

9-1-2022

Measuring transcription factor binding and gene expression using barcoded self-reporting transposon calling cards and transcriptomes

Matthew Lalli

Washington University School of Medicine in St. Louis

Allen Yen

Washington University School of Medicine in St. Louis

Urvashi Thopte

Icahn School of Medicine at Mount Sinai

Fengping Dong

Washington University School of Medicine in St. Louis

Arnav Moudgil

Washington University School of Medicine in St. Louis

See next page for additional authors

Follow this and additional works at: https://digitalcommons.wustl.edu/oa_4

Recommended Citation

Lalli, Matthew; Yen, Allen; Thopte, Urvashi; Dong, Fengping; Moudgil, Arnav; Chen, Xuhua; Milbrandt, Jeffrey; Dougherty, Joseph D; and Mitra, Robi D, "Measuring transcription factor binding and gene expression using barcoded self-reporting transposon calling cards and transcriptomes." *NAR Genomics and Bioinformatics*. 4, 3. Iqac061 (2022).

https://digitalcommons.wustl.edu/oa_4/628

This Open Access Publication is brought to you for free and open access by the Open Access Publications at Digital Commons@Becker. It has been accepted for inclusion in 2020-Current year OA Pubs by an authorized administrator of Digital Commons@Becker. For more information, please contact vanam@wustl.edu.

Authors

Matthew Lalli, Allen Yen, Urvashi Thopte, Fengping Dong, Arnav Moudgil, Xuhua Chen, Jeffrey Milbrandt, Joseph D Dougherty, and Robi D Mitra

Measuring transcription factor binding and gene expression using barcoded self-reporting transposon calling cards and transcriptomes

Matthew Lalli^{1,2,3}, Allen Yen^{1,4}, Urvashi Thopte³, Fengping Dong^{1,2}, Arnav Moudgil^{1,2}, Xuhua Chen^{1,2}, Jeffrey Milbrandt¹, Joseph D. Dougherty^{1,4} and Robi D. Mitra^{1,2,*}

¹Department of Genetics, School of Medicine, Washington University in St. Louis School of Medicine, Saint Louis, MO 63110, USA, ²Edison Family Center for Genome Sciences and Systems Biology Washington University in St. Louis School of Medicine, Saint Louis, MO 63110, USA, ³Seaver Autism Center for Research and Treatment, Icahn School of Medicine at Mount Sinai, New York, NY 10029, USA and ⁴Department of Psychiatry, Washington University School of Medicine, St. Louis, MO 63110, USA

Received April 25, 2022; Revised July 04, 2022; Editorial Decision July 25, 2022; Accepted August 24, 2022

ABSTRACT

Calling cards technology using self-reporting transposons enables the identification of DNA–protein interactions through RNA sequencing. Although immensely powerful, current implementations of calling cards in bulk experiments on populations of cells are technically cumbersome and require many replicates to identify independent insertions into the same genomic locus. Here, we have drastically reduced the cost and labor requirements of calling card experiments in bulk populations of cells by introducing a DNA barcode into the calling card itself. An additional barcode incorporated during reverse transcription enables simultaneous transcriptome measurement in a facile and affordable protocol. We demonstrate that barcoded self-reporting transposons recover *in vitro* binding sites for four basic helix-loop-helix transcription factors with important roles in cell fate specification: ASCL1, MYOD1, NEUROD2 and NGN1. Further, simultaneous calling cards and transcriptional profiling during transcription factor overexpression identified both binding sites and gene expression changes for two of these factors. Lastly, we demonstrated barcoded calling cards can record binding *in vivo* in the mouse brain. In sum, RNA-based identification of transcription factor binding sites and gene expression through barcoded self-reporting transposon calling cards and transcriptomes is an efficient and powerful method to infer gene regulatory networks in a population of cells.

INTRODUCTION

Calling cards is a uniquely powerful method to genetically record interactions between a protein of interest and the genome (1,2). Briefly, a protein of interest is fused to a transposase which can insert a transposon ‘calling card’ into the genome at sites of DNA–protein interaction such as transcription factor binding sites (TFBS). Early protocols recovered inserted transposons from genomic DNA (1), but a recent technical innovation termed the ‘self-reporting transposon’ (SRT) allows for the facile recovery of calling cards through RNA sequencing (RNA-seq) (3). RNA-mediated mapping of transposon insertions is more efficient than previous DNA-based protocols, and this protocol enables the simultaneous identification of TFBS and changes in gene expression in single cells (3). However, in bulk experiments on populations of cells, the RNA-mediated protocol is technically cumbersome, requiring a large number of replicates to identify independent insertions into the same genomic locus (4). Here, we present two crucial modifications of the SRT technology and protocol to facilitate its use and to enable joint recording of TFBS and gene expression in populations of cells: barcoded SRTs and barcoded transcriptomes.

Current implementations of the mammalian calling card protocol employ a hyper-active *piggyBac* transposase (5). An inherent constraint of this transposase is its requirement for a ‘TTAA’ tetranucleotide sequence for transposon insertion. As a result, multiple independent calling card insertions often occur at the same genomic location in different cells. Since the identification of TF binding sites is based on transposition count rather than read density, if these independent insertions are not distinguished, it limits the dynamic range of bulk calling card experiments. In the DNA-based calling card protocol, we solved this problem by including a barcode between the terminal repeats of the

*To whom correspondence should be addressed. Tel: +1 314 362 2751; Fax: +1 314 362 2157; Email: rmitra@wustl.edu

transposon that could be recovered by inverse PCR. However, this location is not compatible with the more efficient RNA-based protocol due to frequent ‘barcode swapping’ during cDNA amplification (6,7). As a result, current best practices for calling card experiments require a large number of biological replicates (typically 8–12) for each condition to increase the number of insertions that can be detected at a given TTAA location (4). While this improves the quantitative readout of these experiments, experimental cost and labor scale linearly with the number of replicates. Therefore, as an alternative approach, we sought to embed a unique barcode within the terminal repeat (TR) of the self-reporting transposon, the best location to enable reliable recovery without barcode swapping. Doing so is challenging, however, because all published sequences of the *piggyBac* transposon TRs are completely invariant, indicating strong sequence constraints on TR function which might preclude barcode insertion (8–12).

Here, we performed targeted mutagenesis of the *piggyBac* terminal repeat sequence to identify sites that could accommodate barcodes in calling card experiments. We discovered at least four consecutive nucleotides within the TR that were tolerant of a range of mutations without major reductions in transposition efficiency. As a resource to the scientific community, we have developed a set of barcoded *piggyBac* SRT plasmids and modified the calling card analysis software to utilize these barcodes. We demonstrate that barcoded SRT calling cards can map the genomic binding sites of transcription factors (TFs) involved in cell fate specification and transdifferentiation *in vitro*. Additionally, we combined barcoded SRT calling cards with bulk RNA barcoding and sequencing (BRB-seq) (13). This enables us to simultaneously identify TFBS and accompanying transcriptional changes from multiple TFs in an easy and affordable protocol. Lastly, we demonstrate that barcoded SRTs facilitate *in vivo* calling card experiments in the mouse brain, reducing labor by 10-fold. These innovations simplify bulk SRT calling card experiments, enable barcoding of experimental conditions, and allow for pooled library preparations that substantially reduce cost and labor. This simple protocol for simultaneously measuring transcription factor binding and gene expression changes will facilitate the inference of gene regulatory networks for TFs involved in development, cellular reprogramming, and disease.

MATERIALS AND METHODS

Transposon mutagenesis

PCR mutagenesis was performed in a 50 μ l reaction containing: 25 μ l 2 \times Kapa HiFi HotStart ReadyMix, 1 μ l of 10 μ M SRT Mutagenesis Forward Primer (either puro or tdTomato version), 1 μ l of 10 μ M SRT Mutagenesis Reverse Primer, 100 ng of SRT DNA (either PB-SRT-puro or PB-SRT-tdTomato), and 22 μ l of ddH₂O. PCR reactions were performed following thermocycling parameters: 95°C for 3 min, 10 cycles of: 98°C for 20 s, 60°C for 30 s, 72°C for 2 min, then 72°C for 10 min and 4°C forever.

PCR reactions were performed in duplicate. Each pool of mutant amplicons was purified with NucleoSpin Gel and PCR Clean-up (Macherey-Nagel). Products were trans-

ferred into separate wells of HEK293T cells to minimize any artifacts.

```
>SRT Mutagenesis Reverse Primer
tgcattctcaggagctcttaaccNNNNaaagatagctctgcgtaaaattgac
>SRT Mutagenesis Forward (puroR)
GCGGAAGGCCGTC AAGGCC
>SRT Mutagenesis Forward (tdTomato)
CACGAGACTAGCCTCGAtcaaggcgcatttaacctagaa
agataa
```

Cell culture

HEK293T cells and Neuro-2a cells were maintained in Dulbecco’s modified Eagle’s media (DMEM) supplemented with 10% fetal bovine serum and 1% penicillin–streptomycin. Cells were passaged every 3–4 days by enzymatic dissociation using trypsin.

Cloning

ASCL1, MYOD1, NEUROD2 and NGN1 were amplified from lentiviral cDNA expression vectors using 2 \times Kapa HiFi HotStart ReadyMix. A nuclear localization sequence was added to the 5’ end of each gene, and an L3 linker (amino acid sequence KLGGGAPAVGGG-PKAADK) (14) was inserted between the TF and hyperactive *piggyBac*.

EF1a_ASCL1, MYOD1 and NEUROG1_P2A_Hygro_Barcode were gifts from Prashant Mali (Addgene plasmid #120427, #120464 and #120467). pHND2-N174 was a gift from Jerry Crabtree (Addgene plasmid #31822).

Animals

All animal practices and procedures were approved by the Washington University in St. Louis Institutional Animal Care and Use Committee.

Calling card experiments

Calling card experiments are performed as described with minor modifications (4,15,16). Twenty-four hours before transfection, 250 000 HEK293T or Neuro-2a cells are plated per well in a 12-well plate. The next day, cells are transfected using PEI (Polysciences) with 1 μ g of total DNA comprising 500 ng of *piggyBac* (fused or unfused) and 500 ng of donor SRT (purified PCR product or miniprep DNA). Medium is changed 24 h after transfection. Three days after transfection, each well is trypsinized and replated into a T25 flask. For puromycin-resistance SRTs, puromycin is added 24 h later (2 μ g/ml). Three days after puromycin selection, total RNA is harvested using Direct-zol RNA MiniPrep kit (Zymo Research). RNA concentration and integrity was assessed using Nanodrop.

Virus generation and injections

Transposase and donor transposon constructs (barcoded or wildtype) were cloned into independent AAV transfer

vectors and used for *in vitro* transfection or viral packaging. Plasmids were packaged into AAV9 by the Hope Center Viral Vectors Core at Washington University School of Medicine. For *in vivo* experiments, intracranial injections into the cortex of wildtype C57BL6/J P0-1 mice of both sexes were performed as previously described (15). 1 μ l of viral mix was delivered across three sites per hemisphere for a total of 6 μ l per brain. Viral titers (viral genomes [vg] per milliliter) were $\sim 1.0 \times 10^{13}$. Animals were euthanized at P21 for analysis, cortices were dissected, and RNA extraction was conducted as described (15).

