

Figure S1

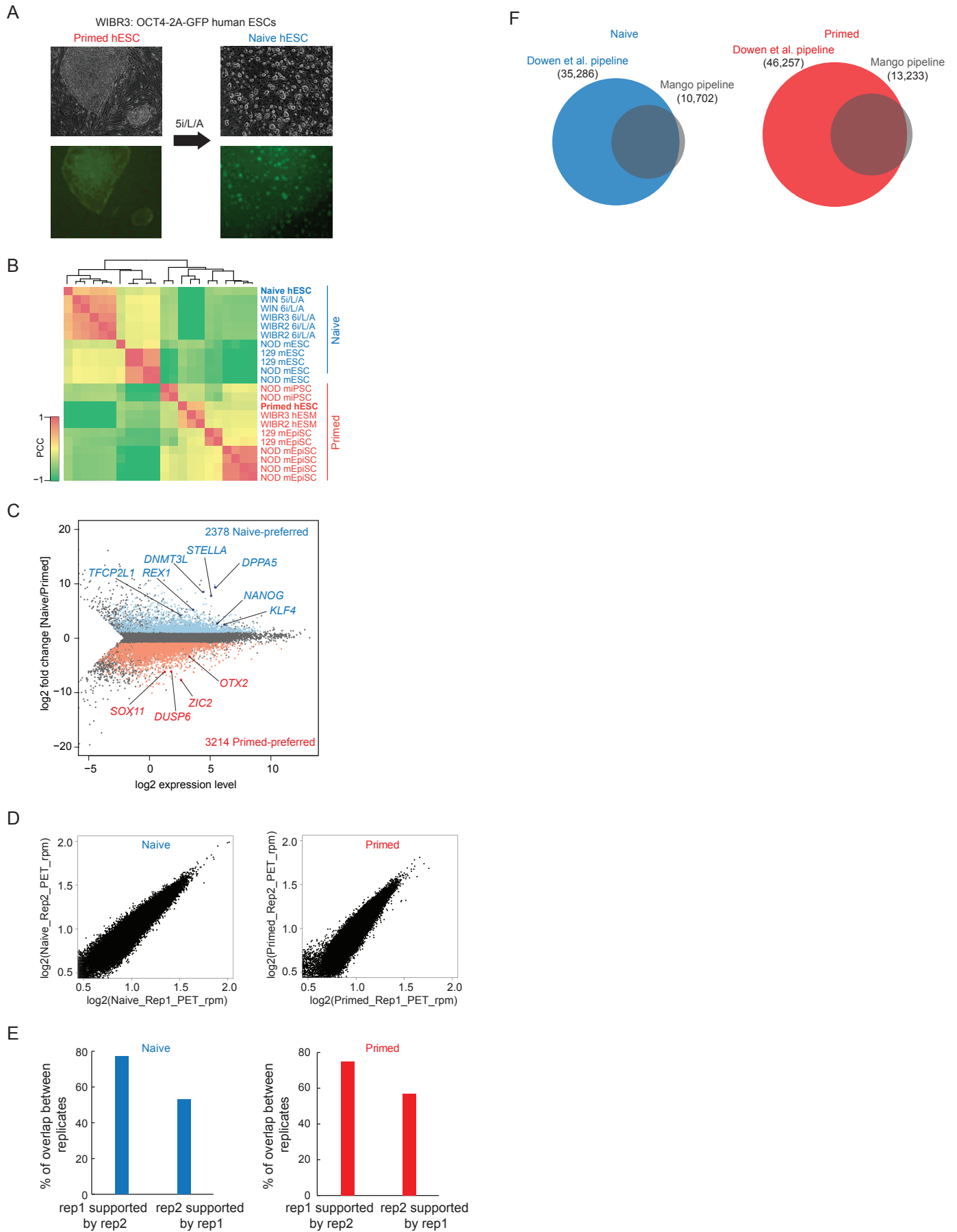


Figure S2

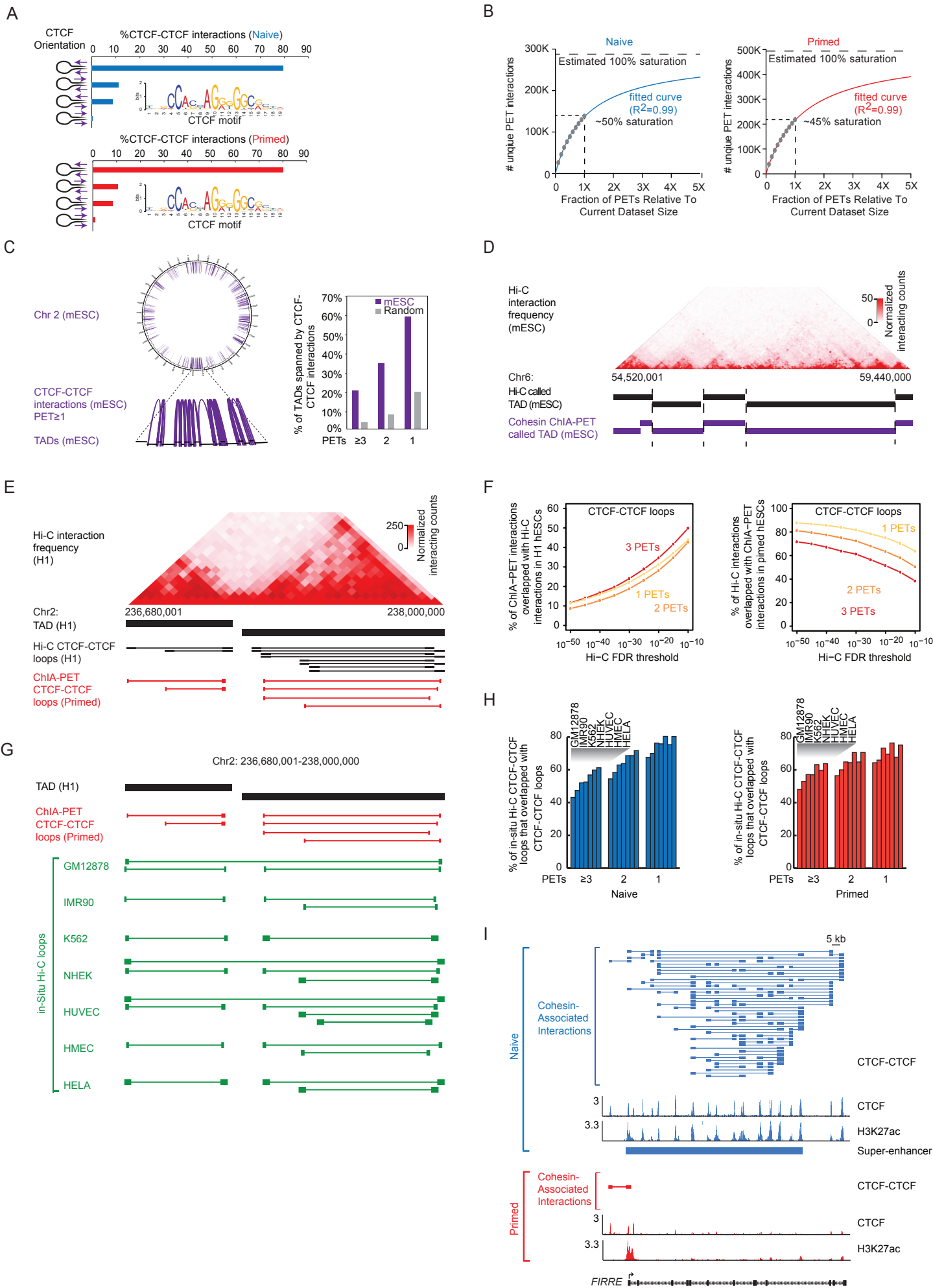


Figure S3

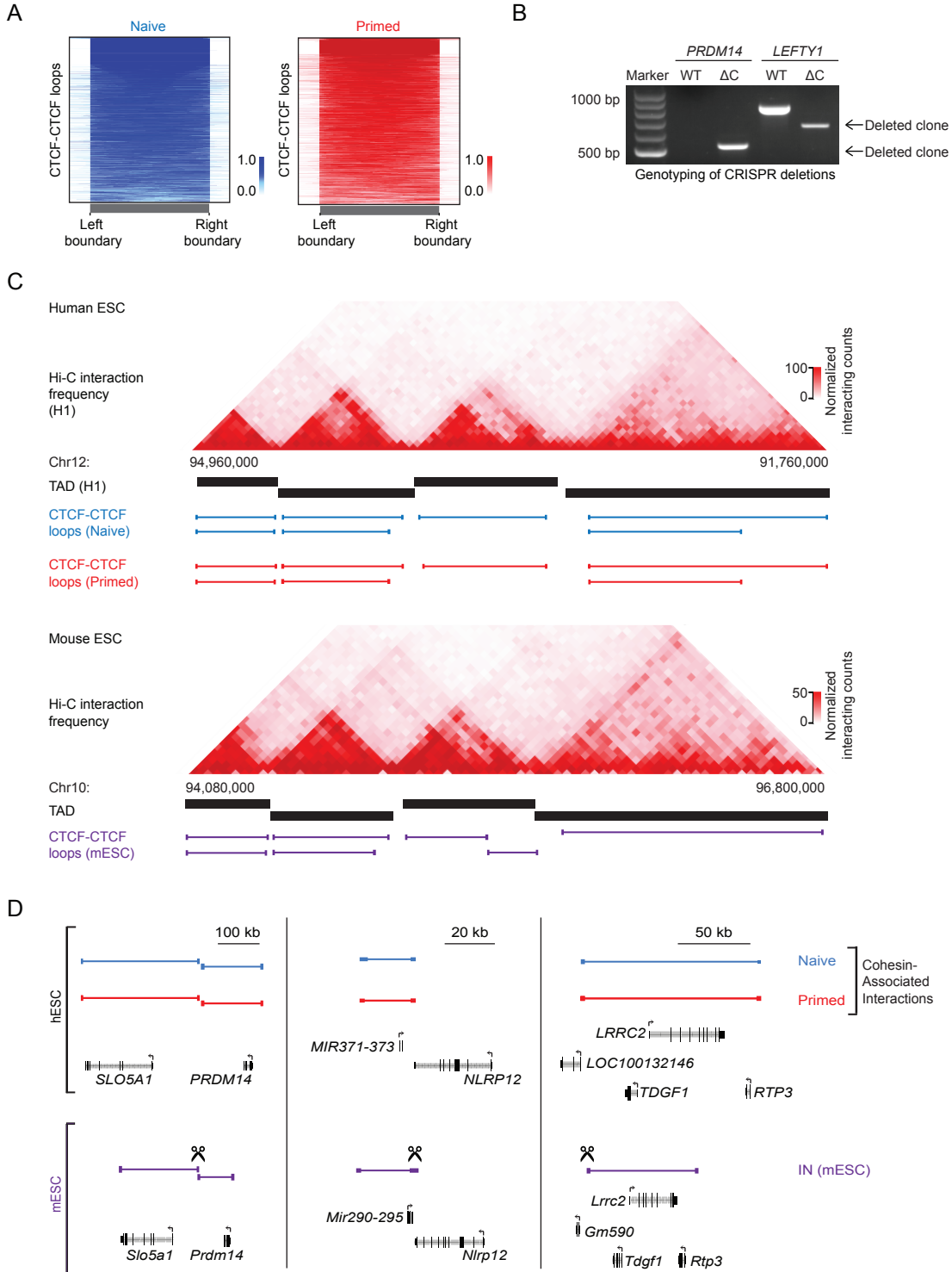


