

1 **Massively multiplex single-cell Hi-C**

2 Vijay Ramani¹, Xinxian Deng², Kevin L Gunderson³, Frank J Steemers³, Christine M
3 Disteche^{2,4}, William S Noble¹, Zhijun Duan^{5,6#}, Jay Shendure^{1,7#}

4

5 ¹ Department of Genome Sciences, University of Washington, Seattle, WA

6 ² Department of Pathology, University of Washington, Seattle, WA

7 ³ Illumina Inc., Advanced Research Group, San Diego, CA

8 ⁴ Department of Medicine, University of Washington, Seattle, WA

9 ⁵ Division of Hematology, University of Washington School of Medicine, Seattle, WA

10 ⁶ Institute for Stem Cell and Regenerative Medicine, University of Washington, Seattle,

11 WA

12 ⁷ Howard Hughes Medical Institute, Seattle, WA

13

14 #Correspondence to Zhijun Duan (zjduan@uw.edu) and Jay Shendure
15 (shendure@uw.edu)

16

17 **Keywords:** Nuclear architecture; Hi-C; single-cell

18 **Working title:** Single-cell Hi-C by combinatorial cellular indexing

19

20 **Abstract**

21 We present combinatorial single cell Hi-C, a novel method that leverages combinatorial
22 cellular indexing to measure chromosome conformation in large numbers of single cells.

23 In this proof-of-concept, we generate and sequence combinatorial single cell Hi-C

24 libraries for two mouse and four human cell types, comprising a total of 9,316 single cells

25 across 5 experiments. We demonstrate the utility of single-cell Hi-C data in separating

26 different cell types, identify previously uncharacterized cell-to-cell heterogeneity in the

27 conformational properties of mammalian chromosomes, and demonstrate that

28 combinatorial indexing is a generalizable molecular strategy for single-cell genomics.

1 **Main Text**

2 Our understanding of genome architecture has largely progressed through the successive
3 development of new technologies¹. Advances in microscopy revealed the presence of
4 “chromosome territories”—nuclear regions that preferentially self-associate in a manner
5 correlated with transcriptional activity². The invention of Chromosome Conformation
6 Capture (3C) and its derivatives³ resulted in a proliferation of data measuring genome
7 architecture and its relation to other aspects of nuclear biology at increasing resolution.

8
9 3C assays rely on the concept of proximity ligation, a technique that has been used to
10 measure local protein-protein⁴, RNA-RNA⁵, and DNA-DNA interactions⁶. By coupling
11 an “all-vs-all” 3C assay with massively parallel sequencing^{7,8} (*e.g.* “Hi-C”), one is able to
12 query relative contact probabilities genome-wide. However, contact probabilities
13 generated by these assays represent ensemble averages of the respective conformations of
14 the millions of nuclei used as input, and scalable techniques characterizing the variance
15 underlying these population averages remain largely underdeveloped. A pioneering study
16 in 2013 demonstrated proof-of-concept that Hi-C could be performed on single isolated
17 mouse nuclei, but relied on the physical separation and processing of single murine cells
18 in independent reaction volumes, with consequent low-throughput⁹.

19
20 The repertoire of high-throughput single-cell techniques for other biochemical assays has
21 expanded rapidly as of late¹⁰⁻¹³. Single-cell RNA-seq (scRNA-seq) was recently paired
22 with droplet-based microfluidics to markedly increase its throughput^{11,12}. Orthogonally,
23 we introduced the concept of combinatorial cellular indexing¹⁰, a method that eschews

1 microfluidic manipulation and instead tags the DNA within intact nuclei with successive
2 (combinatorial) rounds of nucleic acid barcodes, to measure chromatin accessibility
3 (scATAC-seq) in thousands of single cells without physically isolating each single cell.
4 However, such throughput-boosting strategies have yet to be successfully adapted to
5 single-cell chromosome conformation analysis.

6
7 To address this gap, we sought to develop a high-throughput, easy-to-implement single-
8 cell Hi-C protocol (**Figure 1a**), based on the concept of combinatorial indexing and also
9 building on recent improvements to the Hi-C protocol^{14,15}. A population of cells is fixed,
10 lysed to generate nuclei, and restriction digested *in situ* with the enzyme DpnII. Nuclei
11 are then distributed to 96 wells, wherein the first barcode is introduced through ligation
12 of barcoded biotinylated double-stranded bridge-adaptors. Intact nuclei are then pooled
13 and proximity ligated all together, followed by dilution and redistribution to a second 96-
14 well plate. Importantly, this dilution is carried out such that each well in this second plate
15 contains at most 25 nuclei. Following lysis, a second barcode is introduced through
16 ligation of barcoded Y-adaptors.

17
18 As the number of barcode combinations (96 x 96) exceeds the number of nuclei (96 x
19 25), the vast majority of single nuclei are tagged by a unique combination of barcodes.
20 All material is once again pooled, and biotinylated junctions are purified with
21 streptavidin beads, restriction digested, and further processed to Illumina sequencing
22 libraries. Sequencing these molecules with relatively long paired-end reads (*i.e.* 2 x 250
23 base pair (bp)) allows one to identify not only the genome-derived fragments of

1 conventional Hi-C, but also external and internal barcodes (each combination of which is
2 hereafter referred to as a ‘cellular index’) which enable decomposition of the Hi-C data
3 into single-cell contact probability maps (**Figure 1b**). Like scATAC-seq with
4 combinatorial cellular indexing¹⁰, this protocol can process hundreds to thousands of cells
5 per experiment without requiring the physical isolation of each cell.

6
7 As a proof-of-concept, we applied combinatorial single cell Hi-C to synthetic mixtures of
8 cell lines derived from mouse (primary mouse embryonic fibroblasts (MEFs), and the
9 ‘Patski’ embryonic fibroblast line) and human (HeLa S3, the HAP1 cell line, K562, and
10 GM12878; all five experiments and sequenced libraries are summarized in **Table 1**,
11 although we focus on ML1 and ML2 biological replicates in the text). All experiments
12 were carried out such that subsets of cell types received specific barcodes during the first
13 round of barcoding (*e.g.* in ML1 and ML2, each well during the first round of barcoding
14 contained either HeLa S3 + Patski cells or HAP1 + MEF cells; see **Methods**).

15
16 Before deconvolving the resulting data to single cells, we examined the overall
17 distribution of ligation junctions (*i.e.* contacts). Encouragingly, there were very few
18 contacts between mouse and human (ML1: 0.006%; ML2: 0.008%), demonstrating
19 minimal cross-talk between cellular indices, and that nuclei remain intact through all
20 ligation steps (confirmed through phase-contrast microscopy; **Supplementary Figure 1**).
21 We also examined the *cis:trans* ratio, defined here as the ratio of long-range (*i.e.* >20 kb)
22 intrachromosomal contacts to interchromosomal contacts (**Figure 1c**), and found it to be
23 on par with expectation for high-quality Hi-C datasets (ML1: 4.41; ML2: 4.38).

1

2 We next split the Hi-C data by cellular index and characterized the number of unique
3 read-pairs associated with each, the vast majority of which should correspond to single
4 cells. When examining a histogram of unique index occurrences as a function of read
5 depth, we noted a bimodal distribution, reminiscent of patterns seen in scATAC-seq
6 datasets¹⁰, where low-coverage indices likely represent ‘noise’ consequent to tags from
7 free DNA in solution (**Supplementary Figure 2**). After discarding these, we infer 1,081
8 cellular indices in ML1, with a median of 9,274 unique read-pairs per index (ML2: 841
9 cellular indices; median of 8,335 unique read-pairs per index). Importantly, we also
10 observe minimal barcode bias across replicate experiments (**Supplementary Figure 3**),
11 as well as similar median *cis:trans* ratios per cell (ML1: 4.43 with median absolute
12 deviation (MAD) of 1.66; ML2: 4.34 with MAD of 1.66) (**Figure 1d, Supplementary**
13 **Figure 4**).

14

15 The only previously published example of single-cell Hi-C data suggests that high single
16 cell *cis:trans* ratios are a hallmark of high-quality single-cell data⁹. The high *cis:trans*
17 ratios that we observe are comparable to those of the 10 single-cell maps generated in
18 that study, which reported a median value of 6.26 (MAD = 0.74), calculated as the ratio
19 of *all* intrachromosomal contacts to interchromosomal contacts (*i.e.* with no cutoff for
20 minimal intrachromosomal distance). Reanalyzing our own data using this more liberal
21 criterion yielded similar ratios of 6.17 (ML1; MAD = 1.99) and 5.96 (ML2; MAD =
22 1.94). Of note, our ratios are calculated over 1,922 cellular indices (ML1 and ML2
23 combined), 857 of which have more than 10,000 unique contacts, compared to the 10

1 previously reported single cells each with at least 10,000 unique contacts. This
2 comparison illustrates the scalability of combinatorial methods, as compared with
3 methods relying on the physical isolation and serial processing of each single cell.

4

5 We designed our experiments to facilitate validation of the single-cell origin of each
6 cellular index. Uniquely tagged cells should be associated with species-specific cellular
7 indices in mixture experiments, with a collision rate broadly defined by a formulation of
8 the “birthday problem”¹⁰. Consistent with the expected collision rate, we observed that
9 4.53% of all ML1 cellular indices (4.40% in ML2) were “collisions” (*i.e.* had less than
10 95% of reads mapping to either the mouse or human genome) (**Figure 2a,b**). For further
11 analyses we filtered out any cellular indices failing this criterion, while accepting that we
12 remain blind to “within species” collisions. We also filtered out indices where the
13 associated *cis:trans* ratio was less than 1 (1.94% of indices in ML1; 1.62% in ML2),
14 which could suggest broken nuclei.

15

16 Before continuing, we combined filtered data from ML1 and ML2 with equivalently
17 filtered data from secondary experiments (PL1 and PL2) (**Table 1, Supplementary**
18 **Figure 5**). We then employed a conservative genotype filter¹⁶ which removed 20.4% of
19 human cellular indices (**Supplementary Figure 6**), leaving us with a combined dataset of
20 3,609 human single cell Hi-C maps. Together with mouse data (which were filtered for
21 coverage, *cis:trans* ratio, and species purity), a total of 8,141 single cell Hi-C maps were
22 generated across these four experiments.

