

1 A single-cell massively parallel reporter assay detects cell type 2 specific cis-regulatory activity

3 Siqi Zhao^{1,2}, Clarice KY Hong^{1,2}, David M Granas^{1,2}, Barak A Cohen^{1,2}

4 1 Edison Family Center for Systems Biology and Genome Sciences, 2 Department of Genetics,
5 Washington University School of Medicine, USA.

6
7 **We developed a single-cell massively parallel reporter assay (scMPRA) to measure the**
8 **activity of libraries of cis-regulatory sequences (CRSs) across multiple cell-types**
9 **simultaneously. As a proof of concept, we assayed a library of core promoters in a mixture**
10 **of HEK293 and K562 cells and showed that scMPRA is a reproducible, highly parallel,**
11 **single-cell reporter gene assay. Our results show that housekeeping promoters and CpG**
12 **island promoters have lower activity in K562 cells relative to HEK293, which likely reflects**
13 **developmental differences between the cell lines. Within K562 cells, scMPRA identified a**
14 **subset of developmental promoters that are upregulated in the CD34⁺/CD38⁻ sub-state,**
15 **confirming this state as more “stem-like.” Finally, we deconvolved the intrinsic and**
16 **extrinsic components of promoter cell-to-cell variability and found that developmental**
17 **promoters have a higher proportion of extrinsic noise compared to housekeeping**
18 **promoters, which may reflect the responsiveness of developmental promoters to the**
19 **cellular environment. We anticipate scMPRA will be widely applicable for studying the role**
20 **of CRSs across diverse cell types.**

21 Introduction

22 The majority of heritable variation for human diseases maps to the non-coding portions of the
23 genome¹⁻⁶. This observation has led to the hypothesis that genetic variation in the *cis*-regulatory
24 sequences (CRSs) that control gene expression underlies a large fraction of disease burden⁷⁻¹⁰.
25 Because many CRSs function only in specific cell types¹¹, there is intense interest in high-
26 throughput assays that can measure the effects of cell-type-specific CRSs and their genetic
27 variants.

28 Massively Parallel Reporter Assays (MPRAs) are one family of techniques that allow
29 investigators to assay libraries of CRSs and their non-coding variants *en masse*¹²⁻¹⁸. In an MPRA
30 experiment, every CRS drives a reporter gene carrying a unique DNA barcode in its 3' UTR, which
31 allows investigators to quantify the activity of each CRS by the ratio of its barcode abundances in
32 the output RNA and input DNA. This approach allows investigators to identify new CRSs, assay
33 the effects of non-coding variants, and discover general rules governing the functions of
34 CRSs^{12,19-23}. One limitation of MPRAs is that they are generally performed in monocultures, or as
35 bulk assays across the cell types of a tissue. Performing cell-type specific MPRAs in tissues will
36 require methods to simultaneously readout reporter gene activities and cell type information in
37 heterogeneous pools of cells.

38 To address this problem, we developed scMPRA, a procedure that combines single-cell
39 RNA sequencing with MPRA. scMPRA simultaneously measures the activities of reporter genes
40 in single cells and the identities of those cells using their single-cell transcriptomes. The key

41 component of scMPRA is a two-level barcoding scheme that allows us to measure the copy
42 number of all reporter genes present in a single cell from mRNA alone. A specific barcode marks
43 each CRS of interest (CRS barcode, “cBC”) and a second random barcode (rBC) acts as a proxy
44 for DNA copy number of reporter genes in single cells (**Fig. 1a**). The critical aspect of the rBC is
45 that it is complex enough to ensure that the probability of the same cBC-rBC appearing in the
46 same cell more than once is vanishingly small. In this regime, the number of different cBC-rBC
47 pairs in a single cell becomes an effective proxy for the copy number of a CRS in that cell. Even
48 if a cell carries reporter genes for multiple different CRS, and each of those reporter genes is at
49 a different copy number, it is still possible to normalize each reporter gene in each individual cell
50 to its plasmid copy number. With this barcoding scheme, we can measure the activity of many
51 CRSs with different input abundances in single cells.

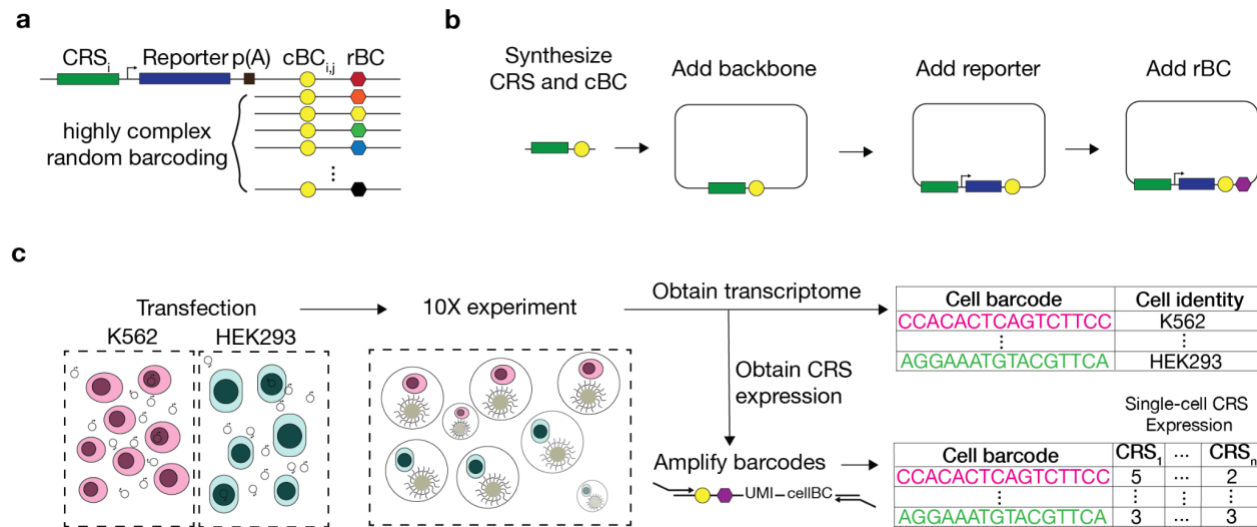
52 Results

53 scMPRA enables single-cell measurement of CRS activity

54 As a proof of principle, we used scMPRA to test whether different classes of core promoters show
55 different activities in different cell types. Core promoters are the non-coding sequences that
56 surround transcription start sites, where general cofactors interact with RNA polymerase II^{24,25}.
57 Core promoters are divided into different classes by the functions of their host genes
58 (housekeeping vs developmental), as well as by the sequence motifs they contain (TATA-box,
59 downstream promoter element (DPE), and CpG islands). We selected 676 core promoters that
60 we previously tested²⁴ and cloned them into a double-barcoded MPRA library (**Supplementary**
61 **Table 1**). In the first stage of library construction each core promoter reporter gene was
62 represented by 10 unique cBCs. We then added rBCs to the library by cloning a 25 nt random
63 oligonucleotide (oligo) directly downstream of the cBCs. The library contains $\sim 1.4 \times 10^7$ unique
64 cBC-rBC pairs (**Methods, Fig. 1b**). Using this complexity, we calculated that the probability of
65 plasmids with the same cBC-rBC pair occurring in the same cell is less than 2×10^{-3} with our
66 transfection protocols (**Methods**). Given this low likelihood, the number of rBC per cBC in a cell
67 represents the copy number of a CRS in that cell. Knowing the copy number of CRSs in single
68 cells allows us to normalize reporter gene expression from each CRS to its copy number in
69 individual cells.

70 We performed a cell mixing experiment to test whether scMPRA could measure cell type
71 specific expression of reporter genes. We transfected K562 and HEK293 cells (**Methods**), and
72 performed scMPRA on a 1:1 mixture of those cell lines (**Fig. 1c**). We harvested cells and prepared
73 them for sequencing using the 10X Chromium™ platform. The mRNA from single cells was
74 captured, converted to cDNA, and pooled together. We then split the samples, with a quarter of
75 the amplified cDNA library used for amplifying the cBC-rBC pairs and three-quarters used to
76 amplify the transcriptome. The resulting reporter barcode abundances and transcriptome of each
77 single cell are linked by their shared 10X cell barcode (**Methods**).

78



79
80
81
82
83
84
85
86
87
88
89

Figure 1 scMPRA measures CRS at single-cell resolution. (a) Each CRS reporter construct is barcoded with a cBC that encodes the identity of the CRS, as well as a highly complex rBC. The complexity of the cBC-rBC pair ensures that the probability of identical plasmids being introduced into the same cell is extremely low. (b) Cloning strategy for the double barcoded library. CRSs and their corresponding cBCs are synthesized together and cloned into an appropriate backbone. 25 nt rBCs are introduced to the plasmids with Hifi assembly. (c) Experimental overview for scMPRA using mixed cell experiment as an example. K562 cells and HEK293 cells are transfected with the double-barcoded core promoter library. After 24 hours, cells were harvested and mixed for 10X scRNA-seq. Cell identities were obtained through measuring the single transcriptome, and single-cell expression of CRSs was obtained by quantifying the barcodes. The cell identity and CRSs expression were linked by the shared 10X barcodes.