Figure S4

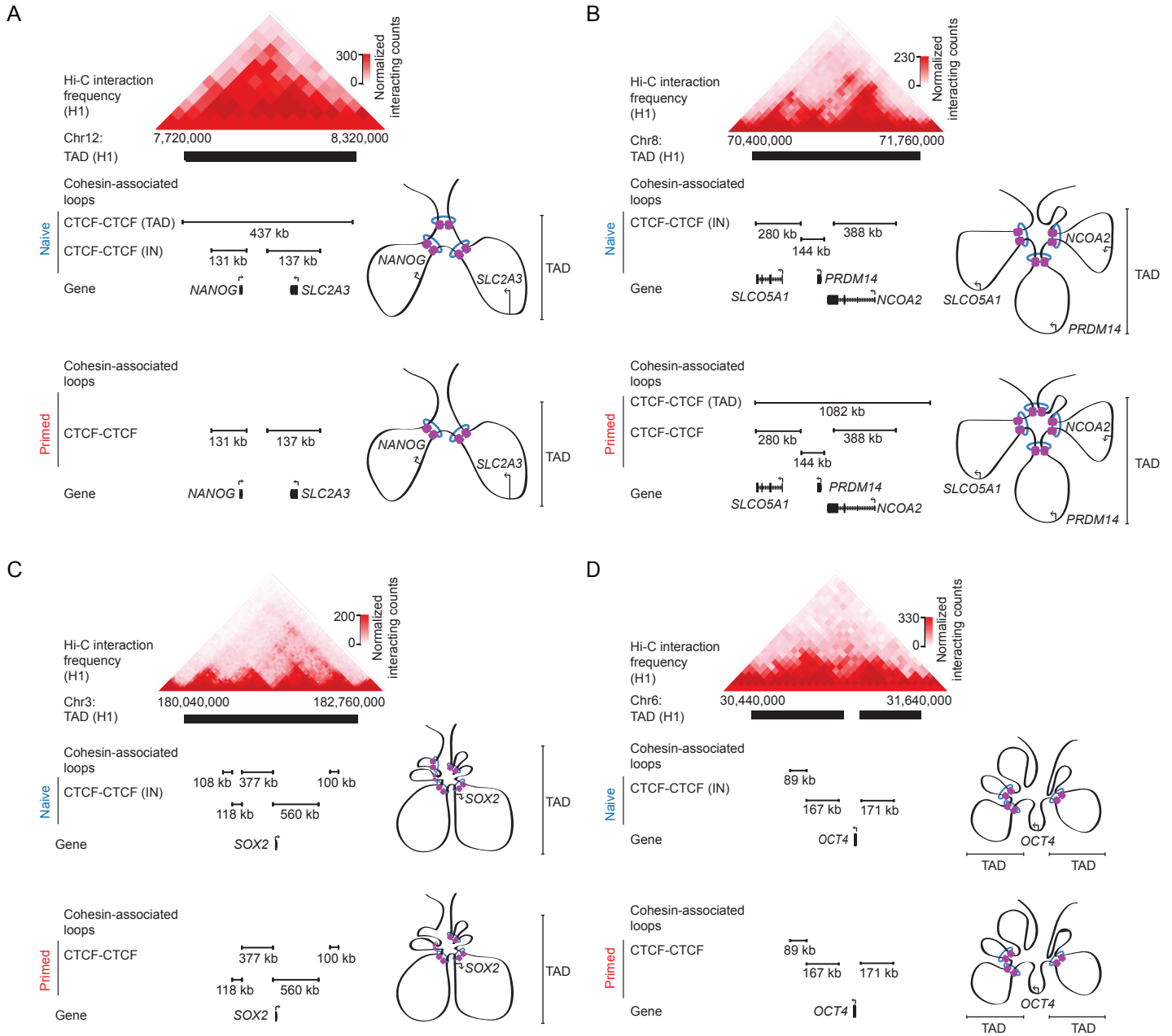
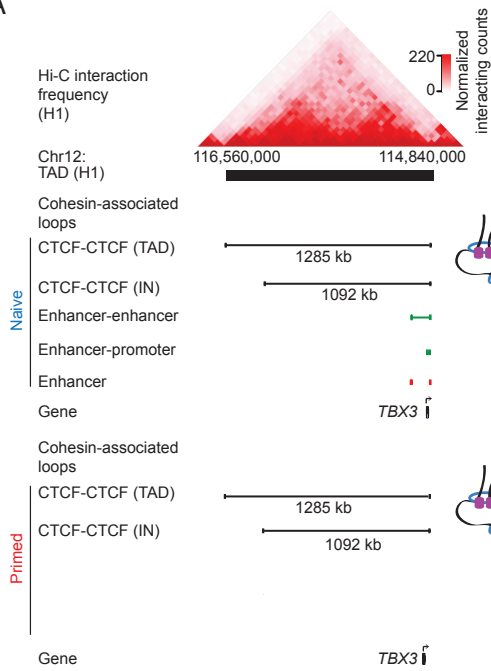
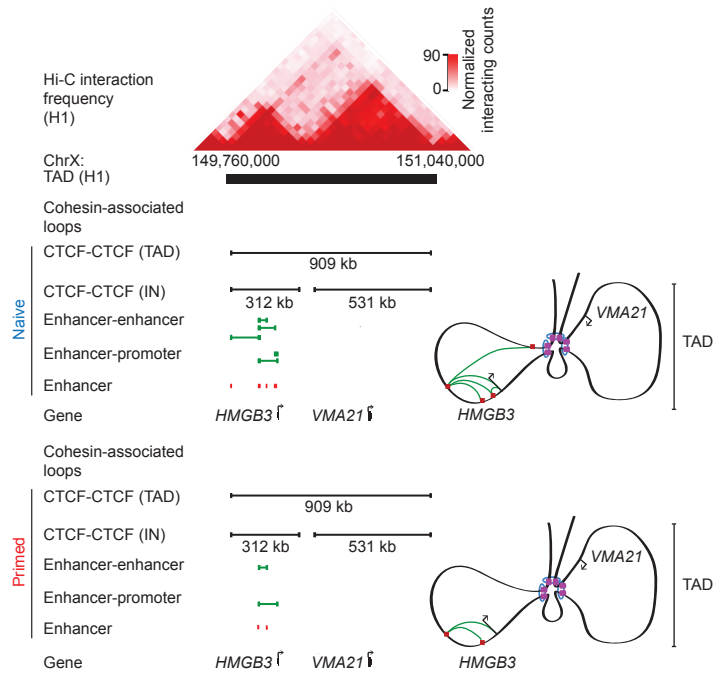


