

Supplemental Materials (Beagan et al.)

Supplemental Figures

Figure S1 (related to Figure 1): Progression of 5C data through analysis pipeline.

Figure S2 (related to Figure 1): Progression of 5C data through alternative 5C analysis approaches.

Figure S3 (related to Figure 3): Methodology for identification of significant 3-D interaction classes.

Figure S4 (related to Figures 2,4,5,6): RNA-seq library normalization and quality control.

Figure S5 (related to Figures 4,5): The Klf4 gene engages in both ES-iPS (purple class) and NPC-iPS (blue class) 3-D interactions.

Figure S6 (related to Figure 5): NPC-specific genes and enhancers are enriched in NPC only (green class) interactions.

Figure S7 (related to Figure 6): The Mis18 and Urb1 genes engage in ES only (red class) 3-D interactions linked to inadequately reprogrammed, ES-specific CTCF binding.

Supplemental Figure Captions

Supplemental Tables

Table S1: Summary of paired-end read alignments for 5C libraries, related to Experimental Procedures and Supplemental Experimental Procedures.

Table S2: Spearman's rank correlation coefficients calculated for distance-corrected interaction frequencies of pairs of biological replicates, related to Experimental Procedures.

Table S3: Summary of paired-end read alignments for RNAseq libraries, related to Experimental Procedures and Supplemental Experimental Procedures.

Table S4: Summary of external ChIP-seq libraries analyzed in this study, related to Experimental Procedures and Supplemental Experimental Procedures.

Table S5: ES-specific genes, Related to Figures 4, 5, 6 and Supplemental Experimental Procedures.

Table S6: NPC-specific genes, Related to Figures 4, 5, 6 and Supplemental Experimental Procedures.

Table S7. Interactions selected for representative interaction score barplots, Related to Figures 4, S5, 5, 6, S7.

Supplemental Experimental Procedures

ES cell culture

Primary Neural Progenitor Cell isolation

iPS cell culture

Culture of pluripotent cells in 2i media

3C template generation and characterization

5C primer design

5C library generation and sequencing

iPS cell transgene integration detection by 5C primers

RNA-seq library preparation

RNA-seq data processing

CTCF binding enrichment by ChIP-qPCR

5C data processing pipeline

Paired-end read mapping and counting

Low count primer removal

Raw contact matrix visualization

Quantile normalization

Primer correction

Low count fragment-fragment pair removal

Contact matrix binning

Pseudo-fragment level 5C mapping resolution

Identification of bad primer gaps

Distance-dependence normalization

Probabilistic model fitting and distance-corrected interaction scores

GC content bias investigation

Comparison of 5C analysis pipeline to alternative approaches

Principal component analysis

Classification of cell type-specific 3-D interactions

Empirical false discovery rate calculation

Justification of strategy

Model generation – mean parameter estimation

Model generation – estimating the mean-variance relationship

Model generation – variance parameter estimation

Simulations

Monte Carlo, p-value calculation, classification

Computing the false discovery rates for each 3-D interaction class

6 Sample vs 10 Sample 5C data processing

Interaction adjacency clustering

ChIPseq peakcalling

Parsing ES-specific and NPC-specific genes

Parsing ES-specific and NPC-specific enhancers

Parsing ES-specific and NPC-specific CTCF sites

Computing enrichments

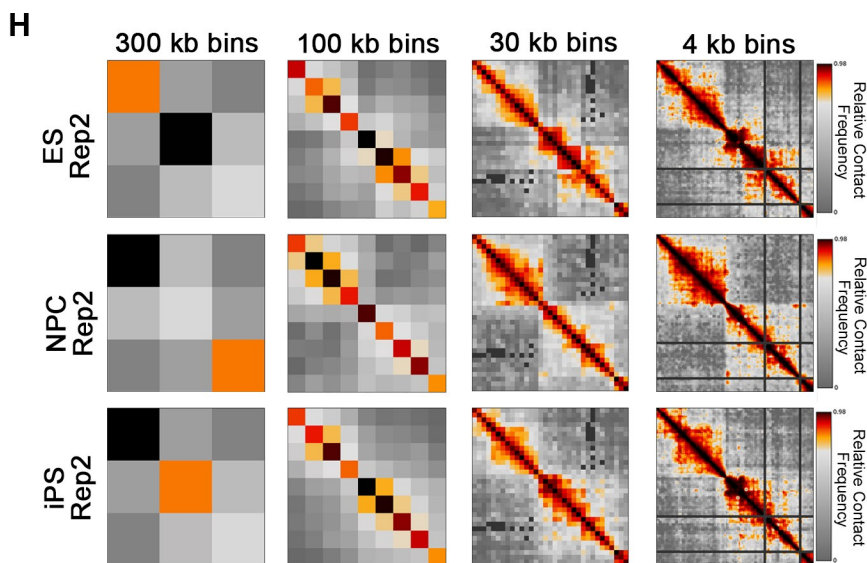
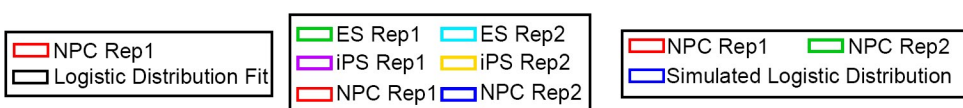
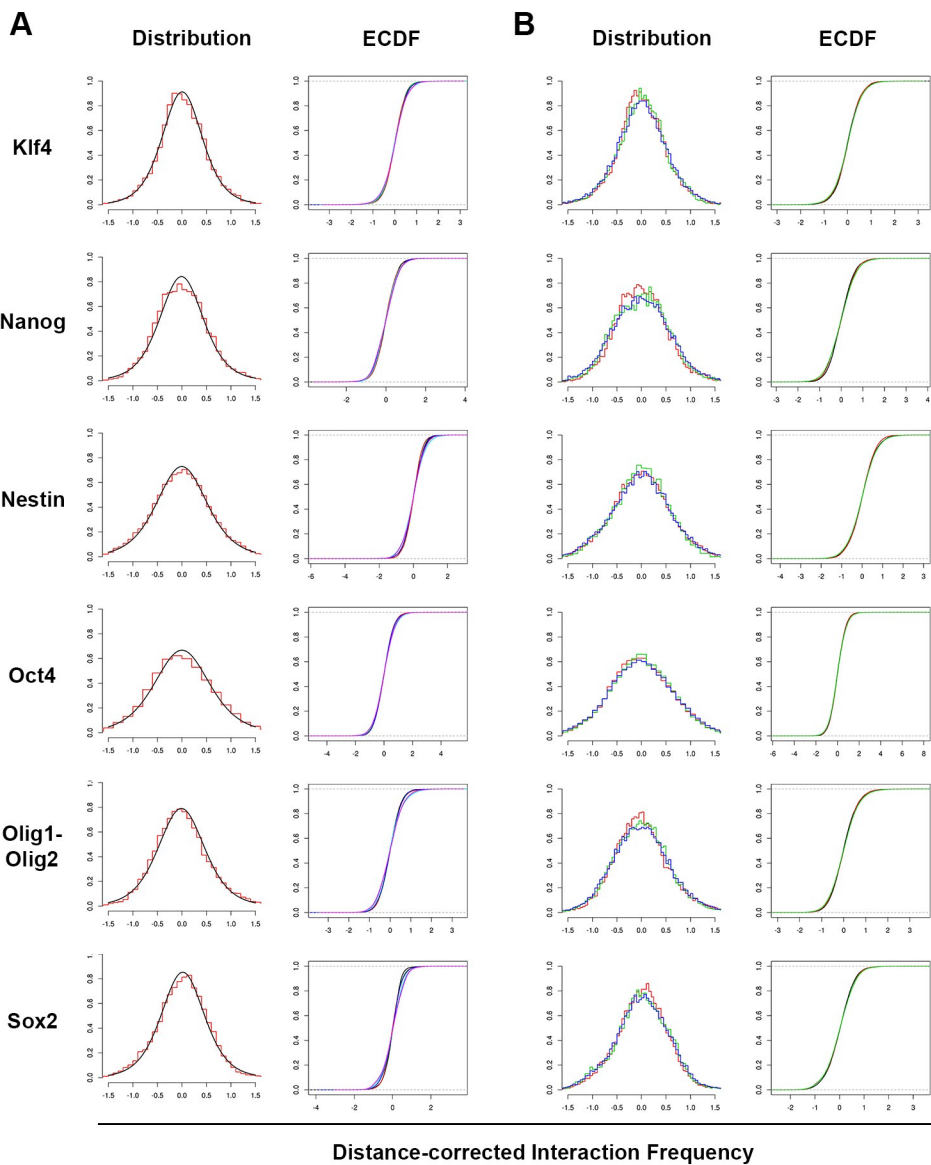
- Annotation intersections