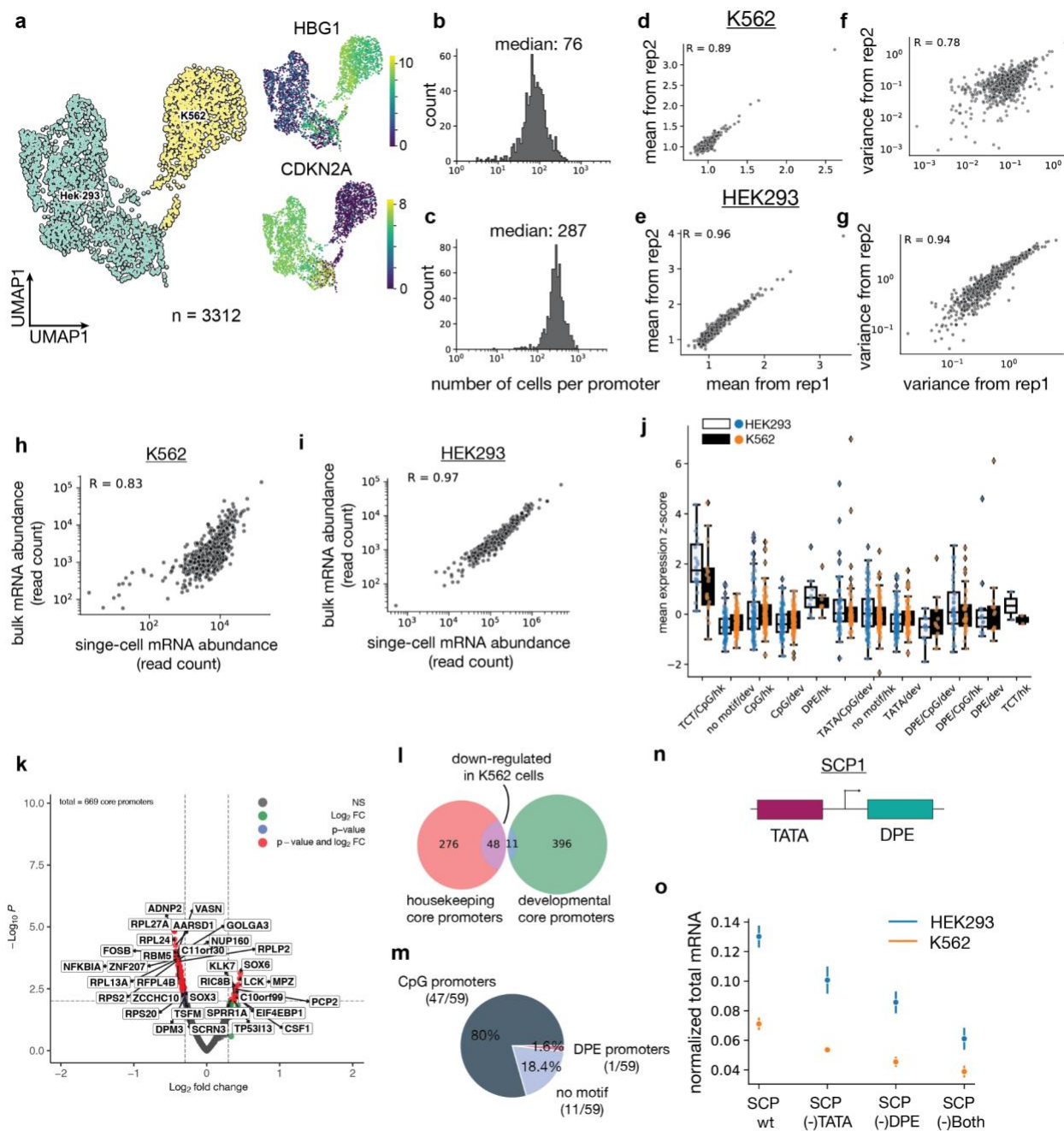
We recovered a total of 3112 cells (1524 in replicate 1 and 1588 in replicate 2) that are unambiguously assigned to one of the two cell types (**Fig. 2a, Supplementary Figs S1 a,b**). We determined the efficiency of our method by calculating the recovery rate of our input promoters. We then calculated the core promoter expression by taking the average of the cBC expression for the same promoter. We found that scMPRA recovered 99.5% (673 out of 676 core promoters) of the input library for K562 cells and 100% (676 out of 676 core promoters) for HEK293 cells, highlighting the efficiency of our method for recovering input elements.

We next calculated the number of individual cells in which each core promoter is measured. We found that the empirical distribution of the number of cells per core promoter is log normal, with a median of 76 cells per core promoter for K562 cells and 287 cells per core promoter for HEK293 cells (**Fig. 2b,c**). Given that the number of pBC-rBC pair is effectively the number of plasmids per cell, we also calculated the number of plasmid per cell, and found that fewer number of plasmids were incorporated into K562 cells compare to HEK293 cells (median plasmid number in K562 cells: 164, median plasmid number in HEK293 cells: 341. **Supplementary Fig. 1c,d**). The difference in transfection efficiency between these cell types with the same input likely reflects global cellular differences between them, and is representative of the condition when performing scMPRA in different cell types.

We calculated the biological reproducibility and found that scMPRA is highly reproducible in both cell types for measurements of mean expression (K562: Pearson R = 0.89, HEK293: Pearson R = 0.96) and cell-to-cell variance (K562: Pearson R = 0.78, HEK293: Pearson R = 0.94, **Fig 2 d-g**). To validate the measurements, we conducted bulk RNA-seq for the core promoter library in the two cell types separately, and found the bulk measurements correlate well with the aggregated single-cell measurements (**Fig. 2 h,i, Supplementary Fig. 1e,f**). This analysis shows

107
108
109
110
111
112

113 that single-cell measurements of library members in as few as 70 individual cells still correlate
 114 well with bulk measurements, highlighting the sensitivity of our method.
 115



116 **Figure 2. scMPRA detects cell type specific CRS activity.** (a) UMAP of the transcriptome from the mixed-cell scMPRA
 117 experiment. 3312 out of 3417 cells are assigned to either K562 or HEK293 cells. Cell-type specific genes were used to identify the
 118 cell clusters (HBG1 for K562 cells and CDKN2A for HEK293 cells). Cells are labeled by their cell type. (b,c) Histogram of the
 119 number of cells per core promoter for HEK293 and K562 cells. (d-g) Reproducibility for expression mean and cell-to-cell variance for
 120 both K562 and HEK293 cells. (h,i) Scatterplot of reproducibility of scMPRA mean expression with bulk MPRA measurement using
 121 read count normalization. (j) Boxplot of mean expression from different categories of core promoters in K562 (orange) and HEK293
 122 (blue) cells. (k) Volcano plot for differential expression (DE) of the core promoters in K562 and HEK293 cells (Significant DE
 123 reporters have p -value < 0.01 and \log_2 fold change greater than 0.3). (l) A Venn diagram of the functional characterization
 124 (housekeeping vs developmental) of down-regulated reporters in K562 cells. Housekeeping promoters are enriched (p -value =
 125 1.08×10^{-11} from hypergeometric test). (m) Pie chart of the sequence features (CpG, DPE, TATA) of down-regulated reporter genes.
 126

127 CpG promoters are enriched ($p=2.18 \times 10^{-6}$, from hypergeometric test). (n) Schematic of SCP1 binding sites. (o) Expression of wild-
128 type and mutated (TATA⁻, DPE⁻, and Both) versions of SCP1 core promoter (error bar: 1 s.d.)

129

130

131 **scMPRA detects cell type specific CRS activity and non-coding variant effect**

132 We asked whether the data allowed us to detect core promoters with differential activity between
133 K562 and HEK293 cells. While different classes of core promoters had similar activities in both
134 cell lines (**Fig. 2j**), our differential analysis using DEseq2²⁶ identified a small number of promoters
135 (11 out of 669) that are upregulated in K562 cells, and 59 promoters that are downregulated in
136 K562 cells (adjusted $p < 0.01$, log2 fold change > 0.3 , **Fig. 2k, Supplementary Table 2**). Among
137 the down-regulated promoters, 48 out of 59 core promoters belong to housekeeping genes ($p =$
138 1.08×10^{-11} , **Fig. 2l**), and 46 out of 59 core promoters are CpG-island-containing core promoters
139 ($p=2.18 \times 10^{-6}$, **Fig. 2m**). This down-regulation might be explained by the fact that the K562 cell
140 line is a cancer derived cell line, and a hallmark regulatory change in cancer cells is the
141 hypermethylation of CpG promoters²⁷. These results demonstrate the ability of scMPRA to detect
142 CRSs with cell-type specific activities.

143 Another application of scMPRA is to detect cell type specific effects of non-coding
144 variants. To test whether our method can detect the effects of mutations in a given CRS, we
145 included an artificial core promoter SCP1²⁸ along with mutated versions without a TATA Box or
146 DPE motif in our library (**Fig. 2n**). We first computed the total number of captured reporter gene
147 transcripts, since it is the closest proxy to the bulk expression measurement. We found that
148 deletions of the TATA motif or DPE motif both reduced expression (**Fig. 2o**) and we observed a
149 similar trend in the bulk data (**Supplementary Fig. 1g**). When we directly calculated the mean
150 of the single-cell expression distribution instead of total number of captured reporter gene
151 transcripts, we found that the deletion of the DPE motif has a stronger effect in K562 cells than
152 in HEK293 cells (40% reduction vs 20% reduction) (**Methods, Supplementary Fig. 1 h,i**). We
153 hypothesized that the differential expression of transcription factors between K562 and HEK293
154 cells leads to differential sensitivity to the TATA and DPE motifs. We examined the single-cell
155 transcriptome and found that TAF9, which recognizes the DPE motif²⁹, is more highly expressed
156 in K562 cells compared to HEK293 Cells (**Supplementary Fig. 1j**, Wilcoxon $p=4.27 \times 10^{-94}$). This
157 observation likely explains why the deletion of the DPE motif has a stronger effect in K562 cells.
158 Our results demonstrate that scMPRA can identify and explain cell-type specific effects of non-
159 coding variants.

160

161 **scMPRA detects cell sub-state specific CRS activity**

162 Single-cell studies have revealed heterogeneity in cell states even within isogenic cell types^{30–33}.
163 Therefore, we asked if scMPRA can identify CRSs with cell-state specific activity. We repeated
164 scMPRA on K562 cells alone and obtained a total of 5141 cells from two biological replicates.
165 Measurements of the mean and variance of each library member were again highly correlated
166 between replicates and agree well with independent bulk measurement (**Supplementary Fig. 2**
167 **a-d**).