23

1 We next explored whether cell types could be separated *in silico* on the basis of single-
2 cell Hi-C signal. We generated matrices where rows represent single cells, and columns
3 represent the number of contacts between pairs of chromosomes (**Supplementary Figure**
4 **7a**). Principal components analysis (PCA) on this matrix resulted in separation of single
5 HeLa S3 and HAP1 cells (**Figure 2c**), which was validated by our programmed barcode
6 associations. Principal component 1 (PC1), which strongly correlated with coverage
7 (**Supplementary Figure 8**), accounted for the majority of the variance (52.1%), while
8 the combination of PC1 and principal component 2 (PC2; 1.07% of the variance)
9 separated HeLa S3 and HAP1 cells. We then analyzed the “loadings” of our features in
10 PC2, the axis separating HeLa S3 and HAP1 cells, and found that the strongest loadings
11 recapitulated known translocations specific to HAP1¹⁷ (namely, translocations between
12 chromosomes 15 and 19, and between chromosomes 9 and 22), while other strong
13 loadings corresponded to documented HeLa S3 translocations^{16,18} (**Figure 2d**). Repeating
14 these analyses by i.) removing specific interactions from the matrices and repeating PCA
15 (**Supplementary Figure 9**) ii.) using an alternate feature set (interacting 10 Mb
16 intrachromosomal windows; **Supplementary Figures 7b, 10**), iii.) separating cells by
17 replicate (**Supplementary Figure 11**), and iv.) sequencing 908 additional human cells
18 (K562 and GM12878; Library ML3 containing 1,175 cells total; **Supplementary Figure**
19 **12**), all recapitulated cell-type separation to varying degrees, demonstrating that PCA can
20 potentially be used to separate cell types on the basis of Hi-C signal (with the caveat that
21 the separations observed here may be driven by karyotype differences between these cell
22 types).
23

1 We next examined the heterogeneity present in single cell Hi-C maps in terms of polymer
2 conformation. We plotted contact probability as a function of genomic distance for 769
3 single cells, each with at least 10,000 unique contacts (**Figure 3a**), finding that the
4 pattern of scaling observed for single cells was markedly more disperse when compared
5 to a shuffled control where the assignment of cellular indices to reads are randomized,
6 regardless of species analyzed. We then examined the relationship between single-cell
7 power-law scaling coefficients (**Figure 3b**), calculated between distances of 50 kb and 8
8 Mb^{19,20}, and single-cell *cis:trans* ratios, noting a correlation across four out of five
9 experiments (**Figure 3c, Supplementary Figure 13**) between high *cis:trans* ratios and
10 shallow scaling coefficients. Although beyond the scope of our methodological proof-of-
11 concept, these empirical observations of cell-to-cell heterogeneity in contact probability
12 distributions are likely to be highly useful in constraining computational models of
13 mammalian chromosome conformation.

14

15 In summary, we present a novel method for single-cell Hi-C that relies on the concept of
16 combinatorial cellular indexing for rapid scaling to large numbers of cells. For this proof-
17 of-concept, we applied this method to generate single-cell Hi-C maps for 9,316 cells with
18 at least 1,000 unique contacts. This dataset is two orders of magnitude larger than the
19 only published single-cell Hi-C dataset, with 2,563 filtered cells containing more than
20 10,000 unique contacts, compared to the 10 existing single-cell maps defined using a
21 similar coverage cutoff. Looking forward, an important technical goal is to further
22 increase the number of unique contacts obtained per single cell, as well as to increase the
23 number of single cells processed per experiment. Importantly, our combinatorial

1 approach is internally controlled in the sense that key steps are carried out in a “single
2 pot”, thus mitigating technical confounders of conventional (serial) replicates of single
3 cell or bulk experiments.

4

5 Given the generally similar workflow of our method and traditional bulk Hi-C, it may be
6 possible to incorporate into routine practice, thus adding a ‘single cell’ dimension to Hi-C
7 data production and a means of obtaining single-cell and bulk measurement at once (the
8 latter generated by summing single cells). Furthermore, our demonstration that thousands
9 of single-cell Hi-C maps can be generated in a single workflow, without the need to
10 isolate each cell, demonstrates the power of combinatorial indexing for large-scale single
11 cell biology. Combinatorial indexing may thus be generalizable to additional aspects of
12 single cell or even intracellular biology where DNA barcodes can be incorporated *in situ*.

13

14 **Methods**

15 **Cell Culture**

16 HeLa S3 (CCL2.2), primary MEFs, and Patski cells were cultured at 37°C, 5% CO₂ in
17 DMEM supplemented with 1X Pen-Strep (Gibco), and 10% FBS (Gibco). HAP1 cells
18 were cultured were cultured at 37°C, 5% CO₂ in IMDM supplemented with 1X Pen-Strep
19 and 10% FBS. K562 cells were cultured at 37°C, 5% CO₂ in RPMI-1640 supplemented
20 with 1X Pen-Strep and 10% FBS. GM12878 cells were cultured at 37°C, 5% CO₂ in
21 RPMI-1640 supplemented with 1X Pen-Strep and 15% FBS.

22

23 **Cell Fixation**

1 Adherent cells (*i.e.* HeLa S3, HAP1, Patski, MEF) were washed once with 1X PBS (Life
2 Technologies), trypsinized (0.25% Trypsin-EDTA, Life Technologies), spun down at
3 500xg for 5 min., and resuspended in 20 mL serum-free DMEM (IMDM for HAP1).
4 Cells were crosslinked by adding 1.12 mL (2% final concentration, for HeLa S3, HAP1,
5 and MEF) or 1.4 mL (2.5% final concentration, for Patski) 37% formaldehyde (Alcon)
6 and incubated at RT (25°C) for 10 min., after which crosslinking was quenched using 1
7 mL 2.5M glycine. Quenched reactions were incubated on ice for 15 min., spun down at
8 800xg for 5 min., resuspended in 1X PBS, aliquoted into 10E6 cell aliquots, pelleted once
9 again at 800xg for 5 min, decanted, snap frozen in liquid nitrogen, and finally stored
10 indefinitely at -80°C.

11

12 Suspension cells (*i.e.* K562, GM1878) were spun down at 500xg for 5 min., resuspended
13 in 20 mL serum-free RPMI-1640, crosslinked with a final concentration of 2%
14 formaldehyde, and processed as above.

15

16 **Combinatorial Single Cell Hi-C**

17 For the step-by-step combinatorial single cell Hi-C protocol, see **Supplementary**
18 **Protocol**. Like the recently published scDNase-seq protocol²¹, combinatorial single cell
19 Hi-C uses carrier plasmid to prevent DNA losses during steps of the protocol where small
20 amounts of DNA are handled.

21

22 All oligonucleotide sequences used in this study were obtained from IDT Technologies,
23 and are detailed in **Supplementary File 1**. All libraries were sequenced on a HiSeq 2500.

1

2 *Barcode Programming*

3 Our primary datasets (Library ML1 and biological replicate library ML2), used HeLa S3,
4 HAP1, Patski, and MEFs, with subsets of human and mouse cell types in distinct wells
5 during the first round of barcoding (HeLa S3 + Patski in half of wells; HAP1 + MEFs in
6 half of wells). Our secondary datasets (Library PL1 and biological replicate PL2) were
7 generated using the same cell types, but a subtly different programming scheme
8 (illustrated in **Supplementary Figure 14**), wherein each well contained only a single cell
9 type during the first round of barcoding. Finally, we generated and lightly sequenced a 5th
10 library (Library ML3), mixing the same murine cell types as before with two new human
11 cell types—GM12878 and K562—in a similar manner to Libraries ML1 and ML2
12 (GM12878 + Patski in half of wells; K562 + MEFs in half of wells).

13

14 **Bridge Adaptor Barcode Design**

15 Bridge adaptor barcodes were drawn from randomly generated 8-mers, such that the
16 following criteria were met: i.) all adaptors must have a minimum pairwise Levenshtein
17 distance of 3; ii.) adaptors must not contain the sequences TTAA or AAGCTT; iii.)
18 adaptors must contain >60% GC content; iv.) adaptors must not contain homopolymers
19 \geq length 3; and v.) adaptors must not be palindromic.

20

21 **Processing Combinatorial Single Cell Hi-C Data**

22 All code used for combinatorial single-cell Hi-C data analysis will be available (along
23 with all data) upon publication at <https://github.com/VRam142/combinatorialHiC>.

1 Below, we describe in detail the analytical pipeline used to process combinatorial single-
2 cell Hi-C data. The analytical steps broadly fall under three categories: i.) Barcode
3 Identification & Read Trimming, ii.) Read Alignment, Read Pairing, & Barcode
4 Association, and iii.) Cellular Demultiplexing & Quality Analysis.

5

6 *Barcode Association & Read Trimming*

7 First, to obtain round 2 (*i.e.* terminal) barcodes, we use a custom Python script to iterate
8 through both mates, compare the first 8 bases of each read against the 96 known barcode
9 sequences, and then assign barcodes to each mate using a Levenshtein distance cutoff of
10 2. Reads “split” in this way are output such that the first 11 bases of each read, which
11 derive from the custom barcoded Y adaptors, are removed. Mates where either terminal
12 barcode went unidentified, or where the terminal barcodes did not match, are discarded.

13

14 For each resulting “split” pair of reads, the two reads are then scanned using a custom
15 Python script to find the common portion of the bridge adaptor sequence. The 8 bases
16 immediately 5’ of this sequence are isolated and compared against the 96 known bridge
17 adaptor barcodes, again using a Levenshtein distance cutoff of 2. There are cases where
18 the entire bridge adaptor, including both barcodes flanking the ligation junction, is
19 encountered in one mate, and not the other. To account for these cases, we also isolate the
20 8 bases flanking the 3’ end of the common bridge adaptor sequence (when it is
21 encountered within a read), reverse complement it, and compare the resulting 8-mer
22 against the 96 known bridge adaptor barcodes. Output reads are then clipped to remove
23 the bridge adaptor and all 3’ sequence. Barcodes flanking the ligation junction should

1 match; again, mates where barcodes do not match, or where a barcode is not found are
2 discarded.

3

4 The result of this processing module are three files: filtered reads 1 and 2, and an
5 “associations” file—a tab-delimited file where the name of each read passing the above
6 filters and their associated barcode combination are listed.

7

8 *Read Alignment, Read Pairing, & Barcode Association*

9 As is standard for Hi-C reads, the resulting processed and filtered reads 1 and 2 were
10 aligned separately using bowtie2/2.2.3 to a Burrows-Wheeler Index of the concatenated
11 mouse (mm10) and human (hg19) genomes. Individual SAM files were then converted to
12 BED format and filtered for alignments with MAPQ ≥ 30 using a combination of
13 samtools, bedtools, and awk. Using bedtools closest along with a BED file of all DpnII
14 sites in both genomes (generated using HiC-Pro²²), the closest DpnII site to each read
15 was determined, after which BED files were concatenated, sorted on read ID using UNIX
16 sort, and then processed using a custom Python script to generate a BEDPE format file
17 where 5' mates always precede 3' mates, and where a simple Python dictionary is used to
18 associate barcode combinations contained in the “associations” file with each pair of
19 reads. Reads were then sorted by barcode, read 1 chromosome, start, end, read 2
20 chromosome, start, and end using UNIX sort, and deduplicated using a custom Python
21 script on the following criteria: reads were considered to be PCR duplicates if they were
22 associated with the same cellular index, and if they comprised a ligation between the
23 same two restriction sites as defined using bedtools closest.