Figure S5

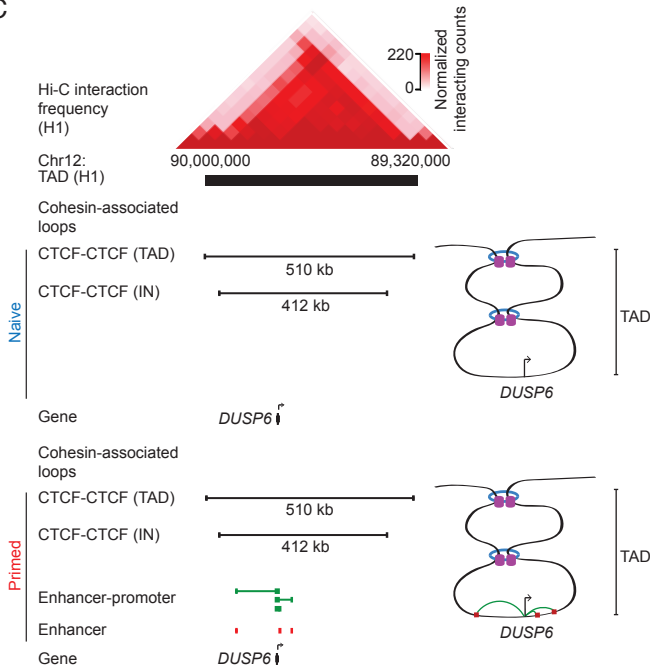
A



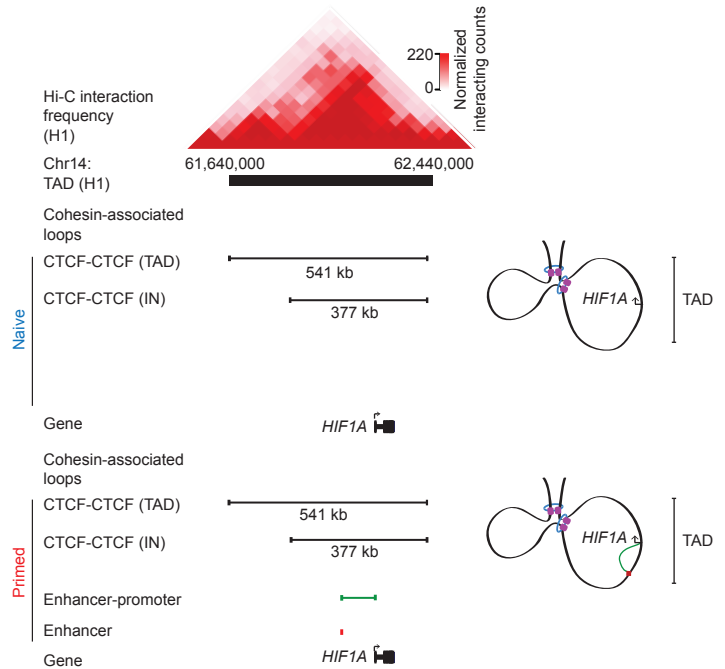
B



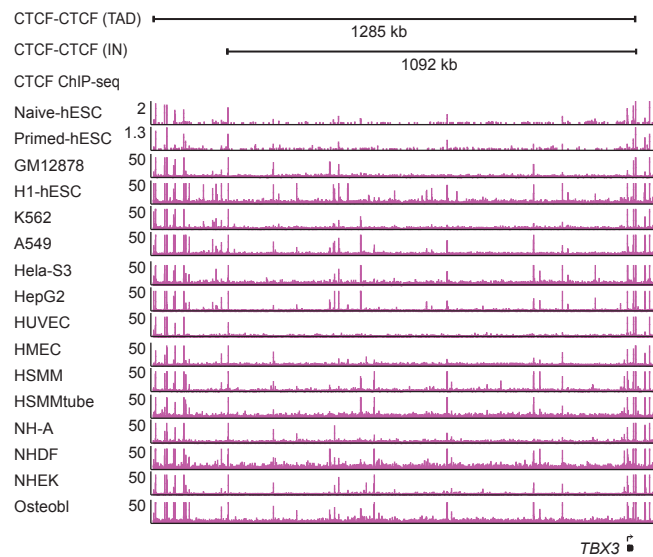
C



D



E



F

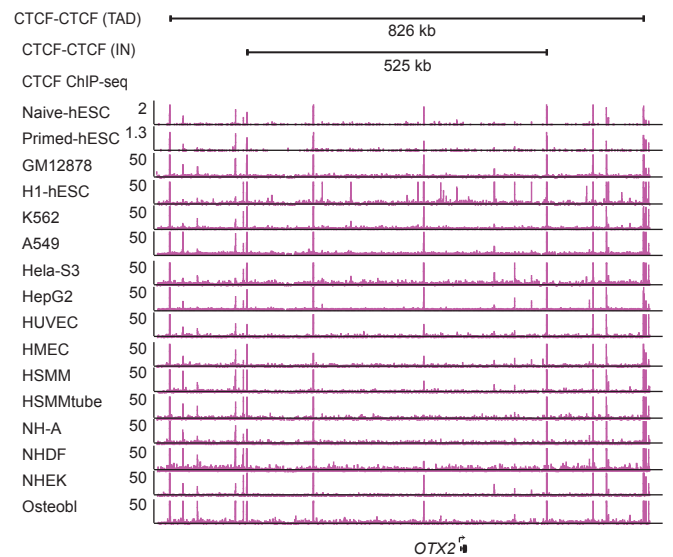
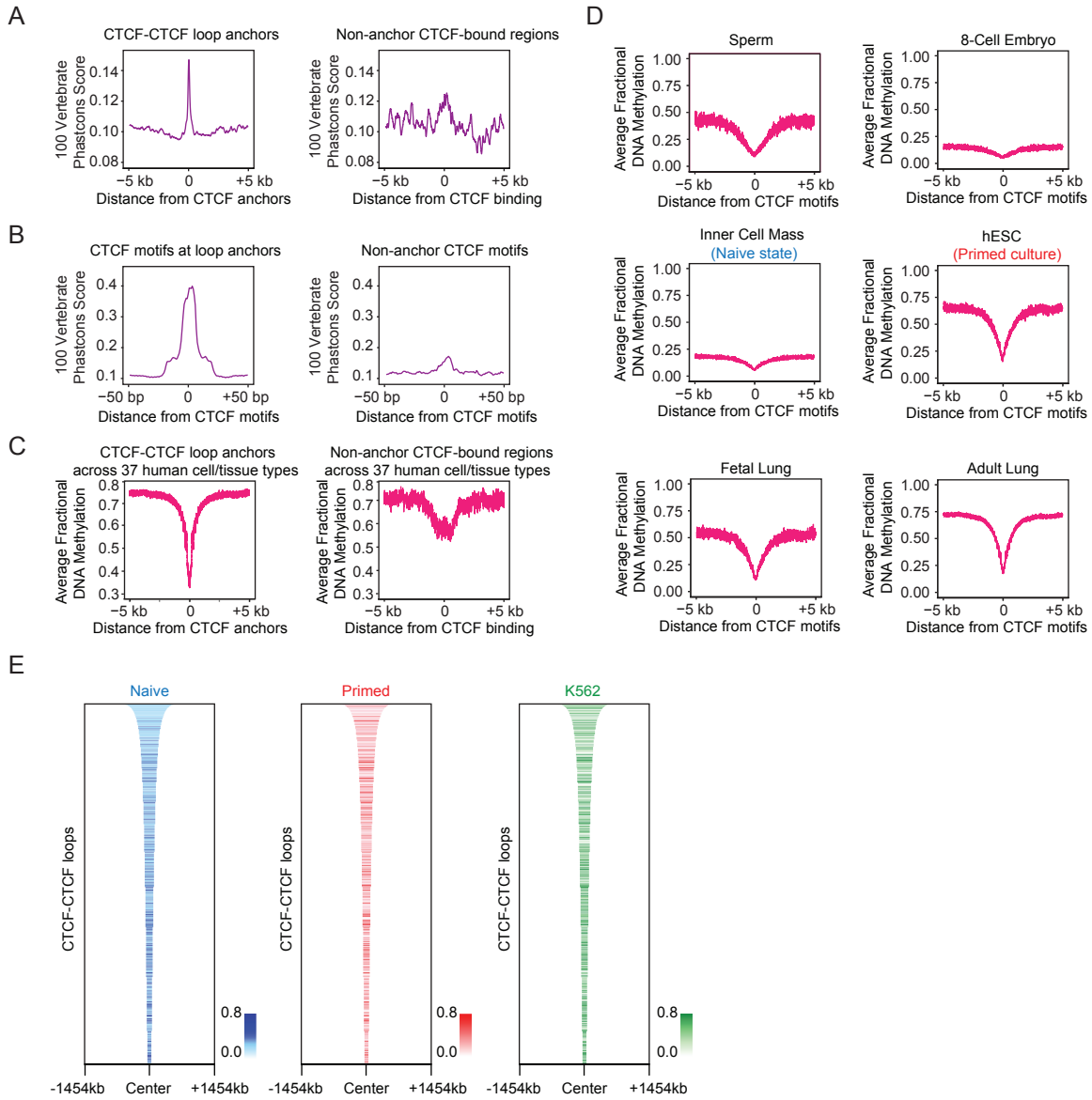


Figure S6



Supplemental Figure Legends

Figure S1. Human ESCs, expression analysis and ChIA-PET data

(A) Phase and fluorescence images of primed hESCs (endogenous OCT4-2A-GFP) and emerging naive colonies induced by treating these primed hESCs with 5i/L/A medium for 10 days. 40x magnification.

(B) Cross-species hierarchical clustering of expression datasets from naive and primed pluripotent cells in both mouse and human highlights the similarity of our datasets to the existing datasets for these cell states in human and mouse samples.

(C) Comparison between the transcriptomes of naive and primed hESCs reveals common and differentially expressed genes.

(D) Correlation analysis for two replicates of cohesin ChIA-PET dataset were displayed by scatter plot.

(E) Percentage of cohesin ChIA-PET interactions that overlap in replicates in naive and primed hESCs.

(F) The overlap of ChIA-PET interactions called by the Downen et al., 2014 pipeline and the Mango pipeline in naive and primed hESCs.

Related to Figure 1

Figure S2. Cohesin-associated interactions are largely responsible for the organization of TADs

(A) CTCF motif orientation analysis of CTCF-CTCF loops. The percentage of each type of CTCF motif orientation is shown in a bar graph.

(B) Saturation analysis for the cohesin ChIA-PET datasets in naive (left panel) and primed (right panel) hESCs.

(C) CTCF-CTCF loops span many TADs identified using Hi-C data in mESCs. Chromosome 2 is displayed as a circos plot in mESCs, with a zoomed in region below. CTCF-CTCF loops (≥ 1 PETs) are indicated as purple arcs. The bar graphs show percentages of TADs spanned by CTCF-CTCF loops when various confidence thresholds (1, 2, ≥ 3 PETs) were used. As a background control, we used random shuffling of TAD locations (100 iterations).

(D) Cohesin ChIA-PET data can be used to discover TADs in mESCs. A comparison of TADs derived with the same algorithm from Hi-C data in mESCs (Dixon et al., 2012) and cohesin ChIA-PET data (mESCs) for a portion of chromosome 6.

(E) Comparison of CTCF-CTCF loops from H1 hESC Hi-C dataset to primed hESC ChIA-PET CTCF-CTCF loops at a locus on chromosome 2. Normalized Hi-C interaction frequencies in H1 hESCs are displayed as a two dimensional heatmap. CTCF-CTCF loops derived from H1 hESC Hi-C dataset are colored in black, ChIA-PET CTCF-CTCF loops are colored in red.

(F) The percent of ChIA-PET CTCF-CTCF loops in primed hESCs present in the Hi-C CTCF loops in H1 hESCs (or vice versa) are plotted as a function of significance thresholds (false discovery rate–FDR) for calling Hi-C interactions displayed as line plots for ChIA-PET CTCF-CTCF loops when various thresholds (1, 2, ≥ 3 PETs) were used.

(G) Comparison of primed hESC ChIA-PET CTCF-CTCF loops to CTCF-CTCF interactions (green lines) derived from in-situ Hi-C in GM12878, IMR90, K562, NHEK, HUVEC, HMEC, and HeLa cells (Rao et al., 2014).

(H) Bar plot indicating the percent of in situ Hi-C CTCF-CTCF loops that overlap with ChIA-PET CTCF-CTCF loops in naive and primed hESCs when various thresholds (1, 2, ≥ 3 PETs) were used.

(I) Cohesin-associated interactions at the *FIRRE* locus are shown. The cohesin-associated interactions are shown as blue lines (naive) and red lines (primed). The

ChIP-seq binding for CTCF and H3K27ac are shown. The blue bar indicates a super-enhancer in naive hESCs.

Related to Figure 2

Figure S3. Cohesin ChIA-PET interactions

(A) Heatmap showing that cohesin ChIA-PET interactions occur predominantly within CTCF-CTCF loops that define putative insulated neighborhoods in hESCs. See the section entitled “Heatmap Representation of High-confidence ChIA-PET Interactions” for details. The color bar indicates normalized high-confidence interactions per loop.

(B) Genotyping for CRISPR-mediated deletions of anchors of CTCF-CTCF loops constraining the super-enhancer associated genes *PRDM14* and *LEFTY1*.

(C) CTCF-CTCF loops tend to be preserved in syntenic regions of human and mouse ESCs. Heatmaps of Hi-C interaction frequencies in H1 hESCs (upper panel) or mESCs (lower panel) are displayed to illustrate a syntenic region (human chr12: 91,760,000-94,960,000, mouse chr10: 94,080,000-96,800,000). Shared CTCF-CTCF loops are indicated as blue lines (naive hESCs) and red lines (primed hESCs). Mouse CTCF-CTCF loops are shown below in purple.

(D) Multiple loops forming insulated neighborhoods (IN) in mESCs whose CTCF boundaries were previously shown to be necessary for insulator function are preserved in human ESCs. The scissor-marked regions were deleted by CRISPR/Cas9 editing in mESCs, which caused local mis-regulation of gene expression (Downen et al., 2014).

Related to Figure 3

Figure S4. 3D structures of TADs containing key pluripotency genes in naive and primed hESCs

(A-D) Schematics of 3D structure for TADs containing *NANOG*, *PRDM14*, *SOX2* and *OCT4* in naive and primed hESCs. For each TAD, Hi-C interaction data (Dixon et al., 2015) is shown together with cohesin-associated loop data for TAD-spanning CTCF loops and insulated neighborhood spanning CTCF loops. A subset of CTCF-CTCF loops was selected for display based on a directionality index (Extended Experimental Procedures) and a subset of genes present in these loops is shown for simplicity. These schematics represent one potential conformation of TADs, but because the underlying data originates with a population of cells, additional conformations are possible.

Related to Figure 4

Figure S5. Differential regulated genes occur in 3D regulatory structures of TADs in naive and primed hESCs

(A-D) Schematics of 3D structure for TADs containing *TBX3*, *HMGB3*, *DUSP6* and *HIF1A* in naive and primed hESCs. For each TAD, Hi-C interaction data (Dixon et al., 2015) is shown together with cohesin-associated loop data for TAD-spanning CTCF loops, insulated neighborhood spanning CTCF loops, enhancer-enhancer loops and enhancer-promoter loops. A subset of CTCF-CTCF loops was selected for display based on a directionality index (Extended Experimental Procedures) and a subset of genes present in these loops is shown for simplicity. These schematics represent one potential conformation of TADs, but because the underlying data originates with a population of cells, additional conformations are possible.

(E) CTCF binding to the TAD and putative insulated neighborhood (IN) anchor sites is preserved in a broad spectrum of human cell types in the domain containing *TBX3*.

(F) CTCF binding to the TAD and putative insulated neighborhood (IN) anchor sites is preserved in a broad spectrum of human cell types in the domain containing *OTX2*.

Related to Figure 5

Figure S6. Conservation of hESC CTCF loop anchors

(A) DNA sequence in anchor regions of CTCF-CTCF loops in hESCs is more conserved in vertebrates than DNA sequence in hESC regions bound by CTCF that do not serve as loop anchors.

(B) The CTCF sequence motif at sites used to anchor DNA loops in hESCs is more conserved in vertebrates than that motif at sites that do not serve as loop anchors in hESCs.

(C) Anchor regions of hESC CTCF-CTCF loops are hypomethylated relative to regions bound by CTCF that do not serve as anchors.

(D) DNA hypomethylation at CTCF-CTCF loop anchors is constitutive throughout the life cycle of humans.

(E) CTCF loops are largely preserved between normal cells (hESCs) and cancer cells (K562). The 9,344 CTCF-CTCF loops that define putative insulated neighborhoods in naive hESCs were ranked by size and shown. The color bar indicates normalized PET-signal at these CTCF-CTCF loops.

Related to Figure 6

Tables

Table S1. RNA-seq gene expression in naive and primed hESCs. Related to Figure 1.

Table S2. SMC1 ChIA-PET, H3K27ac ChIP-seq, CTCF ChIP-seq peaks for hESCs.

Related to Figure 1.

Table S3. High confidence SMC1 ChIA-PET interactions for naive and primed hESCs.

Related to Figure 1, 2.

Table S4. Differential Super-enhancers and differential CTCF-CTCF loops between naive and primed hESCs. Related to Figure 5

Table S5. Cancer mutations identified at CTCF motifs within CTCF-CTCF loop anchors.

Related to Figure 6

Table S6. Mango-called high confidence SMC1 ChIA-PET interactions for naive and primed hESCs. Related to the experimental procedures.