- Computing percentage incidence, fold-enrichment above background, and p-values

- Visualizing enrichments

Computing connectivity

Supplementary References



I Spearman's ρ , 300 kb bins

ES Rep2	1		
NPC Rep2	0.556	1	
iPS Rep2	0.874	0.678	1
	ES Rep2	NPC Rep2	iPS Rep2

Spearman's ρ , 100 kb bins

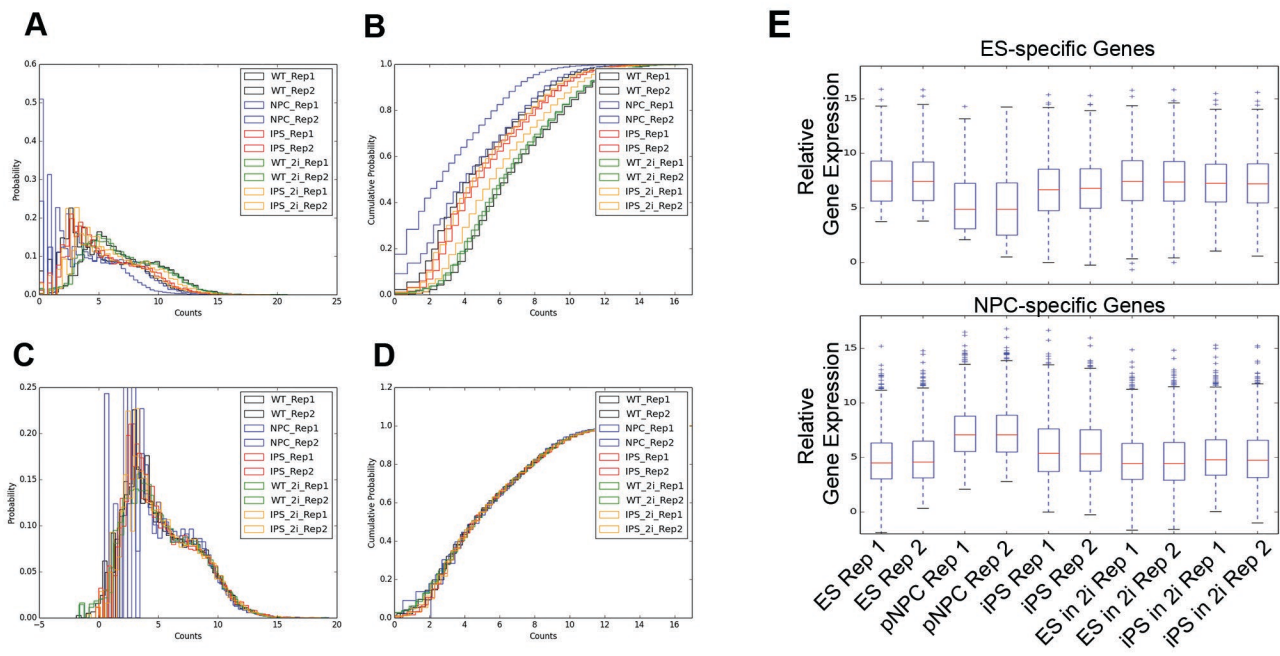
ES Rep2	1		
NPC Rep2	0.423	1	
iPS Rep2	0.748	0.595	1
	ES Rep2	NPC Rep2	iPS Rep2