1

2 *Cellular Demultiplexing & Quality Analysis*

3 When demultiplexing cells, we run two custom Python scripts. First, we generate a
4 “percentages” file that includes the species purity of each cellular index, the coverage of
5 each index, and the number of times a particular restriction fragment is observed once,
6 twice, thrice, and four times. We also include the *cis:trans* ratio described above, and, if
7 applicable, the fraction of homozygous alternate HeLa alleles observed. We use these
8 percentages files to filter BEDPE files (see below) and generate, at any desired
9 resolution, single cell matrices in long format (*i.e.* BIN1-BIN2-COUNT), with only the
10 “upper diagonal” of the matrix included to reduce storage footprint. These matrices are
11 then converted to numpy matrices for visualization and further analysis.

12

13 *Filtration of Cellular Indices*

14 We applied several filters to our resulting cellular indices to arrive at the cells analyzed in
15 this study. We first removed all cellular indices with fewer than 1000 unique reads. We
16 next filtered out all indices where the *cis:trans* ratio was lower than 1. Finally, for all
17 experiments we removed cellular indices where less than 95% of reads aligned uniquely
18 to either the mouse (mm10) or human (hg19) genomes. For all human cells from HAP1
19 and HeLa S3 mixing experiments (Libraries ML1, ML2, PL1, and PL2) further filtration
20 by genotype was performed. For each cellular index, we examined all reads overlapping
21 with known alternate homozygous sites in the HeLa S3 genome and computed the
22 fraction of sites where the alternate allele is observed. We then drew cutoffs to filter out
23 all cells where this fraction fell between 56% and 99%.

1

2 We do acknowledge that particular applications (*e.g.* structural modeling) may require
3 more stringent filtration for cellular indices covering single cells. As such, we provide
4 with the raw data a supplementary file specifying the “species purity” of each barcode
5 combination in each sequenced library, along with the number of times DpnII restriction
6 fragments are observed in a cell once, twice, thrice, or four times, with the expectation
7 that given some tolerable noise level, one should only observe restriction fragment copy
8 numbers equal to or less than the copy number of that fragment for that cell type.
9 Relatedly, we note that further inspection of the HAP1 cells used in this study revealed
10 that they were not entirely haploid. HAP1 cells, an engineered haploid line, have faster
11 doubling times compared to HeLa S3, and have been described as having a relatively
12 large frequency of diploid cells²³. FACS analysis (data not shown) of the stock used for
13 these experiments showed that ~40% of cells analyzed harbored $2n$ nucleic acid content,
14 indicating haploid cells in G2 or reverted diploid cells in G1.

15

16 **Data Analysis**

17 *PCA of Combinatorial Single-Cell Hi-C Data*

18 Single-cell matrices at interchromosomal contact resolution (\log_{10} of contact counts) and
19 10 Mb resolution (binarized; 0 if absent, 1 if present) were vectorized and concatenated
20 using custom Python scripts. Concatenation was performed such that redundant entries of
21 each contact matrix (*i.e.* C_{ij} and C_{ji}) were only represented once. Resulting matrices,
22 where rows represent single-cells and columns represent observed contacts, were then
23 decomposed using the PCA function in scikit-learn. For interchromosomal matrices,

1 entries for intrachromosomal contacts (*i.e.* the diagonal) were set to 0. For 10 Mb
2 intrachromosomal matrices, all interchromosomal contacts were ignored and all entries
3 C_{ij} where $|i - j| < 3$ were set to zero.

4

5 *Calculation of Contact Probabilities in Single Cells*

6 Methods to calculate the scaling probability within single cells were adapted from
7 Imakaev, Fudenberg *et al*¹⁹ and Sanborn, Rao *et al*²⁰. A histogram of contact distances
8 normalized by bin size was generated using logarithmically increasing bins (increasing
9 by powers of 1.12ⁿ). We obtained the scaling coefficient by calculating the line of best fit
10 for the log-log plot of this histogram between distances of 50 kb and 8 Mb. Shuffled
11 controls were generated by randomly reassigning all cellular indices and repeating the
12 above analysis; this importantly maintains the coverage distribution of the new set of
13 simulated “single cells.”

14

15 All plots were generated in R using ggplot2 (<http://ggplot2.org/>).

16

17 **References**

- 18 1. Ramani, V., Shendure, J. & Duan, Z. Understanding Spatial Genome
19 Organization: Methods and Insights. *Genomics Proteomics Bioinformatics* **14**, 7–
20 20 (2016).
- 21 2. Cremer, T. & Cremer, C. Chromosome territories, nuclear architecture and gene
22 regulation in mammalian cells. *Nat Rev Genet* **2**, 292–301 (2001).
- 23 3. van Steensel, B. & Dekker, J. Genomics tools for unraveling chromosome
24 architecture. *Nat Biotechnol* **28**, 1089–1095 (2010).
- 25 4. Soderberg, O. *et al.* Direct observation of individual endogenous protein
26 complexes in situ by proximity ligation. **3**, 995–1000 (2006).
- 27 5. Ramani, V., Qiu, R. & Shendure, J. High-throughput determination of RNA
28 structure by proximity ligation. *Nat Biotechnol* (2015). doi:10.1038/nbt.3289
- 29 6. Dekker, J., Rippe, K., Dekker, M. & Kleckner, N. Capturing chromosome

- 1 conformation. *Science* **295**, 1306–1311 (2002).
- 2 7. Lieberman-Aiden, E. *et al.* Comprehensive Mapping of Long-Range Interactions
3 Reveals Folding Principles of the Human Genome. *Science* **326**, 289–293 (2009).
- 4 8. Duan, Z. *et al.* A three-dimensional model of the yeast genome. *Nature* **465**, 363–
5 367 (2010).
- 6 9. Nagano, T. *et al.* Single-cell Hi-C reveals cell-to-cell variability in chromosome
7 structure. *Nature* **502**, 59–64 (2013).
- 8 10. Cusanovich, D. A. *et al.* Epigenetics. Multiplex single-cell profiling of chromatin
9 accessibility by combinatorial cellular indexing. *Science* **348**, 910–914 (2015).
- 10 11. Klein, A. M. *et al.* Droplet barcoding for single-cell transcriptomics applied to
11 embryonic stem cells. *Cell* **161**, 1187–1201 (2015).
- 12 12. Macosko, E. Z. *et al.* Highly Parallel Genome-wide Expression Profiling of
13 Individual Cells Using Nanoliter Droplets. *Cell* **161**, 1202–1214 (2015).
- 14 13. Rotem, A. *et al.* Single-cell ChIP-seq reveals cell subpopulations defined by
15 chromatin state. *Nat Biotechnol* **33**, 1165–1172 (2015).
- 16 14. Rao, S. S. P. *et al.* A 3D map of the human genome at kilobase resolution reveals
17 principles of chromatin looping. *Cell* **159**, 1665–1680 (2014).
- 18 15. Deng, X. *et al.* Bipartite structure of the inactive mouse X chromosome. *Genome*
19 *Biol.* **16**, 152 (2015).
- 20 16. Adey, A. *et al.* The haplotype-resolved genome and epigenome of the aneuploid
21 HeLa cancer cell line. *Nature* **500**, 207–211 (2013).
- 22 17. Essletzbichler, P. *et al.* Megabase-scale deletion using CRISPR/Cas9 to generate a
23 fully haploid human cell line. *Genome Research* **24**, 2059–2065 (2014).
- 24 18. Naumova, N. *et al.* Organization of the mitotic chromosome. *Science* **342**, 948–
25 953 (2013).
- 26 19. Imakaev, M. *et al.* Iterative correction of Hi-C data reveals hallmarks of
27 chromosome organization. *Nat Meth* **9**, 999–1003 (2012).
- 28 20. Sanborn, A. L. *et al.* Chromatin extrusion explains key features of loop and
29 domain formation in wild-type and engineered genomes. *Proc. Natl. Acad. Sci.*
30 *U.S.A.* **112**, E6456–65 (2015).
- 31 21. Jin, W. *et al.* Genome-wide detection of DNase I hypersensitive sites in single
32 cells and FFPE tissue samples. *Nature* **528**, 142–146 (2015).
- 33 22. Servant, N. *et al.* HiC-Pro: an optimized and flexible pipeline for Hi-C data
34 processing. *Genome Biol.* **16**, 259 (2015).
- 35 23. Carette, J. E. *et al.* Ebola virus entry requires the cholesterol transporter Niemann-
36 Pick C1. *Nature* **477**, 340–343 (2011).
- 37
38

39 **Acknowledgements**

40 The authors thank S. Kasinathan, members of the UW Center for Nuclear Organization
41 and Function, and members of the Shendure lab (particularly M. Kircher), for helpful
42 discussions. HeLa S3 cells were used as part of this study. Henrietta Lacks, and the HeLa

1 cell line that was established from her tumor cells in 1951, have made significant
2 contributions to scientific progress and advances in human health. We are grateful to
3 Henrietta Lacks, now deceased, and to her surviving family members for their
4 contributions to biomedical research. Primary MEF aliquots were a gift from Carol Ware.
5 This work was funded by grants from the NIH (5T32HG000035 to VR, DP1HG007811
6 to JS and U54DK107979 to XD, CMD, WSN, ZD and JS). JS is an Investigator of the
7 Howard Hughes Medical Institute.

