

## Genomic context sensitivity of insulator function

André M. Ribeiro-dos-Santos<sup>1,3</sup>, Megan S. Hogan<sup>1,3</sup>, Raven Luther<sup>1</sup>, Matthew T. Maurano<sup>1,2,4</sup>

<sup>1</sup> Institute for Systems Genetics, NYU Grossman School of Medicine, New York, NY 10016, USA.

<sup>2</sup> Department of Pathology, NYU Grossman School of Medicine, New York, NY 10016, USA.

<sup>3</sup> These authors contributed equally to this work.

<sup>4</sup> Corresponding author: [maurano@nyu.edu](mailto:maurano@nyu.edu)

## Abstract

Compartmentalization of interactions between genomic regulatory elements and potential target genes is influenced by the binding of insulator proteins such as CTCF, which act as potent enhancer blockers when interposed between an enhancer and a promoter in a reporter assay. But only a minority of CTCF sites genome-wide function as bona fide insulators, depending on cellular and genomic context. To dissect the influence of genomic context on enhancer blocker activity, we integrated reporter constructs with promoter-only, promoter and enhancer, and enhancer blocker configurations at hundreds of thousands of genomic sites using the Sleeping Beauty transposase. Deconvolution of reporter activity by genomic position revealed strikingly different patterns of reporter function, including a compartment of enhancer blocker reporter integration sites with robust expression. The high density of integration sites permits quantitative delineation of characteristic genomic context sensitivity profiles, and their decomposition into sensitivity to both local and distant DNaseI hypersensitive sites. Furthermore, a single-cell expression approach permits direct linkage of reporters integrated into the same clonal lineage with differential endogenous gene expression, we observe that CTCF insulator activity does not completely abrogate reporter effects on endogenous gene expression. Collectively, our results lend new insight to genomic regulatory compartmentalization and its influence on the determinants of promoter-enhancer specificity.

## Introduction

Insulators are a class of genomic regulatory elements that block interaction of enhancers with their cognate promoters (Phillips and Corces 2009). The insulator hypothesis offers an attractive paradigm for understanding regulatory specificity in mammalian genomes through the delineation of regulatory domains. Historically, insulator function has been primarily defined in an ectopic or reporter context, although insulator function has been identified at various endogenous sites such as the *Igf2/H19* locus (Bell and Felsenfeld 2000). Enhancer blocker activity is canonically defined by a reporter assay which interposes a candidate insulator element between a weak promoter and an enhancer (Chung et al. 1993), while barrier insulators protect transgenes from silencing due to spreading of heterochromatin (West et al. 2002). Insulators have also been employed to counter genotoxicity from transgene enhancer activation of endogenous oncogenes (Li et al. 2009; Liu et al. 2015). Known insulators such as the canonical chicken  $\beta$ -globin hypersensitive site 4 element are composite elements with enhancer blocker, barrier, and other activities (Dickson et al. 2010), and often have secondary functions, such as silencers (Qi et al. 2015).

The architectural protein CTCF is the only known vertebrate insulator protein and its binding can confer a potent enhancer blocking effect (Phillips and Corces 2009). Additionally, binding sites for CTCF co-localize with genomic features such as topologically associated domain boundaries, but direct functional analysis of these sites is impeded by the difficulty of genome engineering at relevant scales. While binding affinity and recognition sequence orientation appear to confer some specificity to CTCF sites involved in domain organization, this remains inadequate to explain the activity of ~100,000 cell-type specific CTCF sites genome-wide and to what extent specificity is conferred by nearby binding sites for other factors (Maurano et al. 2015; Guo et al. 2015), resulting in a major gap in our understanding of the sequence determinants of genomic regulatory architecture. Stably integrated reporter assays have shed light on the mechanics of insulator function, such methods typically do not assess interaction with the surrounding endogenous genomic elements (Walters et al. 1999). In contrast, integrated barcoded reporter assays offer the potential to directly assess reporter response to genomic landscape (Akhtar et al. 2013; Moudgil et al. 2020).

Here we describe a high-throughput, randomly integrated barcoded reporter platform to analyze insulator activity in varied genomic contexts. We developed an enhancer blocker construct interposing potent CTCF insulator elements (Liu et al. 2015) between a  $\beta$ -globin HS2 enhancer (HS2) and  $\gamma$ -globin promoter. Barcoded reporters with or without insulator elements were randomly integrated into the genome of cultured K562 cells using the Sleeping Beauty transposase system, and subsequently mapped to enable barcode-specific readout of genomic context effects. We find that reporters with and without insulator elements are distinguished by characteristic response signatures to genomic context. Finally, we employ single cell RNA-seq to link cells deriving from the same initial clone, assess the potential for interference by nearby reporter insertions, and link specific integrations to perturbations on target genes.

## Results

### Characterization of enhancer blocker reporter

We developed and characterized a series of reporter constructs based on well-characterized genomic regulatory elements including the murine  $\gamma$ -globin promoter and murine  $\beta$ -globin locus control region (LCR) hypersensitive site 2 (HS2) enhancer. A potent insulator element (A1 or C1) previously identified through an analysis of highly occupied CTCF sites in the human genome (Liu et al. 2015) (**Supplemental Fig. S1**) was interposed between the promoter and enhancer in an enhancer blocker position (**Fig. 1**). Reporter expression drove a PuroGFP fusion protein to enable selection and/or measurement of transcriptional activity on a cellular level. The reporter was flanked by Sleeping Beauty (Mátés et al. 2009) inverted terminal repeats (ITRs) to enable transposition into the genome. Reporter plasmids were transiently co-transfected with a plasmid expressing SB100X, a highly active variant of the Sleeping Beauty transposase (Mátés et al. 2009).

We first characterized the activity of several different classes of reporters based on this scaffold, including GGlo (promoter-only), GGlo+HS2 (promoter and enhancer), and Ins+GGlo+Ins+HS2 (enhancer blocker) reporters. K562 erythroleukemia cells were transfected with reporter plasmid and SB100X transposase plasmid (**Supplemental Table S1**). Reporter activity was characterized using flow cytometry to measure the proportion of GFP<sup>+</sup> cells (**Fig. 1**). Ins+GGlo+Ins+HS2 reporters showed low expression, comparable to promoter-only GGlo constructs, confirming that the CTCF site acts as an enhancer blocker. An insulator element truncated to just the core 54 bp of the CTCF recognition sequence showed similar enhancer blocker activity (**Fig. 1, Supplemental Fig. S1**). Insulator elements showed no capacity to augment transcription on their own without promoter or enhancer, suggesting that the ITRs do not interfere with reporter function. Finally, CTCF effect on reporter activity was orientation-independent. These results confirm the readout of our enhancer blocker reporter assay.

### Reporter activity in genomic context

As flow cytometry assesses single-cell GFP activity representing the sum of all reporters integrated in that cell, it does not reflect the activity of individual insertion sites. We developed a strategy to deconvolute the activity of individual reporters in mixed culture using unique 16 nt reporter barcode (BC) sequences. We generated three types of libraries based on this reporter BC strategy (**Fig. 2A, Supplemental Fig. S2, Supplemental Fig. S3**): the genomic location of

reporter integration sites were mapped using inverse PCR (iPCR), their representation was determined using DNA libraries, and their expression was assessed with RNA libraries. Sequencing libraries were constructed using a two-stage nested PCR to add Illumina adapters. We incorporated a 8-12 nt unique molecular identifier (UMI) (Jee et al. 2016) to permit targeted single-molecule counting.

We performed a series of 4 independent experiments using the GGlo, GGlo+HS2, and Ins+GGlo+Ins+HS2 constructs (**Supplemental Table S1**). After growth under puromycin selection for 8-11 days, 2-4 replicate DNA, RNA, and iPCR libraries were generated for each experiment (**Supplemental Table S2**) and sequenced to saturation. Replicate libraries exhibited high consistency. Activity averaged across all insertion sites recapitulated cellular activity (**Fig. 1, Supplemental Fig. S4A**). Each experiment averaged 26,765 insertion sites analyzed after quality control (**Fig. 2B, Supplemental Table S3**). These data yielded high-resolution maps of reporter activity, with an average distance between reporters of 23-48 kb, and an average distance from DHS to reporter 9-19 kb. Libraries were merged and analyzed together, yielding 308,664 insertion sites (**Table 1; Supplemental Table S3**).

Examination of the  $\beta$ -globin (**Fig. 2C**) and *MYC* (**Fig. 2D**) loci demonstrated notable differences in patterns of reporter activity. GGlo exhibited variable activity that was highly responsive to local genomic context: at the  $\beta$ -globin locus, its activity was concentrated tightly around the endogenous genes; at *MYC*, activity localized to two separate regions around the gene body and distal ALL enhancers. GGlo+HS2 showed more variable insertion location and site-specific activity. We observed peak activity at a subset of regions, but insertions were depleted over several regions, including the DHS cluster immediately downstream of *MYC*, suggesting that the GGlo+HS2 construct does not support expression at these regions, or that those insertions have a negative effect on growth. Ins+GGlo+Ins+HS2 instead exhibited more uniform activity regardless of surrounding genomic context and showed reduced position preference throughout the window.

To systematically assess the length scale of sensitivity of different reporters to genomic context, we computed the correlation in activity between insertions at adjacent sites (**Fig. 3A, Supplemental Fig. S4B**). GGlo+HS2 showed the lowest correlation across all size ranges, suggesting the least influence of genomic context; GGlo and Ins+GGlo+Ins+HS2 showed high correlation at short-range (<5 kb) but diverged at longer range: GGlo correlation dropped to nearly zero while

Ins+GGlo+Ins+HS2 remained high beyond 50 kb. To provide an easily computed metric reflecting the contribution of genomic context at different distance scales, we computed the number of DHSs within 5 kb and 100 kb of the reporter/insertion site. We then used a linear model to systematically quantify the effect of these indicators of genomic context on activity of reporters of different classes (**Fig. 3B**). Genomic context offered distinctively lower predictive power for GGlo+HS2, and all three reporters showed distinct contributions of short and long-range genomic context (**Fig. 3C-D**).

### **Clonal analysis using integrated barcodes**

Droplet-based single cell RNA-seq (scRNA-seq) approaches provide the compartmentalization needed to associate reporter BCs integrated in the same cell. We adapted our integrated reporter assay to the 10x Genomics scRNA-seq platform and performed a pilot experiment using Ins+GGlo+Ins+HS2 (**Fig. 4A**; **Supplemental Table S4**). We generated transcriptomic scRNA-seq libraries, as well as amplicon-targeted libraries to enrich for reporter transcripts in individual cells. This showed that, when an individual reporter was present in multiple cells, its activity was highly reproducible; instead control comparisons with different reporters in the same cells or with permuted data showed little correlation (**Fig. 4B**).