168 As the phases of the cell cycle represent distinct cell-states, we asked whether scMPRA
169 could identify reporter genes with differential activity through the cell cycle. We assigned cell cycle
170 phases to each cell using their single cell transcriptome data (**Fig. 3a**) and then calculated the
171 mean expression of each reporter gene in different cell cycle phases. We found that most core

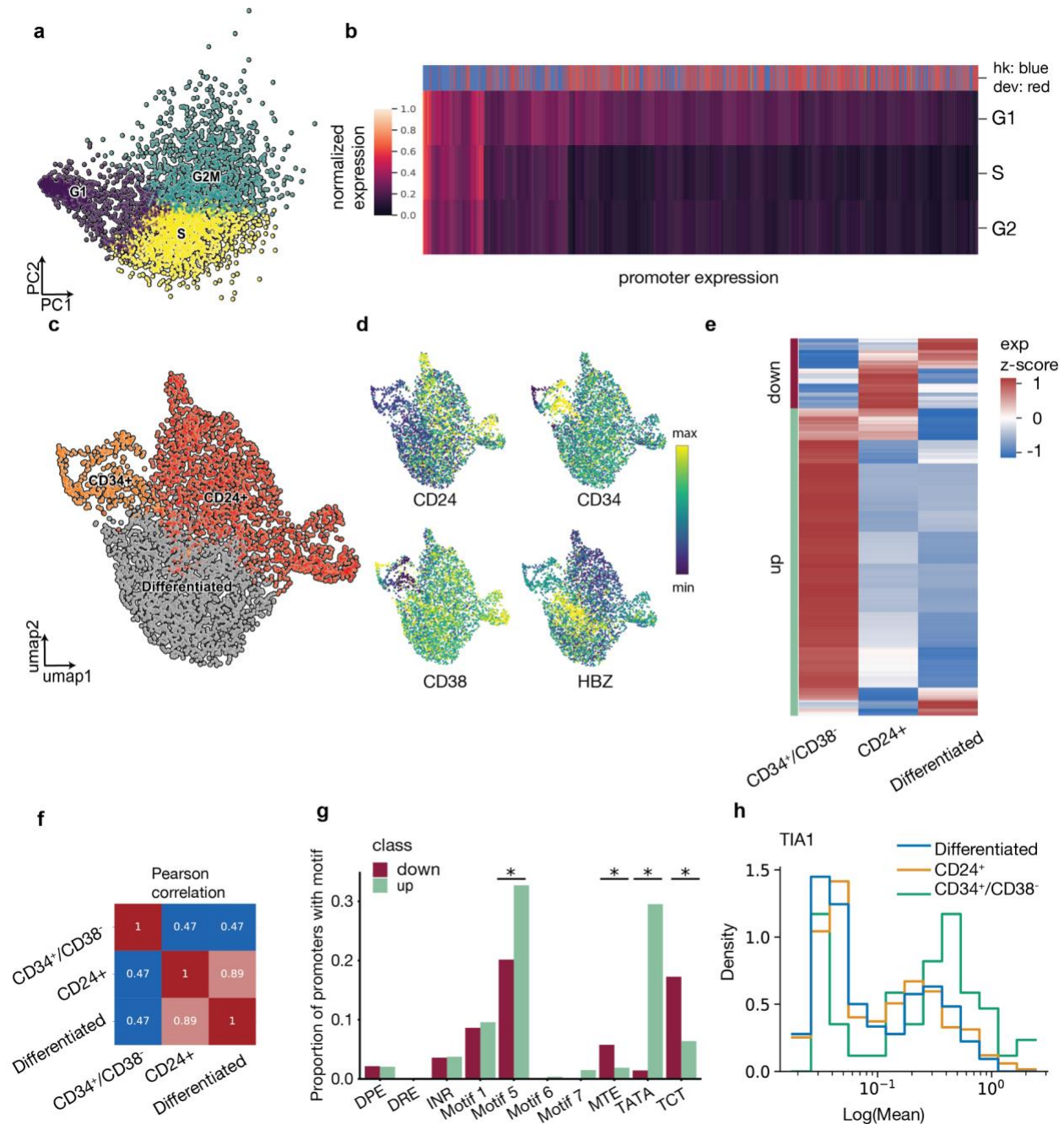
172 promoters in our library are upregulated in the G1 phase of the cell cycle, and some housekeeping
173 promoters are highly expressed through all cell cycle phases (**Fig. 3b**). We also identified core
174 promoters with different expression dynamics through the cell cycle. For example, we found the
175 core promoter for UBA52 remains highly expressed in the S phase, whereas the core promoter
176 for CXCL10 is lowly expressed throughout (**Supplementary Fig. 2e**). This analysis illustrates the
177 ability of scMPRA to identify CRSs whose expression naturally fluctuates with cellular dynamics.

178 We then asked whether scMPRA could detect reporter genes with activities that were
179 specific to other cell-states in K562 cells, after normalizing for cell cycle effects. We focused on
180 two specific sub-states that have been reported and experimentally validated for high proliferation
181 rates in K562 cells^{34,35}. The first is the CD34⁺/CD38⁻ sub-state that has been identified as a
182 leukemia stem-cell subpopulation, and the second is the CD24⁺ sub-state that is linked to
183 selective activation of proliferation genes by bromodomain transcription factors^{31,32}. To identify
184 these sub-states in our single-cell transcriptome data, we first regressed out the cell cycle effects
185 and confirmed that the single cell transcriptome data no longer clustered by cell cycle phase
186 (**Supplementary Fig. 2f**). We then identified clusters within K562 cells that have the CD34⁺/CD38⁻
187 expression signature, or the CD24⁺ signature (**Fig. 3 c,d**). Although the CD34⁺/CD38⁻ cells
188 represent only 9.3% of the cells in a K562 culture, scMPRA revealed two distinct classes of core
189 promoters that are upregulated and downregulated in these cells respectively (**Fig 3e**).
190 Conversely, the expression patterns of promoters are similar between the CD24⁺ cluster and cells
191 in the “differentiated” cluster (**Fig. 3e, f**). Motif analysis of the up/down regulated classes of
192 promoters in CD34⁺/CD38⁻ cells showed that different core promoter motifs are enriched in each
193 class, with the TATA box and Motif 5 being enriched in the upregulated class and MTE and TCT
194 motifs being enriched in downregulated class (**Fig. 3g, Methods**). This result suggests that
195 differences in core promoter usage might be driving the differences between CD34⁺/CD38⁻ and
196 the other clusters. Because the TATA box is mostly found in developmental core promoters, the
197 CD34⁺/CD38⁻ subpopulation likely reflects a more “stem-like” cellular environment in these cells.
198 Our analysis highlights the ability of scMPRA to identify CRSs with differential activity in rare cell
199 populations.

200 With the single-cell expression data, we asked how certain promoters achieve higher
201 expression in the CD34⁺/CD38⁻ state. We asked whether the single-cell expression distribution
202 for the CD34⁺/CD38⁻ state is shifted higher than for the other states, or if the range of expression
203 is the same for each sub-state, with only the proportion of cells with high expression changing in
204 each state. To answer this question, we calculated the proportion of cells in each sub-state
205 belonging to the 90th percentile of the total single cell expression distribution. For the majority of
206 promoters, the CD34⁺/CD38⁻ cluster has a much higher proportion of cells in the 90th percentile
207 (**Supplementary Fig 3a**). At the same time, there is no difference in the maximum expression of
208 cells in different sub-states, and this maximum level is mainly set by the promoter identity
209 (**Supplementary Fig 3b**). Even for the most differentially expressed promoter in the CD34⁺/CD38⁻
210 subpopulation, TIA1, the expression distributions for cells in the three sub-states cover the same
211 range, but the proportion of cells in the right-tail of the distribution is higher for CD34⁺/CD38⁻ cells
212 (**Fig. 3h**). This result suggests that the “stem-like” cellular environment of the CD34⁺/CD38⁻
213 subpopulation increases the probability of certain promoters having higher expression, without
214 shifting the maximum expression those promoters achieve. Taken together, these analyses

215 highlight how the joint transcriptome and CRS measurements in scMPRA can be used to
 216 understand differences in behavior in cellular sub-states.

217
 218



219
 220 **Figure 3. scMPRA detects cell sub-state-specific CRS activity.** (a) PCA plot of K562 cells classified based on the cell cycle
 221 score. (b) Heatmap of reporter expression in different cell cycle phases (Color bar indicates housekeeping (blue) vs developmental
 222 (red) promoters). (c) Representative expression dynamics of reporter genes through cell cycle for UBA52, CSF1, and CXCL10. (d)
 223 UMAP embedding of K562 cells with high proliferation sub-states (CD34⁺/CD38⁻ and CD24⁺). (e) Marker gene expression signifies
 224 different cell sub-states in K562 cells. CD34, CD38 marks the “leukemia stem cell” sub-state; CD24 marks a high proliferation sub-
 225 state, and HBZ marks the differentiated leukemia sub-state; left color bar: hierarchical clustering showing 2 clusters based on
 226 expression pattern in the three substates. (f) Heatmap showing the correlation matrix of core promoter expression in three substates
 227 (CD34⁺/CD38⁻, CD24⁺, and Differentiated). (g) Proportion of promoters in each cluster that contains the indicated core promoter

228 motif. * represents significant enrichment in one cluster over the other ($p < 0.05$, Fisher's exact test). (h) Histogram of single-cell
229 expression of TIA1 promoter in three substates.