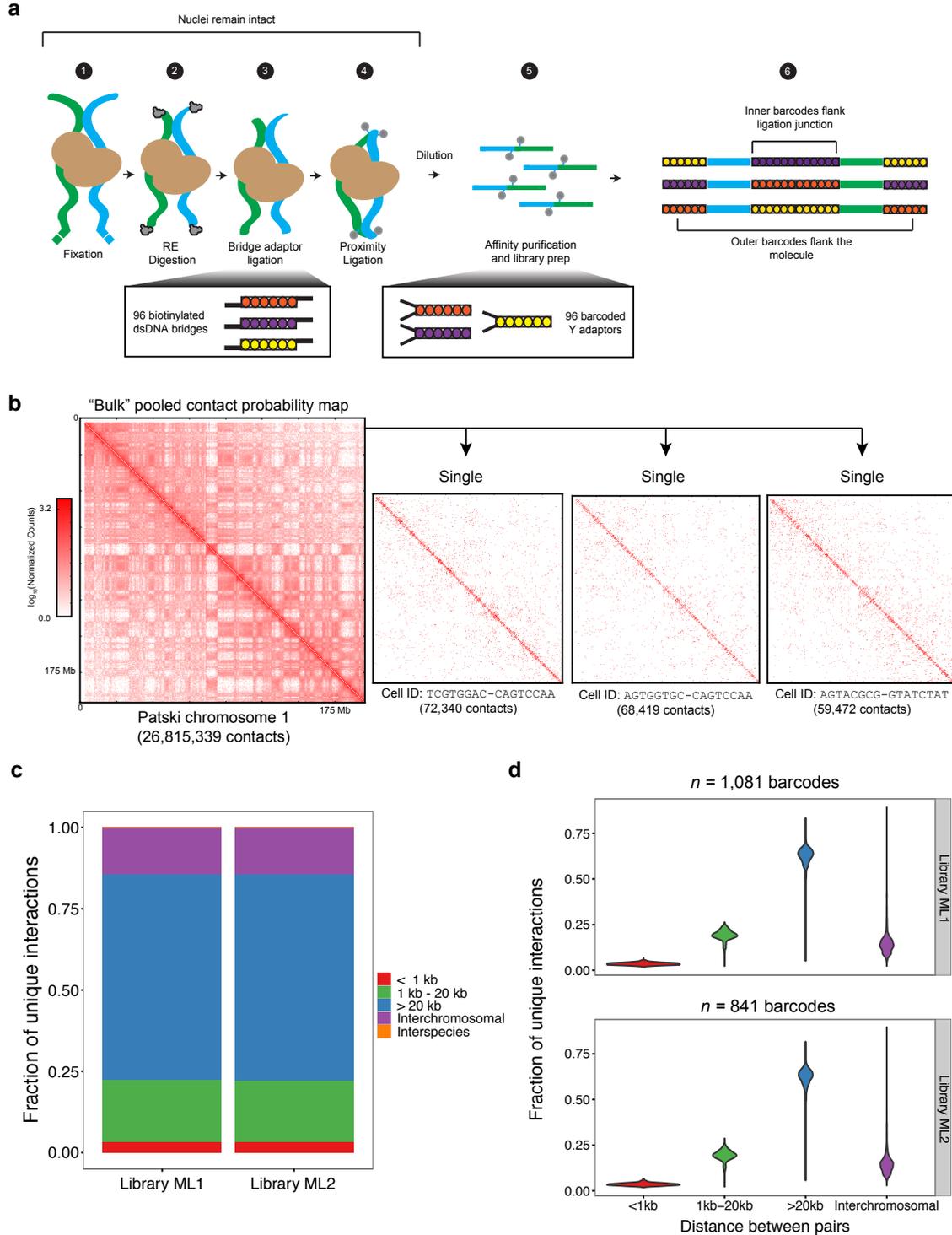
8

9 **Competing Financial Interests**

10 K.G. and F.S. are employees of Illumina Inc.

11

1 Figures

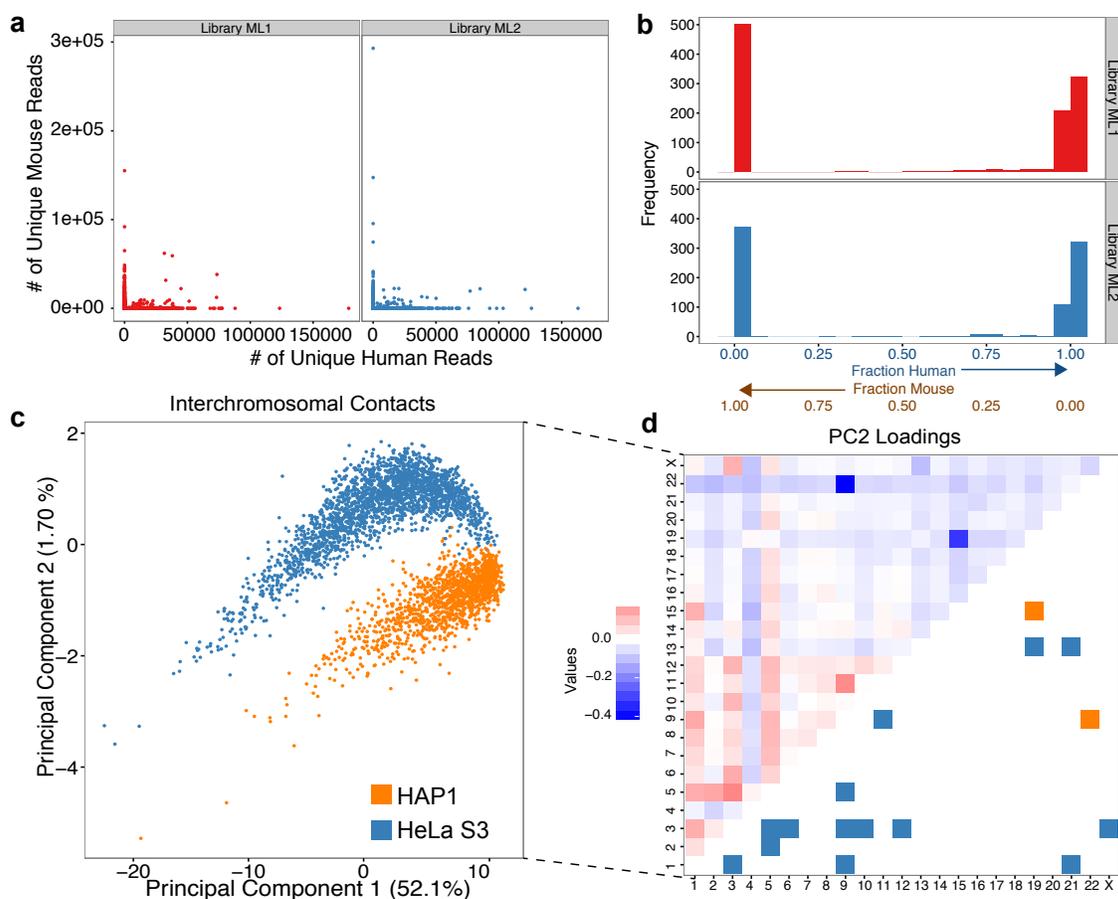


2

3 **Figure 1: Combinatorial single cell Hi-C integrates the *in situ* Hi-C protocol with**

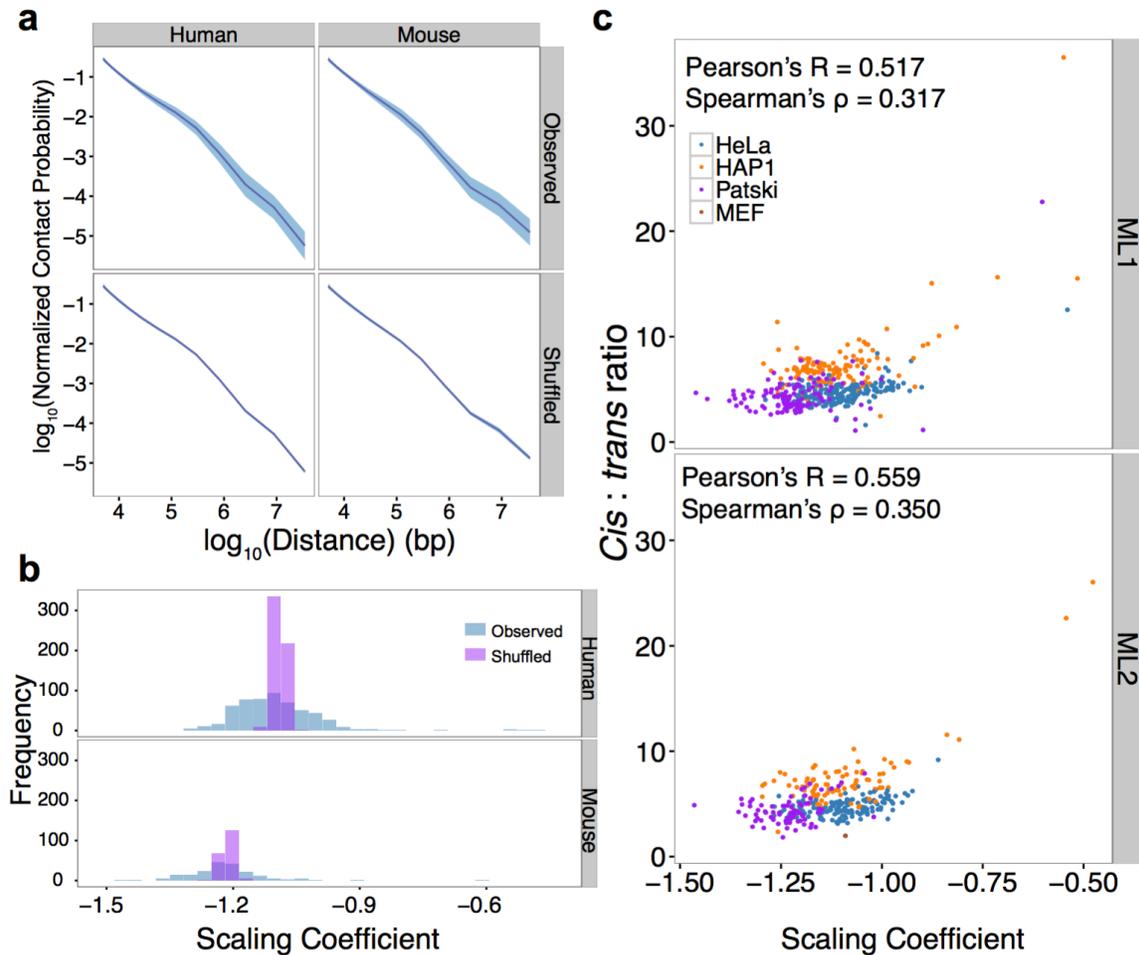
4 **combinatorial cellular indexing to generate signal-rich bulk Hi-C maps that can be**

1 **decomposed into single cell Hi-C maps.** a.) Combinatorial single cell Hi-C follows the
2 traditional paradigm of fixation, digestion, and re-ligation shared by all Hi-C assays
3 (Steps 1 – 4), but importantly uses a biotinylated bridge adaptor to incorporate a first
4 round of barcodes prior to proximity ligation (Step 3), and custom barcoded Illumina Y-
5 adaptors (Step 5) to incorporate a second round of barcodes prior to affinity purification
6 and library amplification (Steps 5 – 6). b.) Bulk data generated by this protocol can be
7 decomposed to single cell Hi-C maps. Shown are combinatorial single cell Hi-C reads for
8 mouse chromosome 1, corresponding to data from the Patski cell line. This data can then
9 be separated on the basis of cellular indices into many single cell contact maps, shown
10 here for three single Patski cells. c.) Our Hi-C libraries demonstrate a high *cis:trans* ratio,
11 measured as the ratio of intrachromosomal contacts > 20 kb apart to interchromosomal
12 contacts. d.) The high *cis:trans* ratio observed in bulk data is maintained after libraries are
13 decomposed to ~1800 single cell Hi-C maps. All indices tagging fewer than 1,000 reads
14 are ignored in this analysis.