Spearman's ρ , 30 kb bins

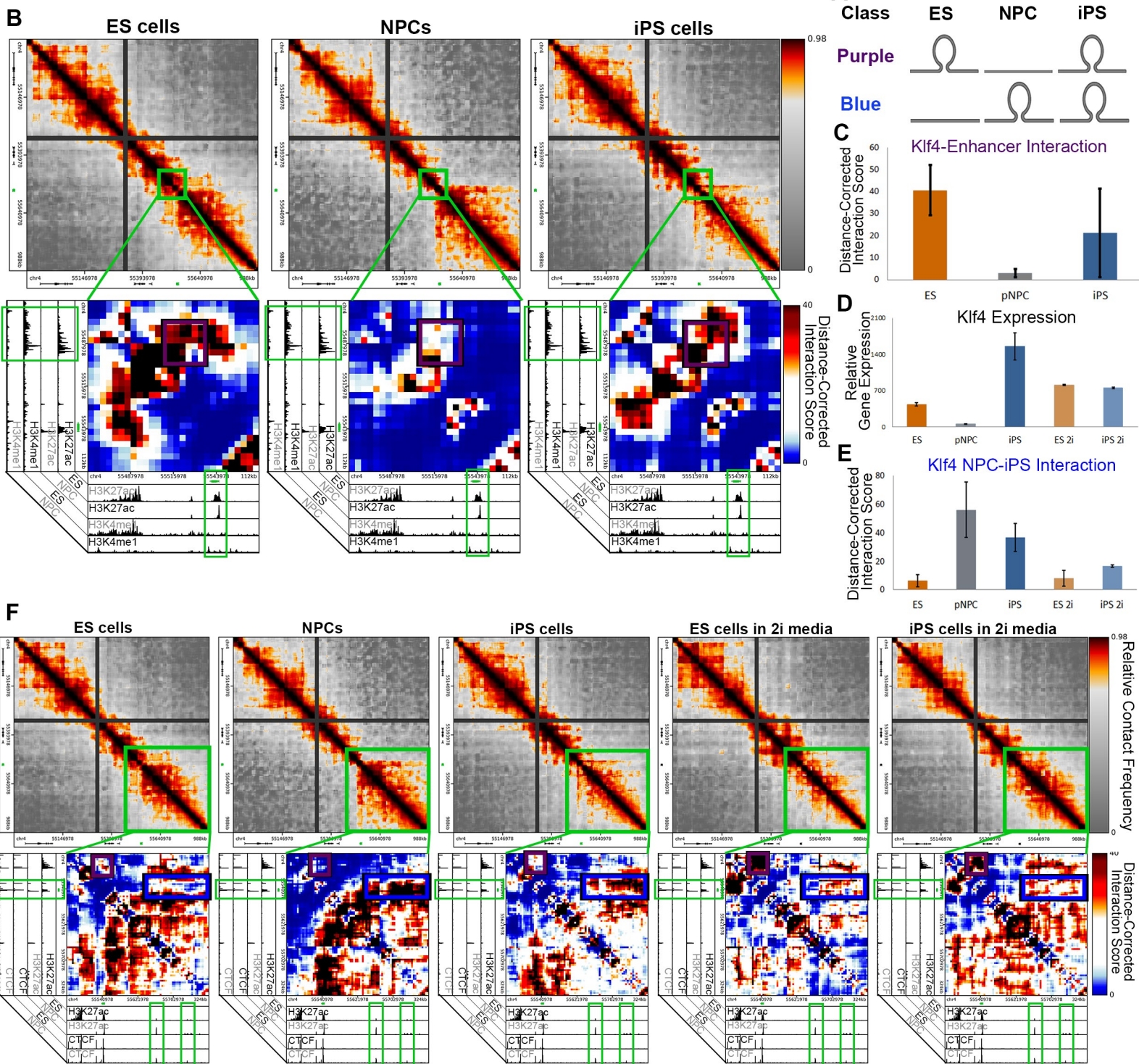
ES Rep2	1		
NPC Rep2	0.332	1	
iPS Rep2	0.521	0.526	1
	ES Rep2	NPC Rep2	iPS Rep2