Calling card library preparation

A detailed protocol for calling card library preparation is available at protocols.io.

Calling card libraries were prepared as described with minor modifications (4,15,16). We performed first-strand reverse transcription reactions in 20 μ l total volume using 2 μ g of RNA from each *in vitro* sample and 4 μ g from each *in vivo* sample. RNA mixed with water and dNTPs was hybridized to oligo-dT primers (1 μ l of 50 μ M SMART_dTVN) by incubation at 65°C for 5 min and immediately transferred onto ice. 0.5 μ l of Maxima H Minus Reverse Transcriptase, RNasin RNase inhibitor, and 5 \times RT buffer were added and samples are incubated at 50°C for 60 min for reverse transcription.

Barcoded calling card and transcriptome library preparation

Bulk RNA Barcoding and sequencing (BRB-seq) was performed with minor modifications (13,17). We performed first-strand reverse transcription reactions in 20 μ l total volume using 2 μ g of RNA from each sample. Barcoded BRB-seq_dT30VN primers were modified to mimic the 10 Genomics v2 chemistry. RNA mixed with water and dNTPs was hybridized to barcoded oligo-dT primers (2 μ l of 25 μ M stock) by incubation at 65°C for 5 min and immediately transferred onto ice. 1 μ l of template switch oligo (TSO_SMART), 0.5 μ l of Maxima H Minus Reverse Transcriptase, RNasin RNase inhibitor, and 5 \times RT buffer were added and samples were incubated at 50°C for 60 min for reverse transcription. In the current version of the SRT library protocol, we use the BRB-seq_dT30VN oligos for reverse transcription.

For barcoded SRT and transcriptome experiments studying ASCL1, MYOD1 and unfused *piggyBac*, 3 μ l of barcoded reverse-transcription product from four replicates for each factor were pooled together for transcriptome analysis. Four replicates of each factor (5 μ l per replicate) were pooled and purified in parallel for calling card library preparation. Pooled samples were purified using NucleoSpin Gel and PCR Clean-up (Macherey and Nagel) and eluted with 30 μ l of elution buffer. We designed barcoded oligoDT-VN oligos to mimic the configuration of 10x Genomics v2 chemistry: partial seq1, 16 bp cell barcode extracted randomly from the 10x Genomics safelist, and a 10 bp UMI (5N + 5V). Barcoded primer sequences are provided in Supplementary Table S5.

For transcriptome profiling, 4 μ l of the purified pool of barcoded first-strand reactions were mixed with 19 μ l water, 1 μ l partial seq1 primer (25 μ M), 1 μ l SMART primer (25 μ M), and 25 μ l 2 \times KAPA HiFi HotStart ReadyMix (Roche). 10 cycles of PCR with a long extension time (98° 20 s, 60° 30 s, 72° 6 min) were performed. Amplified cDNA was purified with 0.6 \times AMPure XP (Beckman Coulter) magnetic beads. DNA was eluted with 20 μ l water and concentration was measured using the TapeStation D5000 ScreenTape (Agilent). 600 pg of product were tagged and amplified with barcoded N7 primers and P5-index-seq1 primers using the Nextera XT kit (Illumina). BRB-seq libraries were sequenced on a Novaseq 6000 paired-end with 28 \times 91 reads.

SRT libraries were prepared similarly. 4 μ l of the purified pool of barcoded first-strand reactions were mixed with 19 μ l water, 1 μ l partial seq1 primer (25 μ M), 1 μ l SRT_PAC_F1 primer (25 μ M) and 25 μ l 2 \times KAPA HiFi HotStart ReadyMix (Roche). Twenty cycles of PCR with a long extension time (98° 20 s, 65° 30 s, 72° 5 min) were performed. Amplified cDNA was purified with 0.6 \times AMPure XP (Beckman Coulter) magnetic beads. DNA was eluted with 20 μ l water and concentration was measured using the TapeStation D5000 ScreenTape (Agilent). 600 pg of product were tagged and amplified with barcoded N7 primers and P5_BC_SRT primers using the Nextera XT kit (Illumina).

Because SRT libraries have low diversity on Read1, we designed a set of six P5 SRT primers with stagger regions of different lengths to introduce sequence diversity. We recommend sequencing at least four SRT libraries on the same flow cell and using 20% PhiX DNA spike-in. Barcoded *piggyBac* primers with stagger regions are listed in Supplementary Table S6.

Primers

```
>SMART_dT18VN
AAGCAGTGGTATCAACGCAGAGTACGTTTTTTT
TTTTTTTTTTTTTTTTTTTTTTTTTTTTVN
>BRB-seq_dT30VN, e.g. pSeq1-BC1-UMI-dtVN
CTACACGACGCTCTCCGAT
CTCTGATAGCATGGTCATNNNNNVVVVTT
TTTTTTTTTTTTTTTTTTTTTTTTTTTTVN
>SRT_PAC_F1
CAACCTCCCCTTCTACG*A*G*C
Asterisks indicate phosphorothioate bond substitutions
>SRT_tdTomato_F1
TCCTGTACGGCATGGAC*G*A*G
> SMART_TSO
AAGCAGTGGTATCAACGCAGAGTACrGrGrG
>SMART
AAGCAGTGGTATCAACGCAG*A*G*T
>Partial Seq1
CTACACGACGCTCTCCGA*T*C*T
> P5_BC_SRT_STAGGER1 (XXXXX is i5 index, stag-
ger region is underlined)
AATGATACGGCGACCACCGAGATCTACACXX
XXXXACACTCTTTCCCTACACGACGCTCTCCGA
TCTTGCCTCAATTTTACGCAGACTATCTTT
```

```

>P5-index1-Seq1 (index sequence is underlined)
AATGATACGGCGACCACCGAGATCTACA
CAGGACAACACTCTTTCCCTACACGACGCTCTT
CCGATCT
>Nextera_N70X (index sequence is underlined)
CAAGCAGAAGACGGCATAACGAG
ATTTCGCCTTAGTCTCGTGGGCTCGG

```

Library prep and sequencing

Purified PCR product was measured using a D5000 Screen-Tape (Agilent). cDNA samples were diluted to 600 pg/ μ l and 2 μ l of this was used for tagmentation with Nextera XT kit.

RNA-seq analysis

Sequencing data corresponding to barcoded bulk RNA transcriptomes were processed using the 10x Genomics software package Cell Ranger (v 2.1.0). The output filtered gene expression matrices were imported into R (v 3.5.1) for further analysis (18). Gene counts were used directly in edgeR for standard bulk RNA-seq analysis (19).

Calling card analysis

Sequencing and analysis:

Bulk barcoded RNA calling card libraries were sequenced and analyzed as described with modifications to utilize the SRT barcode (4). Calling card reads begin with a stagger region that serves as a library barcode, the barcoded transposon TR, the insertion motif TTAA, then the genome at the site of insertion. Reads are checked for the library barcode, TR sequence and TTAA and these sequences are trimmed. SRT barcodes are extracted by UMI-tools (20), and appended as a sequence tag to the read. Any remaining Nextera adaptors are trimmed before mapping the reads to the human genome (hg38) using NovoAlign or STAR. Aligned reads are validated as insertions if adjacent to a TTAA site in the genome. Bona fide insertions are then converted to qBED format (née .ccf) (21). SRT barcodes were incorporated into the barcode column of the qBED. If non-overlapping barcode sets were used to define experiments, qBED files can be demultiplexed by this barcode field.

Peak calling: Calling card peaks were called as described (1,22) using in-house peak calling software. Specifically, peaks were called using the call_peaks_mac3 python script, which follows the algorithm used by MACS to call ChIP-Seq peaks (23) modified for the analysis of calling card data. The main peak calling function is passed an experiment frame, a background frame, and an TTAA_frame, all in qBED/ccf format (21). It then builds interval trees containing all of the background and experiment hops (insertion events) and all of the TTAA's. Next, it scans the genome with a window of window_size and step size of step_size and looks for regions that have significantly more experimental hops than background hops (poisson w/ pvalue_cutoff). It merges consecutively enriched windows and computes the center of the peak. Next it computes lambda, the number of insertions per TTAA expected from the background distribution by taking the max of lambda_bg, lambda_1, lambda_5,

lambda_10. It then computes a *p*-value based on the expected number of hops = lambda \times number of TTAA's in peak \times number of hops in peak. Finally, it returns a frame that has Chr, Start, End, Center, Experiment Hops, Fraction Experiment, Background Hops, Fraction Background, Poisson *P*-value as columns. We used parameters: -pc 0.001 -peak_finder_pvalue 0.01 -window 1000 -step 500 -pseudocounts 0.2 for peak calling.

RESULTS

Identifying candidate regions for barcode insertion in piggyBac terminal repeat

The SRT consists of a promoter driving a reporter (e.g. fluorescent protein or puromycin resistance cassette) flanked by the transposon terminal repeat sequences (TR), the part of a transposon that is recognized by its cognate transposase. Importantly, there is no polyadenylation (poly(A)) signal sequence after the reporter gene, so gene transcription proceeds through the TR and into the genome. This design allows the SRT to report its genomic location in cellular RNA (Figure 1A) (3). To maximize compatibility with calling cards library preparation and minimize template switching, the ideal barcode location would be as close to the genomic insertion site as possible. Because sequences outside of the TRs are not inserted into the genome, barcodes cannot be introduced there (Figure 1A, site 1). A barcode inserted between the reporter gene and the TR, as implemented in our DNA-based calling cards protocol (1), would be \sim 300 bp away from the informative transposon-genome junction (Figure 1A, site 2). This would retain a long stretch of shared sequence present in all amplicons that would lead to extensive barcode swapping during the SRT amplification PCR step in library preparation (6,7).

Therefore, we sought to introduce a barcode into the TR itself (Figure 1B, site 3), directly adjacent to the TR-genome junction. Such a strategy has two major advantages compared to other approaches. First, a barcode in this position could be captured in the same sequencing read as the transposon-genome junction, simplifying the protocol. Second, by eliminating as much constant intervening sequence as possible, there is little risk of introducing aberrant chimeric PCR products during sequencing library preparation (6,7). Whereas modifications to TRs from other transposases such as *Sleeping Beauty* have been successfully engineered (9), similar efforts have revealed extensive sequence constraints on *piggyBac* TRs for efficient transposition (10,11). Nevertheless, we sought to identify candidate regions within the TR that might accommodate a DNA barcode.

The minimal *piggyBac* TR consists of a 19-bp internal repeat (IR), a 3-bp spacer, and a 13-bp terminal invert repeat (12) (Figure 1B). These sequences are critical for *piggyBac* recognition, cleavage, and transposition. Notably, all published sequences of the 13-bp terminal invert repeat in the *piggyBac* TR are completely invariant. DNase I footprinting of *piggyBac* binding to its TRs revealed strong binding across much of this region (8), yet a few bases were less protected and therefore might be a candidate region for inserting a barcode (Figure 1B, underlined nucleotides, gold).