The integrated reporter BCs provide a unique combinatorial genetic identifier for cells derived from a given clone during transfection (Bidy et al. 2018). To identify cells deriving from the same clone and impute reporter presence and flanking gene expression levels across cells, we developed a clonal inference approach (**Fig. 4C**).

Given this success, we performed a scaled-up experiment using 3 different classes of reporter constructs (**Supplemental Table S4**), including GGlo+HS2, Ins+GGlo+Ins+HS2 (in replicate), and Ins+GGlo+HS2+Ins (intended to test the activity of a reporter fully flanked by insulator elements). Given that each transfection is distinguished by a distinct set of reporter BCs, we pooled the cells from four independent transfections for generation of scRNA-seq libraries and super-loaded to maximize power. Reporter BCs specific to each individual transfection were identified using the DNA, RNA, and iPCR libraries from bulk cells. Deconvolution without using transfection labels yielded only 4.75% of clones harboring BCs from two independent transfections (**Fig. 4D**). Reporter BC labels were used to prune conflicting cells and clones (**Fig. 4E**). Clones contained a

median of 2 cells (**Fig. 4F**) and a median of 4 reporter BCs (**Fig. 4G**). This showed that the distance between insertions in a given cell was sufficient to enable independent readout of hundreds or thousands of reporters (**Fig. 4H**).

### **scRNA-seq integration analysis**

The compartmentalization provided by the single-cell readout enables direct linkage of specific insertion events to their effect on adjacent genes (Gasperini et al. 2019). To permit direct correlation between reporter activity and effect on the local genomic landscape, we developed an analysis approach accounting for the sparsity of single-cell data to identify nearby genes with significant differential expression (**Fig. 5A, Methods**).

We conducted a power simulation to estimate the ability of our statistical framework to detect expression perturbations. We simulated a dataset using a negative binomial distribution based on gene expression values observed on the actual experimental results (**Methods**). These results showed good power to detect expression perturbations with  $\text{abs}(\text{fold change}) > 2$  and average expression above 1 UMI (**Supplemental Fig. S5**). To ensure that tests were well powered to detect expression changes, we ignored reporters that perturbed fewer than 3 cells, genes with overall average expression inferior to 0.05 UMIs, and genes with average count less than 10 UMIs (experiment 4) or less than 5 UMIs (experiment 5) among perturbed or unperturbed cells.

We then assessed significant differential expression linked to reporter insertions. We pooled analysis results from all Ins+GGlo+Ins+HS2 experiments (T0190, T0221, T0222). Ins+GGlo+Ins+HS2 and Ins+GGlo+HS2+Ins showed a concentration of significant tests at close range to the TSS (**Fig. 5C**). We observed that insertions of all three reporter classes were more likely to affect gene expression of TSSs in the same TAD (**Fig. 5D**) and when the reporter was inserted in the gene body itself (**Fig. 5E**). Collectively, these results underscore the role for genomic context in dictating the effect of ectopically delivered regulatory elements on endogenous gene expression.

### **Discussion**

A key challenge in recognizing functional genomic regulatory variation is the specificity of enhancer-promoter interactions. Indeed, genomic context is a key predictive feature of models for recognizing functional regulatory variation (Halow et al. 2021), yet there remains no systematic

and mechanistic approach to incorporating context as a feature. Our work suggests that assessment of genomic regulatory element function can consider multiple orthogonal axes, including expression (i) level, (ii) consistency (stochastic vs. deterministic), and (iii) sensitivity to the local and/or long-range regulatory landscape. We show how these properties differ when assayed in a cellular vs. site-specific context. Our approach provides a ready platform for large-scale characterization of classes of genomic regulatory elements along these lines, and their incorporation into existing models of functional regulatory variation (Halow et al. 2021).

Although our enhancer blocking insulator reporters demonstrate a strong effect on a cellular level, we show that these reporters actually demonstrate a range of expression levels depending on genomic context. Indeed, the high reporter expression at some genomic integration sites implies a total abrogation of insulator function. A strict definition of insulator function is difficult to reconcile with these results, or with the abundance of CTCF sites in the genome. Instead, our results are more consistent with a model of insulator function which moderates but does not eliminate sensitivity to genomic context.

While our assay shows high technical reproducibility, correlation of reporters independently inserted at close range (<500 bp) reaches only  $R=0.6$  (**Fig. 3A**). This suggests that, after reporter sequence content and genomic context, additional epigenetic and stochastic factors play a strong secondary role. Transcriptional enhancement is an inherently stochastic process, and single-cell approaches show enhancers increase the frequency of a given cell undergoing transcription rather than augmenting the transcriptional rate of a given cell (Weintraub 1988; Walters et al. 1996). Consistent with this, ectopic CTCF sites often bind CTCF but do not always form loops as measured by 4C (Redolfi et al. 2019). It is possible that more complex, composite regulatory elements at key genomic loci might contain additional functional elements that would reduce their variability of expression even in an ectopic context. For example, multiple tandem CTCF sites have been shown to increase the durability of insulation (Huang et al. 2020). We expect that our approach will enable further dissection of the interplay between sequence, genomic context, and single-cell behavior in the future.



## Methods

### Plasmid cloning and barcoding

pCMV(CAT)T7-SB100 (SB100X) and pT2/LTR7-GFP were gifts from Zsuzsanna Izsvak (Addgene plasmids #34879 and #62541, respectively) (Mátés et al. 2009).

The  $\gamma$ -lobin promoter, mouse  $\beta$ -globin hypersensitive site 2 (HS2) enhancer, A1 insulator, A2 insulator, C1 insulator, A1 Core, and C1 Core DNA fragments (**Supplemental Table S5**) were synthesized by Genscript USA (Piscataway, NJ). All plasmids used in this study are listed in **Supplemental Table S6**.

Sleeping Beauty reporter constructs used in this study were barcoded using a Gibson Assembly approach prior to introduction into K562 cells (**Supplemental Fig. S2**). The plasmid backbone to be barcoded was PCR amplified using pTR-GibsonBC-FW and pTR-GibsonBC-RV primers, and the correct length fragment was purified from a 1% agarose gel. Next, Gibson Assembly was performed using the amplified plasmid backbone and a synthesized DNA fragment “GibsonBC4” according to the manufacturer’s protocol (NEB cat# E2611L). Barcoded plasmid library DNA was purified using the Zymo Clean and Concentrate-5 (Zymo Research cat# D4014) protocol prior to transformation. Purified barcoded plasmid DNA was transformed into electrocompetent MegaX DH10B-T1 bacteria (Fisher cat# C640003) using an Eppendorf 2510 electroporator set to 1800 V. After recovering for 1 h at 37 °C, transformation reactions were transferred to 50 mL LB Media with 100  $\mu$ g/mL Ampicillin and incubated at 37 °C for 16 h shaking at 220 RPM. Barcoded plasmid library DNA was purified using the ZymoPure II Plasmid Maxiprep kit protocol and quantified on a Nanodrop.

### Cell culture and transfection

K562 cells were obtained from ATCC (ATCC CCL-243) and cultured in RPMI 1640 medium with glutamine (Fisher cat# MT10040CV) supplemented with 10% FBS (Gemini Bio-Products cat# 100-106), 1 mM sodium pyruvate, and 10 U/mL penicillin-streptomycin. Cultures were maintained at 37 °C and 5% CO<sub>2</sub>, and were subcultured once cultures reached a density of 5x10<sup>5</sup> cells/mL.

1x10<sup>6</sup> K562 cells were transfected using the ThermoFisher Neon Transfection System 100  $\mu$ L Kit according to the manufacturer's instructions with varying amounts of transposon and transposase (**Supplemental Table S1**). Cells were transfected with 4  $\mu$ L TE to use as a negative control for puromycin selection. Transfected cells were selected with puromycin (2.5  $\mu$ g/mL). K562 media with puromycin was replaced every 2 days. Cell counts were performed either using PrestoBlue (Life Technologies cat# A13261) and fluorescence detection with the Synergy H1 Multi-Mode Microplate Reader, or were stained with trypan blue and counted on a hemocytometer.

### Flow cytometry of GFP Expression Assays

On day 8 after transfection, GFP expression was measured using the SONY SH800S Cell Sorter. For each experiment, a 100  $\mu$ M chip and the Optical Filter Pattern 2 were used, the 405 nm, 488

nm and 561 nm lasers were enabled, automatic color compensation was turned off, and sensor gain settings were set to the following values: forward scatter (FSC) = 3, back scatter (BSC) = 30.5%, and FL2 (GFP) = 36.5%.

Using FlowJo software (v10.7.2), single live cells were gated using side scatter (SSC) and FSC values from a TE (mock) transfected cell sample. GFP expression data were plotted on a histogram of unit area vs. GFP fluorescence (525±50 nm). The GFP-negative cell population was defined using the GFP expression of TE (mock) transfected single-cells. The GFP-positive cell population was defined by cells that had a greater GFP fluorescence than TE (mock) transfected cells.

### **Genomic DNA Purification**

11-14 days post-transfection, cell pellets containing  $3 \times 10^6$  to  $4 \times 10^6$  cells each were snap frozen in LN<sub>2</sub> and stored at -80°C until DNA extraction. Cell pellets were allowed to warm to room temperature, and then were resuspended in 385  $\mu$ L DNA Quick Extract (Lucigen cat# QE09050) and transferred to a 1.5 mL tube. Cells were incubated at 65 °C for 15 min, followed by 98 °C for 5 min. After cooling briefly, 10  $\mu$ L Proteinase K (Sigma-Aldrich cat# P4850-5ML) was added, and cell lysate was incubated at 55 °C overnight. The following day, 5  $\mu$ L RNase A (Sigma-Aldrich cat# R4642-50MG) was added, and the cell lysate was incubated at 37 °C for 30 min. Genomic DNA was precipitated by adding 4  $\mu$ L Glycoblue (Fisher cat# AM9515), 40  $\mu$ L 3M Sodium Acetate, and 1 mL ice-cold 100% ethanol. After incubating at -80 °C for 1 h, DNA was pelleted by centrifugation at 20,000 g for 30 min at 4 °C. The DNA pellet was washed twice with 70% ethanol, and then resuspended in 200  $\mu$ L Buffer EB (Qiagen cat# 19086).