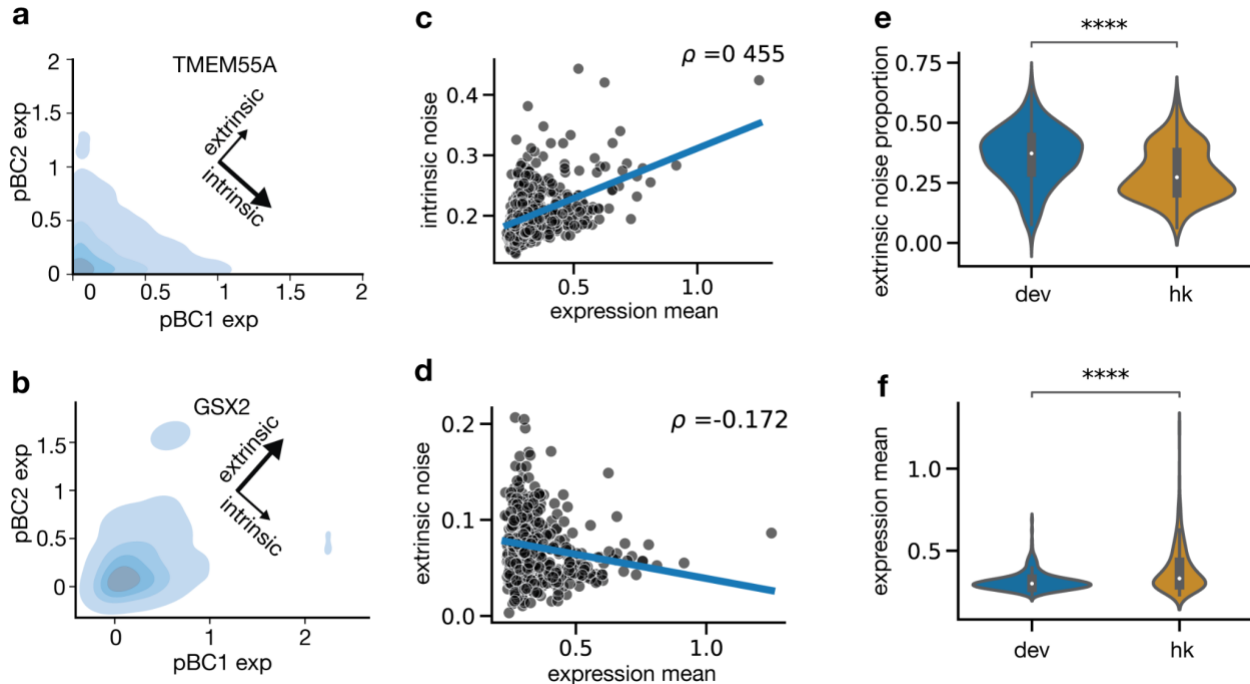
230

231

232 **scMPRA enables decomposition of intrinsic and extrinsic noise**

233 Finally, we analyzed the cell-to-cell variability of reporter genes across K562 cells. Cell-to-cell
234 variability, or expression noise, is the phenomenon where gene expression varies among the cells
235 of an isogenic population. Expression noise has important roles in development³⁶, rare-cell cancer
236 resistance^{30,37}, and its origin is a central question in single-cell biology. A common framework for
237 studying expression noise is to decompose it into its intrinsic component, which arises from the
238 thermal fluctuations of macromolecular interactions, and its extrinsic component, which results
239 from fluctuations in the global cellular environment³⁸⁻⁴². Intrinsic and extrinsic noise can be
240 decomposed using dual-reporter experiments, where two identical reporter genes are measured
241 across the same single-cells³⁹. High covariance of the two reporter genes indicates high extrinsic
242 noise and low intrinsic noise, while independent variation of the two reporters suggests high
243 intrinsic noise and low extrinsic noise. In scMPRA, plasmids with the same CRS but different
244 barcodes are sometimes incorporated into the same cells, effectively serving as a dual-reporter
245 experiment. We extracted pair-wise expression for the same core promoter labeled with different
246 cBCs from our scMPRA data, and computed intrinsic noise and extrinsic noise using a previously
247 developed statistical framework⁴³ (**Methods**). We found that different core promoters have distinct
248 intrinsic and extrinsic noise profiles (**Fig 4 a,b**). Globally, we found that intrinsic noise correlates
249 with mean expression levels (Pearson $\rho = 0.455$), while extrinsic noise is not correlated with mean
250 expression (Pearson $\rho = -0.172$, **Fig. 4 c,d**). This result agrees with the notion that intrinsic noise
251 arises from the thermodynamics of transcription at different promoters, whereas many sources
252 for extrinsic noise are independent of the specific promoters. We also found that developmental
253 promoters have a higher proportion of noise that is extrinsic, reflecting their role in driving
254 developmental promoters that respond to extrinsic cues during development (**Fig. 4 e,f**). This
255 analysis suggests that scMPRA could be a powerful tool to study the mechanistic origin of cell-
256 to-cell variability in a high throughput manner.

257



258
259
260
261
262
263
264
265
266

Figure 4. scMPRA deconvolves intrinsic and extrinsic cell-to-cell variability. (a, b) Density plots for single-cell expression of paired cBC expression for the same promoter. TMEM55A has high intrinsic noise, and GSX2 has high extrinsic noise. (c) Scatterplot of expression against intrinsic noise. Blue line shows the linear regression (Pearson $\rho = 0.455$) (d) Scatterplot of expression against extrinsic noise. Blue line shows the linear regression (Pearson $\rho = -0.172$) (e) Violin plot of extrinsic noise proportion for housekeeping and developmental promoters (Mann-Whitney U test. Stars indicate significance: **** : $p < 1 \times 10^{-4}$) (f) Violin plot of expression mean for housekeeping and developmental promoters (Mann-Whitney U test. Stars indicate significance: **** : $p < 1 \times 10^{-4}$)

267 Conclusions

268 We have presented a method to measure the cell-type and cell-state specific effects of CRSs by
269 devising a barcoding scheme to read out input copy number with mRNA. We demonstrated that
270 scMPRA detects cell-type specific reporter gene activity in a mixed population of cells, and cell-
271 state specific activity in an isogenic population. We also demonstrated that scMPRA can be a
272 powerful tool to study how different CRS control cell-to-cell variability. The assay is reproducible
273 and reports accurate mean levels of reporter gene activity in as few as 70 cells. The primary
274 limitation of scMPRA is that it relies on mRNA counts of the rBC to estimate plasmid DNA
275 abundance, and therefore it cannot accurately measure CRSs that are truly silent in a given cell
276 type. The inclusion of a separate constitutive promoter on each plasmid driving expression of the
277 rBCs would allow us to quantify plasmid copy number independent of the expression of the
278 reporter gene.

279 A future direction is to perform scMPRA in complex tissues to measure the cell type
280 specific effects of genetic variation in CRSs. With the burgeoning of Adeno-associated viral
281 delivery systems with distinct tropisms⁴⁴⁻⁴⁷, we anticipate that scMPRA will be widely used to
282 study cis-regulatory effects in a variety of complex tissues.

283
284

285

286 Methods

287 Cell culture

288 K562 cells were cultured using a medium consisting of Iscove's Modified Dulbecco's Medium
289 (IMDM) + 10% Fetal Bovine Serum (FBS) + 1% non-essential amino acids + 1% pen/strep at 37
290 C with 5% of CO₂. HEK293 cells were cultured using a medium consisting of Eagle's Minimum
291 Essential Medium (EMEM) + 10% Fetal Bovine Serum (FBS) + 1% pen/strep at 37 C with 5% of
292 CO₂.
293

294 Cloning Strategy

295 We developed a two-level barcoding technology to enable single-cell normalization for plasmid
296 copy number. We applied this strategy to a promoter library we previously tested in bulk
297 assays²⁴. The original library contains 676 core promoters with a length of 133bp. Each core
298 promoter has 10 promoter barcodes to provide redundancy in the measurements. We then
299 synthesized a single-stranded 90 bp DNA oligonucleotide containing a 25 bp random sequence,
300 a restriction site, and 30 bp homology on each side of the barcode region.
301

302 We used Hifi Assembly™ to add the random barcodes to the plasmid library. 4 µg of the
303 plasmid library were split into 4 reactions and digested with 2µl of Sall for 1.5 hours at 37°C. The
304 digested products were run at 100V for 2 hours on a 0.7% agarose gel. The correct size band
305 was cut and purified with the Monarch Gel Extraction Kit (New England BioLabs T1020L). The
306 insert single-stranded DNA was diluted in TE to a stock concentration of 100 uM. The insert was
307 then further diluted to 1 uM with ddH₂O. Three assembly reactions were pooled together, each
308 reaction containing 100 ng of digested library backbone, 1 uM of insert DNA, 1µl of NEBuffer 2,
309 10 µl of 2X Hifi assembly mix, and H₂O up to 20 ul, The reaction was incubated at 50°C for 1
310 hour. The assembled product was purified with the Monarch PCR&DNA Cleanup kit (New
311 England BioLabs T1030L) and eluted in 12 µl of H₂O.
312

313 The assembled plasmid was transformed using Gene Pulser Xcell Electroporation Systems by
314 electroporation (BIO-RAD 1652661), 50 µl of ElectroMax DH10B electrocompetent cells
315 (Invitrogen 18290015) with 1 µl of hifi assembled product at 2 kV, 2000 Ω, 25 nF, with 1 mm
316 gap. 950 µl of SOC medium (Invitrogen 15544034) was added to the cuvette and then
317 transferred to a 15 ml Falcon tube. Two transformations were performed, and each tube was
318 incubated at 37 °C for 1 hour on a rotator with 300 rpm. The culture was then added to pre-
319 warmed 150 µl LB/Amp medium and grown overnight at 37 °C. 1 µl of the culture was also
320 diluted 1:100 and 50 µl of the diluted cultured was plated on a LB agar plate to check the
321 transformation efficiency. For the core promoter library, we obtained more than 4X10⁸ colonies,
322 large enough to cover a complex library.

323 Estimating Library Complexity

324 To estimate the library complexity, we sequenced the DNA library using a nested PCR-based
325 Illumina library preparation protocol. Briefly, we first used Q5 polymerase (New England
326 BioLabs M0515) to amplify the region containing the two barcodes with SCARED P17 (5'-
327 GACGAGCTCTATAAGTAATCTAGA-3') and SCARED P18 (5'-TTTTCTAGGTCTCTGGTCTCGA-
328 3'). The total reaction volume is 50 μ l with 50ng of plasmids with 2.5 μ l of 10uM primer each. The
329 annealing temperature is 61°C with an extension time of 10s. 25 cycles of amplification were
330 done. The product was then purified with the Monarch PCR&DNA Cleanup kit (New England
331 BioLabs T1030L), and eluted with 20 μ l of ddH₂O. For the second PCR (SCARED P19: 5'-
332 GGACGAGCTCTATAAGTAATCTAGA-3', SCARED P20: 5'-
333 GTGACTGGAGTTCAGACGTGTGCTCTTCCGATCT-3'), a 25 μ l reaction was set up with 0.25
334 μ l product from the previous step, the annealing temperature is 61°C, and the extension time is
335 10s, a total of 10 cycles was done. The PCR product was cleaned up using the Monarch
336 PCR&DNA Cleanup kit. For the last PCR to add the P5 and P7 Illumina adapters (P5: 5'-
337 AATGATACGGCGACCACCGAGATCTACACACCCGCACACTCTTCCCTACACGACGCT-3',
338 P7:5'-CAAGCAGAAGACGGCATAACGAGATAAGTTGACAGTACTGGAGTTCAGACGTG-3'), a
339 reaction with 25 μ l of total volume was set up with 2 μ l of cleaned product from PCR2, a total of
340 10 cycles of PCR was done.