Sequences used in this study

Gene	Sequence	Application
<i>PRDM14</i> sgRNA 1	GTGACACTGTGCAGACCACT	sgRNA target sequence
<i>PRDM14</i> sgRNA 2	ATAAGAAGGGTGGCCGGGCG	sgRNA target sequence
<i>LEFTY1</i> sgRNA 1	AAGGTGGGTCTCACAGGATT	sgRNA target sequence
<i>LEFTY1</i> sgRNA 2	GAAATAGGTAACCTTTTAA	sgRNA target sequence
TAD L sgRNA 1	GGGGAGGTGCTCCGTA CTTC	sgRNA target sequence
TAD L sgRNA 2	AAACAGCTGACAACATCGAA	sgRNA target sequence
TAD R sgRNA 1	GAGCCATCCGGTGGTAGATT	sgRNA target sequence
TAD R sgRNA 2	CAGAGTTGGTGACTCCGTAA	sgRNA target sequence
<i>PRDM14</i> F	CCTGACATCTCAGTGCACGT	Genotype PCR

<i>PRDM14</i> R	CCTTGCTCTATCGCCCAGTC	Genotype PCR
<i>LEFTY1</i> F	AGCGGAAAACAACAGCAAAT	Genotype PCR
<i>LEFTY1</i> R	GCAACTGAAGTGAGTGCATGA	Genotype PCR
TAD L F	TGGCACTAGATATTTGAGAGAAATTG	Genotype PCR
TAD L R	TCTTCCAGGTTCAACGCTCT	Genotype PCR
TAD R F	CAAGTCCTGGGTTCTCATCC	Genotype PCR
TAD R R	TTGAGATCCCAGGAGTGAGG	Genotype PCR
<i>GAPDH</i> F	CGAGATCCCTCCAAAATCAA	RT-qPCR
<i>GAPDH</i> R	ATCCACAGTCTTCTGGGTGG	RT-qPCR
<i>PRDM14</i> F	ACACGCCTTTCCCGTCCTA	RT-qPCR
<i>PRDM14</i> R	GGGCAGATCGTAGAGAGGCT	RT-qPCR
<i>SLCO5A1</i> F	ACCTCAGCAAACCTTCTCGG	RT-qPCR
<i>SLCO5A1</i> R	GAGACCATTAACGCCTGGATG	RT-qPCR
<i>LEFTY1</i> F	TGATCGTCAGCATCAAGGAG	RT-qPCR
<i>LEFTY1</i> R	GAGCACAGAGCATTTGTCCA	RT-qPCR
<i>SDE2</i> F	AGGATTCCGTCCTCAAAGGT	RT-qPCR
<i>SDE2</i> R	TGGACCCTTCTGCAGTCTCT	RT-qPCR
<i>KLF4</i> F	GATGGGGTCTGTGACTGGAT	RT-qPCR
<i>KLF4</i> R	CCCCCAACTCACGGATATAA	RT-qPCR
<i>NANOG</i> F	GCAGAAGGCCTCAGCACCTA	RT-qPCR
<i>NANOG</i> R	AGGTTCCCAGTCGGGTTCA	RT-qPCR
<i>OCT4</i> F	GCTCGAGAAGGATGTGGTCC	RT-qPCR
<i>OCT4</i> R	CGTTGTGCATAGTCGCTGCT	RT-qPCR
<i>OTX2</i> F	CAAAGTGAGACCTGCCAAAAGA	RT-qPCR
<i>OTX2</i> R	TGGACAAGGGATCTGACAGTG	RT-qPCR
BAC1 Probe	RP11-487J21	3D DNA FISH
BAC2 Probe	RP11-13714	3D DNA FISH

Extended Experimental Procedures:

Cell Culture

Primed and naive hESCs were cultured as previously described (Theunissen et al., 2014). Primed hESCs were maintained on mitomycin C-inactivated MEF feeder layers and passaged every 7-10 days. When passaging primed hESCs, clumps of cells were partially dissociated with collagenase type IV (GIBCO, 17104-019), and then subjected to two sedimentation steps in stationary 50 cm tubes for 10 minutes at room temperature in primed hESC medium to remove single cells. Primed hESC medium (500 ml) consisted of 400 ml of Dulbecco's Modified Eagle Medium: Nutrient Mixture F-12 (DMEM/F12, Invitrogen, 11320), 75 ml Fetal Bovine Serum (FBS, Hyclone, SH30071.03HI), 25 ml KnockOut™ Serum Replacement (KSR, Invitrogen, 10828-028),

supplemented with 1 mM glutamine (Invitrogen, 25030-024), 1% nonessential amino acids (Invitrogen, 11140-050), penicillin-streptomycin (Invitrogen, 15140-122), 0.1 mM β -mercaptoethanol (Sigma, M6250-100ML), and 4 ng/ml FGF2 (R&D systems, 233-FB-025).

For the induction of naive hESCs, primed hESCs were cultured for 24 hr in the primed hESC medium described above, further supplemented with 10 μ M ROCK inhibitor Y-27632 (Stemgent, 04-0012). Colonies were then trypsinized to form a single cell suspension and cells were plated onto a MEF feeder layer in the primed hESC medium + ROCK inhibitor described above. 24 hr later, the medium was switched to 5i/L/A naive hESC medium. The 5i/L/A naive hESC medium (500 ml) used for induction and maintenance of naive hESCs was made up of 240 ml DMEM/F12, 240 ml Neurobasal (Invitrogen, 21103), 5 ml N2 supplement (Invitrogen, 17502048) and 10 ml B27 supplement (Invitrogen, 17504044), supplemented with 10 μ g recombinant human LIF (purified in-lab from *E. coli*), 1 mM glutamine, 1% nonessential amino acids, 0.1 mM β -mercaptoethanol, penicillin-streptomycin, 50 μ g/ml BSA (Sigma, A4737-25G), and the following small molecules and cytokines: 1 μ M PD0325901 (Stemgent, 04-0006), 1 μ M IM-12 (Enzo, BML-WN102-0005), 0.5 μ M SB590885 (R&D systems, 2650/10), 1 μ M WH-4-023 (A Chemtek) 10 μ M Y-27632 (Stemgent, 04-0012), and 10 ng/ml Activin A (Peprotech, 120-14). Following an initial wave of widespread cell death, dome-shaped naive hESC colonies appeared within 10 days and could be expanded and maintained in 5i/L/A naive hESC medium.

Naive hESCs were maintained on mitomycin C-inactivated MEF feeder cells and passaged every 5-7 days. The naive hESCs were passaged by dissociating cells with accutase (GIBCO, A1110501), and then centrifuging cells at 1000 rpm for 5 minutes at room temperature in neutralization medium (DMEM supplemented with 10% FBS, 1 mM glutamine, 1% nonessential amino acids, penicillin-streptomycin, and 0.1 mM β -mercaptoethanol). To harvest cells for downstream experiments, primed and naive hESCs were trypsinized and subsequently pre-plated on gelatin-coated dishes to deplete MEF feeder cells. All cell culture experiments were performed under physiological oxygen conditions (5% O₂, 3% CO₂).

Genome Editing

The CRISPR/Cas9 system was used to create hESCs with CTCF site deletions. For each experiment two target-specific oligonucleotides (sgRNA) flanking the proposed deletion were cloned into plasmids carrying a codon-optimized version of Cas9 (pX330, Addgene: 42230) that had been further engineered with either a GFP or mCherry fluorescent reporter. WIBR3 primed hESCs were cultured in 10 μ M ROCK inhibitor (Stemgent; Y-27632) 24 hr prior to electroporation. Two confluent six well plates of cells were harvested using 0.25% trypsin/EDTA (Invitrogen) and resuspended in phosphate buffered saline (PBS). Cells were electroporated with 20 μ g of pX330-sgRNA-GFP and 20 μ g pX330-sgRNA-mCherry targeting up and downstream of the intended CTCF site deletion. Cells were subsequently plated in MEF feeder layers in primed hESC medium supplemented with 10 μ M ROCK inhibitor. 48 hr post electroporation, cells were harvested using 0.25% trypsin/EDTA and double positive GFP+/mCherry+ cells were isolated by Fluorescent Activated Cell Sorting (FACS). After sorting, GFP+/mCherry+ cells were plated on MEF feeder layers in primed hESC medium supplemented with 10 μ M ROCK inhibitor. 8-12 days later individual colonies were picked, expanded and genotyped by PCR.

shRNA knockdown

VSVG coated lentiviruses were generated in HEK-293 cells. Viral containing supernatant was collected 48 and 72 hr post-transfection. Viral supernatant was filtered through a 0.45 mm filter. 24 hr prior to infection primed human ESCs were treated with 10 μ M ROCK inhibitor. On the day of infection naive and primed human ESCs were single cell disassociated with Accutase and 0.25% trypsin respectively. Cells were then resuspended in lentiviral supernatant w/polybrene in ultra-low attachment plates and spun in a centrifuge at 2000 rpm. This spin infection was conducted for 1.5 hr. Cells were then replated on DR4 MEF feeder layer (primed cells were supplemented w/ ROCKi). Medium was changed after 20 hr. 48 hr post-infection medium was supplemented with puromycin (0.5 μ g/ml) to select for proviral integration, and doxycycline (2 μ g/ml) to induce expression of the shRNA. Seven days later, RNA was extracted and gene expression level was measured by RT-qPCR.

Gene Expression Analysis

RNA was isolated using Trizol reagent (Invitrogen, 15596-026), and reverse transcribed using oligo-dT primers and SuperScript III reverse transcriptase (Invitrogen, 18080044) according to the manufacturer's instructions. Quantitative real-time PCR was performed on a 7000 ABI Detection System with FAST SYBR Green Master Mix (Applied Biosystems, 4309155). Gene expression was normalized to GAPDH.

3D DNA FISH

3D DNA FISH was performed as previously described (Bolland et al., 2013). Briefly, cells were attached to slides using a Cytospin at 500 rpm, 3 min, then fixed with 4% paraformaldehyde (PFA) for 10 min at room temperature, then quenched in 0.1 M Tris-HCl, pH 7.4 for 10 min at room temperature. Next, cells were permeabilized in 0.1% saponin/0.1% Triton X-100 in PBS for 10 min at room temperature. Slides were then washed twice in PBS for 5 min at room temperature and subsequently incubated for at least 20 min in 20% glycerol/ PBS at room temperature. Slides were freeze/ thawed in liquid nitrogen three times, and washed twice in PBS for 5 min at room temperature. Slides were incubated in 0.1 M HCl for 30 min at room temperature, washed in PBS for 5 min at room temperature and permeabilized in 0.5% saponin/0.5% Triton X-100/PBS for 30 min at room temperature. This was followed by two washes in PBS for 5 min at room temperature before equilibration in 50% formamide/2x SSC for at least 10 min at room temperature. BAC probes (Empire Genomics) were pipetted onto a coverslip. FISH slides were air dried and heated to 78 °C for precisely 2 min on a hot plate. Coverslip w/probe was mounted in slides and sealed with rubber cement. Slides were incubated overnight at 37 °C in a dark humidified chamber. The next day, rubber cement was removed and slides were placed in 2x SSC until coverslips detached. Slides were washed in 50% formamide/2x SSC for 15 min at 45 °C, washed in 0.2x SSC for 15 min at 63 °C, washed in 2x SSC for 5 min at 45 °C and washed in 2x SSC for 5 min at room temperature. Subsequently slides were washed in PBS for 5 min at room temperature and stained with DAPI (5 μ g/ml in 2x SSC) for 2 min at room temperature. Finally, slides were destained in PBS for 5 min at room temperature and coverslips were mounted on slides. Imaging was carried out by confocal microscopy in the Whitehead Keck imaging facility. Spatial distance between FISH probes was quantified using ImageJ (FIJI).

BACs

High-Throughput 3D DNA FISH BAC clones were purchased from BACPAC (CHORI) and were used to generate fluorescently labeled probes as described (Shachar et al., Cell 2015). Probes were as follows: RP11-261P1 (chr5_TAD1), RP11-810B19

(chr5_TAD2), RP11-1029M14 (chr5_equidist_con), RP11-258M5 (chr11_TAD1), RP11-52J19 (chr11_TAD2), RP11-2L5 (chr11_equidist_con).

High-Throughput 3D DNA FISH

High-throughput 3D DNA FISH in naive human embryonic stem cells was done as previously described in (Shachar et al., 2015). Briefly, 10^5 cells per well were plated in 96-well plates, fixed in 4% PFA in PBS for 15 min at room temperature, permeabilized and denatured as described (Shachar et al., 2015). A mix containing 300 ng of each fluorescently labeled probe was ethanol precipitated and re-suspended in 25 μ l of hybridization buffer. Cells were denatured with probe mix at 90 °C for 8 min and left to hybridize overnight. Images were acquired using a Perkin Elmer Opera automated imaging system at >100 randomly sampled fields in multiple wells using a 40X water objective. Image analysis was carried out as described in (Burman et al., 2015). Briefly, FISH spot coordinates were detected in individual nuclei that contained an equal number of spots in each channel. The minimal distance in pixels between each combination of spot pairs (TAD border 1, TSD border 2, control) was calculated using Acapella 2.0 (PerkinElmer) and R (<http://www.R-project.org/>). For statistical analysis, the Mann-Whitney test was used to compare the closest distance between TAD borders to an equidistant control locus.

RNA-seq

RNA-seq was performed for naive and primed hESCs. 6 million cells were used for each RNA extraction. Total RNA was purified using the mirVana™ miRNA Isolation Kit (Life Technologies, AM1560) following the manufacturer's instructions. 1 μ g of total RNA was used for the RNA-seq library construction. A technical replicate was performed for both naive and primed hESCs. Polyadenylated RNA-seq libraries were prepared using the TruSeq Stranded mRNA Library Prep Kit (Illumina, RS-122-2101). The RNA-seq libraries were sequenced on the Illumina HiSeq 2000.