### **RNA Purification**

11-14 days post-transfection, cell pellets containing  $1 \times 10^6$  cells each were resuspended in 350  $\mu$ L Trizol solution (Fisher cat# 15596026) and stored at -80 °C until RNA extraction. Frozen samples were allowed to warm to room temperature, and then 350  $\mu$ L cell solution was transferred to a Phase-Lock gel tube (Fisher cat# NC1093153). 70  $\mu$ L chloroform was added to each Phase-lock tube and shaken vigorously, followed by a 2 min incubation at room temperature. Tubes were centrifuged at 12,000 x g for 10 min at 4 °C. Following centrifugation, the aqueous phase was decanted from each tube and transferred to a new tube. 350  $\mu$ L 70% ethanol was added and mixed well, and the solution was transferred to a Qiagen RNeasy-mini spin column. Samples were centrifuged at 13,000 x g for 15 s, and the flow-through was discarded. 350  $\mu$ L Buffer RW1 was added to each column, samples were centrifuged at 13,000 x g for 15 s, and the flow-through was discarded. This Buffer RW1 wash was repeated once more for a total of 2 washes. Next, 500  $\mu$ L Buffer RPE was added to each column, samples were centrifuged at 13,000 x g for 15 s, and the flow-through was discarded. This Buffer RPE wash was repeated once more for a total of 2 washes. After the last RPE wash, the column was centrifuged for an additional 2 min at 13,000 x g to remove residual ethanol. Samples were eluted in 40  $\mu$ L RNase-free H<sub>2</sub>O.

To ensure that the RNA preparation was DNA-free, we utilized the Ambion TURBO DNA-free kit protocol (Thermo Fisher Scientific cat# AM1907). Following DNase treatment, RNA was transferred to a fresh tube, and the concentration was quantified on the Nanodrop.

### **Amplicon Library Preparation**

For DNA libraries, unique Molecular Identifiers (UMIs) and the inner portion of the P5 sequencing adapter were added. 20  $\mu\text{g}$  of genomic DNA was digested with PstI (NEB cat# R3140L) for 1 h at 37 °C, and then purified using the Zymo Clean and Concentrate-25 (Zymo Research cat# D4034) protocol. One cycle of PCR was performed with the following conditions: 8 replicate 50  $\mu\text{L}$  reactions were prepared, each containing 500 ng PstI digested DNA, 1x Phusion Hot Start Flex Mastermix (NEB cat# M0536L), and 200 nM of the primer P5\_Plasmid\_8N/9N/10N, and incubated at 98 °C for 5 min, 60 °C for 1 min, and 72 °C for 10 min. Replicate reactions were combined and then purified using the Zymo Clean and Concentrate-5 (Zymo Research cat# D4014) protocol, eluting the DNA in 20  $\mu\text{L}$ .

For RNA libraries, cDNA was synthesized using the Superscript IV First Strand Synthesis kit (Invitrogen), with 5  $\mu\text{g}$  RNA template and 2  $\mu\text{M}$  primer P5\_barcode\_0N/1N/2N (containing a truncated sequencing adapter) in two replicate reactions per sample. RNA was first incubated with primers and dNTPs at 60 °C for 10 min, then placed on ice for one min. The remaining RT reagents were added, and samples were incubated at 55 °C for 10 min, 80 °C for 10 min, and then cooled to 4 °C. Next, 1  $\mu\text{L}$  RNaseH was added to each reaction, and incubated at 37 °C for 20 min. Single-stranded cDNA was purified using the Zymo Clean and Concentrate-5 protocol (Zymo Research cat# D4014), using 7 volumes of DNA binding buffer and eluting in 10  $\mu\text{L}$  Zymo DNA elution buffer. Unique Molecular Identifiers (UMIs) and the inner portion of the P7 sequencing adapter were added to each single-stranded cDNA molecule using 1 cycle of PCR with the following conditions: 2 replicate 50  $\mu\text{L}$  reactions were prepared, each containing 5  $\mu\text{L}$  cDNA, 1x Phusion Hot Start Flex Mastermix (NEB cat# M0536L), and 200 nM of the primer P7\_Plasmid\_8N/9N/10N, and incubated at 98 °C for 5 min, 64 °C for 5 min, and 72 °C for 5 min. Replicate reactions were combined and then purified using the Zymo Clean and Concentrate-5 (Zymo Research cat# D4014) protocol, eluting the DNA in 20  $\mu\text{L}$ .

For inverse PCR (iPCR) libraries, 40  $\mu\text{g}$  genomic DNA was digested with DpnII for 2 h at 37 °C. Digested DNA was purified using the Zymo Clean & Concentrate-25 column protocol, and digestion was verified by running 100 ng DpnII digested DNA out on a 1% agarose gel. Intramolecular DpnII ligation was performed using DpnII digested DNA at a concentration of 5  $\mu\text{g}/\text{mL}$ , and T4 DNA ligase at a concentration of 10,000 U/mL. Ligation reactions were incubated overnight at 4 °C, and ligation products were purified using the Zymo Clean & Concentrate-25 column protocol.

DNA, RNA, and iPCR libraries then amplified using a nested PCR approach to add full Illumina sequencing adapters in two stages. To add the inner P5 and P7 sequencing adapters (DNA and RNA samples already had P5 or P7 added, respectively), samples were amplified for 20-30 PCR cycles. 8 replicate 50  $\mu\text{L}$  reactions were prepared, each containing 2  $\mu\text{L}$  DNA, 1x Phusion Hot

Start Flex Mastermix (NEB cat# M0536L), 200 nM of the appropriate P5 and P7 primers for each library type (**Supplemental Table S5**), and incubated 1 cycle at 98 °C for 5 min; 20-30 cycles (sample dependent) at 98 °C for 15 s, 55 °C for 15 s, and 72 °C for 30 s; and 1 cycle of 72 °C for 10 min. Replicate reactions were combined and then purified using the Zymo Clean and Concentrate-5 (Zymo Research cat# D4014) protocol, eluting the DNA in 20  $\mu$ L.

The remaining (outer) adapter sequences with indexing barcodes were added to each library using 10 cycles of PCR with the following conditions: one 50  $\mu$ L reaction was prepared per library, each containing 1  $\mu$ L DNA purified from the previous round of PCR, 1x Phusion Hot Start Flex Mastermix (NEB cat# M0536L), 200 nM of each indexed P5 and P7 primers (e.g. P5\_amplicon\_S502 and P7\_Amplicon\_N704, **Supplemental Table S5**), and incubated 1 cycle at 98 °C for 5 min, 10 cycles at 98 °C for 15 s, 71 °C for 15 s, and 72 °C for 30 s, and 1 cycle of 72 °C for 10 min. Final DNA libraries were purified using the Zymo Clean and Concentrate-5 (Zymo Research cat# D4014) protocol, eluting the DNA in 20  $\mu$ L. Completed libraries were quantified using the Qubit dsDNA HS (Fisher cat# Q32851) kit protocol.

### Single Cell RNA-seq (scRNA-seq) Library Preparation

Cells were transfected as described above (**Supplemental Table S1**). To enrich for cells that received the transposase construct, mKate-positive cells were sorted into new plates 24 h after transfection using a Sony SH800 cell sorter as described above, except that the 665 $\pm$ 30 nm optical filter was used to gate for mKate fluorescence, and expanded.

scRNA-seq expression libraries were generated using the 10x Chromium NextGem Single Cell 3' workflow (10X Genomics cat# 1000128). For experiment 4, an additional sort for GFP-positive cells was performed on day 4, and 5700 cells were collected for scRNA-seq (**Supplemental Table S4**) while the remaining cells were expanded in culture for an additional 7 d, at which point cell pellets were collected for RNA, DNA, and iPCR libraries.

For experiment 5, cell pellets were collected for RNA, DNA, and iPCR libraries and cells were frozen 14 days post-transfection. After thawing cells and expanding for 2 days, cells were collected from each of 4 separate transfections, and pooled for scRNA-seq libraries (**Supplemental Table S4**). Additional pellets were collected post-thaw for additional RNA, DNA, and iPCR libraries.

Separate libraries enriched for reporter transcripts were generated from the cDNA produced in the Post GEM-RT Cleanup & cDNA Amplification step of the 10X Chromium Single Cell 3' (v3.1) library protocol using PCR with the following conditions: 8 replicate 25  $\mu$ L reactions were prepared, each containing 1  $\mu$ L amplified 10X Chromium Single Cell cDNA, 1x Phusion Hot Start Flex Mastermix (NEB cat# M0536L), 500 nM of the primer P5\_Halfsite, and 500 nM of the primer P7\_10xSBbarcodeV2\_0N, and incubated 1 cycle at 98 °C for 5 min; 8 cycles (sample dependent) at 98 °C for 15 s, 66 °C for 15 s, and 72 °C for 30 s; and 1 cycle of 72 °C for 10 min. Replicate

reactions were combined and then purified using the Zymo Clean and Concentrate-5 (Zymo Research cat# D4014) protocol, eluting the DNA in 20  $\mu$ L. Completed indexed adapter sequences were added to each library during a final 10 cycles of PCR using the conditions described in DNA Library Preparation above. Completed libraries were quantified using the Qubit dsDNA HS (Fisher cat# Q32851) kit protocol.

### Sequencing and analysis

Illumina libraries were generated and sequenced on an Illumina NextSeq 500. Reads were demultiplexed by a standard pipeline using Illumina bcl2fastq v2.20 requiring a perfect match to indexing BC sequences.

DNA, RNA, iPCR, and enriched scRNA-seq libraries were processed by a custom pipeline (**Supplemental Table S2**). Read pairs whose sequence comprised >75% G bases were dropped. PCR primer sequence was removed and UMI and cell barcodes (cellBC) were extracted using UMI-tools v1.0.1 (Smith et al. 2017) including the option `--quality-filter-threshold=30`. Reporter barcodes (BCs) were extracted based on position from read pairs matching the expected template sequence with fewer than 10% mismatched bases. BCs were required to have 2 bases or fewer with base quality score below 30. Reporter BCs, cellBCs, and UMIs were each deduplicated using a directed adjacency approach based on that of UMI-tools (Smith et al. 2017).

iPCR libraries were trimmed to remove plasmid sequences, including the potential for digestion at a secondary DpnII site using cutadapt v2.9 (Martin 2011). Reads were then mapped to a hg38 reference genome augmented with transposon and sleeping beauty sequences using BWA v.0.7.12 (Li and Durbin 2009a). Libraries where read 1 was sequenced to 24 bp or more beyond the end of the plasmid sequence were mapped in paired-end mode using BWA-MEM with -Y option. Otherwise read 2 was mapped in single-end mode using BWA aln and samse. Reads without reporter BCs, aligned with insertions or deletions, with >10% mismatch rate, with mapping quality <10, or with >1 kb between mates (for paired mapping) were excluded from further analysis. The integration insertion site was defined as the 5' mapping site of read 2. Reporter BCs were additionally deduplicated using coordinates to group. Sites with the same reporter BC within  $\pm 5$  bp of the same BC were collapsed. Integrations with <2 reads, representing <1% of the total coverage at a given genomic position, or BCs found at multiple sites were excluded.

Read counts for DNA, RNA, and iPCR libraries were normalized per 1M sequenced reads and merged based on reporter BC. Missing RNA counts were imputed as 0, and only BCs with >10 DNA reads and an integration site were considered. Reporter activity was computed as  $\log_2(\text{RNA}/\text{DNA} + 1)$ .