1
2 **Figure 2: The large number of cellular indices generated through combinatorial**
3 **single cell Hi-C are overwhelmingly species-specific, and can be separated by cell**
4 **type.** a.) Species mixing experiments enable explicit definition of the “collision rate” in a
5 combinatorial single cell Hi-C experiment. In libraries ML1 and ML2, similar levels of
6 collision (defined as any cellular index with at least 1,000 unique reads, but <95% species
7 purity) are observed, and fall within the expected range. In these plots, points close to the
8 *x* and *y* axes likely represent single cells. b.) Species contamination visualized as a
9 histogram of the fraction of reads mapping to the human genome. Cellular indices (here
10 filtering out all indices with fewer than 1,000 unique reads) are largely species specific.
11 c.) Projection onto the first two principal components from PCA analysis of
12 interchromosomal contact matrices results in separation of HeLa S3 and HAP1, two

1 karyotypically different cell lines ($n = 3,609$ cells). Percentages shown are the percentage
2 of variance explained by each plotted PC. d.) Principal component 2 loadings represent
3 the contribution of each feature (interchromosomal contact) to the observed cell type
4 separation. Strongly blue features appear to be HAP1 specific, while red features appear
5 to be HeLa S3 specific. Known translocations for each cell type are mirrored against the
6 loading heatmap.



1

2 **Figure 3: Combinatorial single cell Hi-C captures cell-to-cell heterogeneity masked**

3 **by bulk measurement, a.)** Decay in contact probability for all primary experiment (ML

4 libraries) cells with at least 10,000 unique contacts ($n = 769$ cells). Plotted is the mean

5 contact probability for each bin (purple), along with standard deviation (blue). Shuffled

6 controls where all cellular index assignments have been randomized demonstrate

7 strikingly lower variance compared to observed single cells, for both mouse and human.

8 b) Scaling coefficients calculated for a.), for distances between 50 kb and 8 Mb. Shuffled

9 controls demonstrate a tighter distribution of coefficients compared to the observed single

10 human cells. c.) Single-cell scaling coefficients reproducibly demonstrate positive

11 correlation with single-cell *cis:trans* ratios in both mouse and human cells.

1

	ML1	ML2	PL1	PL2	ML3
Programming During First Barcoding / Cell Type Composition	HeLa S3 + Patski mixed in 4 rows; HAP1 + MEF mixed in 4 rows	HeLa S3 + Patski mixed in 4 rows; HAP1 + MEF mixed in 4 rows	HeLa S3 in 2 rows; Patski in 2 rows; HAP1 in 2 rows; MEF in 2 rows	HeLa S3 in 2 rows; Patski in 2 rows; HAP1 in 2 rows; MEF in 2 rows	GM12878 + Patski mixed in 4 rows; K562 + MEF mixed in 4 rows
# of sequenced reads	128,660,854	124,833,616	107,609,682	181,089,598	20,302,425
# of associated reads	62,566,086	59,931,225	59,054,337	94,661,167	10,192,548
# of mapped, associated, MAPQ >= 30	27,589,249	27,012,845	28,295,707	43,763,077	5,161,279
# of unique mapped, associated, MAPQ >= 30	14,997,809	11,799,762	21,980,282	37,360,138	4,622,296
# of cellular indices	1,081	841	3,184	3,743	1,291
Median read count per cellular index	9,274	8,335	4,137	6,270	2,146
# of cells, post-filtering	975	766	2,900	3,500	1,175
Median read count per filtered cell	9,186	8,390	4,249	6,479	2,138
Median <i>cis:trans</i> Ratio Across Cells	4.43	4.34	3.69	3.42	3.96

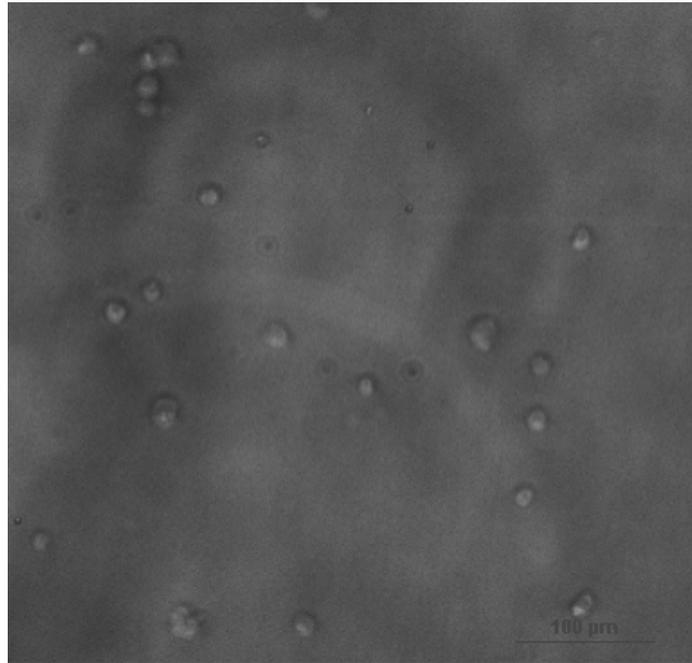
2

3

Table 1: Summary of all 5 sequenced libraries in this study.

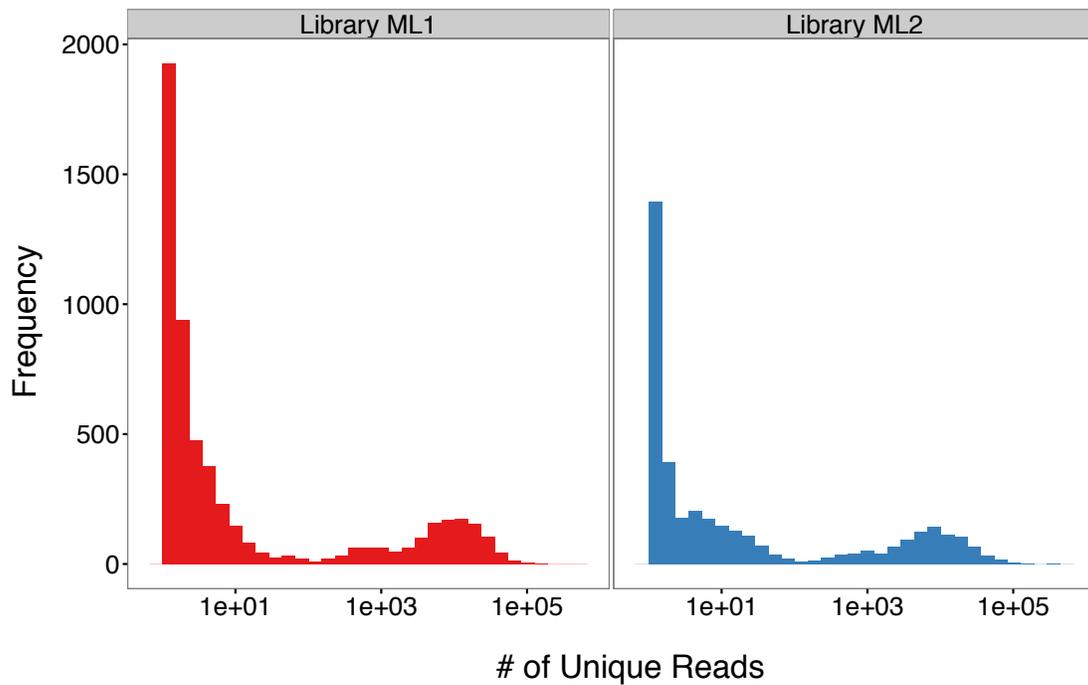
1 **Supplementary Figures**

2



3

4 **Supplementary Figure 1: Nuclei remain intact through proximity ligation in the**
5 **combinatorial single cell Hi-C protocol.** Phase contrast microscopy of HeLa S3 and
6 HAP1 nuclei following proximity ligation and serial dilution shows that nuclei remain
7 intact throughout the combinatorial single cell Hi-C protocol (scale bar = 100 μm).



1

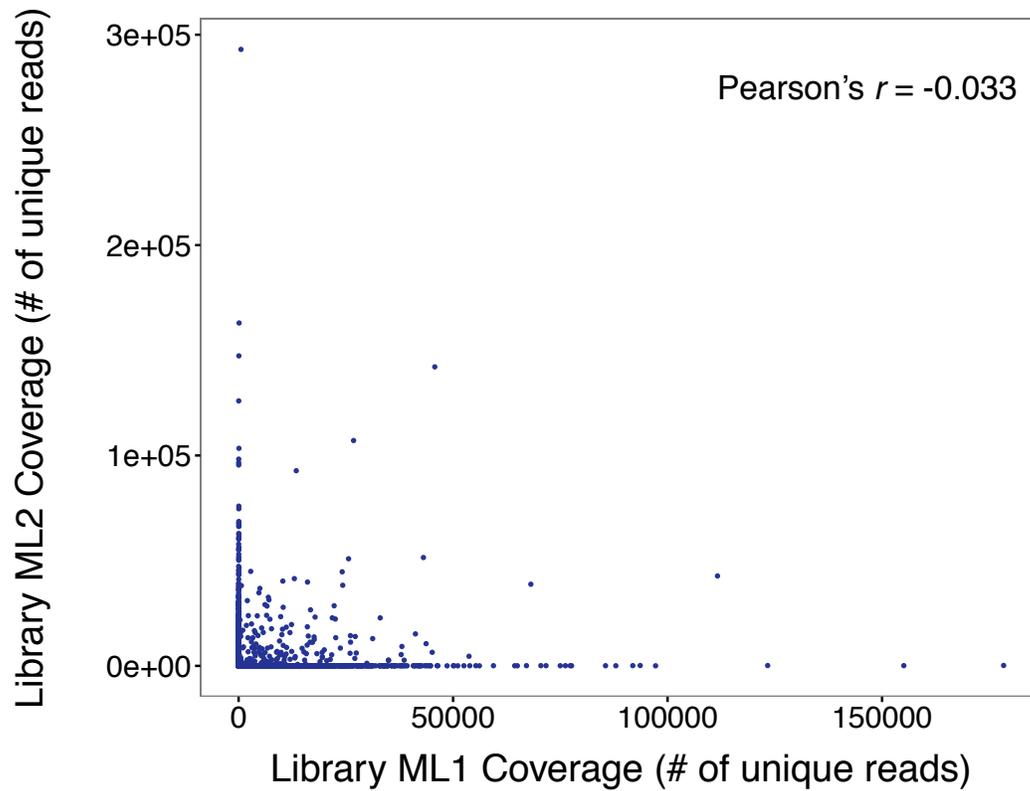
2 **Supplementary Figure 2: Coverage of combinatorial single cell Hi-C cellular indices**

3 **follow a bimodal distribution.** Examining a histogram of the coverage (*i.e.* # of unique

4 reads) of combinatorial single cell Hi-C cellular indices in two replicate libraries reveals