RNA-seq Expression Analysis

RNA-seq alignment and quantification were performed using the TopHat and Cufflinks software tools. RNA-seq reads were first aligned to the human genome (build hg19, GRCh37) using Tophat v2.0.13 (Trapnell et al., 2009) with the parameters: --solexa-quals --no-novel-juncs and using RefSeq gene annotations. The expression levels of RefSeq transcripts were calculated using Cufflinks v2.2.1 (Trapnell et al., 2010). Differentially expressed transcripts were then identified, again using Cufflinks v2.2.1. When multiple transcripts had the same gene name, only the transcript with the highest expression level was kept for further consideration. A gene was considered differentially expressed if it met the following criteria: 1) absolute log₂ fold-change ≥ 1 between the mean expression in the two conditions; 2) false discovery rate q-value ≤ 0.05 .

Three lines of evidence suggested that the RNA-seq datasets were high-quality: 1) ~80% of all reads in all libraries mapped to RefSeq transcript models (hg19), as expected for sequencing of RNA; 2) ~90% of all reads in all libraries mapped to known RefSeq genes (~83% mapped to the exons and ~7% mapped to the introns), as expected for sequencing of poly-A RNA-enriched samples; 3) the replicates of either naive or primed RNA-seq datasets had a Pearson correlation coefficient of expression levels of 0.98 or greater across all RefSeq transcripts.

Cross-Species Gene Expression Analysis

Cross-species gene expression analysis was performed as previously described (Theunissen et al., 2014). For a given gene, the mean expression value for that gene across all human samples was first calculated. Then for each human sample, the expression of that gene in that sample was divided by the mean expression value. The normalization was repeated for all mouse samples. After normalization, all pairwise comparisons of datasets, both intra- and inter-species, were performed using Pearson correlation coefficients (PCCs). The average linkage hierarchical clustering of the Pearson correlation was shown in the heatmap.

ChIP-seq Library Generation and Sequencing

Chromatin immunoprecipitation (ChIP) was performed as previously described (Ji et al., 2015). 50 million naive or primed hESCs were used for each ChIP experiment. The following antibodies were used for ChIP: anti-H3K27ac (Abcam, ab4729), anti-CTCF (Millipore, 07-729), anti-MED1 (Bethyl Labs, A300-793A), anti-OCT4 (Santa Cruz, sc-8628). For each ChIP, 5 µg of antibody and 50 µl protein G Dynabeads (Life Technology, 10004D) were used. The ChIP-seq libraries were prepared using the TruSeq ChIP Sample Prep Kit (Illumina, IP-202-1012), and sequenced on the Illumina HiSeq 2000.

ChIA-PET Library Generation and Sequencing

ChIA-PET was performed using a modified version of a previously described protocol (Downen et al., 2014). 400 million naive or primed hESCs were used for each ChIA-PET library construction. The ChIA-PET libraries were generated in three stages. In the first stage, ChIP was performed using 25 µg anti-SMC1 antibody (Bethyl Labs, A300-055A) and 250 µl protein G Dynabeads (Life technology, 10004D). This stage was the same as the experimental procedure described in the ChIP-seq library generation.

The second stage was proximity ligation of ChIP-DNA fragments, which consists of end blunting and A-tailing to create easily ligated ends, followed by ligation to simultaneously add linker sequences required for later steps and ligate ends of fragments together. The ligation was performed in a large volume to encourage ligation of ends that are in close spatial proximity to each other, ideally from fragments that are co-localized via their interaction with cohesin-bound regions and immunoprecipitation of cohesin. The ChIP-DNA with beads were washed once with TE buffer, then incubated in 1x T4 DNA polymerase buffer (NEBuffer 2.1, New England Biolabs, B7202S), with 7.2 µl T4 DNA polymerase (New England Biolabs, M0203S) and 7 µl of 10 mM dNTPs (Life Technologies, 18427013) in 700 µl total volume at 37 °C for 40 min. The beads were then washed three times with ChIA-PET wash buffer (10 mM Tris-HCl pH7.5, 1 mM EDTA, 500 mM NaCl). The beads were incubated with 1x NEB buffer 2 (New England Biolabs, B7002S) containing 7 µl Klenow fragment (3'-5' exo⁻) (New England Biolabs, M0212S) and 7 µl 10 mM dATP (New England Biolabs, N0440S) in 700 µl total volume at 37 °C for 50 min. The beads were then washed three times with ChIA-PET wash buffer. The beads were then incubated with 1x T4 DNA ligase buffer with 1mM ATP (New England Biolabs, B0202S) containing 42 µl T4 DNA ligase (Life Technologies, 46300018) and 4 µl bridge linker (200 ng/µl including Forward: /5Phos/CGCGATATC/iBiodT/TATCTGACT; Reverse: /5Phos/GTCAGATAAGATATCGCGT) in 14 ml total volume at 16 °C for 22 hr. The beads were then washed three times with ChIA-PET wash buffer. The beads were then incubated with 1x lambda exonuclease buffer (New England Biolabs, M0262S) containing 6 µl lambda exonuclease (New England Biolabs, M0262S), and 6 µl exonuclease I (New England Biolabs, M0293S) in 700 µl total volume at 37 °C for 1 hr.

DNA elution and crosslink reversal were simultaneously performed by incubating the beads at 55 °C overnight. 10 µl of proteinase K (Life Technologies, AM2546) was included during the overnight incubation. The DNA was then purified by phenol-chloroform extraction and ethanol precipitation.

The third stage was the tagmentation of ligated products, purification of the tagmented DNA fragments, amplification of the DNA by PCR, size selection and paired-end sequencing. The ChIA-PET proximity ligation products were tagmented with Tn5 Transposase (5 µl Tn5 transposase (Illumina, FC-121-1030) for 50 ng DNA) at 55 °C for 5 min, then at 10 °C for 10 min. DNA was purified using a Zymo column (VWR, 100554-654) following the manufacturer's instructions. Biotin-labeled DNA was then further affinity purified with M280 streptavidin beads (50 µl for each library, Life Technologies, 11205D), followed by washing five times with 2x SSC/0.5% SDS and then two times with 1x B&W buffer (5 mM Tris-HCl pH7.5, 0.5 mM EDTA, 1 M NaCl). The buffer was discarded and the beads were gently resuspended in 30 µl EB buffer (QIAGEN). 10 µl of the bead slurry was used for PCR amplification. PCR amplification was performed using the Nextera DNA Sample Preparation Kit (Illumina, FC-121-1031) for 10-12 cycles. The DNA was selected for the size range of 300-500 bp and was purified by gel extraction. The ChIA-PET library was subjected to 100 x 100 paired-end sequencing using Illumina HiSeq 2000.

ChIP-seq Data Analysis

All ChIP-Seq datasets were aligned to the human genome (build hg19, GRCh37) using Bowtie (version 0.12.2) (Langmead et al., 2009) with the parameters `-k 1 -m 1 -n 2`. We used the MACS peak finding algorithm, version 1.4.2 (Zhang et al., 2008) to identify regions of ChIP-seq enrichment over input DNA control with the parameters `--no-model --keep-dup=1`. A p-value threshold for enrichment of 1e-09 was used for H3K27ac, H3K27me3 (Theunissen et al., 2014), MED1 and OCT4 datasets, while a p-value of 1e-07 was used for the CTCF dataset. UCSC Genome Browser (Kent et al., 2002) tracks were generated using the MACS wiggle file output option with parameters `-w -S -space=50`. All gene-centric analyses in human ESCs were performed using human (build hg19, GRCh37) RefSeq annotations downloaded from the UCSC genome browser (genome.ucsc.edu).

ChIA-PET Data Processing

All ChIA-PET datasets were processed with a method adapted from a previously published computational pipeline (Dowen et al., 2014; Li et al., 2010). The output of paired-end sequencing is a set of reads, where each read is identified by a read id and consists of two mates that represent sequence from the ends of a DNA fragment. The raw sequences of each mate of each read were analyzed for the presence of the PET linker barcodes and trimmed using Cutadapt with the parameters `-m 17 -a forward=ACGCGATATCTTATCTGACT -a reverse=AGTCAGATAAGATATCGCGT --overlap 10` (Martin, 2011) specifically, we searched for a stretch of at least 10 bp that matched the linker sequence. Once this sequence was identified, the linker sequence and all sequence immediately 3' to this sequence was removed. After removal of linker and 3' sequence, only sequences of at least 17 bp in length were retained. For downstream analysis, all mates from all reads where at least one mate contained the linker sequence were used. Sequences of mates were separately mapped to the hg19 human genome using Bowtie with the parameters `-k 1 -m 1 -v 2 -p 4 --best --strata` (Langmead et al., 2009). These criteria retained only the uniquely mapped mates, with at most two base pair mismatches, for further analysis. Aligned mates were paired using

their respective read ids and now considered PETs (paired-end tags). PETs were filtered for redundancy: PETs with identical genomic coordinates and strand information at both ends were collapsed into a single PET. The PETs were further categorized into intrachromosomal PETs, where the two ends of a PET were on the same chromosome, and interchromosomal PETs, where the two ends were on different chromosomes. The sequences from the ends of all PETs were then analyzed for localized enrichment across the genome using MACS 1.4.2 (Zhang et al., 2008) with the parameters “-p 1e-09 -no-lambda -no-model --keep-dup=2”. Regions identified with MACS were considered PET peaks.

To identify long-range chromatin interactions, we first removed intra-chromosomal PETs of length < 4 kb because these PETs are suspected to originate from self-ligation of DNA ends from a single chromatin fragment in the ChIA-PET procedure (Downen et al., 2014). We next identified PETs that overlapped with PET peaks at both ends by at least 1 bp. Operationally, these PETs were defined as putative interactions. Applying a statistical model based upon the hypergeometric distribution identified high-confidence interactions, representing high-confidence physical linking between the PET peaks. To do this, for each PET peak, we calculated a) the total number of PETs that overlap with the peak and b) the number of PETs that overlap with the peak and also connect to another peak. A hypergeometric distribution was used to determine the probability of seeing at least the observed number of PETs linking the two PET peaks. The correction p-values were calculated using the Benjamini-Hochberg procedure (Benjamini and Hochberg, 1995) to control for multiple hypothesis testing. Operationally, the pairs of interacting sites with three independent PETs and an FDR ≤ 0.01 were defined as high-confidence interactions in the SMC1 ChIA-PET merged dataset and with two independent PETs in the individual SMC1 ChIA-PET replicates. Previously published RAD21 (cohesin) ChIA-PET datasets in K562 (Heidari et al., 2014) were downloaded from ENCODE (<https://www.encodeproject.org/experiments/ENCSTR000FDB/>) and were re-processed exactly as described (Downen et al., 2014; Li et al., 2010).

Additional ChIA-PET interaction analysis

Several additional analyses were conducted to characterize the sets of interactions identified in this work and improve our confidence in these calls. Interactions were compared to those called by a second analysis pipeline to demonstrate robust identification of interactions. Interactions were analyzed for the presence of expected DNA sequence motifs at their ends. Interactions were analyzed for the presence of expected regulatory elements at their ends. Finally, interactions were compared to interactions detected by Hi-C, a second experimental method.

Interactions were compared to interactions called using a second analysis pipeline (called Mango (Phanstiel et al., 2015)). Different analysis pipelines incorporate different biases and assumptions in identifying interactions, and depend on selection of parameters and thresholds, and are thus likely to yield different results. Regardless, one expectation is that a large number of interactions should be identified robustly if bona fide interactions are being found. Thus, comparison of the outputs of the two pipelines provides some measure of the robustness of the identification of interactions. The Mango software (version 1.0.2) was downloaded (<https://github.com/dphansti/mango>), installed, and initially run with default settings. For input, the same data (sequencing reads and genomic regions called enriched for signal) used for the Downen et al. pipeline were applied to the Mango pipeline. 10,702 and 13,233 ChIA-PET interactions were called in naive and primed hESCs, respectively (Table S6) when using the Mango

pipeline. These interactions were then compared to their counterparts derived with the Downen et al. pipeline. Interactions were considered shared if each end of an interaction identified with the Mango pipeline overlapped by at least one base pair with the respective ends of an interaction identified with the Downen et al. pipeline. 86% of the 10,702 interactions in naive hESC were identified in both pipelines. 91% of the 13,233 interactions in primed hESC were identified in both pipelines (Figure S1F). The observation of robust identification of at least a subset of the data using different analysis pipelines generally increases our confidence that a large set of bona fide interactions is being identified.