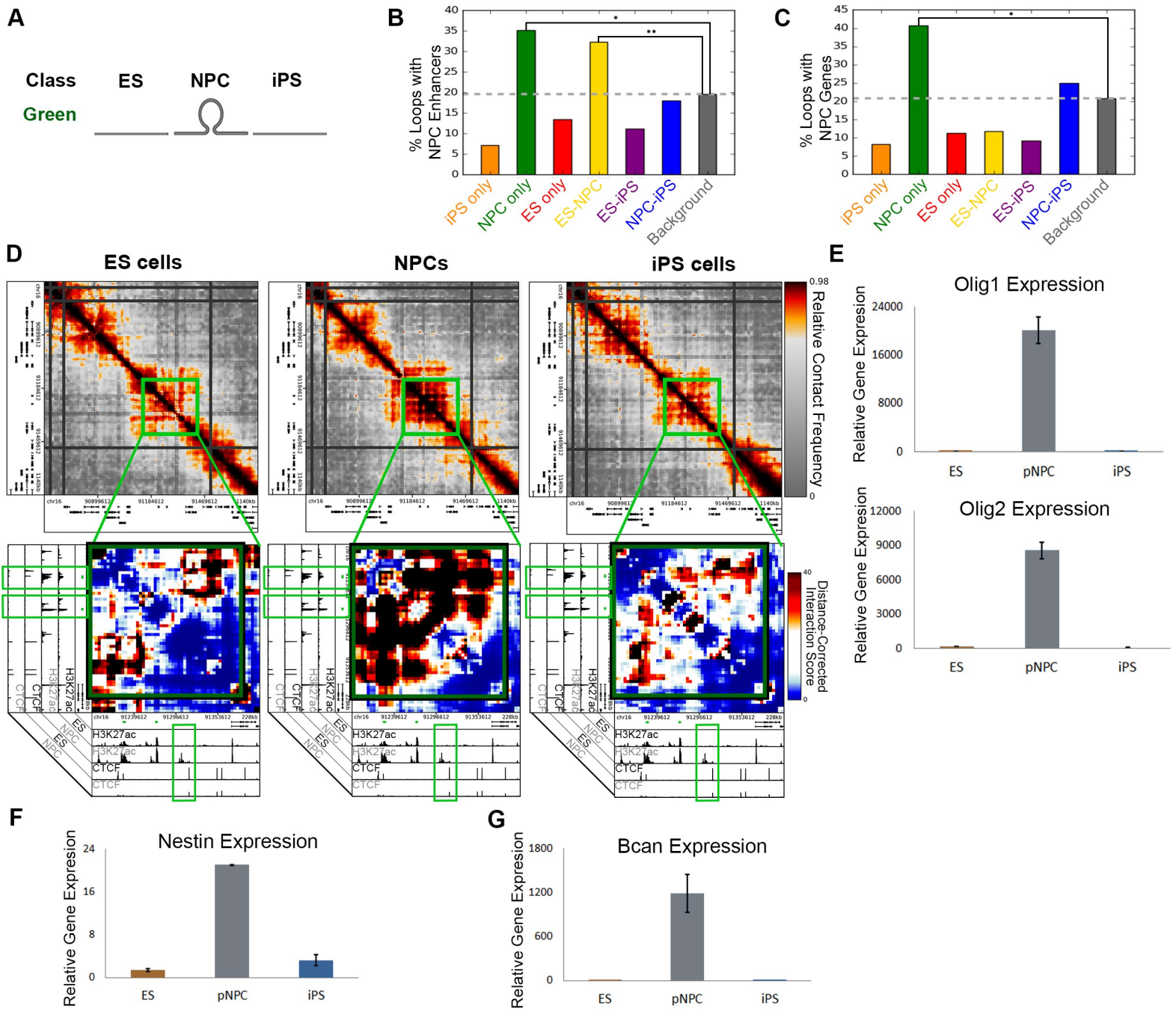
Spearman's ρ , 4 kb bins

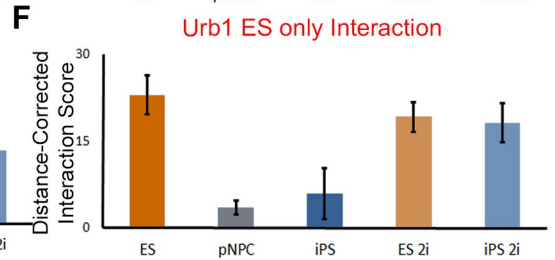
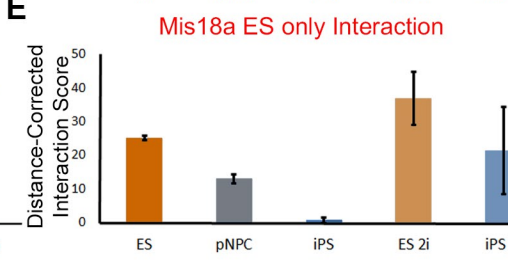
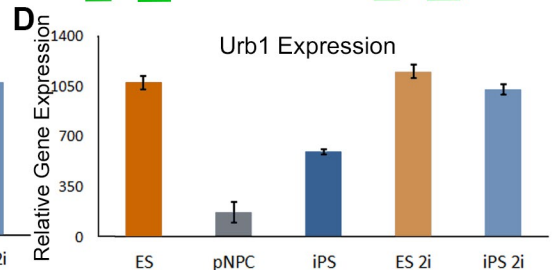
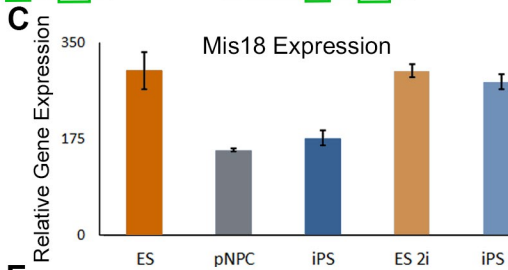
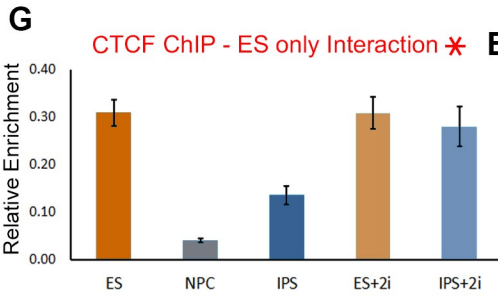
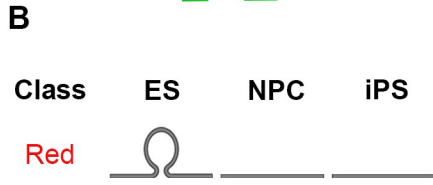
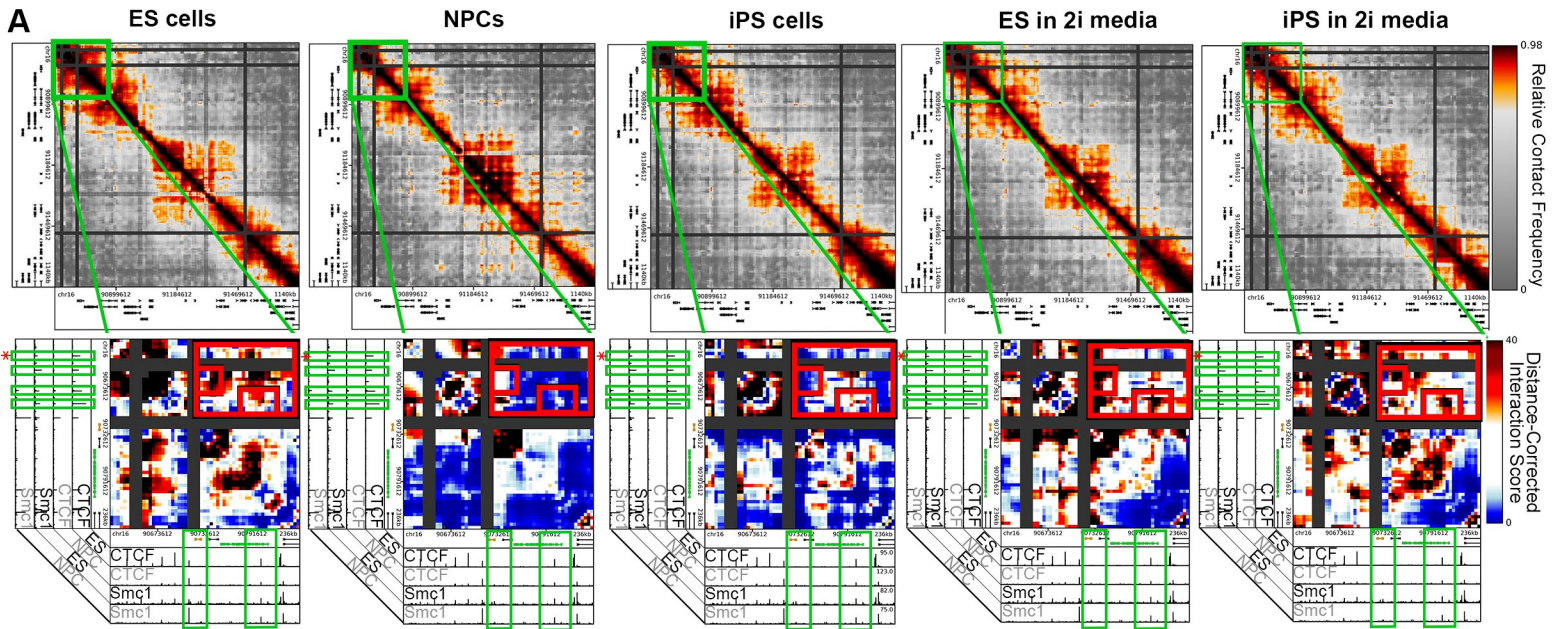
ES Rep2	1		
NPC Rep2	0.268	1	
iPS Rep2	0.427	0.493	1
	ES Rep2	NPC Rep2	iPS Rep2