341
342 The constructed Illumina library was sequenced on an Illumina MiSeq. A total of 1,693,933
343 reads was generated for this library. A filtering strategy was applied to the raw reads, where
344 reads that do not have matching promoter barcodes and wrong-length random barcodes were
345 filtered out. We obtained a total of 1,359,176 reads (80% of the total reads) that contain the
346 correct promoter barcode and correct length random barcode.

347
348 The shallow sequencing of the input plasmid library enabled us to estimate the library
349 complexity and the probability of two identical copies of the plasmid being transfected into the
350 same cell. We first calculated that each random barcode is attached to 1.9 promoter barcodes
351 on average. For a total of 6760 input promoter barcodes, this suggests that a given random
352 barcode is being reused by 3200 different promoters. The reuse of random barcodes is the
353 effective labeling complexity for the double-barcoding. For the Hifi assembly experiment, we
354 used 300 ng input backbone plasmids containing only the promoter barcode (4.5×10^9 total
355 copies and on average 6.65×10^6 copies of plasmids per promoter barcode). Given the effective
356 labeling complexity, the average copy number of the plasmid containing the same promoter
357 barcode-andom barcode pair is at most 2.08×10^3 . For the transfection experiment done in this
358 study, with 2 μ g (6×10^9 copy of plasmids) for cell mixing experiment and 10 μ g (3×10^{10} copy of
359 plasmids) for K562 along experiment, the estimation of the average copy number for an
360 identical plasmid is 4.4×10^2 and 2.2×10^3 respectively.

361
362 After obtaining the average copy number for identical plasmids, we estimate the probability of
363 an identical plasmid being transfected into the same cell. We first define this probability as the
364 collision rate. We note that the transfection of the identical copies of different plasmids are
365 independent, so we could only calculate the collision rate for only one of such plasmids. The

366 calculation of the collision rate for a given library member can be formulated as such: given the
367 number of the identical copies of a plasmid, what is the probability of two or more of the copies
368 being transfected into the same cell? We first write the expectation:

$$n^{-m} \sum_{k=0}^n \binom{n}{k} \sum_{q=0}^{(n-k)} \binom{n-k}{q} \binom{m}{q} q! \left\{ \begin{matrix} m-q \\ n-k-q \end{matrix} \right\}_{n \geq 2} (n-k-q)!(m-q)$$

369 where n denotes the total number of cells, m denotes the total number of identical plasmids, k
370 denotes the number of cells with no plasmid, q denotes the cells with exactly 1 plasmid,
371 parentheses denote binomial coefficient, and brackets denote partition function.
372

373 The above equation was simplified by substituting with the bivariate generating function, and the
374 expected number is:
375

$$m \left(1 - \left(\frac{n-1}{n} \right)^{m-1} \right)$$

376 For a given transfection experiment, we can estimate the effective percentage of plasmid that is
377 successfully transfected into the cell. Given the estimated copy number for identical plasmids is
378 4.4×10^2 and 2.2×10^3 for mixed cell experiment and K562 alone experiment respectively, the
379 expected number of cells having more than 1 identical plasmid can be calculated with the
380 aforementioned equation, and the probability of two copies of an identical plasmid appearing in
381 the same cell is 0.0004 and 0.002 respectively. On a practical note, researchers have
382 suggested that the effective number of the plasmid that are incorporated into the nucleus is
383 about 0.01 - 0.1 of the input amount⁴⁸, hence a library containing around 2.5×10^5 different
384 members transfected to 1 million cells has a theoretical collision rate around 1%.
385
386

387 Transfection

388 K562 cells were transfected using electroporation with the Neon transfection system (Invitrogen
389 MPK5000). 1 million cells were transfected with 2 μ g of plasmid DNA (mixed-cell experiment) or
390 10 μ g of plasmid DNA (K562 sub-state experiment), with 3 pulses of 1450 V for 10 ms. The cells
391 were then plated to pre-warmed K562 medium.
392

393 HEK293 cells were transfected using the Lipofectamine3000 protocol. 4 μ l of p3000 reagent, 4 μ l
394 of Lipofectamine, and OptiMEM were mixed with 2 μ g of plasmid DNA to a volume of 250 μ l.
395 The lipofectamine reagents and plasmid were mixed and incubated at room temp for 15 minutes
396 and then added dropwise to the cells.

397 Bulk RNA extraction and sequencing

398 We determined the optimal harvest time based on plasmid dilution and protein maturation and
399 found the optimal harvest time is between 22 - 28 hours after transfection. The rationale behind
400 the choice of time is to balance the transcription rate and the plasmid dilution during cell
401 replication.
402

403 For both K562 cells and HEK293 cells, we harvested the cells after transfection at 24 hours, and
404 proceeded to extract total mRNA with Qiagen RNeasy kit for K562 cells and Monarch Total RNA
405 miniprep kit for HEK293 cells. The reverse transcription was done with Superscript IV Reverse
406 Transcriptase (Invitrogen 18090010). The final sequencing library was constructed using a
407 nested PCR strategy. Briefly, we first used Q5 (New England BioLabs M0515) polymerase to
408 amplify the region containing the 2 barcodes with SCARED P17 and SCARED P18. The total
409 reaction volume is 50µl with 50ng of backbone with 2.5 µl of 10uM primer each. The annealing
410 temperature is 61°C with an extension time of 10s. 25 cycles of amplification was done. The
411 product was then purified with the Monarch PCR&DNA Cleanup kit (New England BioLabs
412 T1030L), and eluted with 20 µl of ddH₂O. For the second PCR using primers SCARED P19 and
413 SCARED p20, a 25 µl reaction was set up with 0.25 µl product from the previous step, the
414 annealing temperature is 61°C, and the extension time is 10s, a total of 10 cycles was done.
415 The PCR product was cleaned up using the Monarch PCR&DNA Cleanup kit (New England
416 BioLabs T1030L). For the last PCR to add the P5 and P7 Illumina adapters, a reaction with 25
417 µl of total volume was set up with 2 µl of cleaned product from PCR2, a total of 10 cycles of
418 PCR was done. The sequencing library was sequenced on an Illumina Mi-seq machine with
419 other samples pooled in the same lane.

420 10X Experiment for scMPRA

421 We harvested both K562 and HEK293 cells 24 hours after transfection, then followed the cell
422 preparation protocol of 10X genomics. We used the 10X V3.1 chromium kit for our single-cell
423 RNA-seq protocol. All PCRs were performed on an Invitrogen PCR machine. We targeted 2000
424 cells per replicate for each experiment for the mixed cell experiment. We targeted 2500 cells per
425 replicate for the K562 substrate experiment. We followed the 10X protocol
426 (<https://support.10xgenomics.com/single-cell-gene-expression/library-prep/doc/user-guide-chromium-single-cell-3-reagent-kits-user-guide-v31-chemistry>) with 12 cycles of cDNA
427 amplification. To amplify the Capture-Sequence captured reads. 0.25 µl of 100 uM SCARED
428 P32 (5'-GTCAGATGTGTATAAGAGACAG-3') was added to the cDNA amplification mix. For
429 step 2.2, we modify the clean-up protocol by saving both the beads and supernatants and for
430 the supernatants, we use a final concentration of 1.2X beads to pull down the DNA fragments.
431 We then take 25% of both the 0.6X and 1.2X pull down products for the next step of PCR. To
432 construct the illumina sequencing library, we used a 3-step nested PCR strategy. Briefly, we first
433 used Q5 (New England BioLabs M0515) polymerase to amplify the region containing the 2
434 barcodes with SCARED P17 and SCARED P18. We pooled 8 PCR reactions, each with 50 µl of
435 total volume, with 10 cycles to reduce possible jackpotting. The annealing temperature is 61°C
436 with an extension time of 10s. The product was then purified with the Monarch PCR&DNA
437 Cleanup kit (New England BioLabs T1030L), and eluted with 20 µl of ddH₂O. For the second
438 PCR using the following 3 primers (SCARED P21: 5'-
439 GTGACTGGAGTTCAGACGTGTGCTCTTCCGATCTGGACGAGCTCTATAAGTAATCT-3', CAS
440 PC2: 5'-CGAGATCTACACTCGTCGGCAGCGTCAGATGTGTATAAGAGACAG-3', CAS PP2:
441 5'-ATCTACACTCTTTCCCTACACGACGCTCTTC-3'), we pulled 8 PCR reactions, each with 50
442 µl of total volume, with 10 cycles to reduce possible jackpotting , the annealing temperature is
443 61°C, and the extension time is 10s, a total of 10 cycles was done. The PCR product was
444 cleaned up using the Monarch PCR&DNA Cleanup kit (New England BioLabs T1030L). For the
445

446 last PCR to add the P5 and P7 Illumina adapters (CAS P48: 5'-
447 CAAGCAGAAGACGGCATAACGAGATNNNNNNNN[index]GTGACTGGAGTTCAGAC-3', CAS
448 PP4: 5'-AATGATACGGCGACCACCGAGATCTACACTCTTCCCTACA-3', CAS PC4: 5'-
449 AATGATACGGCGACCACCGAGATCTACACTCGTCG-3'), we pulled 8 PCR reactions, each
450 with 50 µl of total volume, with 10 cycles to reduce possible jackpotting, a total of 10 cycles of
451 PCR was done. The transcriptome is generated using the 10X Dual-Index Set TT expression kit
452 ([https://support.10xgenomics.com/single-cell-gene-expression/index/doc/technical-note-](https://support.10xgenomics.com/single-cell-gene-expression/index/doc/technical-note-chromium-next-gem-single-cell-3-v31-dual-index-libraries)
453 [chromium-next-gem-single-cell-3-v31-dual-index-libraries](https://support.10xgenomics.com/single-cell-gene-expression/index/doc/technical-note-chromium-next-gem-single-cell-3-v31-dual-index-libraries)).