DNase-seq data (K562-DS9764) for K562 was downloaded from <https://www.encodeproject.org> and processed using a standard pipeline (<https://github.com/mauranolab/mapping/tree/master/dnase>). DNase I hypersensitive sites were identified using hotspot v1 (John et al. 2011)

hotspot peaks (1% FDR). CTCF binding sites were taken from previously published work (Maurano et al. 2015).

### **Clonal inference**

Cells and reporter BCs deriving from a single initial transfected clone were derived from the enriched scRNA-seq libraries. First, we constructed a bipartite graph whose nodes were cellBC and reporter BCs, connected by edges weighted by the pair's UMI count. Edges with fewer than 2 UMI were dropped.

For Experiment 5, where scRNA-seq data was generated from a superloaded pool of 4 independent transfections, each reporter BC was labelled by its known transfection based on the union of all DNA/RNA/iPCR data. BCs found in more than one transfection were removed from the graph. Edges connecting a cellBC to a reporter BC from a transfection representing <80% of the cell's total UMI were trimmed. Nodes directly connected to two different transfections (i.e. doublets or reporter BC collisions) were dropped.

To reduce the impact of chimeric PCR artifacts, edges representing <10% of total UMI for a given reporter BC or <2% for a given cellBC. Reporter BCs mapping to multiple integration sites or found in multiple transfections were removed. Finally, edges bridging two independent sets of nodes representing <10% of either set's total UMI were pruned to reduce doublets or PCR artifacts. Unconnected nodes were pruned. The remaining connected communities were defined as clones.

### **scRNA-seq analysis**

scRNA-seq 3' libraries were analyzed using Cell Ranger v.4.0.0 (Zheng et al. 2017). A reference was constructed against hg38 and transposon sequences as above using. Ensembl release 93 was used for gene annotations. We obtained the gene expression matrix, whitelist of non-empty cellBCs, blacklist of poor quality cellBCs with few UMIs and too many pSB reads. Only cellBCs contained in the whitelist and absent from these blacklists were considered on further analysis.

### **Reporter effects on gene expression**

We explored the reporter impact on genes whose TSS lay within 250 kb from a reporter. For each reporter and gene, we compared the expression on the set of perturbed cells (those belonging to the reporter's clone) against all other cells. Only genes with Ensembl category of protein\_coding or lincRNA were considered. Only cells included in both single-cell expression and clone assignments were used. To avoid potential confounding from nearby reporter insertions in the same clone, we discarded any reporters with a second reporter within 500 kb. Finally, tests with <3 perturbed cells, average target gene expression in the perturbed or unperturbed cells of <10 UMI (experiment 4) or <5 UMI (experiment 5), or overall average expression <0.05 UMI were excluded from the analysis.

For the differential expression analysis, we modeled reporter effect using a negative binomial (or Gamma-Poisson) regression with regularized dispersion estimate:

$$(1) Y \sim NB(\mu, \theta)$$
$$(2) \log(\mu) = \log(\text{depth}) + \beta_0 + \beta_c X$$

Where  $Y$  are the observed counts for a particular gene across all cells,  $\mu$  is the expected average gene UMIs,  $\theta$  is the gene UMI distribution dispersion,  $\beta_0$  and  $\beta_c$  are the regression coefficients,  $\text{depth}$  is each cell's UMIs, and  $X$  is an indicator vector that is 1 if the cell belongs to the reporter clone (perturbed) or 0 if it does not (unperturbed).

In order to estimate the distribution dispersion ( $\theta$ ) of each gene, we employed the approach of Hafemeister & Satija (Hafemeister and Satija 2019) of fitting a Poisson regression ( $Y \sim \text{Pois}(e^{\beta_0 + \log(\text{depth})})$ ) for a random subset of 2000 genes and estimating  $\theta$  using maximum likelihood. Then, we expanded the estimation to all genes with average expression  $\geq 0.05$  UMIs by fitting a kernel regression of  $\theta$  in relation to the gene average expression.

Perturbation ( $\beta_c$ ) significance was estimated by two-tailed p-value based on Student's t-distribution. Storey's q-value approach was used for multiple testing correction (Storey and Tibshirani 2003) for each transfection individually. Tests with q-value  $< 0.10$  were considered significant.

We performed simulations to estimate detection power over a range of clone sizes ( $n$ ) and effect sizes using equation 1 above with a fixed dispersion ( $\theta$ ) of 100. For each simulation condition, we generated UMI counts for 1000 genes, and the same number of cells as the actual scRNA-seq dataset. Expected UMIs ( $\mu$ ) ranged from the minimum and maximum observed values in the actual scRNA-seq dataset in a logarithm scale. For each gene,  $n$  cells were sampled and their UMI counts altered by the defined effect size. The resulting simulations were evaluated by our analysis algorithm given only the simulated count matrix and cells assignment as perturbed and unperturbed.

Contact matrices were taken from (Rao et al. 2014) GSE63525\_K562\_intrachromosomal\_contact\_matrices.tar.gz), KR normalized, and Armatous v2.2 (Filippova et al. 2014) was used to identify TADs with  $\gamma=0.5.0$  and a resolution of 5 kb and lifted over to hg38 using liftOver.

### Software availability

Code used in sequencing data processing is available at GitHub (<https://github.com/mauranolab/mapping/tree/master/transposon>)

### Competing interest statement

The authors declare no competing interests.

## **Acknowledgements**

This work was partially funded by NIH grant R35GM119703 to M.T.M..



## References

- Akhtar W, de Jong J, Pindyurin AV, Pagie L, Meuleman W, de Ridder J, Berns A, Wessels LFA, van Lohuizen M, van Steensel B. 2013. Chromatin position effects assayed by thousands of reporters integrated in parallel. *Cell* **154**: 914–927.
- Bell AC, Felsenfeld G. 2000. Methylation of a CTCF-dependent boundary controls imprinted expression of the *Igf2* gene. *Nature* **405**: 482–485.
- Biddy BA, Kong W, Kamimoto K, Guo C, Waye SE, Sun T, Morris SA. 2018. Single-cell mapping of lineage and identity in direct reprogramming. *Nature* **564**: 219–224.
- Chung JH, Whiteley M, Felsenfeld G. 1993. A 5' element of the chicken beta-globin domain serves as an insulator in human erythroid cells and protects against position effect in *Drosophila*. *Cell* **74**: 505–514.
- Dickson J, Gowher H, Strogantsev R, Gaszner M, Hair A, Felsenfeld G, West AG. 2010. VEZF1 elements mediate protection from DNA methylation. *PLoS Genetics* **6**: e1000804.
- Filippova D, Patro R, Duggal G, Kingsford C. 2014. Identification of alternative topological domains in chromatin. *Algorithms Mol Biol* **9**: 14.
- Gasperini M, Hill AJ, McFaline-Figueroa JL, Martin B, Kim S, Zhang MD, Jackson D, Leith A, Schreiber J, Noble WS, et al. 2019. A Genome-wide Framework for Mapping Gene Regulation via Cellular Genetic Screens. *Cell* **176**: 1516.
- Guo Y, Xu Q, Canzio D, Shou J, Li J, Gorkin DU, Jung I, Wu H, Zhai Y, Tang Y, et al. 2015. CRISPR Inversion of CTCF Sites Alters Genome Topology and Enhancer/Promoter Function. *Cell* **162**: 900–910.
- Hafemeister C, Satija R. 2019. Normalization and variance stabilization of single-cell RNA-seq data using regularized negative binomial regression. *Genome Biol* **20**: 296.
- Halow JM, Byron R, Hogan MS, Ordoñez R, Groudine M, Bender MA, Stamatoyannopoulos JA, Maurano MT. 2021. Tissue context determines the penetrance of regulatory DNA variation. *Nat Commun* **12**: 2850.
- Huang H, Zhu Q, Jussila A, Han Y, Bintu B, Kern C, Conte M, Zhang Y, Bianco S, Chiariello A, et al. 2020. CTCF Mediates Dosage and Sequence-context-dependent Transcriptional Insulation through Formation of Local Chromatin Domains. *bioRxiv* 2020.07.07.192526.
- Jee J, Rasouly A, Shamovsky I, Akivis Y, Steinman SR, Mishra B, Nudler E. 2016. Rates and mechanisms of bacterial mutagenesis from maximum-depth sequencing. *Nature* **534**: 693–696.
- John S, Sabo PJ, Thurman RE, Sung M-H, Biddie SC, Johnson TA, Hager GL, Stamatoyannopoulos JA. 2011. Chromatin accessibility pre-determines glucocorticoid receptor binding patterns. *Nature Genetics* **43**: 264–268.

- Li CL, Xiong D, Stamatoyannopoulos G, Emery DW. 2009. Genomic and functional assays demonstrate reduced gammaretroviral vector genotoxicity associated with use of the cHS4 chromatin insulator. *Mol Ther* **17**: 716–724.
- Li H, Durbin R. 2009a. Fast and accurate short read alignment with Burrows-Wheeler transform. *Bioinformatics* **25**: 1754–1760.
- Li H, Durbin R. 2009b. Fast and accurate short read alignment with Burrows-Wheeler transform. *Bioinformatics* **25**: 1754–1760.
- Liu M, Maurano MT, Wang H, Qi H, Song C-Z, Navas PA, Emery DW, Stamatoyannopoulos JA, Stamatoyannopoulos G. 2015. Genomic discovery of potent chromatin insulators for human gene therapy. *Nature Biotechnology* **33**: 198–203.
- Martin M. 2011. Cutadapt removes adapter sequences from high-throughput sequencing reads. *EMBnet.journal* **17**: 10–12.
- Mátés L, Chuah MKL, Belay E, Jerchow B, Manoj N, Acosta-Sanchez A, Grzela DP, Schmitt A, Becker K, Matrai J, et al. 2009. Molecular evolution of a novel hyperactive Sleeping Beauty transposase enables robust stable gene transfer in vertebrates. *Nature Genetics* **41**: 753–761.
- Maurano MT, Wang H, John S, Shafer A, Canfield T, Lee K, Stamatoyannopoulos JA. 2015. Role of DNA Methylation in Modulating Transcription Factor Occupancy. *Cell reports* **12**: 1184–1195.
- Moudgil A, Wilkinson MN, Chen X, He J, Cammack AJ, Vasek MJ, Lagunas T, Qi Z, Lalli MA, Guo C, et al. 2020. Self-Reporting Transposons Enable Simultaneous Readout of Gene Expression and Transcription Factor Binding in Single Cells. *Cell* **182**: 992-1008.e21.
- Phillips JE, Corces VG. 2009. CTCF: master weaver of the genome. *Cell* **137**: 1194–1211.
- Qi H, Liu M, Emery DW, Stamatoyannopoulos G. 2015. Functional validation of a constitutive autonomous silencer element. *PLoS One* **10**: e0124588.
- Rao SSP, Huntley MH, Durand NC, Stamenova EK, Bochkov ID, Robinson JT, Sanborn AL, Machol I, Omer AD, Lander ES, et al. 2014. A 3D map of the human genome at kilobase resolution reveals principles of chromatin looping. *Cell* **159**: 1665–1680.
- Redolfi J, Zhan Y, Valdes-Quezada C, Kryzhanovska M, Guerreiro I, Iesmantavicius V, Pollex T, Grand RS, Mulugeta E, Kind J, et al. 2019. DamC reveals principles of chromatin folding in vivo without crosslinking and ligation. *Nat Struct Mol Biol* **26**: 471–480.
- Smith T, Heger A, Sudbery I. 2017. UMI-tools: modeling sequencing errors in Unique Molecular Identifiers to improve quantification accuracy. *Genome Res* **27**: 491–499.
- Storey JD, Tibshirani R. 2003. Statistical significance for genomewide studies. *Proceedings of the National Academy of Sciences of the United States of America* **100**: 9440–9445.