Beagan et al. 2016 - Figure S4







Supplemental Figure Captions

Figure S1 (related to Figure 1). Progression of 5C data through analysis pipeline. (A-F) Grid showing progression of Sox2 region through data processing steps. From top to bottom: **(A)** raw, **(B)** quantile normalized, **(C)** primer corrected, **(D)** binned (4 kb bins; 20 kb smoothing window), **(E)** distance-dependence corrected and **(F)** interaction score computed as $-10 \cdot \log_2(\text{p-value})$ on p-values computed from the distance-dependence normalized data after logistic distribution modeling parameterized for each genomic region. From left to right: (i) contact probability heatmaps for ES Rep1 and NPC Rep1, (ii) boxplots of counts for each primer/bin in the Sox2 region in order of increasing median, (iii) background distance-dependence interaction frequency, showing the mean of the counts at distance scales binned every 40 kb, (iv) kernel density estimates of the counts probability density. **(G)** Boxplots of 'Relative contact frequency' values at 4 kb intervals across the genomic coordinates queried for each 5C region. Plots for the Olig1-Olig2 and Nestin regions of ES Rep 1 are shown. **(H)** Violin plots showing the distribution of log fold enrichment of total cis primer counts over the mean of cis primer counts (x-axis) as a function of each primer's GC content (y-axis). Data for ES Rep 1 is shown at raw, quantile normalization and primer correction stages in the analysis pipeline. **(I)** Heatmaps comparing GC content bias in ES Rep1 in pairwise fragment-to-fragment contacts before and after primer correction. Fold enrichment is computed within each two-sided GC bin as the sum of the counts for all cis primer-primer pairs falling in the GC content range of the bin divided by the expected number of counts for a bin with that many primer-primer pairs in it (see **Supplemental Experimental Procedures**).

Figure S2 (related to Figure 1). Progression of 5C data through alternative 5C analysis approaches. (A-D) Grid showing progression of Sox2 region through our previously published analysis pipeline ([Phillips-Cremins et al., 2013](#)). From top to bottom: **(A)** raw, **(B)** primer corrected, **(C)** distance-dependence normalized via parametric model described in ([Phillips-Cremins et al., 2013](#)) and **(D)** interaction score

computed as $-10 \cdot \log_2(\text{p-value})$ on p-values computed with compound normal-lognormal distribution fits described in (Phillips-Cremins et al., 2013). From left to right: (i) contact probability heatmaps for ES Rep1 and NPC Rep1, (ii) boxplots of counts for each primer/bin in the Sox2 region in order of increasing median, (iii) distance dependence curves, showing the mean of the counts at distance scales binned every 40 kb, (iv) kernel density estimates of the counts probability density. **(E-G)** Grid showing downstream effects of alternative placement of quantile normalization step within the main 5C analysis pipeline. Primer normalized data shown in **(B)** were binned **(E)**, then quantile normalized (in contrast to Figure S1, where quantile normalization is the first step) **(F)**, and finally distance corrected **(G)**.

Figure S3 (related to Figure 3). Methodology for identification of significant 3-D interaction classes. (A-B) Histograms and empirical cumulative distribution functions (ECDF) of distance-corrected interaction frequency values. **(A)** Distributions of NPC Rep 1 (red) superimposed upon a logistic distribution fit with location/scale parameters computed for each region and biological replicate (black). Juxtaposition of models illustrates that our distance-corrected data can be modeled with logistic fits. **(B)** Distributions of the two NPC replicates (red and green) plotted alongside the simulated data distribution (blue). Simulated data closely approximate 5C data, supporting their utility in computing empirical False Discovery Rates. **(C)** Empirical false discovery rates computed from simulated data reported for each classification. FDRs vary slightly depending on which cell-type replicates are used to model parameters of the simulations (see **Supplemental Experimental Procedures**). **(D-G)** Zoomed-in contact density maps for specific **(D)** NPC only interactions (green class), **(E)** iPS only interactions (orange class), **(F)** ES-NPC interactions (yellow class), and **(G)** NPC-iPS interactions (blue class). Classified interaction pixels are outlined in green for each interaction class. **(H)** 5C primer-primer counts data are binned with decreasing bin sizes and displayed as contact density heatmaps. From left to right, heatmaps are shown for bin sizes of 300 kb, 100 kb, 30 kb and finally the 4 kb with a 20 kb smoothing window used in this

study. **(I)** Spearman's rank correlation coefficient was calculated using the distance-corrected interaction frequency data of replicates displayed in **(H)** at each bin size.

Figure S4 (related to Figures 2, 4, 5, 6). RNA-seq library normalization and quality control. (A,C) Frequency histograms of read counts across all genes for each RNA-seq library before **(A)** and after **(C)** normalization. **(B,D)** Cumulative distributions of read counts across all genes for each RNA-seq library before **(B)** and after **(D)** normalization. **(E)** Boxplots of the logged normalized counts of genes parsed as ES-specific or NPC-specific for each replicate.

Figure S5 (related to Figures 4, 5). The Klf4 gene engages in both ES-iPS (purple class) and NPC-iPS (blue class) 3-D interactions. (A) Schematic illustrating the ES-iPS (purple) and NPC-iPS (blue) interaction classes. **(B)** Contact frequency heatmaps (top) and zoomed-in heatmaps of distance-corrected interaction scores (bottom) highlighting a key interaction between Klf4 and an upstream enhancer. Interaction score heatmaps are overlaid on ChIPseq tracks of H3K27ac and H3K4me1 in ES cells and NPCs. **(C)** Distance-corrected interaction score changes among ES, NPC and iPS cells at the Klf4-enhancer ES-iPS (purple class) interaction. Error bars represent standard deviation across two replicates. **(D)** Normalized gene expression for the Klf4 gene is plotted for ES, NPC and iPS cells, as well as ES and iPS cells cultured in 2i media. Error bars represent standard deviation across two replicates. **(E)** Distance-corrected interaction score changes at an NPC-iPS interaction around the Klf4 gene among ES, NPC and iPS cells. Error bars represent standard deviation across two replicates. **(F)** Contact frequency heatmaps (top) and zoomed-in heatmaps of distance-corrected interaction scores (bottom) highlighting the NPC-iPS interaction between the Klf4 gene and a downstream NPC-specific enhancer. Plotted similar to **(B)**.