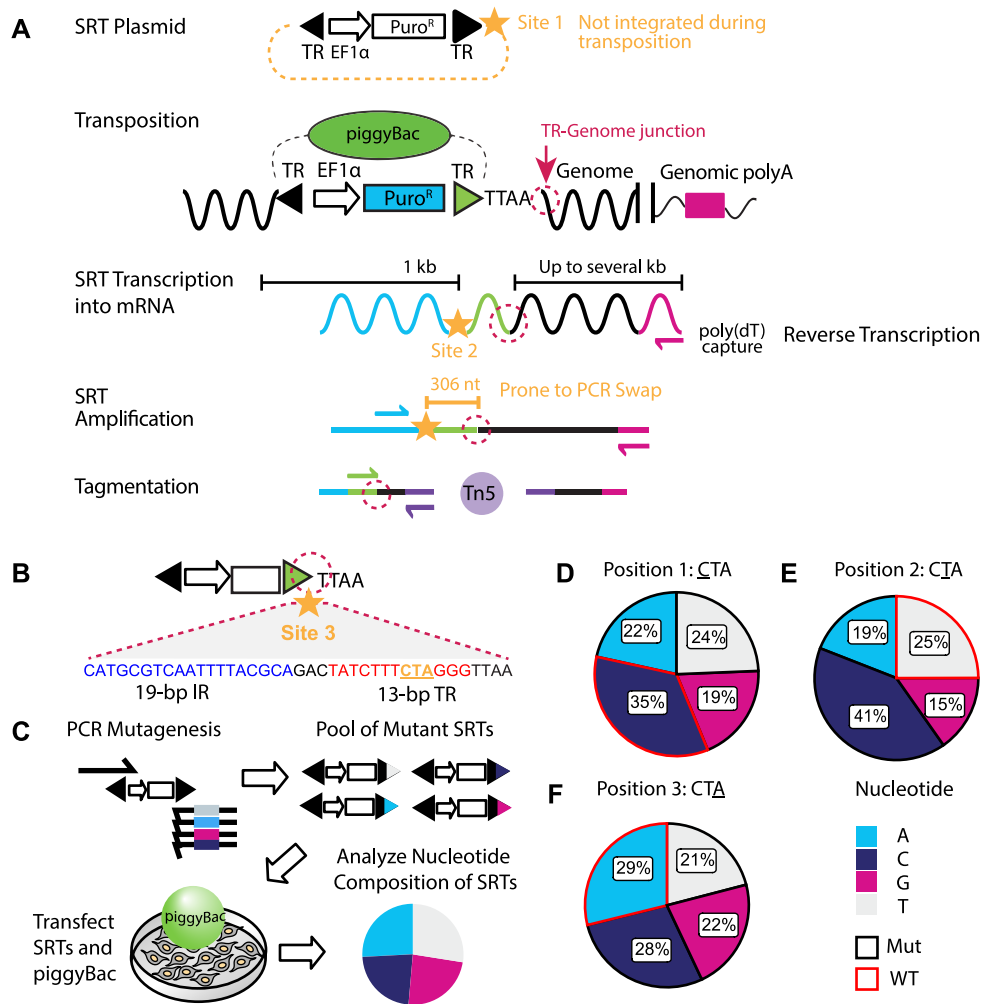


Figure 1. Barcoding the self-reporting transposon. (A) Schematic overview of the SRT construct, Calling Card method, and sequencing library preparation. Candidate sites for barcode insertions are indicated with gold stars. The TR-Genome junction, used to map transposon insertions, is circled in dotted magenta line. (B) Barcode site 3 is within the *piggyBac* TR sequence, immediately adjacent to the TR-Genome junction. Underlined nucleotides in the 13-bp terminal inverted repeat region ('CTA', gold) were targeted for mutagenesis by mutagenic PCR. (C) Overview of calling card rapid mutagenesis scheme. Mutant amplicons were transfected into cells with *piggyBac* transposase and integrated calling cards were collected. Nucleotide frequency for each mutagenized position of integrated SRTs were calculated. Nucleotide frequency at (D) position 1, (E) position 2 and (F) position 3 of integrated mutant SRTs. Wild-type sequences are outlined in red. All four possible nucleotides were well-represented at all three mutated positions. IR: internal repeat. TR: terminal repeat. EF1 α : eukaryotic translation elongation factor 1 α promoter. SRT: self-reporting transposon. nt: nucleotide. kb: kilobase. Puro^R: puromycin resistance cassette. WT: wild-type. Mut: mutant.

Targeted mutagenesis generates mutant SRTs with high transposition efficiency

We developed a simple and rapid screening protocol to generate and identify mutant *piggyBac* TR sequences capable of successful transposition (Figure 1C). We designed primer sequences to introduce single point mutations into our candidate region using PCR. Purified PCR products encoding puromycin-resistance SRTs flanked by mutated TRs were directly transfected into HEK293T cells along with unfused hyper-active *piggyBac*. If mutated amplicons are compatible with transposition, they will be inserted into the genome and confer puromycin resistance. We selected for transposition events after 4 days by adding puromycin. We extracted RNA 3 days after selection, and prepared bulk SRT libraries according to established protocols with modifications described (4) (Methods).

We sequenced calling card libraries using RNA-seq and mapped genomic transposition events from at least two independently generated mutant SRT pools for each position. Each library yielded 75 000–150 000 unique insertion sites providing a representative view of genomic insertion efficiency for mutant SRTs. Analysis of transposition events revealed that all three candidate positions within the *piggyBac* TR accommodated mutations without greatly diminishing transposition ability (Figure 1D–F). Each of the three mutagenized positions tolerated all 4 nucleotides at similar frequencies, hence generating at least 12 unique transposon barcodes.

Having obtained successful transposition of SRTs with single mutations, we next tested whether multi-nucleotide mutations within this region could be tolerated. Using PCR, we introduced three consecutive mixed bases (Ns, where N

can be A, C, G or T) into this region to generate a total of 64 barcoded SRTs. We transfected pools of these mutant SRT PCR amplicons into cells and again prepared calling card libraries after puromycin selection. Analysis of hundreds of thousands of transposition events showed that all 64 mutant transposons could be integrated into the genome, albeit at varying degrees of efficiency (Figure 2A). To better understand sequence preferences governing transposition efficiency, we generated a sequence motif from the top 30 most abundantly inserted transposons. Cytosine was slightly favored in the first two positions, and thymine was strongly disfavored from the third position (Supplementary Figure S1). Among mapped transposition events, we also observed the presence of mutations at a fourth nucleotide position immediately adjacent to our targeted bases, leading us to test whether this position could also be modified. Following the same approach, we generated SRTs with mutations in this position and prepared calling card libraries from two independently transfected sets of cells. As with the other single nucleotide SRT mutants, we found that this position could also tolerate all 4 nucleotides (Figure 2B).

Longer barcodes are preferable in sequencing applications as they not only increase the number of unique sequences available but can also have advantageous properties including error detection and error correction (24). DNA barcodes that differ from each other by a given number of nucleotides can be used to detect and correct errors that arise during sequencing. If barcodes differ by two or three nucleotides, then sequencing errors where one of the barcode bases is misread can be detected or corrected, respectively. A three-nucleotide barcode can encode a maximum of 64 unique sequences including a set of 12 error-detecting barcodes and 4 error-correcting barcodes (Figure 2A). A four-nucleotide barcode can encode up to 256 unique sequences including 48 error-detecting and 12 error-correcting barcodes (Figure 2C). To generate a pool of 256 mutant transposons, we introduced 4 consecutive mixed bases (Ns) into the TR using a degenerate primer. We collected and analyzed over 160 000 unique transposition sites in the genome and found all 256 possible mutated transposons were inserted into the genome (Figure 2C). We analyzed the nucleotide composition of the top 100 most abundantly inserted transposons to reveal sequences mediating transposition efficiency (Figure 2D). Overall preferences were modest except for a strong favoring of C/G in the first position and a disfavoring of thymine in the third position. These results suggest that a fixed sequence for these 4 nucleotides is not required for binding and transposition by the *piggyBac* transposase.

Given the compatibility of mutations in this region of the TR with transposition, we tested whether we could insert a single nucleotide in this region to further increase barcode length. We generated mutant SRTs with a single nucleotide insertion and performed calling cards with these. We observed that very few cells survived selection and consequently few transposition events were recovered from this experiment. Among the recovered transposition events, the most prevalent sequence matched the wild-type SRT with no insertion. Of the recovered SRTs that did contain an inserted nucleotide, many of the sequences also contained a nearby 1-nt deletion which suggests a strict TR length

constraint for successful transposition (Supplementary Table S1). The inserted nucleotide may have disrupted any step of *piggyBac* recognition, cleavage, and transposition by changing the sequence, shape, or flexibility of the transposon (8,25). Thus, focusing on just the four nucleotide barcodes, as a resource to the community, we individually cloned the top 24 integration-competent error-detecting barcodes into two self-reporting vectors (Supplementary Table S2). These SRT vectors include an adeno-associated viral (AAV) vector carrying a tdTomato reporter SRT compatible with *in vivo* calling card experiments (15) and a non-AAV SRT vector encoding the puromycin resistance gene.

Using barcoded SRTs to map binding sites of transcription factors involved in cell fate specification

To demonstrate that barcoded SRTs facilitate TFBS recording in cellular populations, we performed calling card experiments for four TFs using this method. We chose to record the binding of four members of the basic helix-loop-helix (bHLH) family: Achaete-scute homolog 1 (ASCL1), Myogenic Differentiation 1 (MYOD1), Neuronal Differentiation 2 (NEUROD2), and Neurogenin 1 (here referred to as NGN1). These TFs are implicated in cell fate specification and cellular reprogramming (26–31). Interestingly, all four TFs recognize the same canonical E-box motif *in vivo*, bind some overlapping and unique sites in the genome, and regulate distinct gene expression programs (32). To perform calling card experiments, we first created mammalian expression vectors containing fusion proteins of each of the four TFs to the N-terminus of hyperactive *piggyBac* separated by an L3 linker (14). We transfected HEK293T cells expressing fused or unfused *piggyBac* with wild-type or barcoded versions of SRTs encoding either tdTomato or puromycin-resistance reporters and harvested RNA after ~1 week. We prepared and sequenced SRT calling card RNA-seq libraries and analyzed the data to identify transposon insertions in the genome. Calling card peaks were called as described (1,22) and analyzed for enriched motifs and neighboring genes using HOMER (33). We then performed Gene Ontology enrichment analysis on sets of genes located near TFBS (34).