Walters MC, Fiering S, Bouhassira EE, Scalzo D, Goeke S, Magis W, Garrick D, Whitelaw E, Martin DI. 1999. The chicken beta-globin 5'HS4 boundary element blocks enhancer-mediated suppression of silencing. *Mol Cell Biol* **19**: 3714–3726.

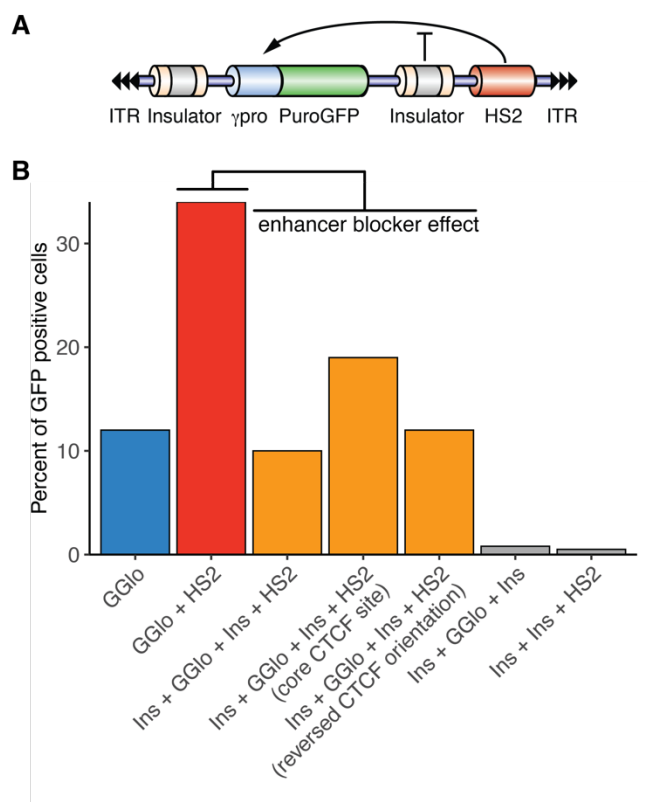
Walters MC, Magis W, Fiering S, Eidemiller J, Scalzo D, Groudine M, Martin DI. 1996. Transcriptional enhancers act in cis to suppress position-effect variegation. *Genes Dev* **10**: 185–195.

Weintraub H. 1988. Formation of stable transcription complexes as assayed by analysis of individual templates. *Proceedings of the National Academy of Sciences of the United States of America* **85**: 5819–5823.

West AG, Gaszner M, Felsenfeld G. 2002. Insulators: many functions, many mechanisms. *Genes Dev* **16**: 271–288.

Zheng GXY, Terry JM, Belgrader P, Ryvkin P, Bent ZW, Wilson R, Ziraldo SB, Wheeler TD, McDermott GP, Zhu J, et al. 2017. Massively parallel digital transcriptional profiling of single cells. *Nat Commun* **8**: 14049.

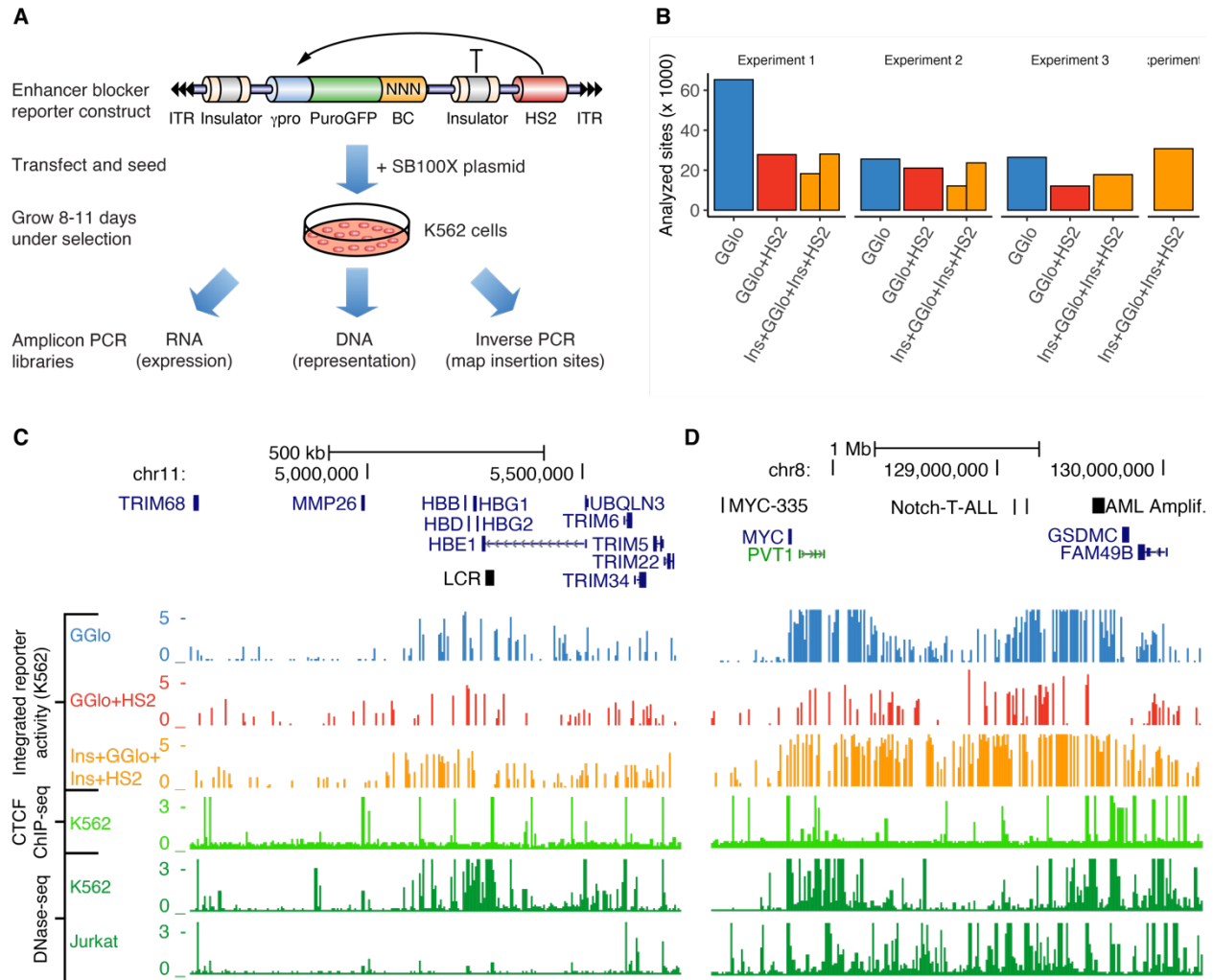
## Figures



**Fig. 1. Cellular activity of enhancer blocker reporter.**

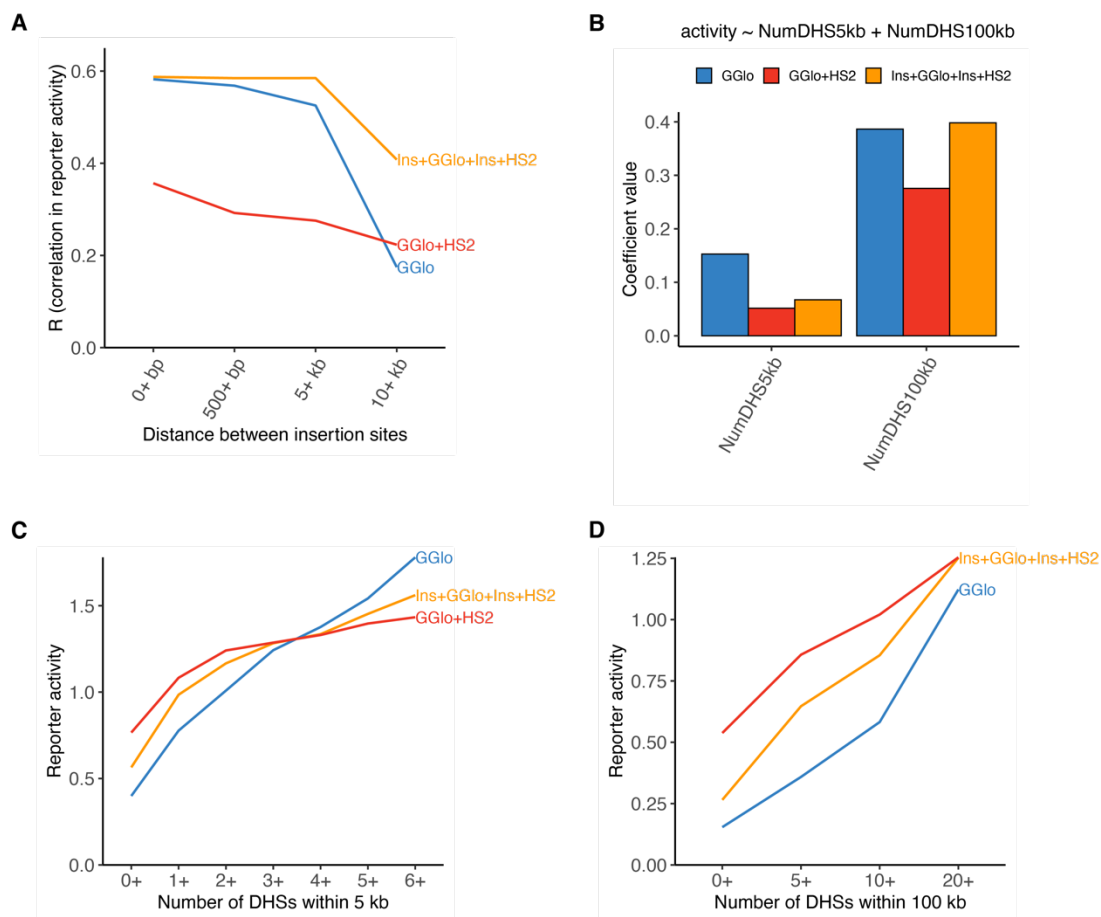
(A) Enhancer blocker reporter scheme consisting of  $\gamma$ -globin promoter driving PuroGFP expression. Reporter expression is reduced when an intervening CTCF site acts as an enhancer blocker to reduce effect of HS2 enhancer.  $\gamma$ pro,  $\gamma$ -globin promoter; HS2,  $\beta$ -globin hypersensitive site 2 enhancer; ITR, Sleeping Beauty inverted terminal repeats.