Figure S6 (related to Figure 5). NPC-specific genes and enhancers are enriched in NPC only (green class) interactions. (A) Schematic illustrating the NPC only (green) interaction class. **(B)** Bar plot displaying the fraction of each looping class containing NPC-specific enhancers compared to the expected background fraction. Fisher's Exact test: *, $P= 3.55182e-58$; **, $P= 0.00063607$. **(C)** Bar plot displaying the fraction of each looping class containing NPC-specific genes compared to the expected background fraction. Fisher's Exact test: *, $P= 1.20143e-86$. **(D)** Zoomed-in heatmaps of distance-corrected interaction scores highlighting key interactions between the Olig1 and Olig2 genes and nearby NPC-active enhancers. Distance-corrected interaction score heatmaps are overlaid on ChIPseq tracks of H3K27ac and CTCF in ES cells and NPCs. **(E-G)** Normalized gene expression for the Olig1 and Olig2 **(E)**, Nestin **(F)** and Bcan **(G)** genes are plotted for ES, NPC and iPS cells. Error bars represent standard deviation across two replicates.

Figure S7 (related to Figure 6). The Mis18 and Urb1 genes engage in ES only (red class) 3-D interactions linked to ES-specific CTCF binding. (A) Contact frequency heatmaps (top) and zoomed-in heatmaps of distance-corrected interaction scores (bottom) highlighting ES only interactions surrounding the Mis18a and Urb1 genes. Interaction score heatmaps are overlaid on ChIPseq tracks of CTCF and Smc1 in ES cells and NPCs. **(B)** Schematic illustrating the ES only (red) class of looping interactions. **(C-D)** Normalized gene expression for the Mis18a **(C)** and Urb1 **(D)** genes are plotted for ES, NPC, iPS cells and ES/iPS cells cultured in 2i media. Error bars represent standard deviation across two replicates. **(E-F)** Distance-corrected interaction score changes at Mis18a **(E)** and Urb1 **(F)** ES-only interactions highlighted on heatmaps with small red boxes in **(A)**. Error bars represent standard deviation across two replicates. **(G)** Relative ChIP-qPCR enrichment of CTCF binding at the ES only interaction displayed in **(A)**. CTCF site queried is denoted by red star in **(A)**.

Supplemental Tables

Table S1: Summary of paired-end read alignments for 5C libraries, related to Experimental Procedures

Library Code	Replicate	Instrument	Lane/Paired End	Total Reads	PE1 Mapped Reads	PE2 Mapped Reads
ES_1	1	Illumina Nextseq 500	L1_P1	33678848	28770023	28505219
			L1_P2	33678848		
			L2_P1	33418892		
			L2_P2	33418892		
			L3_P1	33974768		
			L3_P2	33974768		
			L4_P1	33399920		
			L4_P2	33399920		
ES_2	2	Illumina Nextseq 500	L1_P1	31551080	27875198	27628550
			L1_P2	31551080		
			L2_P1	31299432		
			L2_P2	31299432		
			L3_P1	31772324		
			L3_P2	31772324		
			L4_P1	31309272		
			L4_P2	31309272		
NPC_1	1	Illumina Nextseq 500	L1_P1	27804116	18454027	15832156
			L1_P2	27804116		
			L2_P1	24416680		
			L2_P2	24416680		
			L3_P1	13862024		
			L3_P2	13862024		
			L4_P1	17389664		
			L4_P2	17389664		
NPC_2	2	Illumina Nextseq 500	L1_P1	27793844	18342888	15617223
			L1_P2	27793844		
			L2_P1	24550324		
			L2_P2	24550324		
			L3_P1	13826704		
			L3_P2	13826704		
			L4_P1	17240756		
			L4_P2	17240756		
iPS_1	1	Illumina	L1_P1	23527984	15039775	13005171

		Nextseq 500	L1_P2 L2_P1 L2_P2 L3_P1 L3_P2 L4_P1 L4_P2	23527984 20602800 20602800 11619608 11619608 14506996 14506996		
iPS_2	2	Illumina Nextseq 500	L1_P1 L1_P2 L2_P1 L2_P2 L3_P1 L3_P2 L4_P1 L4_P2	24074808 24074808 21329464 21329464 11963384 11963384 14902364 14902364	15970612	13768584
ES_2i_1	1	Illumina Nextseq 500	L1_P1 L1_P2 L2_P1 L2_P2 L3_P1 L3_P2 L4_P1 L4_P2	22956884 22956884 19862384 19862384 11563004 11563004 14156912 14156912	15065438	12571131
ES_2i_2	2	Illumina Nextseq 500	L1_P1 L1_P2 L2_P1 L2_P2 L3_P1 L3_P2 L4_P1 L4_P2	26479112 26479112 23469892 23469892 13319924 13319924 16661424 16661424	17803279	15151910
iPS_2i_1	1	Illumina Nextseq 500	L1_P1 L1_P2 L2_P1 L2_P2 L3_P1 L3_P2 L4_P1 L4_P2	21352148 21352148 18236676 18236676 10483824 10483824 13062076 13062076	13147449	11301729
iPS_2i_2	2	Illumina Nextseq 500	L1_P1 L1_P2 L2_P1	23812716 23812716 21226860	15400978	12963765

			L2_P2	21226860		
			L3_P1	12105132		
			L3_P2	12105132		
			L4_P1	15124608		
			L4_P2	15124608		

Table S2: Spearman's rank correlation coefficients calculated for distance-dependence corrected interaction frequencies of pairs of biological replicates, related to Experimental Procedures.

ES_Rep_1						
ES_Rep_2	0.830632					
NPC_Rep_1	0.280142	0.243655				
NPC_Rep_2	0.27191	0.267666	0.767335196			
iPS_Rep_1	0.548705	0.581172	0.302233915	0.374198865		
iPS_Rep_2	0.44135	0.426666	0.456490393	0.492875434	0.678932815	
	ES_Rep_1	ES_Rep_2	NPC_Rep_1	NPC_Rep_2	iPS_Rep_1	iPS_Rep_2

Table S3: Summary of paired-end read alignments for RNAseq libraries, related to Experimental Procedures and Supplemental Experimental Procedures.