Interactions were analyzed for the presence of expected DNA sequence motifs at their ends. The subset of cohesin-associated, CTCF-CTCF loops are expected to have convergently oriented CTCF DNA-binding sequence motifs underlying each of the two CTCF binding regions that comprise the ends of the loops. Briefly, the location and orientation of the CTCF motifs at CTCF ChIP-seq peaks were identified using the FIMO software package 4 (Grant et al., 2011; Matys et al., 2006) and searching with the canonical CTCF motif from the Jaspas motif database (ID. MA0139.1). The orientation of CTCF motifs at pairs of CTCF ChIP-seq peaks was next determined. For simplicity, we focused on those CTCF-CTCF loops where CTCF peaks could be unambiguously assigned to a CTCF motif. All pairs of CTCF motifs at the two ends of CTCF-CTCF ChIA-PET interactions were classified into one of the four possible classes of motif orientations: a convergent orientation (forward-reverse), a divergent orientation (reverse-forward), the same direction on the forward strand (forward-forward) or the same direction on the reverse strand (reverse-reverse). Additional details can be found in the section titled “CTCF Motif Orientation Analysis at CTCF-CTCF Loops”. Approximately 80% of the interactions identified here have CTCF sequence motifs in the expected orientation (convergent) at the ends of the interactions (Figure S2A). The observation that CTCF-CTCF loops identified here display the expected orientation of CTCF sequence motifs at their ends increases our confidence that bona fide interactions are being identified.

Interactions were analyzed for the presence of expected regulatory elements at their ends. Cohesin-associated loops are expected to have ends associated with CTCF sites, enhancers and promoters. For this analysis, CTCF sites and enhancers were identified using ChIP-seq data for CTCF and the histone modification H3K27ac, respectively. Briefly, ChIP-Seq datasets were aligned to the human genome to identify regions of ChIP-seq enrichment over input DNA control. A p-value threshold for enrichment of $1e-07$ was used for CTCF, while a p-value of $1e-09$ was used for the H3K27ac dataset. H3K27ac-enriched regions were further filtered for those that were at least 2 kb away from a RefSeq transcription start site to identify the set of enhancer regions. Promoters were defined as the region ± 2 kb around RefSeq transcription start sites. Additional details can be found in the section titled “ChIP-seq Data Analysis”. 75-85% of the interactions identified here in naive or primed hESCs had ends that overlapped with CTCF sites, enhancers or promoters (greater than 1 bp overlap). The large fraction of interactions identified with ends overlapping biologically relevant genomic features increases our confidence that bona fide interactions are being identified.

Interactions were compared to interactions detected by Hi-C, a second experimental method to detect interactions. ChIA-PET detects interactions occurring between sites associated with a specific protein, while Hi-C detects interactions more generally. Thus, the set of ChIA-PET interactions is expected to overlap with the set of Hi-C interactions

A previously published Hi-C dataset from H1 hES cells was downloaded and used to derive a set of interactions (Dixon et al., 2015) Briefly, the raw fragment contact and bias matrices at 40 kb resolution were first obtained using the python hiclib library. The Fit-Hi-C tool (Ay et al., 2014) was then used to call high-confidence DNA interactions (FDR .05). Additional details can be found in the section titled “Calling Hi-C Interactions”. Given the 40 kb bin size, the minimum distance of interactions from the Hi-C data was effectively 80 kb. Thus for comparisons, we compared the Hi-C interactions to ChIA-PET interactions that were 80 kb or greater in length. 74% of the interactions identified using the cohesin ChIA-PET data were also identified in the Hi-C data. The large fraction of the ChIA-PET interactions identified in Hi-C data increases our confidence that bona fide interactions are being identified with the cohesin ChIA-PET data.

Assignment of Interactions to Regulatory Elements

We assigned the PET peaks of interactions to different regulatory elements, including promoters (+/- 2 kb of the Refseq TSS), active enhancers (H3K27ac enriched regions falling outside of promoter regions that are defined as +/- 2 kb of the Refseq TSS), and CTCF ChIP-seq binding sites. Operationally, an interaction was defined as associated with the regulatory element if one of the two PET peaks of the interaction overlapped with the regulatory element by at least 1 base pair. CTCF-CTCF loops were defined as high confidence ChIA-PET loops with CTCF ChIP-seq peaks at both ends of the interactions.

Identification of CTCF-CTCF Loops that Define Putative Insulated Neighborhoods

All CTCF-CTCF loops may potentially form insulated neighborhoods. For this paper putative insulated neighborhoods were defined by incorporating some evidence for loop insulation by measure of directionality index as described below. Briefly, CTCF-CTCF loops were evaluated for putative insulating function by examining the directionality of reads proximal to loop boundaries. One expectation for a loop with insulating function is that, at a loop boundary, interactions originating just upstream of the boundary connect to a distal point located further upstream while interactions originating just downstream of the boundary connect to a distal point located further downstream. Boundaries satisfying these criteria thus have implied functionality in terms of constraining interactions. Adjacent pairs of boundaries satisfying these criteria would thus be candidates for demonstrating insulating function. ChIA-PET interaction directionality preferences were calculated using a method adapted from Hi-C computational analysis (Mizuguchi et al., 2014). Briefly, each chromosome (autosomes and X chromosome) was divided into non-overlapping 40 kb bins. Each intra-chromosomal ChIA-PET interaction (either below or above 4 kb) was then mapped to the matrix comprised of all pairwise combinations of bins. Each end of a ChIA-PET interaction contributed signal to its respective bin, thus generating a matrix of interaction frequencies between bins. ChIA-PET directional preference scores were next calculated from these interaction frequency matrices as the log₂ ratio of upstream to downstream contact frequencies for each region *i* at distances below 400 kb:

$$D_i = \log_2 \left(\frac{\sum_{j=-10}^{j=0} C_{i,i+j}}{\sum_{j=0}^{j=10} C_{i,i+j}} \right),$$

in which C is the ChIA-PET interaction frequency matrix.

Putative insulated neighborhoods were operationally defined as intra-chromosomal CTCF-CTCF interactions where each end of the interaction displayed a change in directional preference. This type of change in interaction preference between upstream and downstream genomic regions was previously used to computationally define topologically associating domains (Dixon et al., 2012; Nora et al., 2012). To improve the robustness of calculating interaction preferences at the CTCF-occupied peaks at CTCF-CTCF interactions, we calculated the average interaction preference at two neighboring bins in the proximity of the CTCF-occupied peaks. Specifically, we first identified the genomic bins where the two ends of CTCF-CTCF interactions were located. For each of the 5' CTCF-occupied PET peaks of these CTCF-CTCF interactions, we selected two bins: one located where the 5' PET peak was located and the other in the immediately neighboring bin in the 3' direction. For each of the 3' CTCF-occupied PET peaks at CTCF-CTCF interactions, we also selected two bins: one located where the 3' PET peak was located and the other in the immediately neighboring bin in the 5' direction. We then filtered for CTCF-CTCF interactions whose mean of interaction directional preference between the two bins at their 5' PET peak was positive (indicating downstream preferences) and mean of interaction directional preference between two bins at their 3' PET peak was negative (indicating upstream preferences). Since the ChIA-PET interaction frequency matrix was calculated using 40 kb bins, this method allowed us to detect putative insulated neighborhoods greater than 80 kb.

TAD schematic construction

For schematics of TAD structures, we show TAD-spanning loops with at least one PET read. We show those putative insulated neighborhoods that pass the directionality index criteria described above. All non-overlapping putative insulated neighborhoods are shown. When overlapping putative insulated neighborhoods are possible, the loop with the most PET reads supporting the interaction was selected for display. When comparing structures encompassing genes with cell type preferred enhancers in naive versus primed hESCs, structures were first identified in the cell type with the cell type preferred enhancer. The second cell type was then examined for the presence of corresponding structures with evidence for CTCF binding. As a default, when enhancer signals were similar, naive hESC structures were first identified and the primed hESCs were then examined for the corresponding structure. For simplicity, a subset of genes is displayed with their associated enhancers. Enhancers were defined as stitched H3K27ac MACS peaks (using the ROSE algorithm). The loop with the highest PET reads supporting each enhancer-promoter or enhancer-enhancer interaction was shown (using $PET \geq 2$).

ChIA-PET Interaction Heatmap at Insulated Neighborhoods

Cohesin ChIA-PET interactions were displayed to examine the similarity of neighborhoods between naive hESCs, primed hESCs, and K562. Insulated neighborhoods for naive hESCs were centered and size-normalized. ChIA-PET PET signal (number of uniquely mapped PETs per million uniquely mapped PETs) was then displayed. For comparison, the ChIA-PET signal from primed hESCs and K562 for the regions with the same coordinates was displayed.

CTCF Motif Orientation Analysis at CTCF-CTCF Loops

The location and orientation of the CTCF motifs at CTCF ChIP-seq peaks were identified using the FIMO software package with a default p value threshold of 10^{-4} (Grant et al., 2011; Matys et al., 2006). In the analysis, the canonical CTCF motif from the Jaspas motif database (ID. MA0139.1) was used. The orientation of CTCF motifs at pairs of

CTCF ChIP-seq peaks were next determined. For simplicity, we focused on those CTCF-CTCF loops where CTCF peaks could be unambiguously assigned to a CTCF motif: each end overlapped a single CTCF ChIP-seq peak by at least 1 base pair and only a single CTCF motif was at the peak. All pairs of CTCF motifs at the two ends of CTCF-CTCF ChIA-PET interactions were classified into one of the four possible classes of motif orientations: a convergent orientation (forward-reverse), a divergent orientation (reverse-forward), the same direction on the forward strand (forward-forward) or the same direction on the reverse strand (reverse-reverse).

Hi-C Interaction Heatmap

To generate a matrix of Hi-C interaction frequencies mapped to a more recent build of the human genome, previously published Hi-C datasets in H1 hESCs (Dixon et al., 2015) were first downloaded from GEO (www.ncbi.nlm.nih.gov/geo/; accession GSM1267196 and GSM1267197). The raw reads from these datasets were mapped to the human genome build hg19 and filtered as previously described (Imakaev et al., 2012). Corrected contact probability matrices at 40 kb resolution were obtained using the python hiclib library (<https://bitbucket.org/mirnylab/hiclib>).

Super-Enhancers in hESCs

Super-enhancers were identified in naive or primed hESCs using ROSE (https://bitbucket.org/young_computation/rose). This code is an implementation of the method used in (Hnisz et al., 2013; Loven et al., 2013). Briefly, regions enriched for H3K27ac signal were identified using MACS. These regions were stitched together if they were within 12.5 kb of each other and enriched regions entirely contained within +/- 2 kb from a TSS were excluded from stitching. Stitched regions were ranked by H3K27ac signal therein. ROSE identified a point at which the two classes of enhancers were separable. Those stitched enhancers falling above this threshold were considered super-enhancers.

SMC1 binding Enrichment Heatmap

The heatmaps show the average ChIP-seq or ChIA-PET read density (r.p.m./bp) of different factors at SMC1 occupied regions. Individual ChIP datasets were processed separately and peaks of enriched signal were identified as described above. For SMC1, the genome was binned into 50 bp bins and read density of signal is shown for the 10 kb region representing +/- 5 kb from the center of each SMC1-enriched region. Similar read density of signal is shown for each other factor at the corresponding regions shown for the SMC1 dataset.

Heatmap Representation of High-confidence ChIA-PET Interactions

ChIA-PET interaction signals relative to the boundaries of CTCF-CTCF loops were mapped in a distance-normalized fashion. For each CTCF-CTCF loop, we demarcated three regions: loop, upstream, and downstream. For the loop region, the region was divided into 50 equally sized bins. For the upstream region, we selected a region extending upstream of the loop itself. The upstream region's length was set at 20% of the length of the corresponding loop. The upstream region was then divided into 10 equally sized bins. Similarly, for the downstream region, we selected a region extending downstream from the loop for a distance corresponding to 20% of the length of the loop itself, and divided the region into 10 equally sized bins.