(B) Specified reporter plasmids were transfected along with a plasmid expressing the Sleeping Beauty SB100X transposase to enable random genomic integration. Activity was measured by flow cytometry. GGlo,  $\gamma$ -globin promoter; Ins, Insulator; HS2,  $\beta$ -globin hypersensitive site 2 enhancer.



## Fig. 2. Site-specific activity of enhancer blocker activity.

(A) Reporter plasmid is co-transfected with plasmid expressing SB100X transposase. PCR-based Illumina library construction enables highly quantitative measurement of hundreds of thousands of barcodes simultaneously. Insertion sites are mapped in multiplex using in-verse PCR; DNA libraries are used to normalize for barcode representation, and RNA libraries to quantify barcode expression.  $\gamma$ pro,  $\gamma$ -globin promoter; BC, unique barcode; HS2, B-globin hypersensitive site 2 enhancer; ITR, Sleeping Beauty inverted terminal repeats. (B) Counts of sites analyzed for 4 experiments of reporters containing promoter only (GGlo), promoter and HS2 enhancer (GGlo+HS2), or with CTCF site interposed between GGlo and HS2 (Ins+GGlo+Ins+HS2). (C-D) Analysis of enhancer-blocker functionality at the *HBB* (C) and *MYC* (D) loci. Top tracks show activity for reporters. Shown are data merged from all experiments. Bottom tracks show CTCF ChIP-seq data for K562 erythroleukemia cells, and DNase-seq data for K562 and Jurkat T-cell leukemia cells. (C) LCR, locus control region (D) Downstream enhancer cluster is highlighted including Notch-T-ALL (Acute Lymphocytic Leukemia) and AML (acute myeloid leukemia) amplified region.

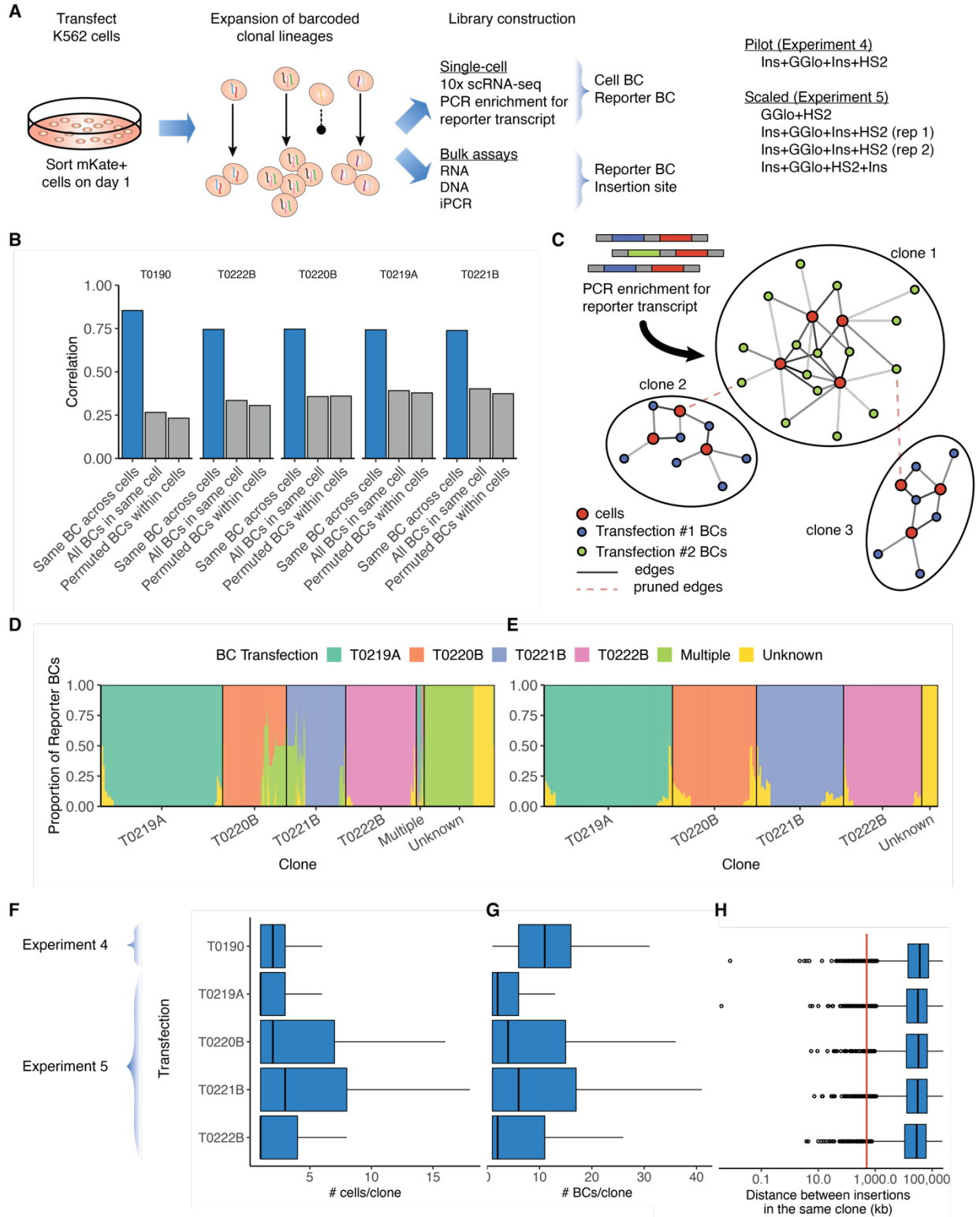


**Fig. 3. Quantitative assessment of genomic context effects on enhancer reporter activity.**

(A) Correlation in activity for nearby insertions by reporter class. Data is merged across all experiments; individual experiments are shown in **Supplemental Fig. S4B**.

(B) Linear regression coefficients for regression of reporter activity on density of local and long-range genomic context.

(C-D) Reporter activity by short-range and long-range genomic context represented by the number of DHSs within 5 kb (C) or 100 kb (D).



**Fig. 4. scRNA-seq inference of clonal relationship of reporter BCs.**

(A) Overview of scRNA-seq experiment.

(B) High correlation of same BCs measured across multiple cells vs. control of unrelated BCs in same cells or permutation analysis.

(C) Graph-based inference of clones from scRNA-seq data using reporter BC to link cells deriving from the same clone.

(D-E) deconvolution of Experiment 5. Clones are shown along the X-axis, grouped by inferred transfection; Multiple, reporter BCs detected in multiple transfusions; None, reporter BCs detected only in scRNA-seq data. Label-free deconvolution shown in (D); final deconvolution removing conflicting clones shown in (E). Excess of Reporter BCs from multiple transfusions in T0220B/T0221B clones results from sub-optimal diversity in the Ins+GGlo+Ins+HS2 plasmid library.

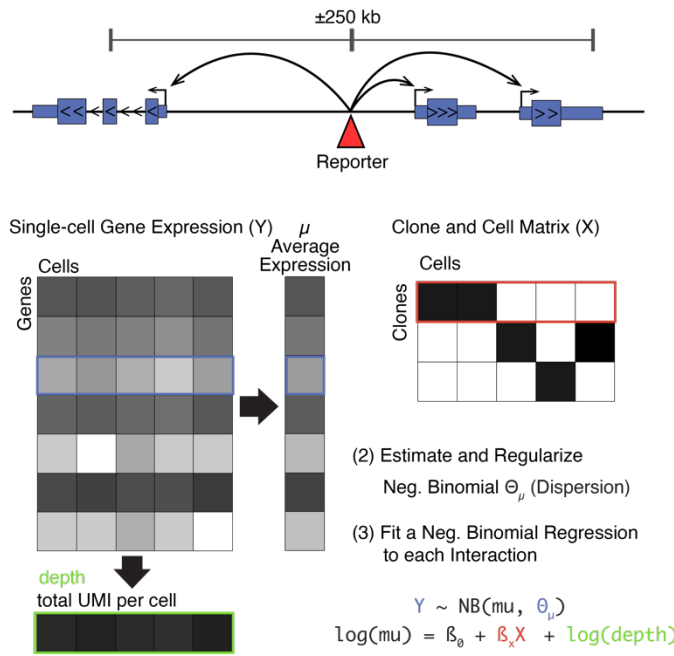
(F) Number of cells per clone

(G) Number of Reporter BCs per clone

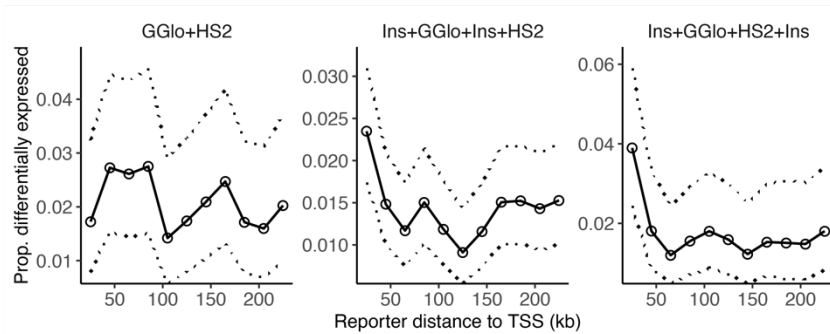
(H) Distance between reporters integrated in cells derived from the same clone. Vertical red line at 500 kb indicates insertions considered far enough to be assumed independent.