Library Code	Replicate	Instrument and Number of Lanes	Lane/Paired End	Total Reads	Alignment Summary
ES_1	1	Illumina Nextseq 500	L1_P1	63385276	Aligned pairs: 47708823 of these: 6941410 (14.5%) have multiple alignments 1525220 (3.2%) are discordant alignments 75.6% concordant pair alignment rate.
			L1_P2	63385276	
			L2_P1	62210488	
			L2_P2	62210488	
			L3_P1	59599184	
			L3_P2	59599184	
			L4_P1	59255860	
			L4_P2	59255860	
ES_2	2	Illumina Nextseq 500	L1_P1	49406568	Aligned pairs: 24184676 of these: 6397683 (26.5%) have multiple alignments 5602245 (23.2%) are discordant alignments 38.8% concordant pair alignment rate.
			L1_P2	49406568	
			L2_P1	48742476	
			L2_P2	48742476	
			L3_P1	46795280	
			L3_P2	46795280	
			L4_P1	46667840	
			L4_P2	46667840	

			L1_P1 L1_P2 L2_P1 L2_P2 L3_P1 L3_P2 L4_P1 L4_P2	28202868 28202868 27832228 27832228 27903044 27903044 27359632 27359632	Aligned pairs: 15612304 of these: 4167968 (26.7%) have multiple alignments 3682919 (23.6%) are discordant alignments 42.9% concordant pair alignment rate.
NPC_1	1	Illumina Nextseq 500	L1_P1 L1_P2 L2_P1 L2_P2 L3_P1 L3_P2 L4_P1 L4_P2	34176084 34176084 33607532 33607532 33839124 33839124 33064788 33064788	Aligned pairs: 16843964 of these: 5246902 (31.2%) have multiple alignments 4743602 (28.2%) are discordant alignments 35.9% concordant pair alignment rate.
NPC_2	2	Illumina Nextseq 500	L1_P1 L1_P2 L2_P1 L2_P2 L3_P1 L3_P2 L4_P1 L4_P2	32832608 32832608 32294456 32294456 32521560 32521560 31787280 31787280	Aligned pairs: 19633261 of these: 5907437 (30.1%) have multiple alignments 2764224 (14.1%) are discordant alignments 52.1% concordant pair alignment rate.
iPS_1	1	Illumina Nextseq 500	L1_P1 L1_P2 L2_P1 L2_P2 L3_P1 L3_P2 L4_P1 L4_P2	66486724 66486724 65682424 65682424 62943540 62943540 62797608 62797608	Aligned pairs: 30717628 of these: 8903977 (29.0%) have multiple alignments 7379227 (24.0%) are discordant alignments 36.2% concordant pair alignment rate.
			L1_P1 L1_P2 L2_P1 L2_P2 L3_P1 L3_P2 L4_P1 L4_P2	23466568 23466568 23169000 23169000 23272384 23272384 22782364 22782364	Aligned pairs: 12293064 of these: 3617560 (29.4%) have multiple alignments 2992072 (24.3%) are discordant alignments 40.1% concordant pair alignment rate.

iPS_2	2	Illumina Nextseq 500	L1_P1	45551664	Aligned pairs: 22993950 of these: 6316523 (27.5%) have multiple alignments 4420192 (19.2%) are discordant alignments 42.1% concordant pair alignment rate.
			L1_P2	45551664	
			L2_P1	44876400	
			L2_P2	44876400	
			L3_P1	43097496	
			L3_P2	43097496	
			L4_P1	42875224	
			L4_P2	42875224	
			L1_P1	29625648	Aligned pairs: 16810920 of these: 4686563 (27.9%) have multiple alignments 3287433 (19.6%) are discordant alignments 46.3% concordant pair alignment rate.
			L1_P2	29625648	
			L2_P1	29151848	
			L2_P2	29151848	
			L3_P1	29348000	
			L3_P2	29348000	
			L4_P1	28673296	
			L4_P2	28673296	
ES_2i_1	1	Illumina Nextseq 500	L1_P1	59127460	Aligned pairs: 42262509 of these: 6635919 (15.7%) have multiple alignments 2569916 (6.1%) are discordant alignments 69.3% concordant pair alignment rate.
			L1_P2	59127460	
			L2_P1	58169908	
			L2_P2	58169908	
			L3_P1	55872492	
			L3_P2	55872492	
			L4_P1	55774136	
			L4_P2	55774136	
				7370792	Aligned pairs: 5861538 of these: 950297 (16.2%) have multiple alignments 370120 (6.3%) are discordant alignments 75.5% concordant pair alignment rate.
				7370792	
				7260416	
				7260416	
				7299252	
				7299252	
				7149924	
				7149924	
ES_2i_2	2	Illumina Nextseq 500		41617840	Aligned pairs: 31055590 of these: 4668288 (15.0%) have multiple alignments 1206653 (3.9%) are discordant alignments 74.2% concordant pair alignment rate.
				41617840	
				40991972	
				40991972	
				39186452	
				39186452	
				39053336	
				39053336	

				12881568 12881568 12701244 12701244 12733836 12733836 12471728 12471728	Aligned pairs: 10705467 of these: 1658953 (15.5%) have multiple alignments 426422 (4.0%) are discordant alignments 81.0% concordant pair alignment rate.
iPS_2i_1	1	Illumina Nextseq 500		43012836 43012836 42257896 42257896 40724336 40724336 40395388 40395388	Aligned pairs: 23098024 of these: 6101858 (26.4%) have multiple alignments 4051057 (17.5%) are discordant alignments 45.8% concordant pair alignment rate.
				7600884 7600884 7470076 7470076 7525056 7525056 7366240 7366240	Aligned pairs: 4591177 of these: 1247317 (27.2%) have multiple alignments 821522 (17.9%) are discordant alignments 50.3% concordant pair alignment rate.
iPS_2i_2	2	Illumina Nextseq 500		45249896 45249896 44710976 44710976 42713772 42713772 42689140 42689140	Aligned pairs: 31164351 of these: 6059475 (19.4%) have multiple alignments 2368280 (7.6%) are discordant alignments 65.7% concordant pair alignment rate.

Table S4: Summary of external ChIP-seq libraries analyzed in this study, related to Experimental Procedures and Supplemental Experimental Procedures.

