1048
1049
1050
1051

is performed with low MOI condition, where most cells have only 1 expressed guide. **i**, An representative estimation results of scMAGeCK-PS and their correlations of ACTB expression.



1052
1053
1054
1055
1056
1057
1058
1059
1060

Figure 2. Additional benchmark results using genome-scale Perturb-seq and ECCITE-seq. a, Benchmark procedure using a genome-scale Perturb-seq and a published, pooled T cell CRISPR screen. **b**, The distribution of unstimulated and stimulated Jurkat cells along the UMAP plot. **c**, The correlation of predicted scores by scMAGeCK-PS and mixscape. **d**, the Receiver-Operating Characteristic (ROC) curve of both methods in separating positive and negative hits. **e-f**, Benchmark using a published ECCITE-seq where PDL1 protein expression is used as gold standard (**e**), and the performance of different methods in terms of predicting PLD1 protein expression (**f**).



1062
1063
1064
1065
1066
1067
1068
1069
1070
1071

Figure 3. Dose-dependent responses of perturbations. **a**, The correlation between a gene’s PS and a phenotype of interest indicates positive (or negative) regulations. **b-c**, The correlation between PDL1 protein expression and the PS of CUL3 (**b**) and STAT1 (**c**). CUL3 is a known negative regulator of PDL1, while STAT1 is a known positive regulator. **d**, The classification of buffered or sensitive genes, based on perturbed gene expression and PS. **e**, The classification of buffered or sensitive genes from published Perturb-seq datasets focusing essential genes in K562²⁶. **f**, The perturbation-expression plot of PSMA3, a buffered gene. **g**, The log fold changes of marker gene expressions (columns) upon perturbing proteasome genes (rows) from the essential gene Perturb-seq dataset.