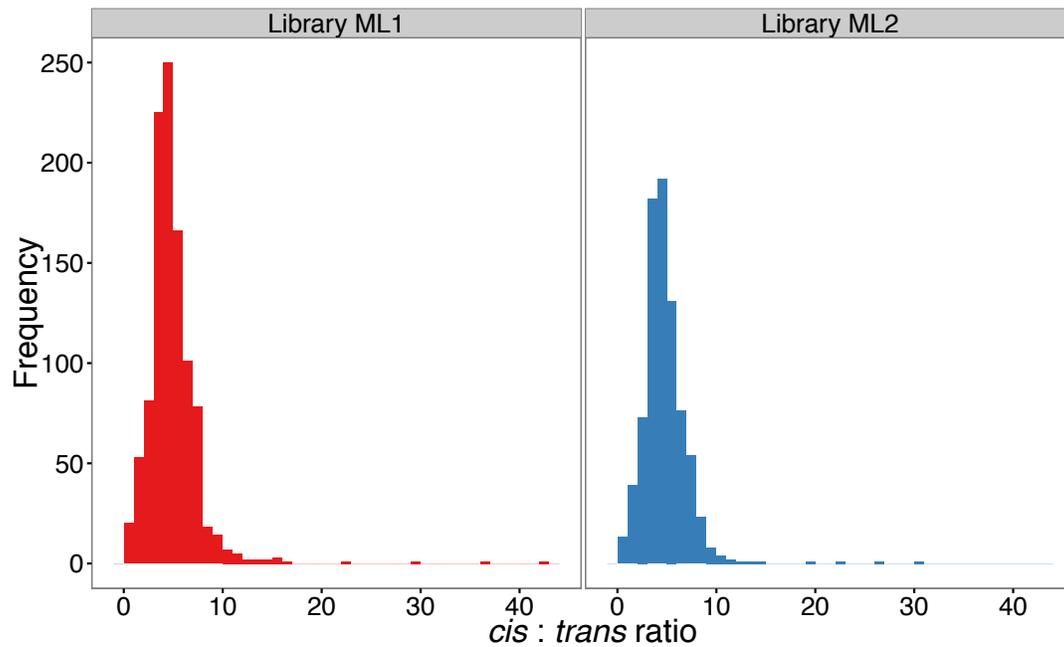
5 a bimodal distribution, where low coverage cellular indices likely represent barcoding of

6 free DNA in solution, rather than intact nuclei.



1

2 **Supplementary Figure 3: Coverage of cellular indices is not correlated between**
3 **replicate experiments.** Scatter plot of coverage per cellular index for all cellular indices
4 with at least 1 unique read in both replicate combinatorial single cell Hi-C libraries. A
5 Pearson's r of -0.03 suggests that there is minimal intrinsic bias (*i.e.* "barcode" effect)
6 biasing coverage of particular cellular indices.



1

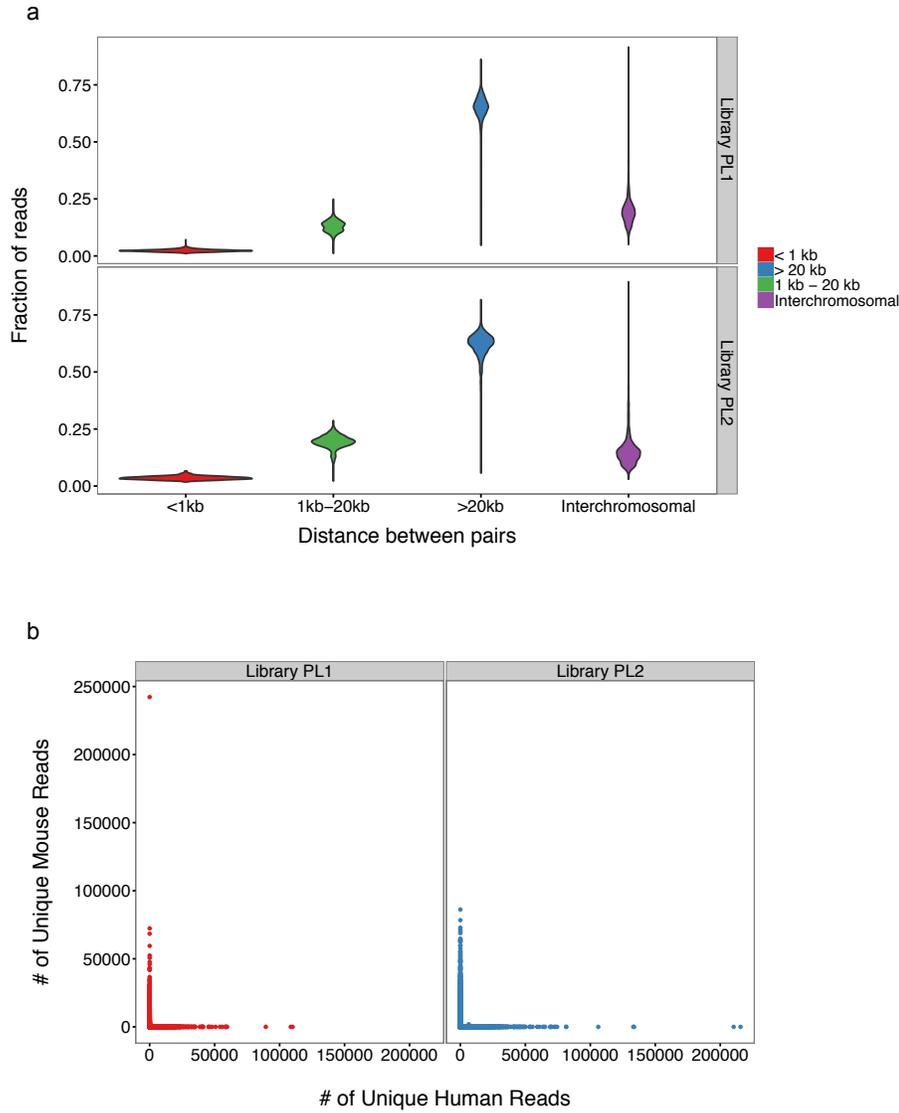
2 **Supplementary Figure 4: Single cellular indices demonstrate high *cis:trans* ratios.**

3 Histogram of the *cis:trans* ratios for cellular indices over two biological replicates. High

4 *cis:trans* ratio suggest that nuclei remain intact during the protocol, and hint at a single-

5 cellular origin for the majority of cellular indices.

6



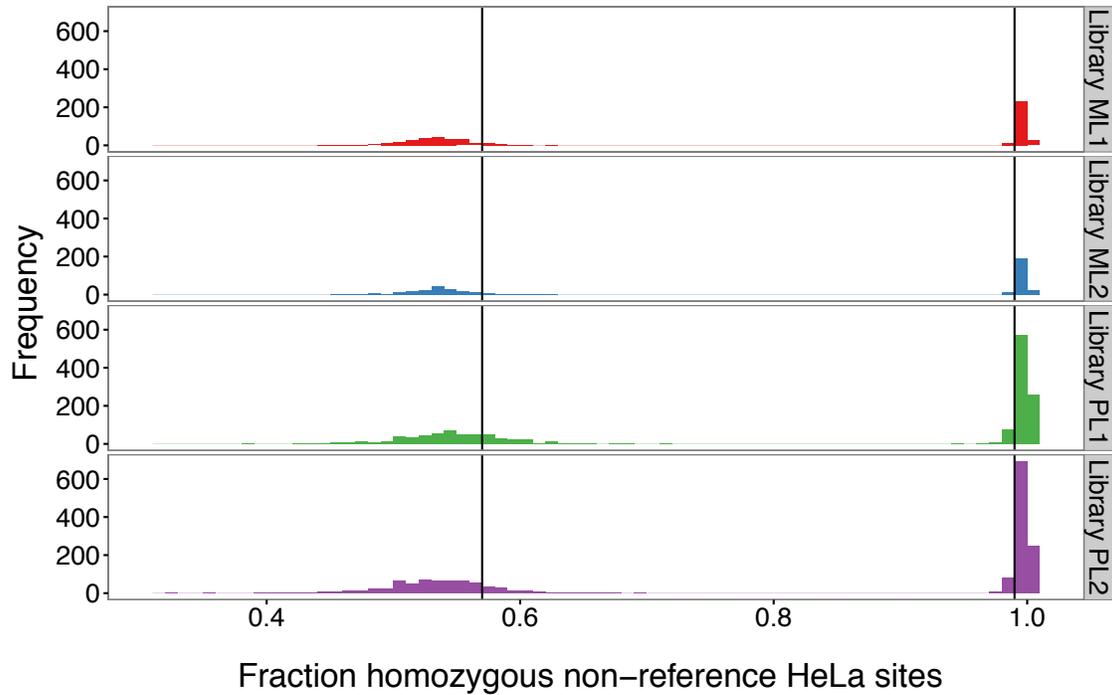
1

2 **Supplementary Figure 5: Quality control statistics for PL1 and PL2 libraries are**

3 **similar to primary experiment libraries.** a.) Violin plots showing the distribution of

4 ligation types across all cellular indices with at least 1,000 reads in libraries PL1 and

5 PL2. b.) Species specificity for both libraries.



1

2 **Supplementary Figure 6: The HeLa genotype enables further filtration of potential**

3 **barcode collisions in combinatorial single cell Hi-C datasets.** We examined all

4 homozygous non-reference sites determined by Adey *et al*¹⁶ and tabulated the fraction of

5 sites where the non-reference allele was found in our sequencing reads, with the

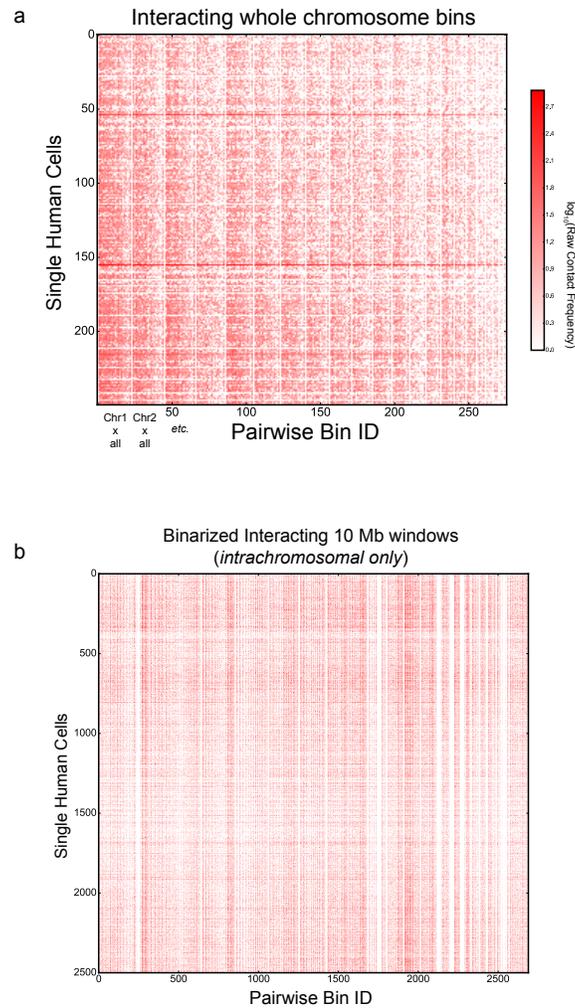
6 expectation that single HeLa cells should have very high (*i.e.* $\geq 99\%$) homozygous non-

7 reference alleles at those sites, with reduced fractions indicating contamination by HAP1.