For each of the four bHLH factors, we recovered hundreds of thousands of genomic insertion events and called thousands of calling card peaks (Supplementary Table S3). Motif enrichment analysis for each factor recovered several enriched bHLH E-box motifs, including the known motifs for Ascl1, MyoD and NeuroD1 (Figure 3A). This motif recovery suggests barcoded calling cards identified bona fide TFBS for these factors. For NEUROD2, the top 3 enriched motifs belonged to specific neuronal bHLHs including NeuroD itself (Figure 3A). Likewise for MYOD1, the top 3 enriched motifs belonged to myogenic bHLHs of the MyoD family (Figure 3A), indicating specificity of the calling card peaks for the TFs of interest. This result supports the interpretation that while the core E-box motif is common to all factors, nucleotides flanking this motif may confer binding specificity (35). For ASCL1, in addition to recovering bHLH motifs, we also observed an enrichment of Jun/Fos and other basic zipper (bZIP) motifs. This might

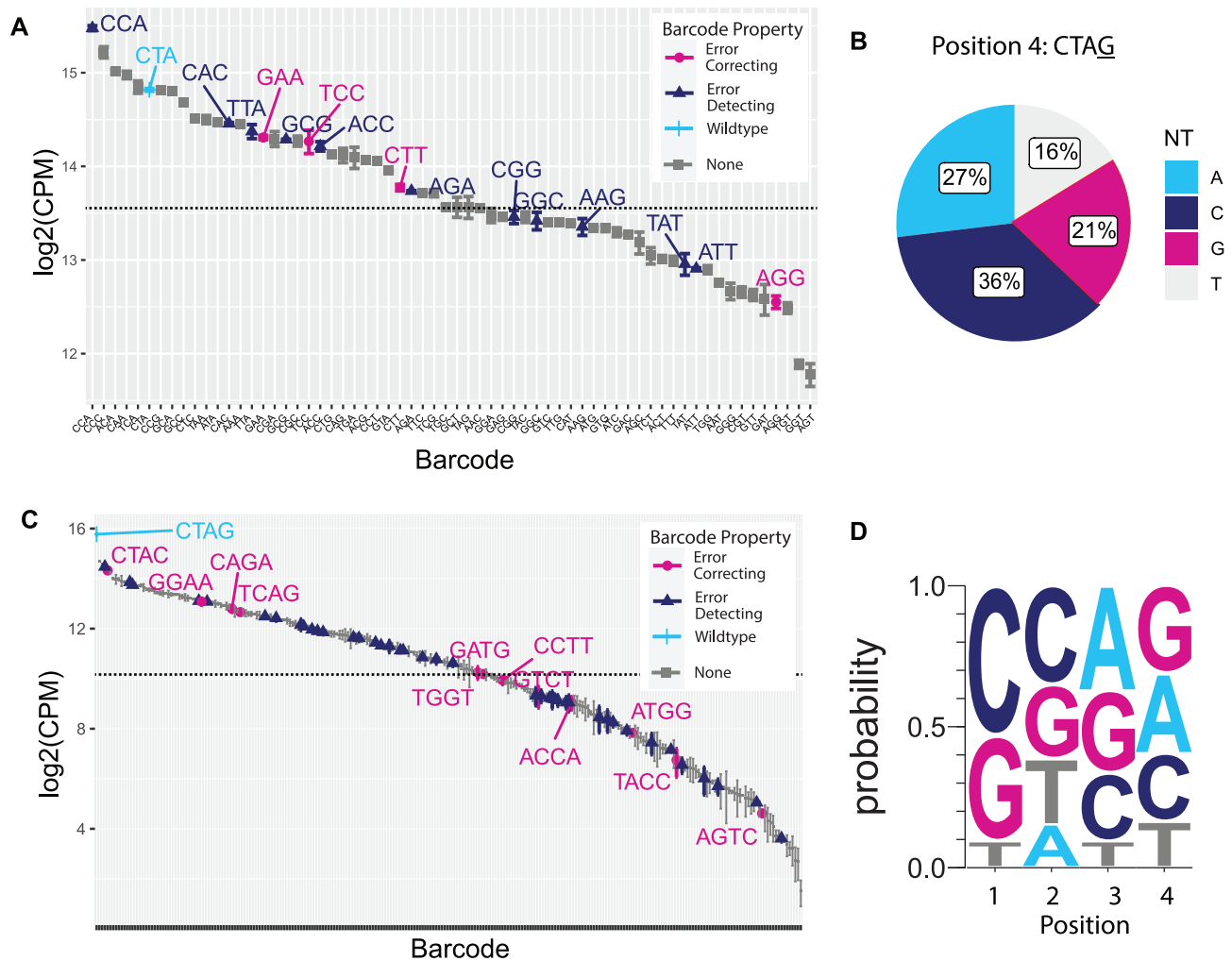


Figure 2. Multi-nucleotide mutagenesis in *piggyBac* terminal repeat discovers integration-competent barcoded SRTs. (A) Normalized counts of integration of events for 64 possible combinations of three nucleotide barcodes at the targeted region are shown (log₂ counts per million (CPM)). All 64 barcoded SRTs could integrate into the genome. Black dotted lines indicate 50th percentile of read counts. Data are plotted as mean and SEM from two independent replicates. (B) Targeted mutagenesis at a fourth position in the terminal repeat identified another site that could tolerate all 4 nucleotide substitutions while retaining integration-competence. Wild-type sequence ('G') is outlined in red. (C) Normalized counts (log₂ CPM) of insertions for 256 combinations of 4-nt barcodes. All 256 barcodes were present at varying degrees of insertional efficiency. Wild-type sequence is colored cerulean. Error-correcting and error-detecting barcodes are colored respectively in magenta and midnight blue. (D) Sequence logo of the top 100 most abundantly inserted 4-nt barcoded SRTs reveals modest sequence preference for integration efficiency. CPM: counts per million sequencing reads.

indicate the binding of bZIP TFs at ASCL1 sites as has been reported for other neuronal bHLH TFs (36).

Next, we identified genes located near each TFBS, and characterized the shared and differential binding of bHLH TFs (37). Consistent with their recognition of a common E-box motif, all four TFs bound near many of the same genes (Figure 3B, Supplementary Figure S2A, B). While many genes were overlapping across TFs, each TF also had its own set of unique genes near TFBS. To gain insight into the regulatory roles of these TFs, we performed Gene Ontology enrichment analysis on sets of genes located near TFBS identified by barcoded SRTs (34). Gene Ontology terms identified for genes proximal to the neurogenic TFs ASCL1, NEUROD2, and NGN1 were enriched for neuronal pathways including axonogenesis and neuron projection development, reflecting their known roles in neuronal reprogramming (Figure 3C) (29,38,39).

MYOD1 binding sites were located near genes strongly enriched for roles in cardiogenesis and muscle development (Figure 3C, Supplementary Figure S2C). Consistent with prior findings of MYOD1 binding some neuronal targets (27), we found some enrichment for binding at genes enriched for neurogenic pathways. The observed enrichment of neuronal and muscle genes is particularly notable given the calling card assay was performed in human embryonic kidney cells which do not natively express any of the assayed TFs. That all factors are able to recognize and bind specific genes enriched for their known functions implies either a permissive binding environment in HEK293T cells or cell-type independent target access by these TFs. This also highlights that subtle differences in nucleotide sequences flanking the common core E-box motif can confer binding specificity at functionally distinct gene sets (35).

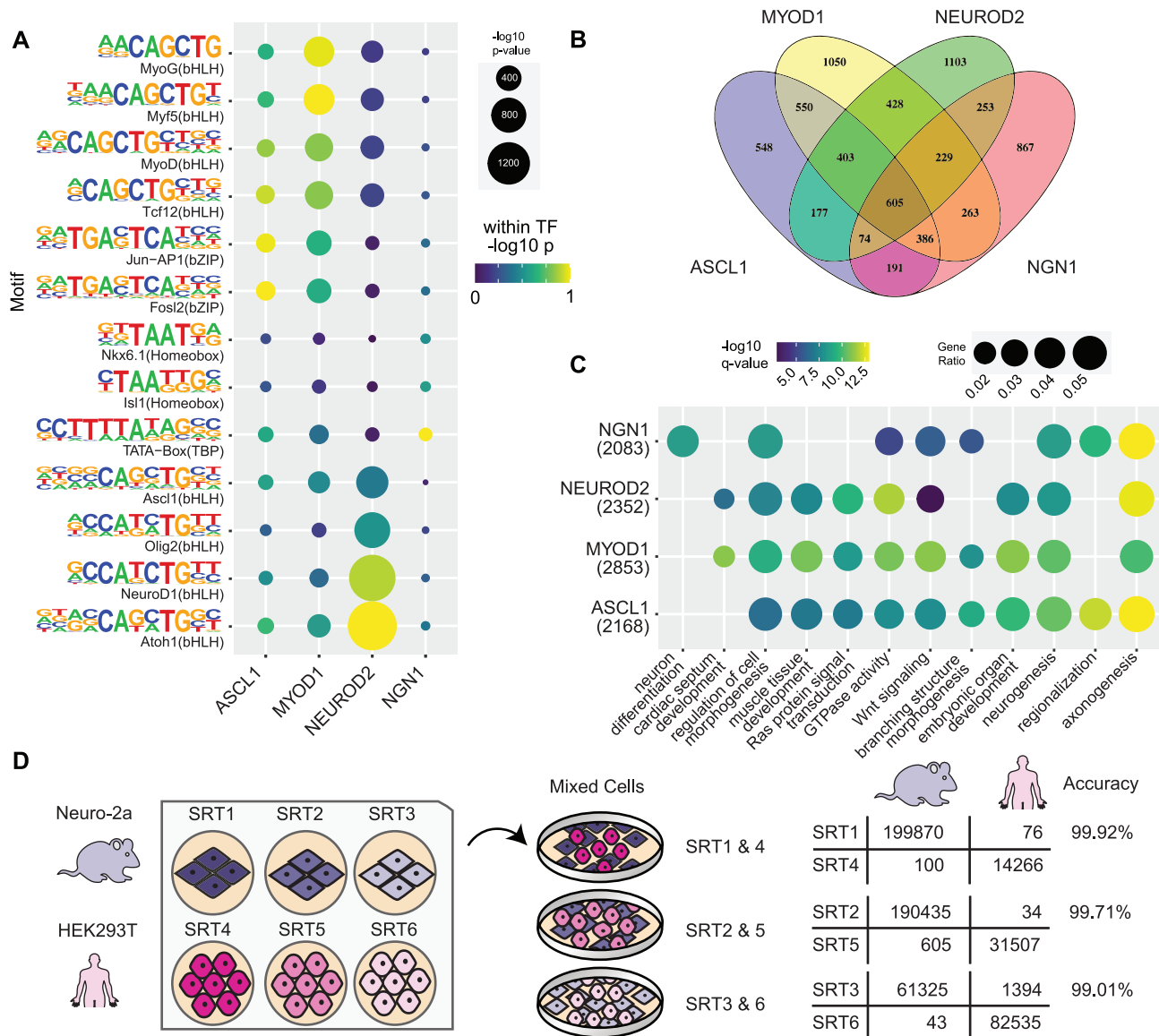


Figure 3. Calling cards using barcoded SRTs recover known binding motifs for bHLH factors near genes related to known TF functions. (A) Top binding motifs for each motif were retrieved from DNA sequences in calling card peaks. These sites are enriched for the canonical E-box motif as well as bHLH TFs including or related to each TF. (B) Venn diagram of genes proximal to called peaks for each TF indicates both shared and distinct binding of these TFs. (C) Gene Ontology enrichment analysis reveals terms related to neurogenesis and myogenesis. (D) Species mixing experiment confirms minimal barcode swapping in SRT library preparation. bHLH: basic helix-loop-helix. bZIP: basic zipper.