To see whether interactions originating within the loop were generally confined within the loop, we first filtered high-confidence interactions in two ways. We required high-

confidence interactions to have at least one end in the interrogated region. This removed interactions where both endpoints of the interaction were anchored outside of the region of interest. We removed interactions that had one end at a domain border PET peak and the other end outside of the domain. This removed interactions that originated at a border and had no end within the domain as we did not consider them to be originating within the domain.

The density of the genomic space covered by ChIA-PET interactions in each bin was next calculated as the number of interactions per bin. Interactions within CTCF-CTCF loops were considered. The density of ChIA-PET interactions was row-normalized to the row maximum for each domain and the normalized frequency was displayed. Interactions connecting enhancers and promoters were considered and displayed. The density of ChIA-PET interactions was row-normalized to the row maximum for each domain and the normalized frequency was displayed.

Differential H3K27ac Signal at Enhancer Clusters Between Naive and Primed hESCs

Enhancer clusters were generated to compare enhancer regions between naive and primed hESCs. We first identified the sets of enhancer clusters in naive and primed hESCs using ROSE (https://bitbucket.org/young_computation/rose). Briefly, regions enriched for H3K27ac signal were identified using MACS. These regions were stitched together if they were within 12.5 kb of each other and enriched regions entirely contained within +/- 2 kb from a TSS were excluded from stitching. Enhancer cluster regions from naive and primed hESCs that overlapped by 1 bp were then merged together to form a representative region that spans the combined genomic region. A total of 24,755 enhancer cluster regions were identified. For each region, the read density in reads per million per base pair (r.p.m./bp) from the replicate data (2 replicate H3K27ac ChIP-seq datasets in naive hESCs and 2 replicate H3K27ac ChIP-seq datasets in primed hESCs) was calculated, and from this the relative read count of each region was obtained by multiplying read density by the length of the region. The edgeR package was used to model technical variation due to noise among duplicate data sets and the biological variation due to differences in signal between naive and primed hESCs (Robinson et al., 2010). Sequencing depth and upper- quartile techniques were used to normalize all 4 datasets together before common and tagwise dispersions were estimated. The statistical significance of differences between naive and primed hESCs was next calculated using an exact test and resulting p values were subjected to Benjamini–Hochberg multiple testing correction (FDR). The final regions with differential H3K27ac signal were required to have the absolute log₂ fold change of normalized H3K27ac signal greater or equal to 2 and FDR less or equal to 0.05.

Fold Change of H3K27ac Signal at Super-Enhancer Clusters

In order to quantify the signal changes of super-enhancers between naive and primed hESCs, H3K27ac ChIP-Seq signal was calculated at the set of all enhancer cluster regions considered as super-enhancers in at least one condition. Sequencing depth and upper-quartile techniques were used to normalize the H3K27ac ChIP-Seq signal at these super-enhancer clusters using normalization factors derived from the total 24,755 enhancer cluster regions described above. The log₂ fold change of normalized H3K27ac signal was displayed.

Saturation Analysis of ChIA-PET Library

To determine the degree of saturation within our ChIA-PET library, we modeled the number of sampled putative interactions, which were defined as PETs that overlapped with two PET peaks at both ends by at least 1 bp, as a function of sequencing depth by a two parameter logistic growth model. Intrachromosomal PETs were subsampled at varying depths, and the number of unique putative interactions that they occupied were counted. Model fitting using non-linear least-squares regression suggested that we sampled approximately 45~50 % of the available intrachromosomal PET space.

Calling Hi-C Interactions

The Fit-Hi-C tool (Ay et al., 2014) was used to call high-confidence DNA interactions from Hi-C datasets in H1 hESCs (Dixon et al., 2015). The raw fragment contact and bias matrices at 40 kb resolution were first obtained using the python hiclib library. The Fit-Hi-C was then used to call high-confidence DNA interactions using the raw fragment contact and bias matrices at 40 kb resolution with the parameters: -L 50,000 -U 5,000,000 -b 200 -p 1 --quiet. The CTCF-CTCF Hi-C interactions were identified by filtering for those Hi-C interactions that have CTCF ChIP-seq peaks within the 40 kb bins at the both ends of the interactions. A Hi-C CTCF-CTCF interaction was classified as “overlapped with a ChIA-PET CTCF-CTCF interaction” if both ends of the Hi-C CTCF-CTCF interaction overlapped with the ends of a ChIA-PET CTCF-CTCF interaction by at least 1 bp. The percentages of Hi-C CTCF-CTCF interactions that overlapped with ChIA-PET CTCF-CTCF interactions (or vice versa) were displayed as line plots.

Comparisons to In-situ Hi-C CTCF-CTCF Interactions

Previously published in-situ Hi-C CTCF-CTCF interactions with CTCF DNA motifs in 7 different cell types (Rao et al., 2014) were first downloaded from GEO (www.ncbi.nlm.nih.gov/geo/; accession GSE63525). They were next compared to the ChIA-PET CTCF-CTCF interactions in naive and primed hESCs. We tested how often these in-situ Hi-C CTCF-CTCF interactions overlapped with the CTCF-CTCF interactions in naive and primed hESCs. Since the in-situ Hi-C CTCF-CTCF interactions were identified by requiring CTCF ChIP-seq peaks within a +/- 15 kb window at the both ends of the interactions in the publication (Rao et al., 2014), an in-situ Hi-C CTCF-CTCF interaction was classified as “overlapped with a ChIA-PET CTCF-CTCF interaction” if both ends of the in-situ Hi-C CTCF-CTCF interaction overlapped with the ends of a ChIA-PET CTCF-CTCF interaction within a +/- 15kb window. The percentages of in-situ Hi-C CTCF-CTCF interactions that overlapped with ChIA-PET CTCF-CTCF interactions were displayed as bar plots.

Calling High-Confidence Cell-Type-Specific CTCF-CTCF Interactions in Naive And Primed hESCs

To identify the naive-specific or primed-specific CTCF-CTCF interactions, we took advantage of the strong signal at the PET peaks to increase the confidence to interpret the ChIA-PET interaction data. This was because the PET counts or the ChIP-seq read counts at PET peaks were frequently an order magnitude higher than the PET count for the number of PETs spanning high-confidence interactions allowing for better dynamic range. Briefly, we applied a negative binomial statistical model from the edgeR package to identify differentially occupied CTCF peaks between naive and primed hESCs using ChIP-seq data (FDR 0.01 and absolute log₂ fold change >= 2) and overlaid these differential ChIP-seq CTCF regions to the CTCF-CTCF ChIA-PET interactions from naive and primed hESCs.

To identify differentially occupied CTCF peaks between naive and primed hESCs, CTCF ChIP-seq peaks from naive and primed hESCs that overlapped by 1 bp were then merged together to form a representative region that spans the combined genomic region. For each region, the read density in reads per million per base pair (r.p.m./bp) from the replicate data (2 replicate CTCF ChIP-seq datasets in naive hESCs and 2 replicate CTCF ChIP-seq datasets in primed hESCs) was calculated, and from this the relative read count of each region was obtained by multiplying read density by the length of the region. The edgeR package was used to model technical variation due to noise among duplicate data sets and the biological variation due to differences in signal between naive and primed hESCs (Robinson et al., 2010). Sequencing depth and upper-quartile techniques were used to normalize all 4 datasets together before common and tagwise dispersions were estimated. The statistical significance of differences between naive and primed hESCs was next calculated using an exact test and resulting p values were subjected to Benjamini–Hochberg multiple testing correction (FDR). The final regions with differential CTCF signal were required to have the absolute log₂ fold change of normalized CTCF signal greater or equal to 2 and FDR less or equal to 0.05. This analysis resulted in 313 naive-specific CTCF peaks and 75 primed-specific CTCF peaks.

We next identified the CTCF-CTCF interactions that were associated with these preferentially occupied CTCF peaks in naive and primed hESCs by requiring at least one end of the CTCF-CTCF interactions overlapped with the preferentially occupied CTCF peaks. To obtain the high-confidence cell-type-specific CTCF-CTCF interactions, we also required the naive-specific CTCF-CTCF interactions that overlapped naive-specific CTCF peaks to have zero PETs in primed hESCs, and primed-specific CTCF-CTCF interactions that overlapped primed-specific CTCF peaks to have zero PETs in naive hESCs. This resulted in only 125 naive-specific CTCF-CTCF interactions and 28 primed-specific CTCF-CTCF interactions.

Topologically Associating Domain (TAD) Calling

TADs were determined from interaction matrices using the method and code previously described in (Dixon et al., 2012). For cohesin ChIA-PET-based TADs, ChIA-PET interactions were used to generate interaction matrices by binning the genome into 40 kb bins and counting the number of PETs connecting any two bins. For H1 hESC Hi-C based TADs, H1 hESC Hi-C data previously generated in (Dixon et al., 2015), was realigned, binned into 40 kb bins, and normalized to generate a Hi-C interaction matrix. Parameters from Dixon et al. were retained (an interaction window of 2 Mb and 40 kb for binning interactions). For human samples, the human reference genome (build hg19, GRCh37) was used and for mouse samples, the mm9 mouse reference genome was used.

Hi-C vs ChIA-PET Interaction Comparison

Hi-C data was examined to see if the Hi-C data supported predicted ChIA-PET interactions. To do this, H1 hESC Hi-C data was first processed to create an interaction matrix as described above. The subset of the Hi-C interaction matrix that could be directly compared to the available ChIA-PET data was then selected. The interaction scores from the Hi-C matrix were then plotted as a box plot. For comparison, a random distribution of Hi-C interactions was generated and also plotted.

TAD Spanning Loops: Percentage and Visualization

TADs derived from Hi-C data from H1 hESCs were examined for the presence of CTCF-CTCF loops that spanned the entire TAD. TADs and Hi-C interactions were derived as described above. For each TAD, we queried if there was at least one CTCF-CTCF loop that connected the upstream and downstream boundaries of the TAD. For this analysis, each boundary was extended by 40 kb both upstream and downstream. A loop was considered spanning if one end was found in the upstream boundary and the other end was found in the downstream boundary. We examined TADs for the number of spanning loops that connected the two boundaries; the percentages of TADs with 1, 2 or 3 spanning loops were reported. For comparison, the analysis was repeated using a set of randomized, shuffled TADs. For the shuffled set, we used the set of H1 Hi-C based TADs but shuffled the chromosome and start site coordinates. Visualization of spanning loops was done using the CRAN-Circlize package (<http://cran.r-project.org/web/packages/circlize/index.html>).

TAD Boundary Overlap

To compare the consistency of TADs called using either Hi-C or ChIA-PET data, we asked if boundaries of Hi-C based TAD were frequently co-localized with boundaries of ChIA-PET based TAD calls. To do this, we examined the overlap of the boundaries of TADs called using ChIA-PET data and Hi-C data. For each boundary, we measured the distance of each Hi-C called TAD boundary to the nearest ChIA-PET called TAD boundary. The distribution of distances was then plotted in a histogram.

Conservation and Disease Analysis

We examined whether the ends of CTCF-CTCF loops overlapped with genomic regions of high sequence conservation or genomic regions associated with disease-causing mutations. We began by identifying the CTCF motifs (as described above) that were within the anchor sites of high confidence CTCF-CTCF ChIA-PET interactions. We considered two sets of regions, the first being CTCF-CTCF anchor sites and the second being CTCF motif sites that are bound by CTCF and within loop anchor sites. For conservation analysis, the 10 kb of sequence around the midpoint of each CTCF-CTCF anchor site (± 5 kb) was used. For each region, for each base pair, the PhastCons score was determined using a 10 way primate multiple alignment (Pollard et al., 2010). This created a vector of PhastCons scores for each region. The vectors for all regions were then averaged and plotted. For association with cancer mutations, the regions described above were overlapped with the coordinates of simple somatic mutations present in cancer from the International Cancer Genome Consortium (ICGC) database (Zhang et al., 2011). For each base pair, the base pair was scored for presence of a mutation. This created a vector of mutation occurrences for each region. The vectors for all regions were then summed and plotted. For association with disease mutations, the regions described above were overlapped with the coordinates of GWAS SNPs (Welter et al., 2014). GWAS SNPs that fell within these regions were reported. All of the analyses were repeated for CTCF motifs. Here, the sequences analyzed included the motif itself plus 200 bp of sequence upstream and downstream.

GWAS Catalog Parsing and Distance Distribution.