8 For this study, we drew conservative cutoffs of 57% and 99% for each species (*i.e.* any

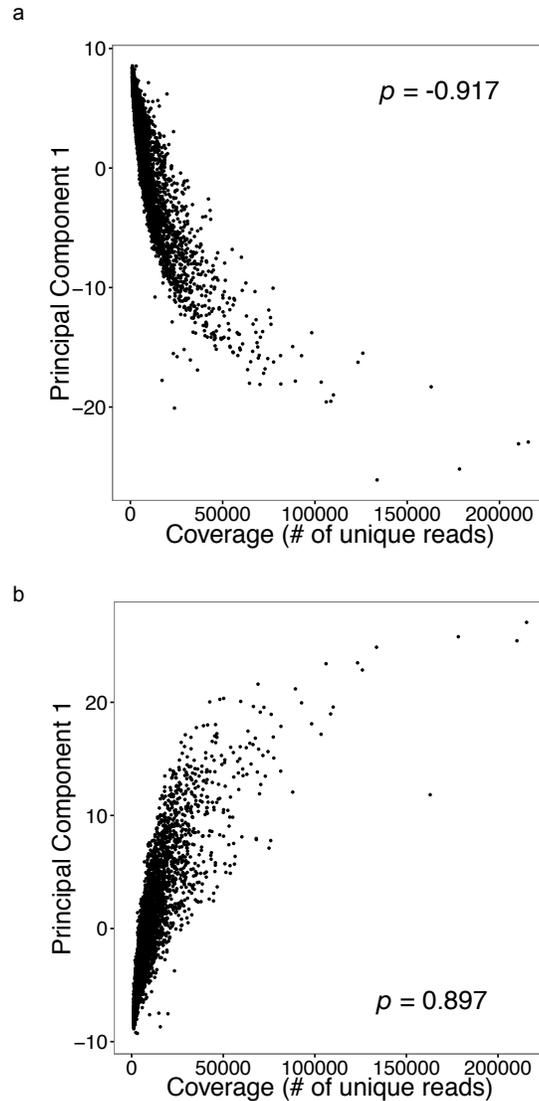
9 cellular indices falling between these values were discarded).

10



1

2 **Supplementary Figure 7: Raw single cell matrices used as input for PCA.** a.) A heat
3 map representation of a portion (250 cells) of the input interchromosomal matrix for
4 PCA. Rows represent single human cells, while columns represent pairwise interactions
5 between two whole chromosomes. For this analysis, raw counts were used, and $n = 3,609$
6 cells. b.) Heat map representation of a portion (2,500 cells) of the input intrachromosomal
7 matrix for PCA. Here, interchromosomal counts were ignored, and interaction
8 frequencies between discrete 10 Mb windows genome-wide were reduced to a binary
9 representation (*i.e.* 1 if present, 0 if absent). Again, $n = 3,609$ cells.



1

2 **Supplementary Figure 8: The first component of PCA using both interchromosomal**

3 **contacts and 10 Mb windowed intrachromosomal contacts strongly correlates with**

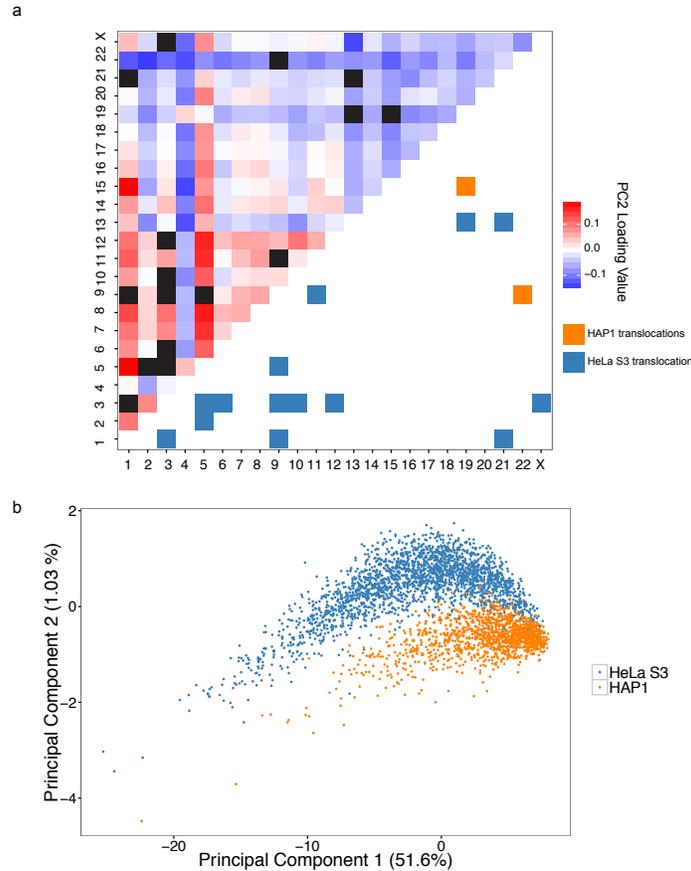
4 **coverage. a.) Correlation between the principal component 1 (PC1) and coverage for**

5 **interchromosomal interactions ($\rho = -0.917$). b.) Correlation between the principal**

6 **component 1 (PC1) and coverage for interacting 10 Mb intrachromosomal windows ($\rho =$**

7 **0.897)**

8



1

2 **Supplementary Figure 9: Analysis of principal component loadings for**

3 **interchromosomal separation experiment reveals that translocations contribute to**

4 **cell type separation in principal component space. a.) Heat map of loadings for**

5 **principal component 2 after all known translocations (blacked out entries) are removed**

6 **from the analysis. b.) After removing all entries corresponding to known translocations**

7 **from the interchromosomal single-cell Hi-C contact matrix, cell-type separation using**

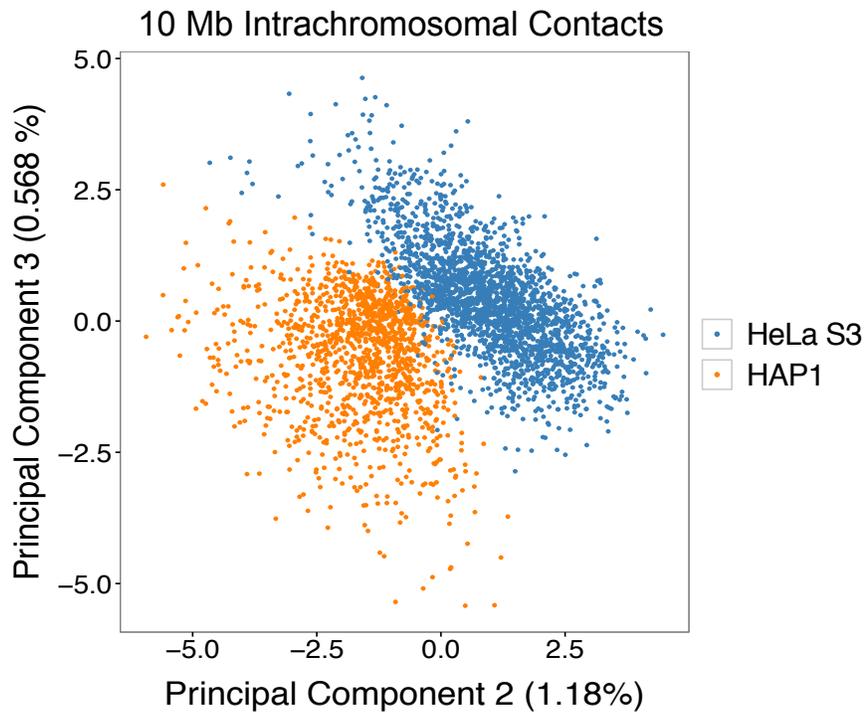
8 **PC1 and PC2 is qualitatively worse but still apparent, suggesting that cell-type specific**

9 **interchromosomal contacts may contribute to the observed separation pattern.**

10 Percentages shown are the percentage of variance explained by each plotted PC.

1

2



3 **Supplementary Figure 10: PCA using an alternate feature set still enables**

4 **separation between HAP1 and K562.** Shown is a projection of principal component 2

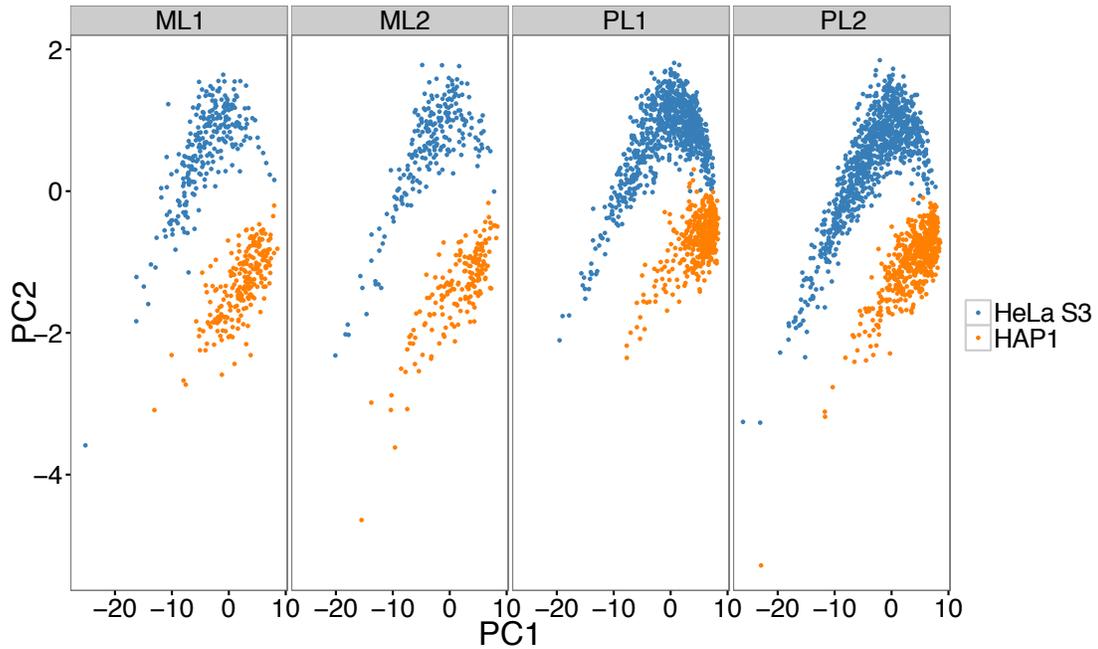
5 and principal component 3 from PCA on the intrachromosomal single cell contact matrix