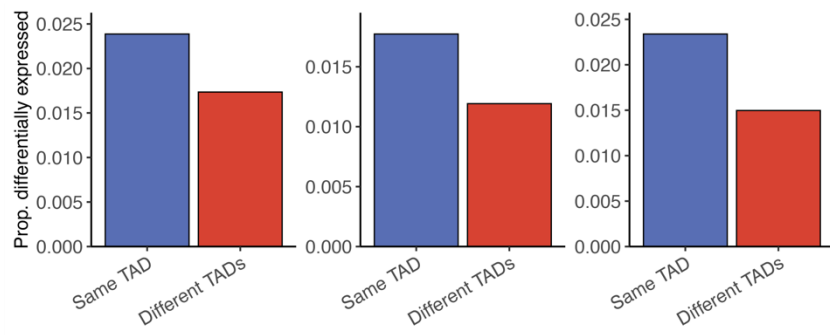
**A** (1) Reporter effect on nearby gene expression



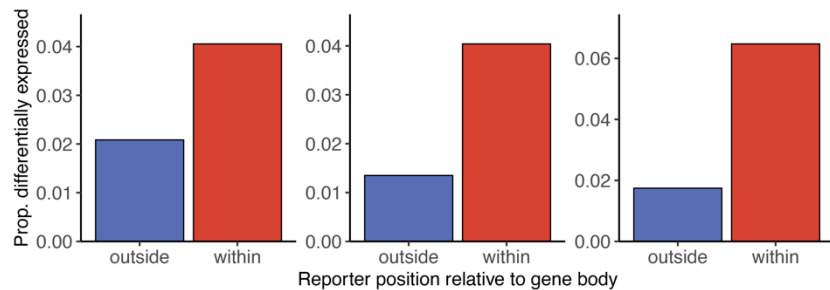
**C**



**D**



**E**



**Fig. 5. Analysis of reporter effect on nearby endogenous gene expression.**

(A) Schematic of analysis framework.

(C-E) Rate of significant effect on gene expression by (C) distance from reporter and TSS (50 kb sliding window and a 25 kb step) (D) whether reporter and gene are in the same TAD, and (E) whether reporter is inside or outside gene body.

## Tables

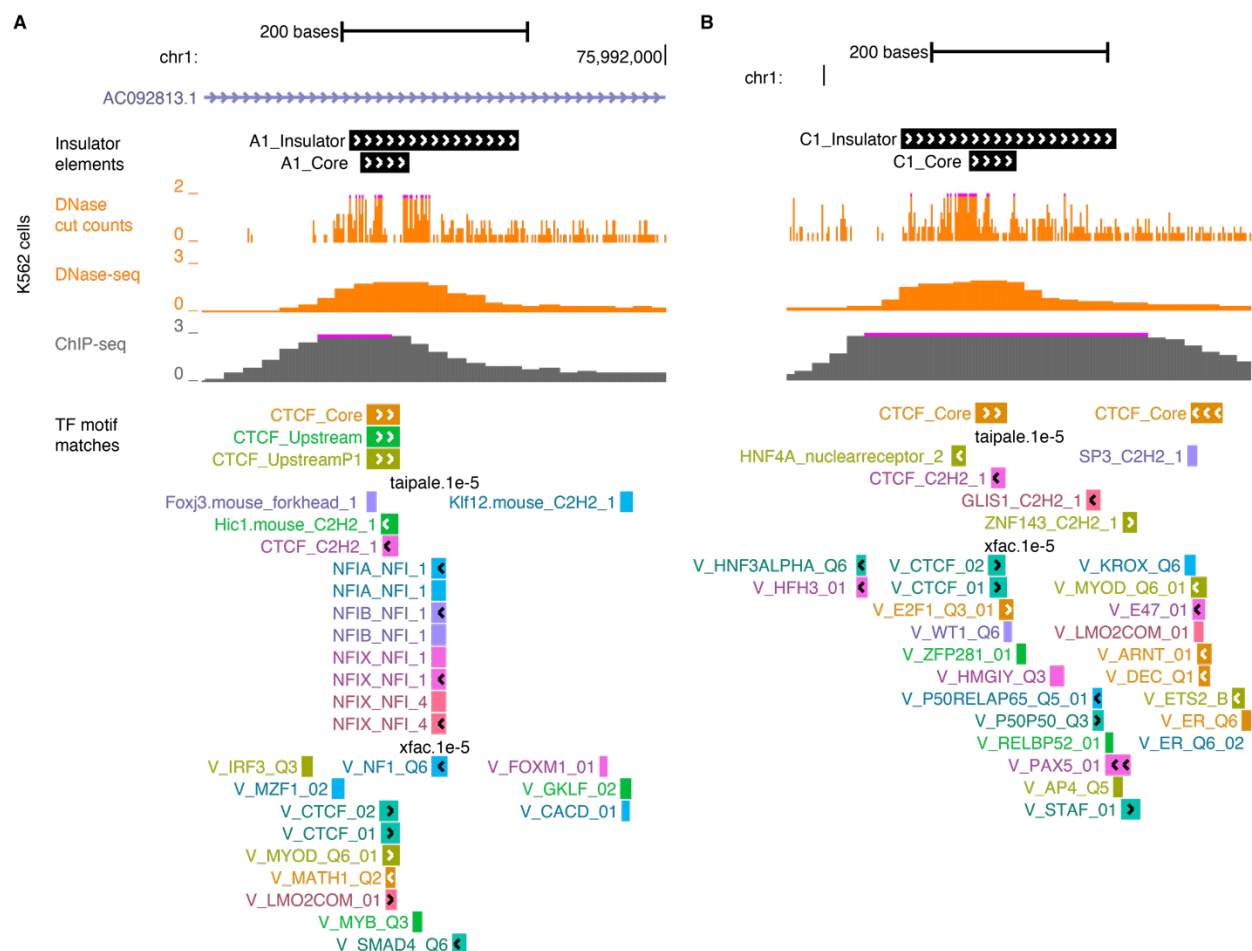
Construct	Experiment 1	Experiment 2	Experiment 3	Experiment 4	Merged
GGlo	65,296	25,588	26,484	-	117,196
GGlo+HS2	27,851	21,060	12,089	-	60,941
Ins+GGlo+Ins+HS2	46,379	35,805	17,796	30,795	130,527

**Table 1. Summary of insertions analyzed per reporter construct and experiment.**

Summary of transfections across four independent experiments. Counts are of insertions passing all QC filters.

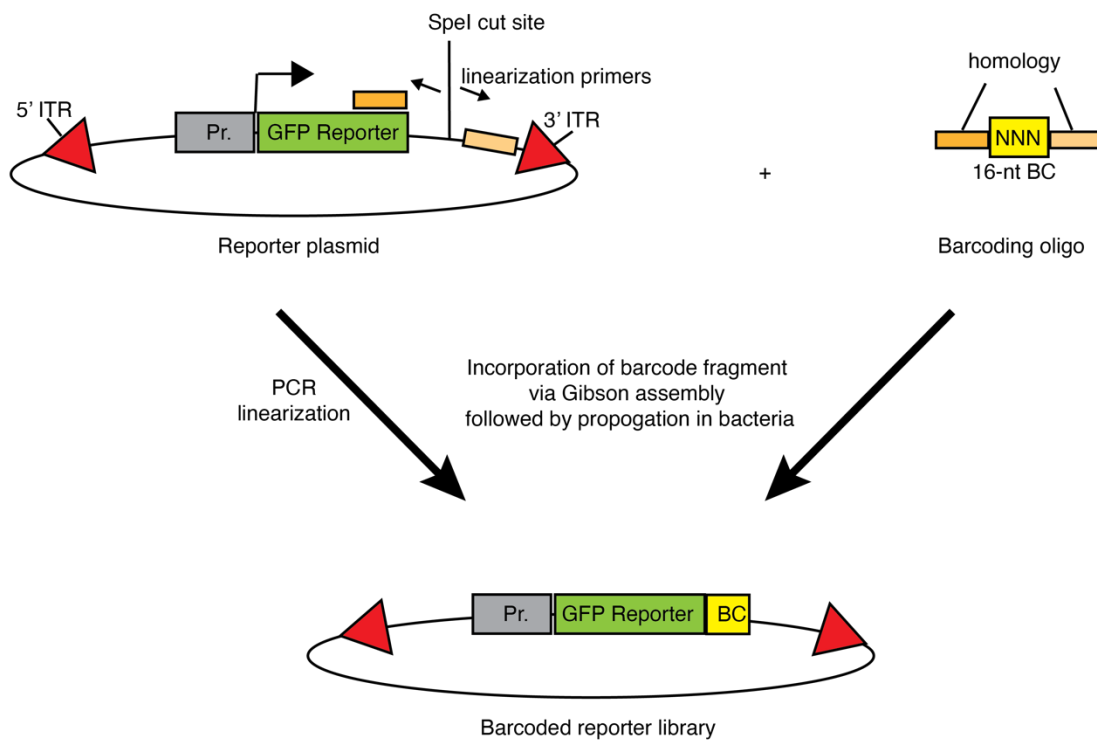
## Supplemental Material

### Supplemental Figures



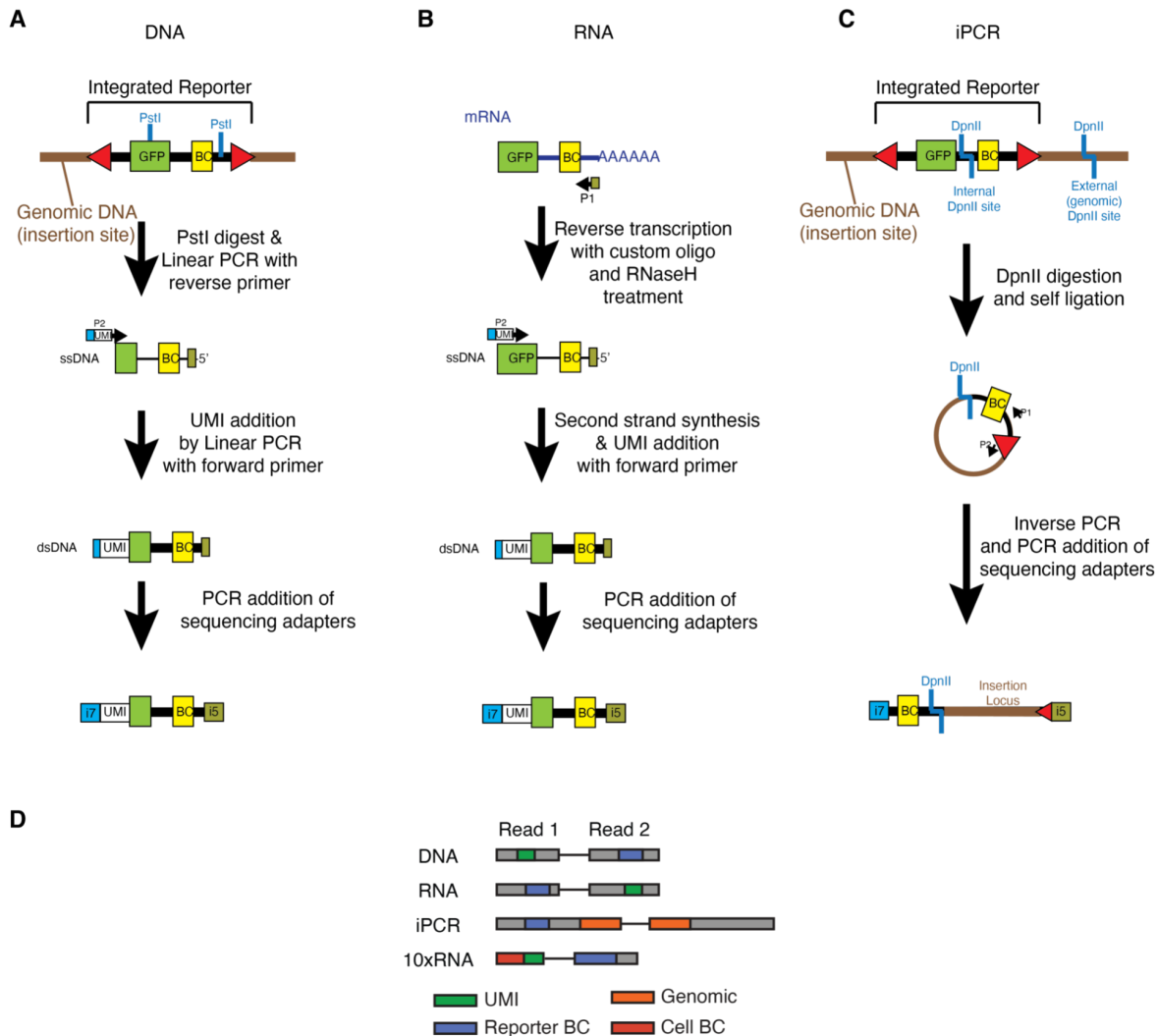
### Supplemental Fig. S1. A1 and C1 genomic insulator elements.