Species mixing experiment confirms low rate of barcode swapping in SRT libraries

To directly assess the rate of barcode swapping during SRT library construction, we performed a species mixing experiment. We transfected individual wells of human (HEK293T) or mouse cells (Neuro-2a) with unique barcoded SRTs and unfused *piggyBac* (Figure 3D, left). After transfection, human and mouse cells were pooled together and grown under selection for SRT insertions. We collected RNA from pooled cells, then constructed and sequenced SRT libraries. We mapped insertions to the mouse and human genomes. In this design, barcode swapping is indicated by mouse transfected SRT barcodes mapping to the human genome and vice versa. Across three biological replicates,

barcoded SRTs mapped to the correct genome in 99.6% of detected SRT insertions (Figure 3D, right), indicating a very low rate of barcode swapping. Such minimal barcode swapping enables multiplexed TF profiling in a single pooled experiment.

Barcoded SRTs and transcriptomes enable simultaneous mapping of TFBS and gene expression

SRTs were specifically invented to enable simultaneous readout of gene expression and transcription factor binding in single cells (3), but can also be used to map TFBS in populations of cells as demonstrated here and previously (3,15). Because SRTs are amplified from poly(A) RNA, we sought to prepare SRT and poly(A) mRNA sequencing libraries in

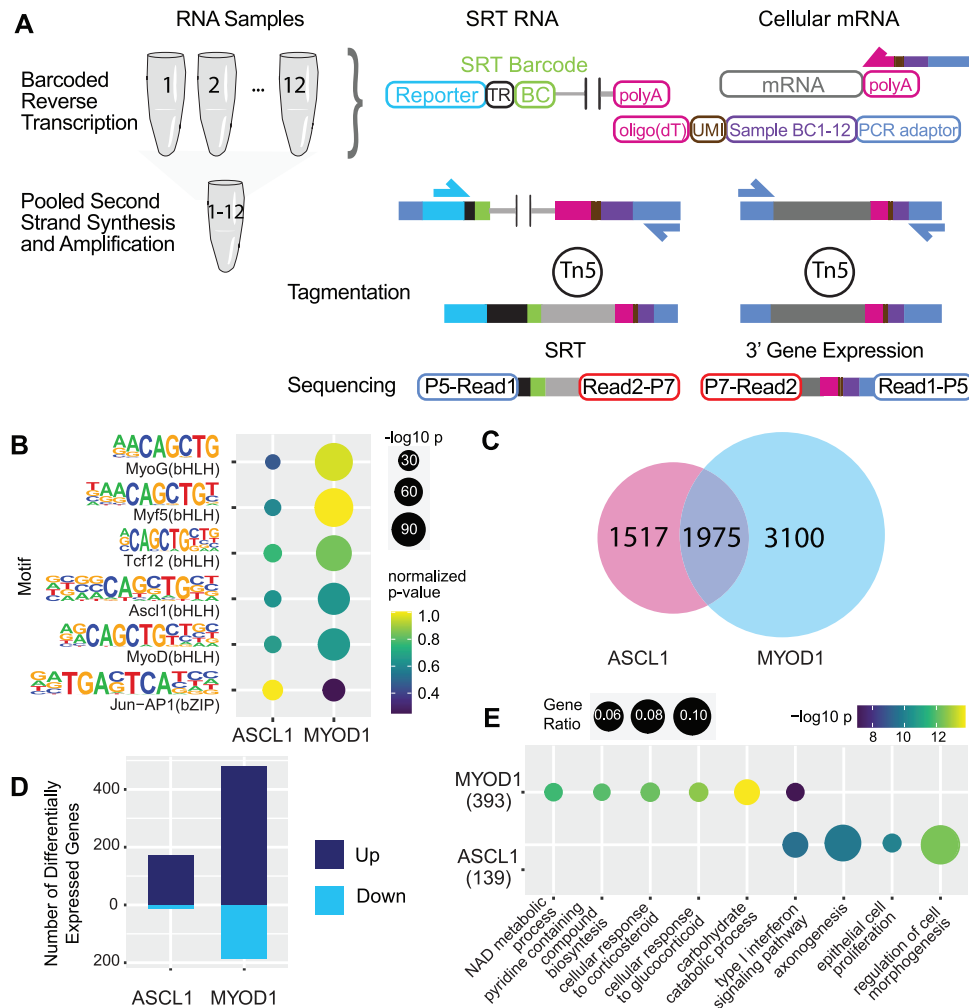


Figure 4. Barcoded SRT calling cards and transcriptomes enables joint measurement of TFBS and gene expression. (A) Schematic overview of barcoded sequencing library preparation. Sample-specific barcode (Sample BC) with unique molecular identifiers (UMI) is introduced during reverse-transcription of poly(A) RNA including SRTs and mRNA. Reverse transcription products (cDNA) can then be pooled for second strand synthesis and amplification. Sequencing libraries are prepared for SRTs and transcriptomes in parallel. (B) Barcoded SRT experiments recover binding motifs for ASCL1 and MYOD1. (C) Venn Diagram showing shared and distinct genes near ASCL1 and MYOD1 binding sites. (D) Transcriptomes profiled by bulk RNA-seq with barcodes revealed differential gene expression for ASCL1 and MYOD1, compared to cells transfected with unfused *piggyBac*. (E) Gene Ontology of differentially expressed genes in ASCL1 and MYOD1 cells.

parallel from the same sample. To reduce cost and labor of library preparation, multiple poly(A) mRNA-seq libraries can be barcoded during reverse-transcription then pooled for library preparation and sequencing (13). We have previously modified this barcoding protocol to employ the 10x Genomics single cell 3' v2 chemistry which enables turnkey analysis of RNA-seq data using Cell Ranger (17,40).

To facilitate simultaneous preparation of SRT calling card and poly(A) 3' RNA-seq libraries, we introduced a sample barcode and unique molecular identifier (UMIs) into the poly(dT) capture oligonucleotide (Figure 4A). Each experimental replicate is reverse transcribed using an oligo(dT) capture oligo with a unique sample barcode, then multiple samples can be pooled for ultra-affordable transcriptomic analysis (13). Calling card experiments could also be designed such that experimental replicates use distinctly barcoded SRTs (individual barcodes or sets of barcodes) so that the same pool of cDNA can then be used to

amplify SRTs and mRNA in parallel reactions. Otherwise, SRT libraries can be amplified individually. Sequencing libraries of amplified products are then prepared by tagmentation (41).

Simultaneous mapping of TFBS and gene expression of pioneer TFs

ASCL1 and MYOD1 belong to a special class of transcription factors called pioneer factors that can access both open and closed chromatin and reprogram cell fate from pluripotent stem cells or fibroblasts to neurons and muscle cells respectively (28,42–44). Chromatin immunoprecipitation followed by sequencing (ChIP-seq) after overexpression of these factors in mouse embryonic fibroblasts revealed a surprising degree of overlapping binding sites between these factors (27). Gene expression profiling of TF overexpression, however, revealed differing transcriptional outcomes

Table 1. Drastic cost and labor reduction of barcoded SRT and transcriptomes compared to original protocol. ‘Original’ calculations use the recommended 12 replicates per TF (4). This experiment assayed 3 TFs (unfused hyper *piggyBac*, ASCL1 and MYOD1). Transfection costs are based on NEON or nucleofector transfection device reactions. Tagmentation costs assume a library is prepared for each of the 12 replicates for both calling cards and transcriptomes. Tapestation costs reflect core facility pricing

	Replicates (<i>n</i>)		Cost (\$USD)	
	Original	Barcoded	Original	Barcoded
Transfections	36	12	720	240
RNA isolation and reverse transcription	36	12	180	60
Amplification	72	2	216	6
Bead Cleanup, Tapestation	72	2	1080	30
Tagmentation	72	2	2160	60
Bead Cleanup, Tapestation	72	2	1080	30
Sequencing			<i>Same</i>	
Total			5436	426

for these two factors (27). As many TF binding events have no or small effects on gene regulation (45–48), integrating TFBS data with mRNA-seq is a powerful method to decipher cis-regulatory modules and identify functional TFBS (47,49–51). Typically, multi-omic measurements are collected from different populations of cells using separate protocols. In contrast, barcoded SRT calling cards and transcriptomes can be collected simultaneously from the same cells which may improve the ability to link TF binding to changes in gene expression.

As a proof-of-principle of this method, we transiently overexpressed unfused hyperactive *piggyBac* or fusions with ASCL1 or MYOD1 in HEK293T cells, then collected RNA after one week. We then prepared SRT calling card and poly(A) 3' RNA-seq libraries in parallel. Cells transfected with ASCL1 and MYOD1 were co-transfected with non-overlapping pools of 12 barcoded SRTs to enable pooled SRT and transcriptome library preparation. Given the low rate of barcode swapping, this design would enable multiplexed TF profiling by pooled library construction. We performed four independent transfections for each factor. Compared to the recommended protocol for the original bulk RNA calling cards method for the same experiment (4), barcoded SRT calling cards and transcriptomes reduces material cost and labor of experiments by over 10-fold (Table 1).

Using the pooled barcoded SRT approach, we recovered hundreds of thousands of genomic insertion sites for each factor (Supplementary Table S4). Compared to unfused *piggyBac* binding sites from previous experiments, barcoded calling card peaks for ASCL1 and MYOD1 were again enriched for bHLH motifs including Ascl1 and MyoD (Figure 4B). Comparing genes near identified TFBS, we again observed ASCL1 and MYOD1 had shared and distinct binding profiles (Figure 4C) consistent with previous studies (26,27). The genomic insertion sites recovered strongly overlapped those from our earlier experiments with unpooled sequencing library preparations, but the total number of sites was lower. This could reflect reduced library complexity after pooling. Future experiments to opti-

mize pooled library complexity would further improve this methodology.

Next, to identify transcriptional consequences of TF overexpression, we analyzed the gene expression profiles that were simultaneously captured with SRTs. Supporting the approach of pooling of barcoded first strands for 3' gene expression library preparation, all 12 samples were well-represented in the sequencing data. Neither the average number of genes detected, nor the total RNA counts differed across factors and samples clustered by experimental condition (Supplementary Figure S3). We performed differential gene expression analysis on transcriptomes of cells transfected with ASCL1 or MYOD1 fusions compared against unfused *piggyBac* and identified 182 and 666 genes differentially expressed respectively. Of the differentially expressed genes, 170 and 480 were upregulated in ASCL1 and MYOD1 transfected cells respectively (Figure 4D), consistent with known roles of these transcription factors as activators of gene expression. Gene Ontology analysis of upregulated genes recapitulated some relevant pathways in ASCL1 transfected cells, but many pathways were not related to neurogenic or myogenic pathways (Figure 4E). Further, while some differentially expressed genes overlapped with genes near TFBS identified by barcoded SRTs, they were not enriched for such overlap. This is consistent with previous studies showing poor correlation between TF binding and gene expression (45–48). Nevertheless, these results demonstrate a novel method to simultaneously collect TFBS and gene expression changes from the same SRT calling card experiment which may facilitate the inference of functional TFBS.