The NHGRI Genome-Wide Association Study (GWAS) database containing SNPs significantly associated with human traits was downloaded 6/19/2015 and parsed as described in (Hnisz et al., 2013). Briefly, trait-associated SNPs with dbSNP identifiers were reproducibly associated with a trait in two independent studies. SNPs were assigned a genomic position using dbSNP build 142. SNPs falling inside RefSeq coding exons were discarded. The distance distribution of trait-associated, noncoding SNPs to

the nearest border of a region in the union of 86 enhancer sets defined in (Hnisz et al., 2013) were shown. The distance distribution of trait-associated noncoding SNPs to the nearest border of a CTCF anchor in the union of the naive and primed anchor sites were shown. SNPs within these regions were assigned to the 0 bin.

Fractional Methylation Analysis at CTCF Sites

We examined the methylation dynamics of CTCF motifs within CTCF-CTCF high-confidence loop anchor sites throughout early embryonic development by using reduced representation bisulfite (RRBS) sequencing from human preimplantation embryos (Smith et al., 2014). The CTCF motifs were flanked by 4991 bp to create a 10 kb region around each CTCF motif. Each 10 kb region was then overlapped with the fractional methylation RRBS data, this generated a set of genomic locations that fell within these 10 kb regions and had fractional methylation values assigned to them. Each of the fractional methylation values was then added to the correct location within the vector of 10000 values relative to the CTCF motif. This created one 10000 value average fractional methylation vector for five of the samples in Figure S6D (sperm, 8-cell embryo, inner cell mass, hESC primed, and fetal lung). This vector was smoothed over using a 10 bp window and plotted as a line plot.

A similar analysis was performed using whole genome bisulfite sequencing from (http://egg2.wustl.edu/roadmap/web_portal/processed_data.html#MethylData). We first averaged the fractional methylation across 37 cell/tissue types. This generated one average fractional methylation value for each base pair in the genome. These values were overlapped as described above with the same 10 kb regions and plotted in the same manner. Also, the fractional methylation plot for the adult lung sample in Figure S6D was generated in the same way as described for the other five samples but used data from WUSTL.

Transcription Factor Motif and Mutation Analysis within CTCF-CTCF Loop Anchors

We determined the average number of mutations found in occurrences of transcription factor motifs that occur in anchor regions. We first downloaded a set of motif instances from (Kheradpour et al., 2013) consisting of sequence motifs, their assignment to transcription factors and their chromosomal location. We next filtered for those motif instances within anchor regions. For each member of the resulting set of motif instances, we counted how many cancer mutations overlapped the motif instance. The counts for all instances assigned to a given factor were summed and divided by the number of instances assigned to that factor. The simple somatic mutations present in cancer were described in the International Cancer Genome Consortium database (Zhang et al., 2011).

ICGC Simple Somatic Mutations within the Loop Anchors

A VCF file of simple somatic mutations was downloaded from the International Cancer Genome Consortium database (Zhang et al., 2011). The file was filtered for mutation calls generated from projects that are not under embargo (project key PACA-AU, PACA-CA, PRAD-CA, BLCA-CN, GACA-CN, LICA-FR, EOPC-DE, MALY-DE, PBCA-DE, LINC-JP, LIRI-JP, LUSC-KR, CLLE-ES, BRCA-UK, CMDI-UK, PRAD-UK, and OV-AU). We first identified all CTCF-CTCF loops for which there are simple mutations overlapping with CTCF motifs within the loop anchors by at least 1 bp for CTCF-CTCF loops in naive and primed hESCs. We next examined the genes that were contained within these CTCF-CTCF loops by requiring that their TSSs be within the loop. The gene

symbols were cross-referenced to Refseq genes, cancer census genes, proto-oncogenes, and tumor suppressor genes. Cancer census genes (Version 73) were downloaded from the COSMIC database (www.cancer.sanger.ac.uk/cosmic). Proto-oncogenes are operationally defined as genes from the cancer census gene list whose mutations result in a dominant phenotype (Bishop et al., 1991). Tumor suppressor genes were downloaded from the TSGene Tumor suppressor gene database (http://bioinfo.mc.vanderbilt.edu/TSGene/Human_716_TSGs.txt) (Zhao et al., 2012).

REFERENCES

- Ay, F., Bailey, T.L., and Noble, W.S. (2014). Statistical confidence estimation for Hi-C data reveals regulatory chromatin contacts. *Genome Res* 24, 999-1011.
- Benjamini, Y., and Hochberg, Y. (1995). Controlling the false discovery rate: a practical and powerful approach to multiple testing. *Journal of the Royal Statistical Society Series B-Methodological* 57, 289-300.
- Bolland, D.J., King, M.R., Reik, W., Corcoran, A.E., and Krueger, C. (2013). Robust 3D DNA FISH using directly labeled probes. *J Vis Exp*.
- Burman, B., Zhang, Z.Z., Pegoraro, G., Lieb, J.D., and Misteli, T. (2015). Histone modifications predispose genome regions to breakage and translocation. *Genes Dev* 29, 1393-1402.
- Dixon, J.R., Jung, I., Selvaraj, S., Shen, Y., Antosiewicz-Bourget, J.E., Lee, A.Y., Ye, Z., Kim, A., Rajagopal, N., Xie, W., et al. (2015). Chromatin architecture reorganization during stem cell differentiation. *Nature* 518, 331-336.
- Dixon, J.R., Selvaraj, S., Yue, F., Kim, A., Li, Y., Shen, Y., Hu, M., Liu, J.S., and Ren, B. (2012). Topological domains in mammalian genomes identified by analysis of chromatin interactions. *Nature* 485, 376-380.
- Dowen, J.M., Fan, Z.P., Hnisz, D., Ren, G., Abraham, B.J., Zhang, L.N., Weintraub, A.S., Schuijers, J., Lee, T.I., Zhao, K., et al. (2014). Control of Cell Identity Genes Occurs in Insulated Neighborhoods in Mammalian Chromosomes. *Cell* 159, 374-387.
- Grant, C.E., Bailey, T.L., and Noble, W.S. (2011). FIMO: scanning for occurrences of a given motif. *Bioinformatics* 27, 1017-1018.
- Heidari, N., Phanstiel, D.H., He, C., Grubert, F., Jahanbani, F., Kasowski, M., Zhang, M.Q., and Snyder, M.P. (2014). Genome-wide map of regulatory interactions in the human genome. *Genome Res* 24, 1905-1917.
- Hnisz, D., Abraham, B.J., Lee, T.I., Lau, A., Saint-Andre, V., Sigova, A.A., Hoke, H.A., and Young, R.A. (2013). Super-Enhancers in the Control of Cell Identity and Disease. *Cell* 155, 934-947.
- Imakaev, M., Fudenberg, G., McCord, R.P., Naumova, N., Goloborodko, A., Lajoie, B.R., Dekker, J., and Mirny, L.A. (2012). Iterative correction of Hi-C data reveals hallmarks of chromosome organization. *Nature Methods* 9, 999-1003.
- Ji, X., Dadon, D.B., Abraham, B.J., Lee, T.I., Jaenisch, R., Bradner, J.E., and Young, R.A. (2015). Chromatin proteomic profiling reveals novel proteins associated with histone-marked genomic regions. *Proceedings of the National Academy of Sciences of the United States of America* 112, 3841-3846.
- Kent, W.J., Sugnet, C.W., Furey, T.S., Roskin, K.M., Pringle, T.H., Zahler, A.M., and Haussler, D. (2002). The human genome browser at UCSC. *Genome Research* 12, 996-1006.
- Kheradpour, P., Ernst, J., Melnikov, A., Rogov, P., Wang, L., Zhang, X., Alston, J., Mikkelsen, T.S., and Kellis, M. (2013). Systematic dissection of regulatory motifs in 2000

predicted human enhancers using a massively parallel reporter assay. *Genome Res* 23, 800-811.

Langmead, B., Trapnell, C., Pop, M., and Salzberg, S.L. (2009). Ultrafast and memory-efficient alignment of short DNA sequences to the human genome. *Genome Biology* 10.

Li, G., Fullwood, M.J., Xu, H., Mulawadi, F.H., Velkov, S., Vega, V., Ariyaratne, P.N., Bin Mohamed, Y., Ooi, H.-S., Tennakoon, C., *et al.* (2010). ChIA-PET tool for comprehensive chromatin interaction analysis with paired-end tag sequencing. *Genome Biology* 11.

Loven, J., Hoke, H.A., Lin, C.Y., Lau, A., Orlando, D.A., Vakoc, C.R., Bradner, J.E., Lee, T.I., and Young, R.A. (2013). Selective inhibition of tumor oncogenes by disruption of super-enhancers. *Cell* 153, 320-334.

Martin, M. (2011). Cutadapt removes adapter sequences from high-throughput sequencing reads. *EMBnet journal* 17, 1-10.

Matys, V., Kel-Margoulis, O.V., Fricke, E., Liebich, I., Land, S., Barre-Dirrie, A., Reuter, I., Chekmenev, D., Krull, M., Hornischer, K., *et al.* (2006). TRANSFAC (R) and its module TRANSCompel (R): transcriptional gene regulation in eukaryotes. *Nucleic Acids Research* 34, D108-110.

Mizuguchi, T., Fudenberg, G., Mehta, S., Belton, J.-M., Taneja, N., Folco, H.D., FitzGerald, P., Dekker, J., Mirny, L., Barrowman, J., *et al.* (2014). Cohesin-dependent globules and heterochromatin shape 3D genome architecture in *S. pombe*. *Nature* 516, 432-435.

Nora, E.P., Lajoie, B.R., Schulz, E.G., Giorgetti, L., Okamoto, I., Servant, N., Piolot, T., van Berkum, N.L., Meisig, J., Sedat, J., *et al.* (2012). Spatial partitioning of the regulatory landscape of the X-inactivation centre. *Nature* 485, 381-385.

Phanstiel, D.H., Boyle, A.P., Heidari, N., and Snyder, M.P. (2015). Mango: a bias-correcting ChIA-PET analysis pipeline. *Bioinformatics* 31, 3092-3098.

Pollard, K.S., Hubisz, M.J., Rosenbloom, K.R., and Siepel, A. (2010). Detection of nonneutral substitution rates on mammalian phylogenies. *Genome Research* 20, 110-121.

Rao, S.S.P., Huntley, M.H., Durand, N.C., Stamenova, E.K., Bochkov, I.D., Robinson, J.T., Sanborn, A.L., Machol, I., Omer, A.D., Lander, E.S., *et al.* (2014). A 3D Map of the Human Genome at Kilobase Resolution Reveals Principles of Chromatin Looping. *Cell* 159, 1665-1680.

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.

Shachar, S., Voss, T.C., Pegoraro, G., Sciascia, N., and Misteli, T. (2015). Identification of Gene Positioning Factors Using High-Throughput Imaging Mapping. *Cell* 162, 911-923.

Theunissen, T.W., Powell, B.E., Wang, H., Mitalipova, M., Faddah, D.A., Reddy, J., Fan, Z.P., Maetzel, D., Ganz, K., Shi, L., *et al.* (2014). Systematic Identification of Culture Conditions for Induction and Maintenance of Naive Human Pluripotency. *Cell Stem Cell* 15, 471-487.

Trapnell, C., Pachter, L., and Salzberg, S.L. (2009). TopHat: discovering splice junctions with RNA-Seq. *Bioinformatics* 25, 1105-1111.

Trapnell, C., Williams, B.A., Pertea, G., Mortazavi, A., Kwan, G., van Baren, M.J., Salzberg, S.L., Wold, B.J., and Pachter, L. (2010). Transcript assembly and quantification by RNA-Seq reveals unannotated transcripts and isoform switching during cell differentiation. *Nature Biotechnology* 28, 511-515.

Welter, D., MacArthur, J., Morales, J., Burdett, T., Hall, P., Junkins, H., Klemm, A., Flicek, P., Manolio, T., Hindorff, L., *et al.* (2014). The NHGRI GWAS Catalog, a curated resource of SNP-trait associations. *Nucleic Acids Research* 42, D1001-1006.

Zhang, J., Baran, J., Cros, A., Guberman, J.M., Haider, S., Hsu, J., Liang, Y., Rivkin, E., Wang, J., Whitty, B., *et al.* (2011). International Cancer Genome Consortium Data Portal-a one-stop shop for cancer genomics data. *Database-the Journal of Biological Databases and Curation*.

Zhang, Y., Liu, T., Meyer, C.A., Eeckhoute, 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 Biology* 9.