(A-B) Shown are full A1 (A) and C1 (B) elements (Liu et al. 2015), and A1 and C1 core elements truncated to just the CTCF footprint. Shown are DNase-seq cut counts and density tracks, and CTCF ChIP-seq tracks, and TF motif matches using FIMO.



**Supplemental Fig. S2. Reporter plasmid barcoding strategy.**

Libraries of barcoded reporter plasmids were generated by PCR linearization followed by incorporation of synthetic oligonucleotide containing a 16-nt random BC sequence using Gibson assembly.

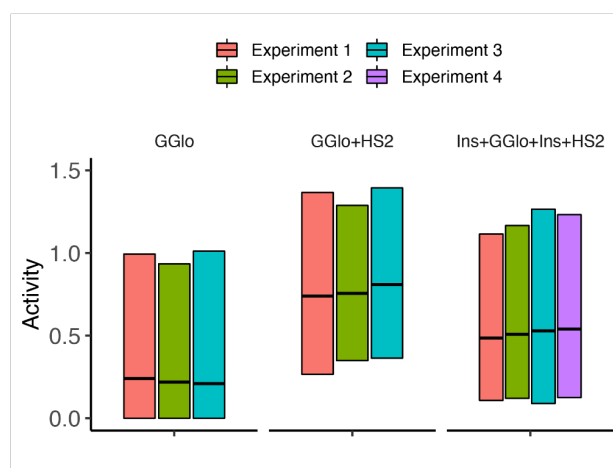


### Supplemental Fig. S3. Amplicon library construction approach.

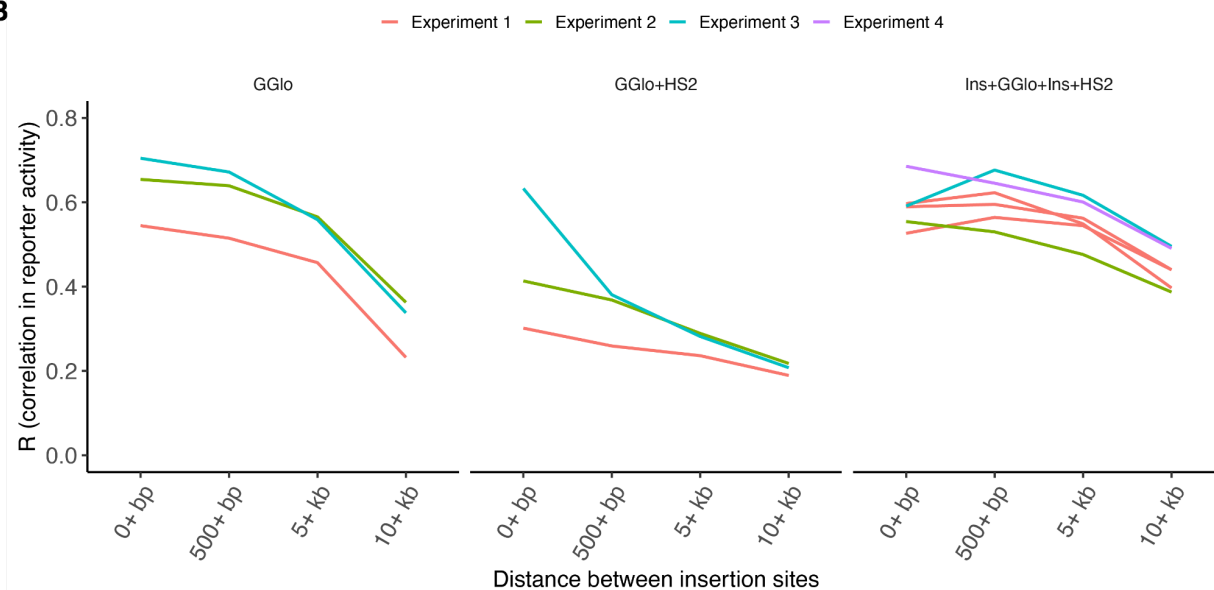
(A-C) Preparation of Illumina sequencing libraries to quantify barcoded reporter representation from DNA, expression from RNA, and integration location using inverse PCR (iPCR). DNA libraries utilized a PstI digest to limit template size, followed by a single linear amplification to add UMIs, and finally an Exo I digest to prevent any amplification from untemplated primers. RNA libraries incorporated UMIs during second strand synthesis.

(D) Schematic of unique molecular identifier (UMI, DNA, RNA, and 10x libraries), Reporter BC, Cell BC (10x libraries), and genomic sequence (iPCR libraries). The number of N nucleotides added to each individual sample varied between 8-12 bp (DNA and RNA) or 0-2 (iPCR) to increase diversity on the sequencing flow cell.

**A**



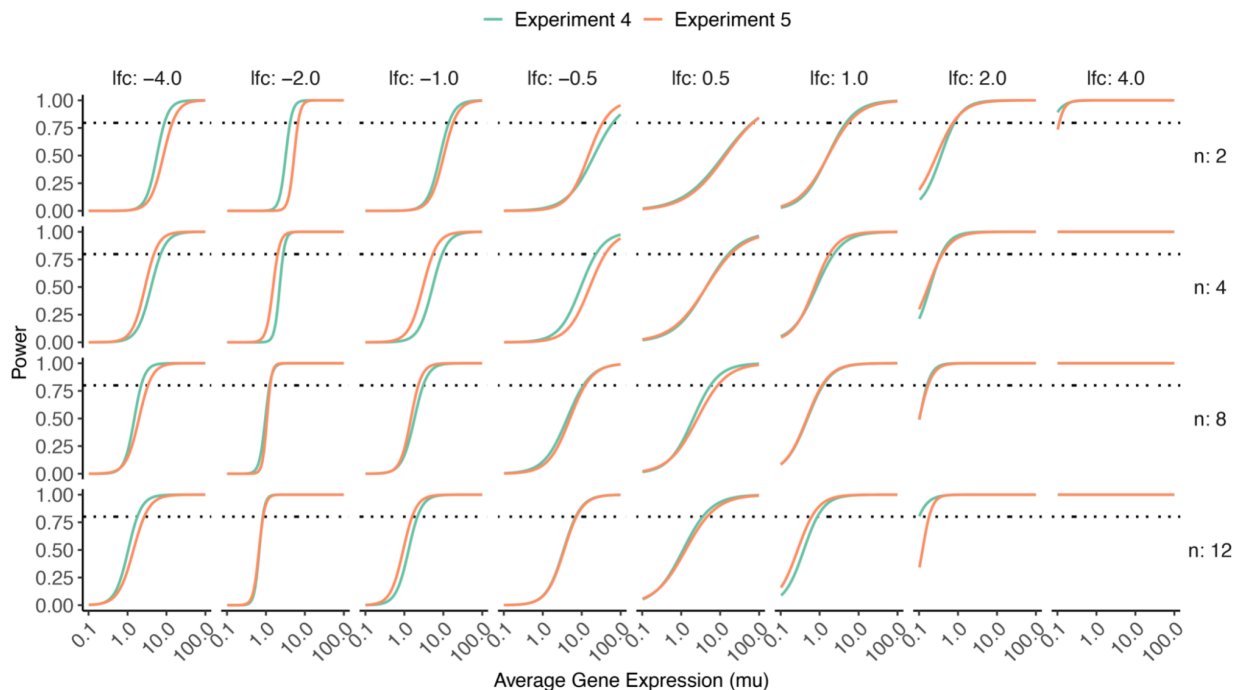
**B**



**Supplemental Fig. S4. Integrated barcoded reporter assay.**

(A) Average activity by reporter class and experiment.

(B) Correlation in activity for nearby insertions by reporter and experiment.



### Supplemental Fig. S5. Reporter impact on gene expression.

(A) Detection power for different circumstances stratified by fold change of the perturbation (columns), number of cells in the clone (rows), and average expression (x-axis). Y-axis shows proportion of significant tests ( $q$ -value  $< 0.10$ ) for each condition. Plotted is a logistic regression fit to the observed response. Horizontal dashed line shows 80% power.



## **Supplemental Tables**

### **Supplemental Table S1. Transfection summaries.**

Summary of transfections and conditions.

### **Supplemental Table S2. Sequencing libraries for DNA/RNA/iPCR/10xRNA experiments.**

Amplicon PCR libraries are listed by experiment. Technical replicates where PCR library construction was repeated on same biological sample are distinguished by letter suffix in “Sample #”. R1 Trim and R2 Trim refer to the number of bps added by the R1 or R2 primer for sequencing diversity or as a UMI.

### **Supplemental Table S3. Reporter experiment summaries.**

Summary of insertions analyzed for individual samples in each experiment. # Cells Seeded refers to the cells seeded on day 0 for Experiments 1-3, and the number of mKate+ cells on day 1 for Experiments 4-5.

### **Supplemental Table S4. Summary of 3’ 10x libraries.**

Shown are 10x scRNA-seq libraries, sequencing statistics, and mapping statistics.

### **Supplemental Table S5. PCR primer and DNA fragment sequences.**

### **Supplemental Table S6. Plasmids.**