Barcoded SRTs facilitate *in vivo* calling card experiments in mouse brain

We have previously demonstrated that SRT calling cards can be used to record TF binding *in vivo* (3,15), though often requiring ~10 technical replicates per biological replicate. To test whether barcoded SRTs can also function *in vivo* and reduce this need for technical replicates, we performed calling card experiments with barcoded and non-barcoded SRTs in the mouse cortex. We packaged tdTomato SRT plasmids with or without barcodes as AAV and delivered them to cortex of mice as described (3,15). Unfused *piggyBac* has an insertion preference at super-enhancers which are a class of enhancers regulating genes linked to cell identity (52,53). Leveraging this property, calling cards have been used to read out these important regulatory elements (3,15,54). To record these sites *in vivo*, we co-transduced mouse cortexes with unfused *piggyBac* and barcoded or non-barcoded SRTs.

After 21 days, we collected similar amounts of brain tissue from mice injected with barcoded or non-barcoded SRTs and prepared calling card libraries (Figure 5A). As with our *in vitro* experiments, all 25 unique barcodes were integrated into the genome and efficiently recovered (Figure 5B). Lower recovery of 2/25 barcodes may reflect imbalances in vector DNA pooling prior to AAV packaging. After normalizing by the total depth of sequencing, we found that use of barcodes improved the recovery of SRTs and yielded around 2-fold more genomic insertions than non-

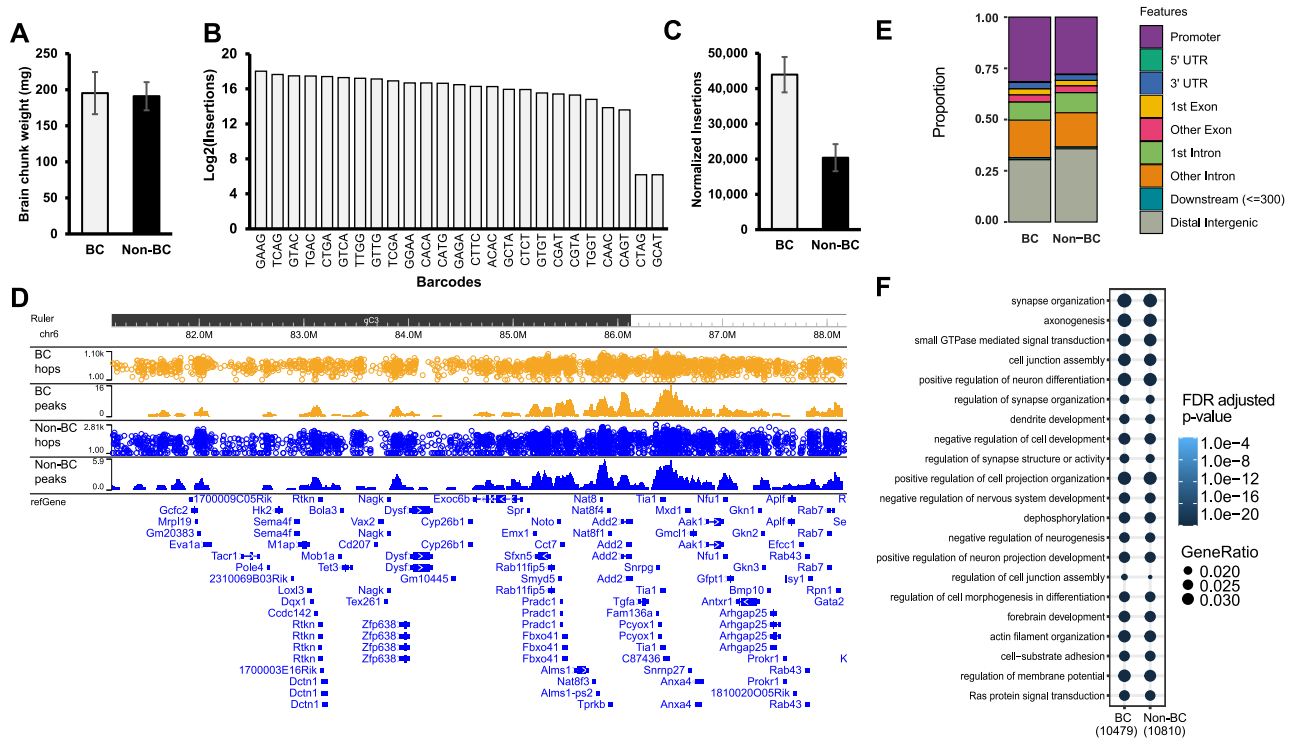


Figure 5. Comparison of barcoded and non-barcoded SRT calling cards *in vivo* in the mouse brain. (A) Equivalent amounts of brain tissue were collected after *in vivo* calling card experiments using a pool of 25 barcoded (BC) or non-barcoded (non-BC) SRT donors delivered by AAV. $n = 4$ for BC and 3 for Non-BC. (B) Number of genomic insertions recovered for each barcoded SRT. (C) Number of genomic insertions recovered at the same depth of sequencing for barcoded and non-barcoded SRTs. (D) Browser view of genomic insertions and called peaks for barcoded and non-barcoded SRTs. (E) Genomic features of peaks called by barcoded and non-barcoded experiments. (F) KEGG pathway enrichment comparison of genes near peaks called by barcoded and non-barcoded experiments.

barcoded counterparts (Figure 5C). Genome-wide, integrations of barcoded and non-barcoded SRTs were highly concordant. Visualizing insertions and called peaks across the genome demonstrates this concordance (Figure 5D). Analysis of genomic features of SRT insertion sites revealed similar insertional preferences (Figure 5E). We recovered more insertions in promoter regions using barcoded SRTs, suggesting the unbarcoded SRTs might have had especially limited dynamic range in these regions. This would be expected as some of these loci are expected to contain strong binding sites or few ‘TTAA’ sequences, limiting the quantification of non-barcoded insertions.

Next, we performed functional enrichment analysis of genes located near insertions. Based on the tropism of AAV9, we expected the vast majority of insertions to be in neuronal cells, with some insertions in astrocytes (3,15). Accordingly, barcoded and non-barcoded insertion sites were located near genes strongly enriched for neurological functions including synapse organization, forebrain development, and axonogenesis (Figure 5F). Functional enrichment was similar for insertions of barcoded and non-barcoded SRTs and consistent with our previous findings (3,15). Altogether, these results demonstrate that barcoded SRTs can recover biologically relevant binding events *in vivo* and outperform non-barcoded SRTs in the number of unique insertions at a fixed sequencing depth, while significantly reducing labor and reagent costs.

DISCUSSION

Understanding where TFs bind in the genome and how they orchestrate gene expression is a central goal in genomics (47,49–51). Calling cards is a powerful functional genomics method to identify the binding sites of TFs and other chromatin-associated factors in mammalian cells both *in vitro* and *in vivo* (1–3,15,22). The recently invented ‘self-reporting transposon’ converts the calling card recordings of TF binding to an RNA readout, enabling simultaneous profiling of gene expression and binding in single cells (3). Here, we present two crucial modifications of the SRT technology and associated protocols to enable parallel recording of TFBS and gene expression in bulk populations of cells: barcoded SRTs and barcoded transcriptomes. Besides enabling transcriptomic measurement, these improvements also drastically reduce the experimental cost and labor of calling card experiments.

First, we performed targeted mutagenesis of the *piggyBac* transposon TR region. We coupled a simple PCR mutagenesis method with SRT calling cards to rapidly screen for positions in the *piggyBac* TR that could be mutated while retaining compatibility with transposition. Through this, we discovered four consecutive nucleotides within the TR that were tolerant of a range of mutations, both singly and in combination, without markedly reducing transposition efficiency. To our knowledge, these are the first reported mu-

tations within the *piggyBac* terminal invert repeat that are compatible with transposition. We note the wild-type TR sequences are inserted at the highest frequency in our *in vitro* experiments compared to any 4-nt barcoded SRT so our targeted mutagenesis approach did not improve overall transposition efficiency. Nevertheless, for applications where the number of integrations is not paramount such as mutagenesis screening (55), cellular lineage tracing (56), and delivery of human gene therapy (57), we anticipate that barcoded *piggyBac* transposons will have broad utility beyond calling card assays.

As a resource to the community, we have individually cloned the top 24 integration-competent error-detecting barcodes into two versatile self-reporting transposon vectors compatible with *in vitro* and *in vivo* experiments. Barcoded error-detecting SRT vectors will allow experimental conditions (e.g. timepoint, drug treatment) to be uniquely barcoded, pooled, and accurately demultiplexed without sample mis-assignment due to sequencing errors. Of note, multiple transcription factors can be assayed simultaneously in a pooled experiment by using non-overlapping sets of barcoded SRTs for each TF.

We demonstrated that barcoded SRT calling cards identified TFBS for four bHLH transcription factors involved in cell fate specification and transdifferentiation: ASCL1, MYOD1, NEUROD2 and NGN1. We identified shared and unique binding sites for each factor and recovered binding motifs that matched known motifs for these factors. Supporting the identification of bona fide TFBS by barcoded SRT calling cards, we found that genes near TFBS were enriched for functions related to known functions of the assayed TFs. Barcoded SRT vectors reduced the experimental cost and labor of the calling card protocol by an order of magnitude which allowed us to easily measure TFBS for these four transcription factors in HEK293T cells. Remarkably, although HEK293T cells do not normally express any of the assayed TFs, all 4 TFs were able to recognize their cognate consensus motifs, and these motifs were located near genes with functions associated with each TF. The successful recovery of TFBS for all four TFs by barcoded calling cards demonstrates its versatility and its promise as an alternative to ChIP-seq, especially in the absence of well-validated antibodies. We have previously demonstrated the strong concordance between calling cards and ChIP-seq (3,15). We also demonstrate that barcoded SRTs improve the recovery of integration events in the mouse cortex. This will facilitate future studies using calling cards to record TFBS *in vivo*. Combined with Cre driver lines, this method is also compatible with cell-type specific recording of TFBS in a mixed populations of cells *in vitro* or *in vivo* (15).