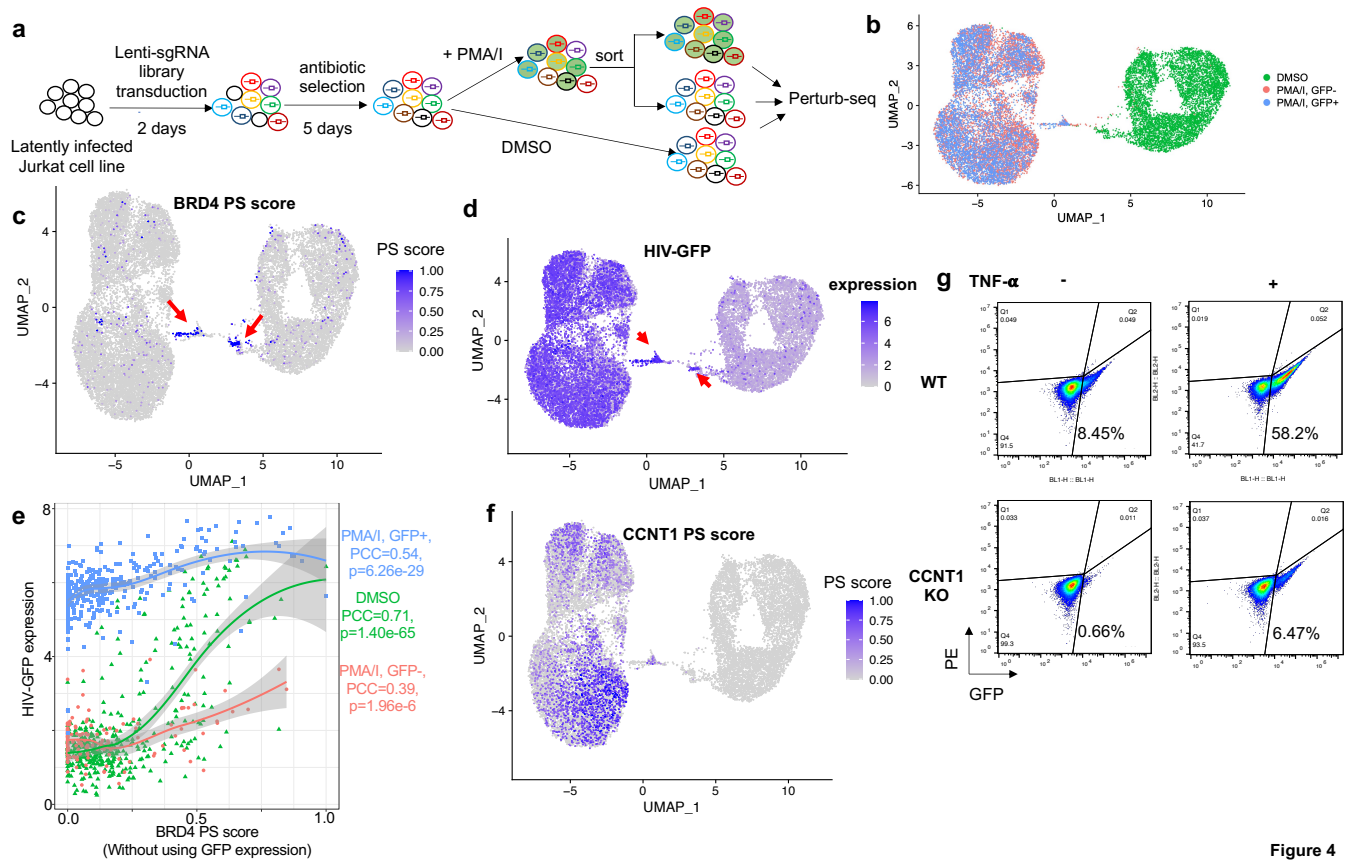


Figure 4

Figure 4. Perturb-seq on HIV latency. **a**, The experimental design of Perturb-seq. **b**, The UMAP plot of single-cell transcriptome profiles. Cells are colored by three different conditions. **c**, The distribution of BRD4 PS. **d**, The expression of HIV-GFP. **e**, The correlations between HIV-GFP expression and BRD4 PS that does not use HIV-GFP as target gene. **f**, The distribution of CCNT1 PS. **g**, The protein expression of HIV-GFP in response to CCNT1 knockout in different cell states (TNF-alpha vs non-stimulated).