454
455 The sequencing was done on an Illumina NextSeq machine. We used 40% of the barcode
456 library, 40% of the balanced scRNA-seq transcriptome, and 20% Phi-X. Sequencing the
457 constructed barcode library with transcriptome and Phi-X is crucial to reduce the sequencing
458 error from the reporter constant sequence. On Read1, only 28 bps contains the 10X cell
459 barcode and UMI was amplified, to avoid sequencing the constant Poly(A) sequence; On
460 Read2, 105 bps was sequenced. For the mixed experiment, we pool reads from a total of 2 runs
461 of NextSeq High Throughput sequencing runs, and for the K562 cells, we pool 3 runs of
462 NextSeq High Throughput runs.

463 scRNA-seq data processing

464 The single-cell RNAseq data were processed using Cellranger 6.0.1
465 (<https://github.com/10XGenomics/cellranger>) and Scanpy 1.8.1⁴⁹
466 (<https://github.com/theislab/scanpy>) following the standard pipeline. Briefly, different sequencing
467 runs from the same biological replicate were pooled together and processed with CellRanger
468 6.1.1; the final output expression matrix was then imported into Scanpy for further normalization.
469 We first removed cells with less than 1000 genes, and genes that were present in less than
470 three cells. We then removed cells with high counts for mitochondrial genes. Next, we
471 normalized the UMI counts to the total cell UMI counts. The normalized expression matrix was
472 used for clustering and visualization with Scanpy. The clustering was done using the Leiden
473 algorithm⁵⁰.

474 scMPRA data processing

475 The relevant script for processing single-cell MPRA reads can be found on a Github repository
476 (<https://github.com/szhao045/scMPRA>). The final sequencing product for scMPRA with Read1
477 contains the cell and molecular information (cellBC and UMI), and Read2 contains the MPRA
478 library information (cBC and rBC). First, we fuzzy-matched the constant sequences before and
479 after both the promoter barcode and random barcode. In this step, we filtered out the reads
480 without correct promoter barcode length, or random barcode length. To increase the speed, we
481 wrote a stand-alone program (https://github.com/szhao045/scMPRA_parsingtools) written with
482 Golang, and can be compiled to work on many operating systems. Second, we filtered out cell
483 barcodes based on the cell barcode list from the CellRanger output barcode list, with error-
484 correction with maximum hamming distance of 1. Third, to mitigate the effect of template-
485 switching during the PCR steps, we plotted the rank read depth for each unique quad of 10X
486 Cell Barcode, UMI, cBC, and rBC. We identified an elbow point with minimum depth of 1 (mixed

487 cell experiment) and 10 (K562 alone experiment), and kept any low-depth unique quad that
488 contains the cBC-rBC pair at most hamming distance of 1 to a high depth pair. Lastly, we
489 remove cells with less than 100 scMPRA-associated UMIs, since the scMPRA reads from those
490 cells were poorly sampled.

491 Cell cycle analysis

492
493 Cell cycle analysis for the scRNA-seq experiment was done with Scanpy 1.8.1 with cell cycle
494 genes⁵¹. The expression profile of each cell was projected onto a PCA plot based on the list of
495 cell cycle genes using Scanpy.

496

497 Motif analysis

498 The core promoters were first clustered according to their expression levels in the different cell
499 sub-state populations by hierarchical clustering. We categorized our data into up/down
500 regulated clusters at the first branching point, aiming to preserve the large structure. We then
501 identified core promoter motifs in each promoter according to the parameters in Zabidi et al⁵².
502 using MAST v4.10.0⁵³ and plotted the proportion of promoters containing each motif in each
503 promoter class.

504 Estimating intrinsic and extrinsic noise

505 Intrinsic and extrinsic noise were estimated using the statistical framework developed for the
506 dual-reporter experiment⁴³. We first extracted the pairwise expression level for cBCs that belong
507 to the same promoter in every single cell. If more than two cBCs are found in the same cell, the
508 pairwise expressions among them are recorded. We then removed promoters with less than
509 100 paired single-cell expression measurements (593 out of 676 promoters passed the filtering
510 step). We then applied the statistical framework developed by Fu and Pachter⁴³. The derivation
511 is abbreviated and can be found in the original publication. Briefly, let C denote the expression
512 for the first pBC in the cell and let Y denote the expression for the second pBC in the cell. Let
513 η_{int} denote the intrinsic noise, and it can be calculated as:

$$\eta_{int} = \frac{1}{a} \left(\sum_{i=1}^n C_i Y_i - n \bar{C} \bar{Y} \right),$$

514

515 where

$$a = (n - 1) \left(1 + \frac{1}{n} \right) + \frac{1}{\rho^2}$$
$$\rho = \frac{\text{Cov}[C, Y]}{\sqrt{\text{Var}[C]} \sqrt{\text{Var}[Y]}}$$

516

517 where n denotes the number of cells.

518

519 Similarly, let η_{ext} denote the extrinsic noise, and it can be calculated as:

$$\eta_{ext} = \frac{1}{2a\bar{C}\bar{Y}} \left(\sum_{i=1}^n (C_i - Y_i)^2 - n(\bar{C} - \bar{Y})^2 \right),$$

520
521 where

$$a = \frac{2n^3 - 7n + 6}{2(n^2 - n)} + \frac{2 - n}{n^2 - n} \frac{\rho}{1 - \rho} + \frac{1}{2(n^2 - n)} \left(\frac{\rho}{1 - \rho} \right)^2$$
$$\rho = \frac{Cov[C, Y]}{\sqrt{Var[C]} \sqrt{Var[Y]}}$$

522
523 where n denotes the number of cells.

524 Statistical Analyses

525 All statistical analyses were done using Python 3.9.6, Numpy 1.12.1⁵⁴, Scipy 1.6.3 and R 4.0.2.
526
527

528 Data and Code Availability

529 Next-generation sequencing data that support the findings of the study are available in the Gene
530 Expression Omnibus using accession code GSE188639.

531
532 The code that supports the findings of this study is available on Github Repository
533 (<https://github.com/szhao045/scMPRA>).

534
535

536 Acknowledgements

537 We thank the members of the Cohen laboratory for their critical feedback on the manuscript. We
538 thank Jess Hoistington-Lopez and MariaLynn Crosby for assistance with high-throughput
539 sequencing. This work is supported by grants to B.A.C from the National Institutes of Health,
540 R01 GM140711 and R01 GM092910.

541 Author Contributions

542 S.Z. and B.A.C. conceived and designed the project, S.Z. performed most of the experiments
543 and analyses with significant technical contributions from C.K.Y.H. and D.M.G., S.Z. and B.A.C.
544 wrote the manuscript with input and feedback from all authors.

545 Competing Interests

546 S.Z. and B.A.C. are inventors on a pending patent filed by Washington University in St. Louis
547 which may encompass the methods, reagents, and data disclosed in this manuscript. B.A.C is
548 on the scientific advisory board of Patch Biosciences.

549
550

551 References

- 552 1. Schaub, M. A., Boyle, A. P., Kundaje, A., Batzoglou, S. & Snyder, M. Linking disease
553 associations with regulatory information in the human genome. *Genome Res.* **22**, 1748–
554 1759 (2012).
- 555 2. Maurano, M. T. *et al.* Systematic localization of common disease-associated variation in
556 regulatory DNA. *Science* **337**, 1190–1195 (2012).
- 557 3. Hindorff, L. A. *et al.* Potential etiologic and functional implications of genome-wide
558 association loci for human diseases and traits. *Proc. Natl. Acad. Sci. U. S. A.* **106**, 9362–
559 9367 (2009).
- 560 4. Yang, J. *et al.* Common SNPs explain a large proportion of the heritability for human height.
561 *Nat. Genet.* **42**, 565–569 (2010).
- 562 5. Vattikuti, S., Guo, J. & Chow, C. C. Heritability and genetic correlations explained by
563 common SNPs for metabolic syndrome traits. *PLoS Genet.* **8**, e1002637 (2012).
- 564 6. Shi, H., Kichaev, G. & Pasaniuc, B. Contrasting the Genetic Architecture of 30 Complex
565 Traits from Summary Association Data. *Am. J. Hum. Genet.* **99**, 139–153 (2016).
- 566 7. Aygün, N. *et al.* Brain-trait-associated variants impact cell-type-specific gene regulation
567 during neurogenesis. *Am. J. Hum. Genet.* **108**, 1647–1668 (2021).
- 568 8. Nott, A. *et al.* Brain cell type–specific enhancer–promoter interactome maps and disease-
569 risk association. *Science* (2019).
- 570 9. Spielmann, M. & Mundlos, S. Looking beyond the genes: the role of non-coding variants in
571 human disease. *Hum. Mol. Genet.* **25**, R157–R165 (2016).