Identifying TFBS is a first step toward understanding gene regulatory networks but many TF binding events have no or small effects on gene regulation (45–48). Integrating multi-omics datasets, such as TFBS and mRNA-seq, is therefore necessary to identify functional TFBS governing gene expression (47,49–51). Often, these multi-omics datasets are generated from different populations of cells using vastly different protocols that can introduce biases and batch effects. Since SRTs are expressed and collected

as RNA, TFBS and gene expression data can be simultaneously generated from the same RNA sample using calling cards. While this method has been demonstrated in single-cell experiments (3), bulk calling cards required modification to allow such joint measurement in bulk experiments. Specifically, we directly barcoded the SRT and introduced an additional barcode during reverse-transcription for barcoding mRNA. Combining barcoded SRT calling cards with bulk RNA barcoding and sequencing (BRB-seq) therefore enabled simultaneous identification of TFBS and gene expression in a protocol with drastically reduced cost and labor (13). In addition to bulk and single-cell calling cards³, emerging multiomic techniques including scDam&T-seq, Paired-Tag and CoTECH also enable the parallel measurement of transcription factor binding sites and gene expression (58–60).

We demonstrated that the combined protocol can jointly recover TFBS and gene expression during TF overexpression of the pioneer factors ASCL1 and MYOD1. Calling cards with barcoded SRTs and transcriptomes is therefore a novel and powerful method to infer functional TFBS in populations of cells. Technical and experimental optimizations of this method may improve its utility in future experiments. Further, using an inducible *piggyBac* system (14,61) would enable temporal measurement of binding and expression changes. That application would be especially powerful during the time course of cellular reprogramming experiments to link TFBS to gene expression changes controlling cell fate specification. Unlike most other methods of profiling DNA-associated proteins, calling cards enables the recovery of binding events at a future time point rather than at the time of binding. This enables the discovery of downstream outcomes of TF binding such as transcriptional changes or cell fate decisions.

Finally, our simple mutagenesis method will be useful for introducing barcodes to other DNA transposons such as *SleepingBeauty* and *Tol2*. Transposons are widely used for transgenics, mutagenesis, and functional genomics experiments (62). As the SRT protocol can easily scale to recover millions of genomic integration sites, insertion preferences for other transposons can be readily ascertained using this method. Each transposon has its own preferences for genomic integration which can have complementary uses. Further, insertion profiles can depend on chromatin state (52) so SRTs can potentially be used to read out chromatin status and histone modifications. For example, we have shown that calling card experiments using unfused *piggyBac* can identify super-enhancers (3,15,54). Joint measurements of *piggyBac* insertions and gene expression with this method may help link super-enhancers to gene regulatory networks.

In conclusion, barcoded SRTs simplify bulk calling cards experiments, enable barcoding of experimental conditions, and allow for pooled library preparations that drastically reduce cost and labor. Incorporating barcoded transcriptomes into the library preparation enables joint measurement of transcription factor binding and gene expression from the same biological sample. This method will facilitate the inference of gene regulatory networks for TFs involved in development, cellular reprogramming, and disease.

DATA AVAILABILITY

Raw and processed sequencing data generated in this study has been submitted to the NCBI Gene Expression Omnibus (GEO; <https://www.ncbi.nlm.nih.gov/geo/>) under accession number GSE195992. All software used to perform the analyses are available at:

SRT Calling card tools: https://github.com/arnavm/calling_cards

Barcoded SRT Calling cards tools: https://github.com/aMattScientist/barcoded_calling_cards

Peak Calling CCF tools: https://gitlab.com/rob.mitra/mammalian_cc_tools

We have implemented a running instance of an example barcoded SRT calling card analysis pipeline on Code Ocean: <https://doi.org/10.24433/CO.6494802.v1> (63).

SUPPLEMENTARY DATA

Supplementary Data are available at NARGAB Online.

ACKNOWLEDGEMENTS

We thank Nancy Craig for useful discussions regarding *piggyBac* mutagenesis. We thank Jessica Hoisington-Lopez and MariaLynn Crosby from the DNA Sequencing Innovation Lab at The Edison Family Center for Genome Sciences and Systems Biology for their sequencing expertise. We thank Mingjie Li at the Hope Center Viral Vectors Core for AAV packaging services.

FUNDING

National Heart, Lung, and Blood Institute [T32HL125241]; National Human Genome Research Institute [T32HG000045 to A.Y.]; National Institute of Child Health and Human Development [U54HD087011]; National Institute of Neurological Disorders and Stroke [R21NS087230-01A1 to R.D.M. and J.M.]; National Institute of Mental Health [RF1MH117070, RF1MH126723 to R.D.M., J.D.D.]; National Institute of General Medical Sciences [R01GM123203 to R.D.M.]; Simons Foundation Autism Research Initiative [SFARI Explorer 500661 to R.D.M., J.D.D.]; National Human Genome Research Institute [R21 HG009750 to R.D.M.]; Eunice Kennedy Shriver National Institute of Child Health & Human Development [P50 HD103525 to J.M., F30 HG009986 to A.M.]; Hope Center Viral Vectors Core and a P30 Neuroscience Blueprint Interdisciplinary Center Core award to Washington University [P30 NS057105]; GTAC@MGI is partially supported by National Institutes of Health (NIH) [P30 CA91842, UL1 TR000448]; M.A.L. is supported by the Seaver Foundation as a Seaver Faculty Scholar.

Conflict of interest statement. None declared.

REFERENCES

- Wang, H., Mayhew, D., Chen, X., Johnston, M. and Mitra, R.D. (2012) 'Calling Cards' for DNA-Binding proteins in mammalian cells. *Genetics*, **190**, 941–949.
- Wang, H., Johnston, M. and Mitra, R.D. (2007) Calling cards for DNA-binding proteins. *Genome Res.*, **17**, 1202–1209.
- Moudgil, A., Wilkinson, M.N., Chen, X., He, J., Cammack, A.J., Vasek, M.J., Lagunas, T., Qi, Z., Lalli, M.A., 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.
- Moudgil, A., Wilkinson, M., Chen, X. and Mitra, R. (2019) Bulk calling cards library preparation. *protocols.io*, <https://doi.org/10.17504/protocols.io.xwhfpb6>.
- Yusa, K., Zhou, L., Li, M.A., Bradley, A. and Craig, N.L. (2011) A hyperactive piggyBac transposase for mammalian applications. *PNAS*, **108**, 1531–1536.
- Omelina, E.S., Ivankin, A.V., Letiagina, A.E. and Pindyurin, A.V. (2019) Optimized PCR conditions minimizing the formation of chimeric DNA molecules from MPRA plasmid libraries. *BMC Genomics*, **20**, 536.
- Kebschull, J.M. and Zador, A.M. (2015) Sources of PCR-induced distortions in high-throughput sequencing data sets. *Nucleic Acids Res.*, **43**, e143–e143.
- Morellet, N., Li, X., Wieninger, S.A., Taylor, J.L., Bischerour, J., Moriau, S., Lescop, E., Bardiaux, B., Mathy, N., Assrir, N. *et al.* (2018) Sequence-specific DNA binding activity of the cross-brace zinc finger motif of the piggyBac transposase. *Nucleic Acids Res.*, **46**, 2660–2677.
- Wang, Y., Pryputniewicz-Dobrinska, D., Nagy, E.É., Kaufman, C.D., Singh, M., Yant, S., Wang, J., Dalda, A., Kay, M.A., Ivics, Z. *et al.* (2017) Regulated complex assembly safeguards the fidelity of sleeping beauty transposition. *Nucleic Acids Res.*, **45**, 311–326.
- Solodushko, V., Bitko, V. and Fouty, B. (2014) Minimal piggyBac vectors for chromatin integration. *Gene Ther.*, **21**, 1–9.
- Li, X., Harrell, R.A., Handler, A.M., Beam, T., Hennessy, K. and Fraser, M.J. (2005) piggyBac internal sequences are necessary for efficient transformation of target genomes. *Insect Mol. Biol.*, **14**, 17–30.
- Li, X., Lobo, N., Bauser, C. and Fraser, M. (2001) The minimum internal and external sequence requirements for transposition of the eukaryotic transformation vector piggyBac. *Mol. Gen. Genomics*, **266**, 190–198.
- Alpern, D., Gardeux, V., Russeil, J., Mangeat, B., Meireles-Filho, A.C.A., Breyse, R., Hacker, D. and Deplancke, B. (2019) BRB-seq: ultra-affordable high-throughput transcriptomics enabled by bulk RNA barcoding and sequencing. *Genome Biol.*, **20**, 71.
- Cadiñanos, J. and Bradley, A. (2007) Generation of an inducible and optimized piggyBac transposon system. *Nucleic Acids Res.*, **35**, e87.
- Cammack, A.J., Moudgil, A., Chen, J., Vasek, M.J., Shabsovich, M., McCullough, K., Yen, A., Lagunas, T., Maloney, S.E., He, J. *et al.* (2020) A viral toolkit for recording transcription factor–DNA interactions in live mouse tissues. *Proc. Natl. Acad. Sci. U.S.A.*, **117**, 10003–10014.
- Moudgil, A., Wilkinson, M., Chen, X. and Mitra, R.D. (2019) Mammalian calling cards quick start guide. *protocols.io*, <https://doi.org/10.17504/protocols.io.xurfnv6>.
- Lalli, M.A., Avey, D., Dougherty, J.D., Milbrandt, J. and Mitra, R.D. (2020) High-throughput single-cell functional elucidation of neurodevelopmental disease-associated genes reveals convergent mechanisms altering neuronal differentiation. *Genome Res.*, **30**, 1317–1331.
- R Core Team (2018) In: *R: A Language and Environment for Statistical Computing R Foundation for Statistical Computing*. Vienna, Austria.
- Robinson, M.D., McCarthy, D.J. and Smyth, G.K. (2010) edgeR: a bioconductor package for differential expression analysis of digital gene expression data. *Bioinformatics*, **26**, 139–140.
- Smith, T.S., Heger, A. and Sudbery, I. (2017) UMI-tools: modelling sequencing errors in unique molecular identifiers to improve quantification accuracy. *Genome Res.*, **27**, 491–499.
- Moudgil, A., Li, D., Hsu, S., Purushotham, D., Wang, T. and Mitra, R.D. (2020) The qBED track: a novel genome browser visualization for point processes. *Bioinformatics*, **37**, 1168–1170.
- Yen, M., Qi, Z., Chen, X., Cooper, J.A., Mitra, R.D. and Onken, M.D. (2018) Transposase mapping identifies the genomic targets of BAP1 in uveal melanoma. *BMC Med Genomics*, **11**, 97.
- Zhang, Y., Liu, T., Meyer, C.A., Eeckhoutte, J., Johnson, D.S., Bernstein, B.E., Nussbaum, C., Myers, R.M., Brown, M., Li, W. *et al.*