1072
1073
1074
1075
1076
1077
1078

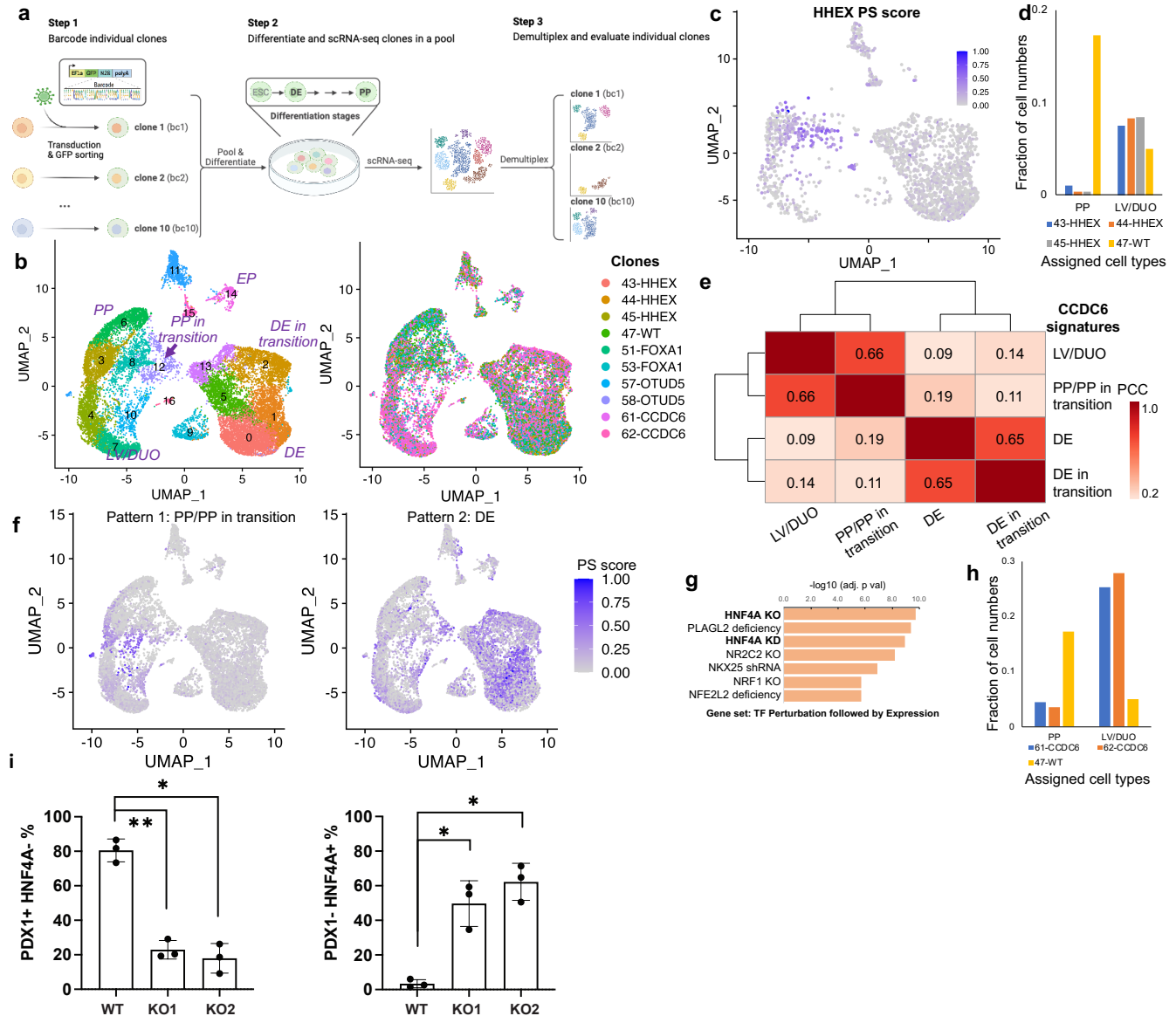


Figure 5

1079
 1080 **Figure 5. Pooled scRNA-seq on pancreatic differentiation.** **a**, Experimental design of multiplexing
 1081 scRNA-seq on the knockout clones of different genes. **b**, The UMAP plot of single-cell transcriptome
 1082 profiles, colored by different clusters (left) or clones (right). **c**, The PS distribution of HHEX. **d**, The
 1083 percentage of cells in PP/LV/DUO cell types from different clones. **e**, The correlations of CCDC6 PSs
 1084 calculated from different HHEX cell types. The Pearson Correlation Coefficient (PCC) is calculated
 1085 from all cells with CCDC6 knockouts and is shown as numbers on the heatmap. **f**, Two different
 1086 distribution patterns of CCDC6 PSs. **g**, The top enriched GO terms of DEGs from PP/PP in transition.
 1087 Enrichr was used to perform enrichment analysis. **h**, The percentage of cells in PP/LV/DUO cell types
 1088 from CCDC6 clones. **i**, The percentage of cells with PDX1+ (a PP marker) or HNF4A+ (a LV marker)
 1089 by flow cytometry sorting. The data is based on two CCDC6 knockouts (KO1, KO2) and one wild-type
 1090 (WT) control. Three independent replicates are performed for each condition. The multiple comparison-
 1091 adjusted p value is calculated by one-way ANOVA test. *p<0.05, **p<0.01.

1094 **Supplementary Figure Legends**

1095

1096 **Supplementary Figure S1. Benchmark different methods using simulated and real datasets. a,**
1097 Steps to generate simulated datasets using scDesign3 from a real scRNA-seq dataset that knocks out
1098 Nelfb gene. **b-e,** The score distribution of scMAGeCK-PS and mixscape using different DEGs and
1099 different values of true efficiencies. **f-g,** Similar with Figure 1f-g, but using a published high MOI
1100 Perturb-seq dataset in the same study. **h-i,** Benchmark results of different methods on another published
1101 CRISPRi-based Perturb-seq, where mismatches are introduced into guides to attenuate perturbation
1102 effects. The Pearson correlation coefficients (PCCs) between the predicted scores of each method and
1103 the expressions of perturbed genes are reported for every perturbed gene (**h**), and between the predicted
1104 scores and predicted sgRNA activities (**i**), using the prediction methods provided in the original study²¹.