6 ($n = 3,609$ cells). For this analysis, only intrachromosomal contacts between 10 Mb

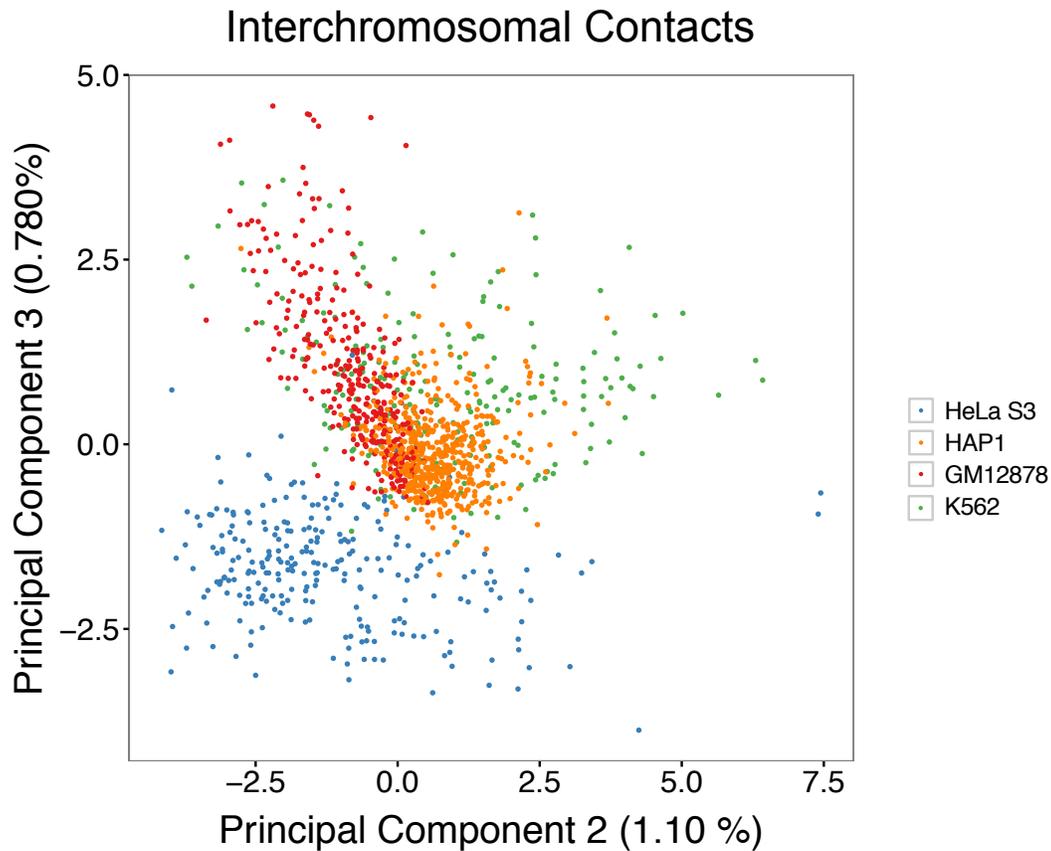
7 windows were used. The matrix used for this computation is shown in **Supplementary**

8 **Figure 7b.** Percentages shown are the percentage of variance explained by each plotted

9 PC.

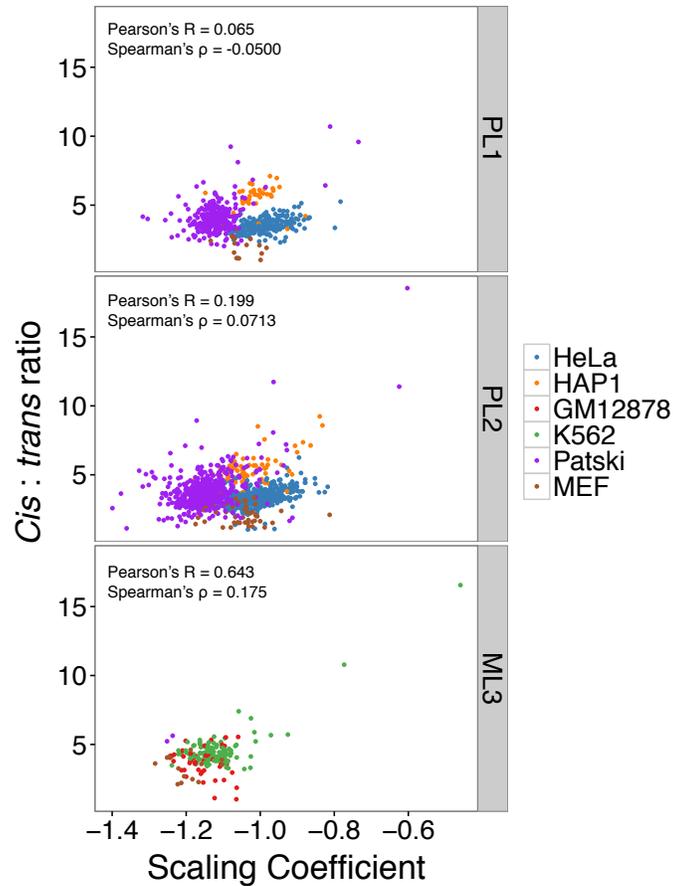


1
2 **Supplementary Figure 11: Separation of cell types by PCA is consistent across**
3 **biological replicate combinatorial single cell Hi-C experiments.** Across 4 different
4 libraries, the separation of single HeLa S3 and HAP1 cells is evident, suggesting that this
5 is not simply a technical artifact or batch effect.



1

2 **Supplementary Figure 12: PCA of single-cell interchromosomal contacts using cells**
3 **from 4 different human cell types results in separation of HeLa S3 from other cell**
4 **lines.** A fifth experiment (Library ML3) containing K562 and GM12878 cells was lightly
5 sequenced and combined with an existing HeLa S3 and HAP1 dataset (Library ML1),
6 resulting in $n = 1,394$ cells. Projection of single cells onto PC2 and PC3 results in
7 separation of HeLa S3 from the remaining three cell types, but weak separation of K562,
8 GM12878, and HAP1. Percentages shown are the percentage of variance explained by
9 each plotted PC.



1

2 **Supplementary Figure 13: Correlation between single cell *cis:trans* ratios and single-**

3 **cell scaling coefficients is reproducible across combinatorial single-cell Hi-C**

4 **experiments.** We observe a correlation between high *cis:trans* ratios and shallow scaling

5 coefficients in both mouse and human cells in both the PL2 (Pearson's R = 0.199;

6 Spearman's ρ = 0.0713) and ML3 (Pearson's R = 0.643; Spearman's ρ = 0.175)

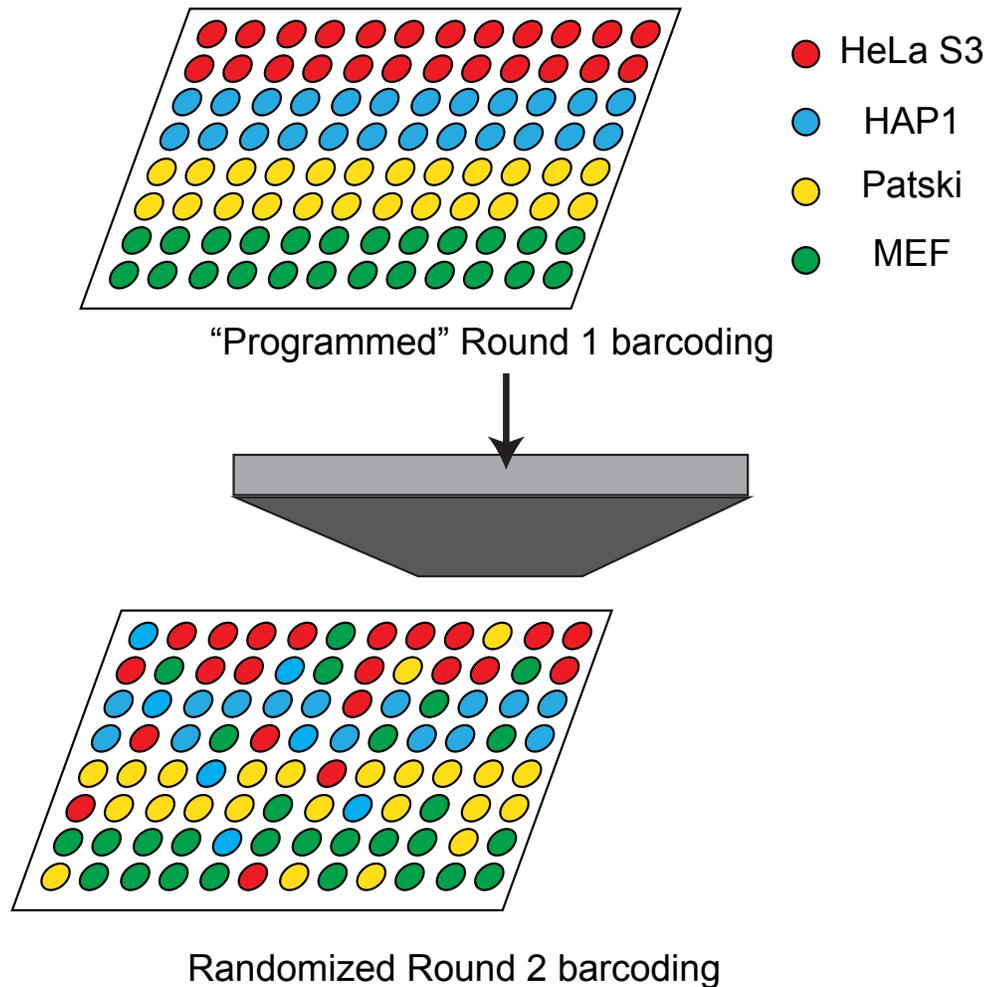
7 experiments. It is possible that the lack of correlation / weaker correlation shown in PL1

8 (Pearson's R = 0.0649; Spearman's ρ = -0.0500) and PL2, respectively, are a result of

9 shallower sequencing, or sampling (*i.e.* perhaps related to the relative abundance of

10 unsynchronized cells in each phase of the cell cycle).

11



1

2 **Supplementary Figure 14: "Programmed" barcoding approaches enable association**
3 **of cell types with unique first round barcodes.** By loading unique cell types into
4 programmed wells during the first round of indexing, we are able to validate cell types *in*
5 *silico*. This schematic shows how libraries PL1 and PL2 were generated, wherein only
6 one cell type was present per cell. By contrast, for ML1, ML2 and ML3, subsets of wells
7 contained mixtures of one human and one mouse cell type.

8