- (2008) Model-based analysis of ChIP-Seq (MACS). *Genome Biol.*, **9**, R137.
24. Buschmann, T. and Bystrykh, L.V. (2013) Levenshtein error-correcting barcodes for multiplexed DNA sequencing. *BMC Bioinf.*, **14**, 272.
 25. Rohs, R., West, S.M., Sosinsky, A., Liu, P., Mann, R.S. and Honig, B. (2009) The role of DNA shape in protein–DNA recognition. *Nature*, **461**, 1248–1253.
 26. Casey, B.H., Kollipara, R.K., Pozo, K. and Johnson, J.E. (2018) Intrinsic DNA binding properties demonstrated for lineage-specifying basic helix-loop-helix transcription factors. *Genome Res.*, **28**, 484–496.
 27. Lee, Q.Y., Mall, M., Chanda, S., Zhou, B., Sharma, K.S., Schaukowitz, K., Adrian-Segarra, J.M., Grieder, S.D., Karetta, M.S., Wapinski, O.L. *et al.* (2020) Neuro-neuronal activity of Myod1 due to promiscuous binding to neuronal genes. *Nat. Cell Biol.*, **22**, 401–411.
 28. Wapinski, O.L., Vierbuchen, T., Qu, K., Lee, Q.Y., Chanda, S., Fuentes, D.R., Giresi, P.G., Ng, Y.H., Marro, S., Neff, N.F. *et al.* (2013) Hierarchical mechanisms for direct reprogramming of fibroblasts to neurons. *Cell*, **155**, 621–635.
 29. Yoo, A.S., Sun, A.X., Li, L., Shcheglovitov, A., Portmann, T., Li, Y., Lee-Messer, C., Dolmetsch, R.E., Tsien, R.W. and Crabtree, G.R. (2011) MicroRNA-mediated conversion of human fibroblasts to neurons. *Nature*, **476**, 228–231.
 30. Blanchard, J.W., Eade, K.T., Szücs, A., Lo Sardo, V., Tsunemoto, R.K., Williams, D., Sanna, P.P. and Baldwin, K.K. (2015) Selective conversion of fibroblasts into peripheral sensory neurons. *Nat. Neurosci.*, **18**, 25–35.
 31. Davis, R.L., Weintraub, H. and Lassar, A.B. (1987) Expression of a single transfected cDNA converts fibroblasts to myoblasts. *Cell*, **51**, 987–1000.
 32. Webb, A.E., Pollina, E.A., Vierbuchen, T., Urbán, N., Ucar, D., Leeman, D.S., Martynoga, B., Sewak, M., Rando, T.A., Guillemot, F. *et al.* (2013) FOXO3 shares common targets with ASCL1 Genome-wide and inhibits ASCL1-Dependent neurogenesis. *Cell Rep.*, **4**, 477–491.
 33. Heinz, S., Benner, C., Spann, N., Bertolino, E., Lin, Y.C., Laslo, P., Cheng, J.X., Murre, C., Singh, H. and Glass, C.K. (2010) Simple combinations of lineage-determining transcription factors prime cis-Regulatory elements required for macrophage and B cell identities. *Mol. Cell*, **38**, 576–589.
 34. Yu, G., Wang, L.-G., Han, Y. and He, Q.-Y. (2012) clusterProfiler: an R package for comparing biological themes among gene clusters. *OMICS*, **16**, 284–287.
 35. Gordán, R., Shen, N., Dror, I., Zhou, T., Horton, J., Rohs, R. and Bulky, M.L. (2013) Genomic regions flanking E-Box binding sites influence DNA binding specificity of bHLH transcription factors through DNA shape. *Cell Rep.*, **3**, 1093–1104.
 36. Brigidi, G.S., Hayes, M.G.B., Delos Santos, N.P., Hartzell, A.L., Texari, L., Lin, P.-A., Bartlett, A., Ecker, J.R., Benner, C., Heinz, S. *et al.* (2019) Genomic decoding of neuronal depolarization by stimulus-specific NPAS4 heterodimers. *Cell*, **179**, 373–391.
 37. Fong, A.P., Yao, Z., Zhong, J.W., Johnson, N.M., Farr, G.H., Maves, L. and Tapscott, S.J. (2015) Conversion of MyoD to a neurogenic factor: binding site specificity determines lineage. *Cell Rep.*, **10**, 1937–1946.
 38. Vierbuchen, T., Ostermeier, A., Pang, Z.P., Kokubu, Y., Südhof, T.C. and Wernig, M. (2010) Direct conversion of fibroblasts to functional neurons by defined factors. *Nature*, **463**, 1035–1041.
 39. Lu, C., Shi, X., Allen, A., Baez-Nieto, D., Nikish, A., Sanjana, N.E. and Pan, J.Q. (2019) Overexpression of NEUROG2 and NEUROG1 in human embryonic stem cells produces a network of excitatory and inhibitory neurons. *FASEB J.*, **33**, 5287–5299.
 40. Zheng, G.X.Y., Terry, J.M., Belgrader, P., Ryvkin, P., Bent, Z.W., Wilson, R., Ziraldo, S.B., Wheeler, T.D., McDermott, G.P., Zhu, J. *et al.* (2017) Massively parallel digital transcriptional profiling of single cells. *Nat. Commun.*, **8**, 14049.
 41. Picelli, S., Björklund, Å.K., Reinius, B., Sagasser, S., Winberg, G. and Sandberg, R. (2014) Tn5 transposase and tagmentation procedures for massively scaled sequencing projects. *Genome Res.*, **24**, 2033–2040.
 42. Tapscott, S.J. (2005) The circuitry of a master switch: myod and the regulation of skeletal muscle gene transcription. *Development*, **132**, 2685–2695.
 43. Iwafuchi-Doi, M. and Zaret, K.S. (2016) Cell fate control by pioneer transcription factors. *Development*, **143**, 1833–1837.
 44. Iwafuchi-Doi, M. and Zaret, K.S. (2014) Pioneer transcription factors in cell reprogramming. *Genes Dev.*, **28**, 2679–2692.
 45. Fisher, W.W., Li, J.J., Hammonds, A.S., Brown, J.B., Pfeiffer, B.D., Weismann, R., MacArthur, S., Thomas, S., Stamatoyannopoulos, J.A., Eisen, M.B. *et al.* (2012) DNA regions bound at low occupancy by transcription factors do not drive patterned reporter gene expression in Drosophila. *Proc. Natl. Acad. Sci. U.S.A.*, **109**, 21330–21335.
 46. Whitfield, T.W., Wang, J., Collins, P.J., Partridge, E.C., Aldred, S.F., Trinklein, N.D., Myers, R.M. and Weng, Z. (2012) Functional analysis of transcription factor binding sites in human promoters. *Genome Biol.*, **13**, R50.
 47. Cusanovich, D.A., Pavlovic, B., Pritchard, J.K. and Gilad, Y. (2014) The functional consequences of variation in transcription factor binding. *PLoS Genet.*, **10**, e1004226.
 48. Paris, M., Kaplan, T., Li, X.Y., Villalta, J.E., Lott, S.E. and Eisen, M.B. (2013) Extensive divergence of transcription factor binding in drosophila embryos with highly conserved gene expression. *PLoS Genet.*, **9**, e1003748.
 49. Dunham, I., Kundaje, A., Aldred, S.F., Collins, P.J., Davis, C.A., Doyle, F., Epstein, C.B., Frietze, S., Harrow, J., Kaul, R. *et al.* (2012) An integrated encyclopedia of DNA elements in the human genome. *Nature*, **489**, 57–74.
 50. Gerstein, M.B., Kundaje, A., Hariharan, M., Landt, S.G., Yan, K.-K., Cheng, C., Mu, X.J., Khurana, E., Rozowsky, J., Alexander, R. *et al.* (2012) Architecture of the human regulatory network derived from ENCODE data. *Nature*, **489**, 91–100.
 51. Kundaje, A., Meuleman, W., Ernst, J., Bilenky, M., Yen, A., Heravi-Moussavi, A., Kheradpour, P., Zhang, Z., Wang, J., Ziller, M.J. *et al.* (2015) Integrative analysis of 111 reference human epigenomes. *Nature*, **518**, 317–330.
 52. Yoshida, J., Akagi, K., Misawa, R., Kokubu, C., Takeda, J. and Horie, K. (2017) Chromatin states shape insertion profiles of the piggyBac, Tol2 and sleeping beauty transposons and murine leukemia virus. *Sci. Rep.*, **7**, 43613.
 53. Whyte, W.A., Orlando, D.A., Hnisz, D., Abraham, B.J., Lin, C.Y., Kagey, M.H., Rahl, P.B., Lee, T.I. and Young, R.A. (2013) Master transcription factors and mediator establish super-enhancers at key cell identity genes. *Cell*, **153**, 307–319.
 54. Kfoury, N., Qi, Z., Prager, B.C., Wilkinson, M.N., Broestl, L., Berrett, K.C., Moudgil, A., Sankaraman, S., Chen, X., Gertz, J. *et al.* (2021) Brd4-bound enhancers drive cell-intrinsic sex differences in glioblastoma. *Proc. Natl. Acad. Sci. U.S.A.*, **118**.
 55. Rad, R., Rad, L., Wang, W., Cadinanos, J., Vassiliou, G., Rice, S., Campos, L.S., Yusa, K., Banerjee, R., Li, M.A. *et al.* (2010) PiggyBac transposon mutagenesis: a tool for cancer gene discovery in mice. *Science*, **330**, 1104–1107.
 56. Siddiqi, F., Chen, F., Aron, A.W., Fiondella, C.G., Patel, K. and LoTurco, J.J. (2014) Fate mapping by piggybac transposase reveals that neocortical GLAST+ progenitors generate more astrocytes than Nestin+ progenitors in rat neocortex. *Cereb. Cortex*, **24**, 508–520.
 57. Li, R., Zhuang, Y., Han, M., Xu, T. and Wu, X. (2013) piggyBac as a high-capacity transgenesis and gene-therapy vector in human cells and mice. *Dis. Models Mech.*, **6**, 828–833.
 58. Rooijers, K., Markodimitraki, C.M., Rang, F.J., de Vries, S.S., Chialastri, A., de Luca, K.L., Mooijman, D., Dey, S.S. and Kind, J. (2019) Simultaneous quantification of protein–DNA contacts and transcriptomes in single cells. *Nat. Biotechnol.*, **37**, 766–772.
 59. Zhu, C., Zhang, Y., Li, Y.E., Lucero, J., Behrens, M.M. and Ren, B. (2021) Joint profiling of histone modifications and transcriptome in single cells from mouse brain. *Nat. Methods*, **18**, 283–292.
 60. Xiong, H., Luo, Y., Wang, Q., Yu, X. and He, A. (2021) Single-cell joint detection of chromatin occupancy and transcriptome enables higher-dimensional epigenomic reconstructions. *Nat. Methods*, **18**, 652–660.
 61. Qi, Z., Wilkinson, M.N., Chen, X., Sankaraman, S., Mayhew, D. and Mitra, R.D. (2017) An optimized, broadly applicable piggyBac transposon induction system. *Nucleic Acids Res.*, **45**, e55.
 62. Kawakami, K., Largaespada, D.A. and Ivics, Z. (2017) Transposons as tools for functional genomics in vertebrate models. *Trends Genet.*, **33**, 784–801.
 63. Lalli, M.A. and Mitra, R. (2022) Barcoded self-reporting transposon calling cards processing pipeline. *Code Ocean*, <https://doi.org/10.24433/CO.6494802.v1>.