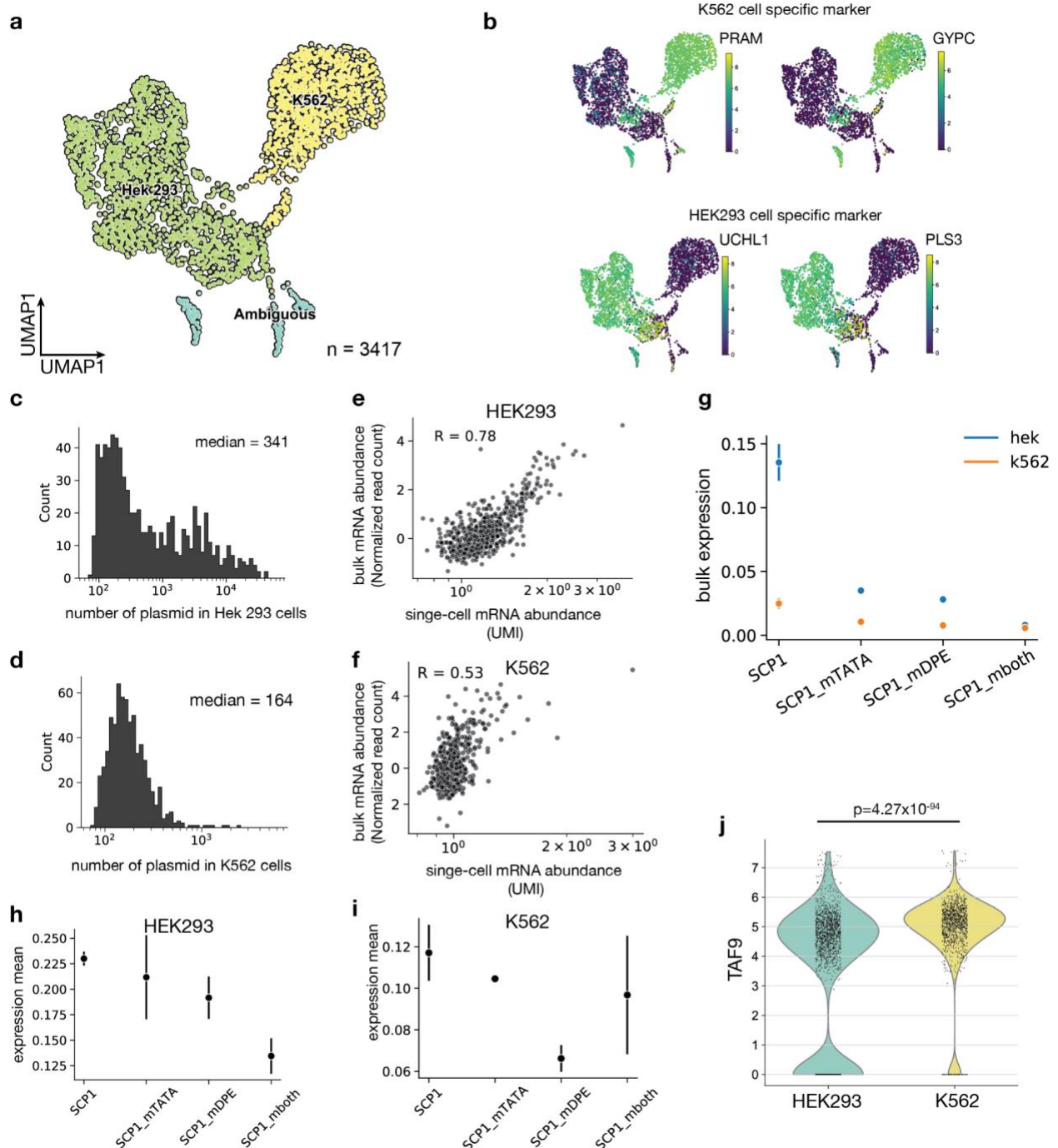
- 572 10. Zhang, F. & Lupski, J. R. Non-coding genetic variants in human disease. *Hum. Mol. Genet.*
573 **24**, R102–10 (2015).
- 574 11. Ong, C.-T. & Corces, V. G. Enhancer function: new insights into the regulation of tissue-
575 specific gene expression. *Nat. Rev. Genet.* **12**, 283–293 (2011).
- 576 12. Arnold, C. D. *et al.* Genome-wide quantitative enhancer activity maps identified by STARR-
577 seq. *Science* **339**, 1074–1077 (2013).
- 578 13. Kwasnieski, J. C., Mogno, I., Myers, C. A., Corbo, J. C. & Cohen, B. A. Complex effects of
579 nucleotide variants in a mammalian cis-regulatory element. *Proc. Natl. Acad. Sci. U. S. A.*
580 **109**, 19498–19503 (2012).
- 581 14. Ireland, W. T. *et al.* Deciphering the regulatory genome of *Escherichia coli*, one hundred
582 promoters at a time. *Elife* **9**, (2020).
- 583 15. Patwardhan, R. P. *et al.* Massively parallel functional dissection of mammalian enhancers
584 in vivo. *Nat. Biotechnol.* **30**, 265–270 (2012).
- 585 16. Sharon, E. *et al.* Inferring gene regulatory logic from high-throughput measurements of
586 thousands of systematically designed promoters. *Nat. Biotechnol.* **30**, 521–530 (2012).
- 587 17. Kinney, J. B., Murugan, A., Callan, C. G., Jr & Cox, E. C. Using deep sequencing to
588 characterize the biophysical mechanism of a transcriptional regulatory sequence. *Proc.*
589 *Natl. Acad. Sci. U. S. A.* **107**, 9158–9163 (2010).
- 590 18. Melnikov, A. *et al.* Systematic dissection and optimization of inducible enhancers in human
591 cells using a massively parallel reporter assay. *Nat. Biotechnol.* **30**, 271–277 (2012).
- 592 19. White, M. A. *et al.* A Simple Grammar Defines Activating and Repressing cis-Regulatory
593 Elements in Photoreceptors. *Cell Rep.* **17**, 1247–1254 (2016).
- 594 20. Kwasnieski, J. C., Fiore, C., Chaudhari, H. G. & Cohen, B. A. High-throughput functional
595 testing of ENCODE segmentation predictions. *Genome Res.* **24**, 1595–1602 (2014).
- 596 21. Chaudhari, H. G. & Cohen, B. A. Local sequence features that influence AP-1 cis-
597 regulatory activity. *Genome Res.* **28**, 171–181 (2018).

- 598 22. Hughes, A. E. O., Myers, C. A. & Corbo, J. C. A massively parallel reporter assay reveals
599 context-dependent activity of homeodomain binding sites in vivo. *Genome Res.* **28**, 1520–
600 1531 (2018).
- 601 23. Tewhey, R. *et al.* Direct Identification of Hundreds of Expression-Modulating Variants using
602 a Multiplexed Reporter Assay. *Cell* **165**, 1519–1529 (2016).
- 603 24. Hong, C. K. Y. & Cohen, B. A. Genomic environments scale the activities of diverse core
604 promoters. *bioRxiv* 2021.03.08.434469 (2021) doi:10.1101/2021.03.08.434469.
- 605 25. Haberle, V. *et al.* Transcriptional cofactors display specificity for distinct types of core
606 promoters. *Nature* **570**, 122–126 (2019).
- 607 26. Love, M. I., Huber, W. & Anders, S. Moderated estimation of fold change and dispersion for
608 RNA-seq data with DESeq2. *Genome Biol.* **15**, 550 (2014).
- 609 27. Amabile, G. *et al.* Dissecting the role of aberrant DNA methylation in human leukaemia.
610 *Nat. Commun.* **6**, 7091 (2015).
- 611 28. Juven-Gershon, T., Cheng, S. & Kadonaga, J. T. Rational design of a super core promoter
612 that enhances gene expression. *Nat. Methods* **3**, 917–922 (2006).
- 613 29. Shao, H. *et al.* Core promoter binding by histone-like TAF complexes. *Mol. Cell. Biol.* **25**,
614 206–219 (2005).
- 615 30. Shaffer, S. M. *et al.* Rare cell variability and drug-induced reprogramming as a mode of
616 cancer drug resistance. *Nature* **546**, 431–435 (2017).
- 617 31. Moudgil, A. *et al.* Self-Reporting Transposons Enable Simultaneous Readout of Gene
618 Expression and Transcription Factor Binding in Single Cells. *Cell* **182**, 992–1008.e21
619 (2020).
- 620 32. Litzenburger, U. M. *et al.* Single-cell epigenomic variability reveals functional cancer
621 heterogeneity. *Genome Biol.* **18**, 15 (2017).
- 622 33. Min, M. & Spencer, S. L. Spontaneously slow-cycling subpopulations of human cells
623 originate from activation of stress-response pathways. *PLoS Biol.* **17**, e3000178 (2019).

- 624 34. Bonnet, D. & Dick, J. E. Human acute myeloid leukemia is organized as a hierarchy that
625 originates from a primitive hematopoietic cell. *Nat. Med.* **3**, 730–737 (1997).
- 626 35. Ishikawa, F. *et al.* Chemotherapy-resistant human AML stem cells home to and engraft
627 within the bone-marrow endosteal region. *Nat. Biotechnol.* **25**, 1315–1321 (2007).
- 628 36. Chang, H. H., Hemberg, M., Barahona, M., Ingber, D. E. & Huang, S. Transcriptome-wide
629 noise controls lineage choice in mammalian progenitor cells. *Nature* **453**, 544–547 (2008).
- 630 37. Emert, B. L. *et al.* Variability within rare cell states enables multiple paths toward drug
631 resistance. *Nat. Biotechnol.* (2021) doi:10.1038/s41587-021-00837-3.
- 632 38. Foreman, R. & Wollman, R. Mammalian gene expression variability is explained by
633 underlying cell state. *Mol. Syst. Biol.* **16**, e9146 (2020).
- 634 39. Elowitz, M. B., Levine, A. J., Siggia, E. D. & Swain, P. S. Stochastic gene expression in a
635 single cell. *Science* **297**, 1183–1186 (2002).
- 636 40. Raj, A., Peskin, C. S., Tranchina, D., Vargas, D. Y. & Tyagi, S. Stochastic mRNA synthesis
637 in mammalian cells. *PLoS Biol.* **4**, e309 (2006).
- 638 41. Hilfinger, A. & Paulsson, J. Separating intrinsic from extrinsic fluctuations in dynamic
639 biological systems. *Proc. Natl. Acad. Sci. U. S. A.* **108**, 12167–12172 (2011).
- 640 42. Sherman, M. S., Lorenz, K., Lanier, M. H. & Cohen, B. A. Cell-to-cell variability in the
641 propensity to transcribe explains correlated fluctuations in gene expression. *Cell Syst* **1**,
642 315–325 (2015).
- 643 43. Fu, A. Q. & Pachter, L. Estimating intrinsic and extrinsic noise from single-cell gene
644 expression measurements. *Stat. Appl. Genet. Mol. Biol.* **15**, 447–471 (2016).
- 645 44. Chan, Y. K. *et al.* Engineering adeno-associated viral vectors to evade innate immune and
646 inflammatory responses. *Sci. Transl. Med.* **13**, (2021).
- 647 45. Byrne, L. C. *et al.* In vivo-directed evolution of adeno-associated virus in the primate retina.
648 *JCI Insight* **5**, (2020).
- 649 46. Wang, D., Tai, P. W. L. & Gao, G. Adeno-associated virus vector as a platform for gene