1105

1106 **Supplementary Figure S2. A genome-scale Perturb-seq. a-b,** The distribution of scMAGeCK-PS and
1107 mixscape predicted scores of top hits including CD247 (**a**) and LCK (**b**) in the pooled screen. **c,** The
1108 correlation between PSs and perturbed gene expression.

1109

1110 **Supplementary Figure S3. Predictions of PDL1 protein expression from a published ECCITE-seq**
1111 **dataset.** The ROC curve, the correlations between scMAGeCK-PS results with PDL1 protein
1112 expression, and the correlations between mixscape results with PDL1 protein expression are reported for
1113 each gene. The correlations are separated by classifications of each single cell: NP (non-perturbed),
1114 defined as mixscape score ≤ 0.5 ; and KO (knockout), defined as mixscape score > 0.5 . For a fair
1115 comparison, we used mixscape classification results to plot PSs (mid panel).

1116

1117 **Supplementary Figure S4. Buffered genes and sensitive genes. a,** RPL4, a buffered gene. **b-c,**
1118 HSPA5 and GATA1, two sensitive genes. **d,** A gene (BRD4) whose expression has no correlation with
1119 PS.

1120

1121 **Supplementary Figure S5. The log fold changes of gene expressions upon perturbing genes within**
1122 **the same protein complex,** including ribosomal subunits (**a**), RNA polymerase (**b**) and mediator
1123 complex (**c**) in essential gene Perturb-seq. (**d**) The log fold changes of proteasome gene expressions
1124 (columns) upon perturbing proteasome genes (rows) from the genome-scale Perturb-seq.

1125

1126 **Supplementary Figure S6. HIV Perturb-seq. a,** The number of genes (nFeature_RNA), UMI counts
1127 (nCount_RNA) and the fraction of mitochondrial RNAs in three different conditions. **b,** Clustering
1128 results. **c,** Enriched Gene Ontology (GO) terms of cluster 8. **d,** The distribution of BRD4-targeting
1129 gRNAs. **e,** The expression distribution of BRD4 signature genes in cluster 8 vs other clusters. Only cells
1130 express BRD4-targeting gRNAs are included. **f,** Differential expression results between BRD4 PS+ cells
1131 vs BRD4 PS- cells.

1132

1133 **Supplementary Figure S7. HIV Perturb-seq. a,** The expressions of CCNT1 (left) and CCNT1-
1134 targeting gRNAs (right). **b,** Differential expression results between CCNT1-targeting cells and non-
1135 targeting control cells in two different cell states. **c,** The expressions of NFKB1. **d,** The quantitative
1136 perturbation-expression relationship between GFP and CCNT1 PS, similar with Figure 4e.

1137

1138 **Supplementary Figure S8. Cell type assignment based on known expression markers of different**
1139 **cell types in pancreatic differentiation scRNA-seq.**

1140
1141 **Supplementary Figure S9. DEG analysis.** a-b, The distribution of FOXA1 PSs across two different
1142 clones. c, The expression pattern of FOXA1. d-e, The DEG analysis results of CCDC6 knockout clones
1143 vs. wild-type clones in different cell types. f, The overlap of statistically significant DEGs in DE and
1144 LV/DUO cell types.
1145
1146 **Supplementary Figure S10. Different CCDC6 functions.** a-b, The two patterns of CCDC6 PSs in
1147 LV/DUO (a) and DE in transition (b) cell types. c-f, Additional enriched terms using Enrichr on DEGs
1148 of CCDC6 knockout.
1149
1150 **Supplementary Figure S11. Flow cytometry analysis of PDX1 and HNF4A expression upon**
1151 **CCDC6 knockout.** One representative plots of three biological replicates are shown.
1152
1153 **Supplementary Table S1. Genome-scale Perturb-seq library design.**
1154
1155 **Supplementary Table S2. HIV Perturb-seq library design.**
1156
1157 **Supplementary Table S3. Sequencing summary of HIV Perturb-seq.**
1158
1159 **Supplementary Table S4. Genotype summary of 10-clone scRNA-seq pancreatic differentiation**
1160 **dataset.**