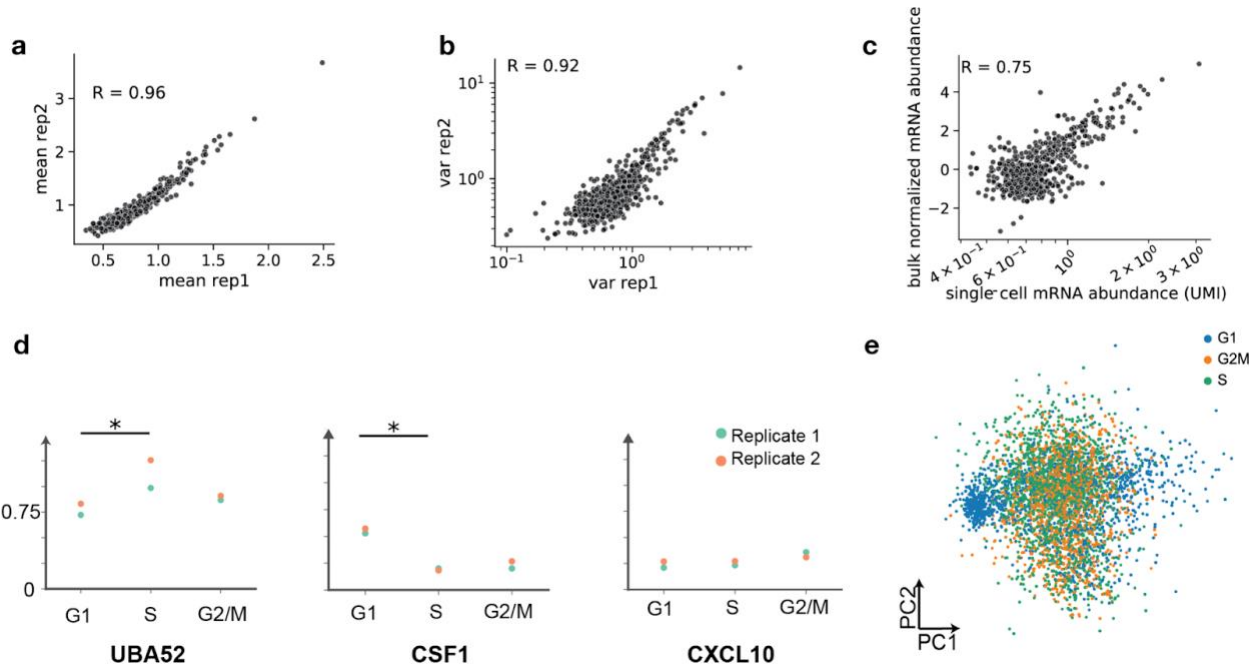
- 650 therapy delivery. *Nat. Rev. Drug Discov.* **18**, 358–378 (2019).
- 651 47. Bryant, D. H. *et al.* Deep diversification of an AAV capsid protein by machine learning. *Nat.*
652 *Biotechnol.* **39**, 691–696 (2021).
- 653 48. Cohen, R. N., van der Aa, M. A. E. M., Macaraeg, N., Lee, A. P. & Szoka, F. C., Jr.
654 Quantification of plasmid DNA copies in the nucleus after lipoplex and polyplex transfection.
655 *J. Control. Release* **135**, 166–174 (2009).
- 656 49. Wolf, F. A., Angerer, P. & Theis, F. J. SCANPY: large-scale single-cell gene expression
657 data analysis. *Genome Biol.* **19**, 15 (2018).
- 658 50. Traag, V. A., Waltman, L. & van Eck, N. J. From Louvain to Leiden: guaranteeing well-
659 connected communities. *Sci. Rep.* **9**, 5233 (2019).
- 660 51. Tirosh, I. *et al.* Dissecting the multicellular ecosystem of metastatic melanoma by single-cell
661 RNA-seq. *Science* **352**, 189–196 (2016).
- 662 52. Zabidi, M. A. *et al.* Enhancer–core-promoter specificity separates developmental and
663 housekeeping gene regulation. *Nature* **518**, 556–559 (2014).
- 664 53. Bailey, T. L. & Gribskov, M. Combining evidence using p-values: application to sequence
665 homology searches. *Bioinformatics* **14**, 48–54 (1998).
- 666 54. Harris, C. R. *et al.* Array programming with NumPy. *Nature* **585**, 357–362 (2020).

667 **Supplementary Figures**



668
669
670
671
672
673
674
675
676

Supplementary Figure 1. scMPRA measures cell-type specific CRS activity. (a) UMAP of the single-cell transcriptome from the mixed-cell experiment. 105 out of 3417 cells (3%) are labeled by both K562 and HEK293 cell genes. (b) UMAP of the mixed-cell experiment with cells marked by other representative markers for K562 and HEK293 cell expression. (c-d) Histogram of the number of plasmids transfected to K562 cells and HEK293 cells. (e-f) Scatterplot of bulk RNA-seq expression against expression mean from scMPRA (Pearson R for K562 cells: 0.53, Pearson R for HEK293 cells: 0.78). (g) Dot plot of the reporter activity of SCP1 and its mutants from bulk RNA-seq data (error bar: 1 s.d.). (h) Dot plot of the mean reporter activity of SCP1 and its mutants from scMPRA experiment for K562 cells. (i) Dot plot of the mean reporter activity of SCP1 and its mutants from scMPRA experiment for HEK293 cells. (j) Violin plot showing the expression distribution of TAF9 in K562 and HEK293 cells. (Wilcoxon rank sum test, $p = 4.27 \times 10^{-94}$).



677

678

679

680

681

682

683

684

685

686

687

688

689

690

691

692

693

694

695

696

697

698

699

700

701

702

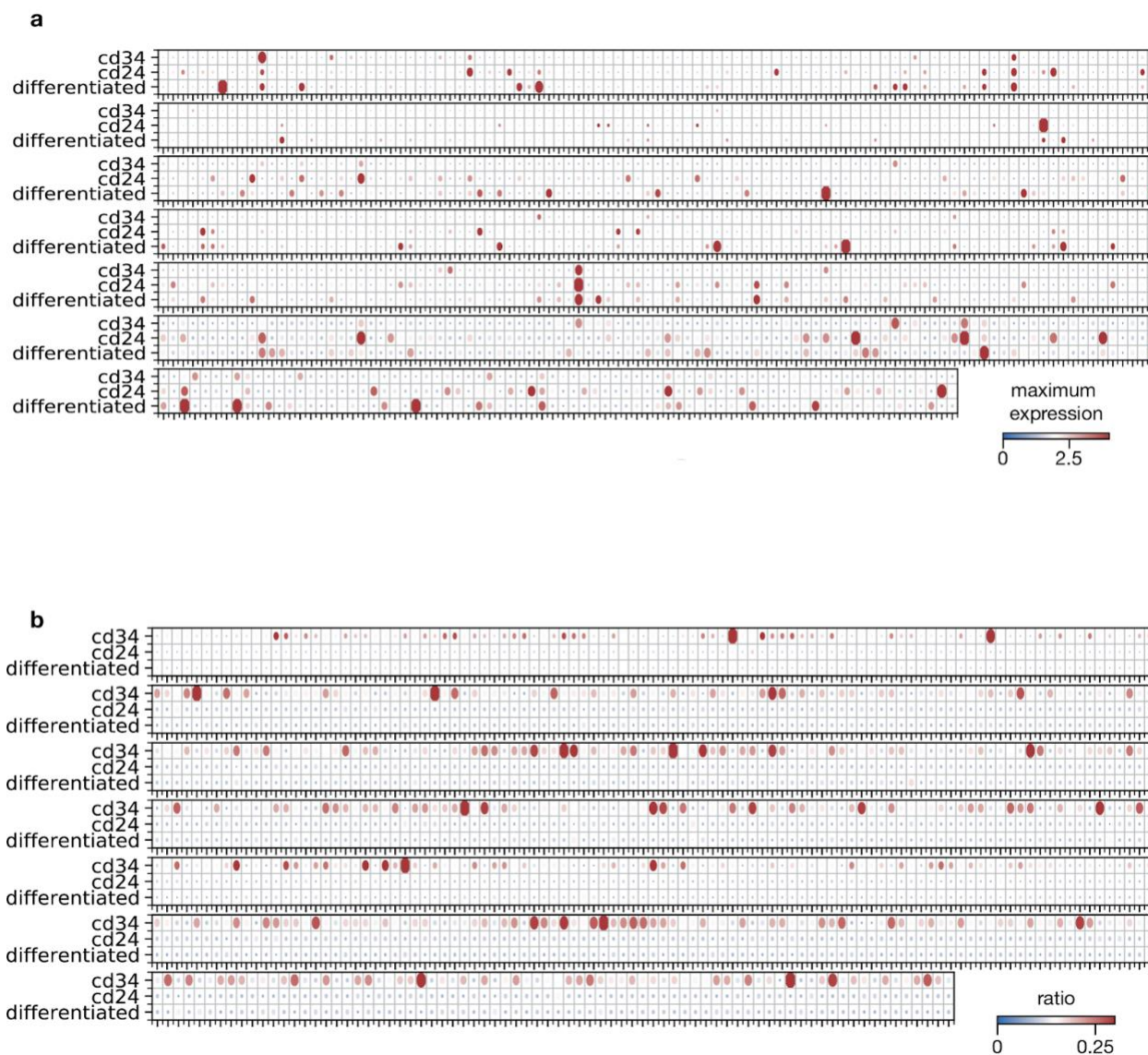
703

704

705

706

Supplementary Figure 2 scMPRA measures CRS activity in K562 cell substates. (a,b) Reproducibility for expression mean and cell-to-cell variance (Pearson Correlation for mean: 0.96, for variance: 0.92). (c) Scatterplot of reproducibility of scMPRA mean expression with bulk MPRA measurement using UMI (Pearson Correlation: 0.75). (d) Different dynamics of expression. For UBA52, the promoter is most highly expressed in S phase; whereas for CSF1, the promoter is most highly expressed in G1 phase. For CXCL10, the promoter is expressed evenly through cell cycle (Stars indicate significance from Wilcoxon rank sum test, *: $p < 0.05$.) (e) Cells no longer cluster together based on cell cycle genes after normalization.



707
708 **Supplementary Figure 3. CD34⁺/CD38⁻ substate changes the probability of cells having higher expression, not the**
709 **maximum expression level. (a)** Dot plot showing the maximum single-cell expression for the core promoter library in
710 CD34⁺/CD38⁻, CD24⁺, and Differentiated clusters. Color and size both indicate the maximum expression change. **(b)** Dot plot
711 showing the percentage of cells in CD34⁺/CD38⁻, CD24⁺, and Differentiated clusters that are in the 90th percentile of expression
712 level per promoter. Color and size both indicate the ratio change.

713
714
715
716
717
718
719