Antibody	Cell Type	Mapped Test ChIP-Seq reads	Test ChIP Reference	Test Sample GEO ID	Control Samples	Mapped Control ChIP-Seq reads	Control Sample GEO ID
CTCF	mES (159-2)	9,562,677	(Stadler et al., 2011)	GSM747534	mES Whole	10,202,630	GSM747545

					Cell Extract		
CTCF	ES-derived NPC	13,641,735	(Phillips-Cremins et al., 2013)	GSM883647	NPC Whole Cell Extract	14,041,323	GSM883648
H3K4me 1	mES (V6.5)	5,707,101	(Meissner et al., 2008)	GSM281695	V6.5 Whole Cell Extract	803,601	GSM307625
H3K4me 1	ES-derived NPC	4,471,210	(Meissner et al., 2008)	GSM281693	NPC Whole Cell Extract	4,369,951	GSM307617
H3K4me 3	mES (V6.5)	6,809,878	(Mikkelsen et al., 2007)	GSM307618	V6.5 Whole Cell Extract	6,008,440	GSM307154 , GSM307155
H3K4me 3	ES-derived NPC	3,397,613	(Mikkelsen et al., 2007)	GSM307613	NPC Whole Cell Extract	4,369,951	GSM307617
H3K27ac	mES (V6.5)	11,128,384	(Creyghton et al., 2010)	GSM594579 Rep2	V6.5 Whole Cell Extract	14,682,811	GSM307154 , GSM307155 , GSM594599
H3K27ac	ES-derived NPC	8,831,628	(Creyghton et al., 2010)	GSM594585	NPC Whole Cell Extract	14,041,323	GSM883648

Table S5: ES Specific Parsed Genes, Related to Figures 4, 5, 6 and Supplemental Experimental Procedures.

Attached as separate excel spreadsheet.

Table S6: NPC Specific Parsed Genes, Related to Figures 4, 5, 6 and Supplemental Experimental Procedures.

Attached as separate excel spreadsheet.

Table S7: Interactions selected for representative interaction score barplots, Related to Figures 4, S5, 5, 6, S7.

Attached as separate excel spreadsheet.

Supplemental Experimental Procedures

ES cell culture

V6.5 ES cells (murine; C57Bl/6 x 129SvJae; male) were purchased from Novus Biologicals. ES cells were expanded on Mitomycin-C inactivated MEF feeder layers in media consisting of DMEM, 15% FBS (Hyclone), 10^3 U/mL leukemia inhibitory factor (Millipore), non-essential amino acids (Lifetech), 0.1 mM 2-mercaptoethanol, 4 mM L-glutamine (Lifetech) and penicillin/streptomycin (Lifetech). Prior to fixation, ES cells were passaged onto gelatin-coated, feeder-free plates to remove feeder layer, and fixed at approximately 70% confluence. Cells were grown to $\sim 7 \times 10^6$ cells per 15 cm dish at the time of fixation.

Primary NPC isolation

Neural progenitor cells were isolated from whole brains of newborn 129SvJae x C57/BL6, Sox2-eGFP mice and cultured as neurospheres in Neural Stem Cell media: DMEM/F12 media (Invitrogen 12100-046 and 21700-075) containing 72 mM glucose, 120 mM Sodium Bicarbonate, 5.6 mM HEPES (Sigma H-0887), 27.5 nM Sodium Selenite (Sigma S-9133), 18 nM progesterone (Sigma P0130), 90 μ g/mL Apotransferrin (Sigma T1428), 23 μ g/mL insulin (Sigma I6634), 100 μ M putrescine (Sigma P-7505), 2 mM L-glutamine (Gibco 25030-081), 1% Pen/Strep (Sigma P0781), 2 μ g/mL heparin, 20 ng/mL rhEGF (R&D Systems) and 10 ng/mL rhFGF (R&D systems). Neurospheres were passaged every 3-4 days to prevent the formation of necrotic cores. After two passages, neurospheres were dissociated with Accutase and plated on Poly-D-Lysine Hydrobromide (100 μ g/mL, Sigma P7280), and laminin (15 μ g/mL, Corning 354232) coated plates at 60,000 cells/cm². Cells were fixed with 1% formaldehyde one day after adherent plating.

iPS cell culture

The iPS cells analyzed in this study were reprogrammed from primary NPCs (pNPCs) as described in (Eminli et al., 2008). Briefly, pNPCs were transduced with lentiviral vectors to ectopically express Oct4,

Klf4, and c-Myc (OKM). iPS cells derived from pNPCs were cultured on irradiated MEFs in medium consisting of Knock-Out DMEM, 15% FBS, Glutamax, non-essential amino acids, penicillin-streptomycin, β -mercaptoethanol and Leukemia Inhibitory Factor (LIF). iPS cells were grown to $\sim 7 \times 10^6$ cells per 15 cm dish at the time of fixation. This iPS clone was extensively characterized for its pluripotent properties as assessed by (i) high expression of endogenous pluripotency markers (Oct4, Sox2, Nanog), (ii) demethylation of Oct4 and Nanog promoters, (iii) in vivo teratoma formation of all three germ layers and (iv) generation of chimeric mice ([Eminli et al., 2008](#)).

Culture of pluripotent cells in 2i media

iPS and ES cells were removed from serum-containing media described above and cultured in 2i serum-free media comprised of 500 mL Knock Out DMEM (Life Technologies # 10829-018), 15% Knockout Serum Replacement (Life Technologies #10828), 5 mL N2 supplement (Life Technologies #17502-048), 5 mL B27 Supplement (Life Technologies #17504-044), 5 mg/mL BSA (Sigma A9418), 1 mM L-Glutamine (Life Technologies # 25030-081), 1% Non-Essential Amino Acids (Millipore #TMS-001-C), 0.1 mM β -Mercaptoethanol (Life Technologies #21985-023), 1% Penicillin-Streptomycin (Sigma #P0781), 10^3 units/mL LIF (Millipore #ESG1107), 3 μ M CHIR99021 (Axon Medchem #1386), and 1 μ M PD0325901 (Axon Medchem #1408) ([Rais et al., 2013](#)). After two passages on feeder cells, ES and iPS cells in 2i media were passaged onto 0.1% gelatin to remove contaminating feeder cells. Cells were grown to $\sim 7 \times 10^6$ cells per 15 cm dish at the time of fixation with 1% formaldehyde before 5C.

3C template generation and characterization

3C templates were produced as previously described ([Dekker et al., 2002](#); [Gheldof et al., 2010](#); [Phillips-Cremins et al., 2013](#); [van Berkum and Dekker, 2009](#)) for ES (n=2), NPC (n=2), iPS (n=2), ES+2i (n=2) and iPS+2i (n=2) pellets. Briefly, cells were fixed in culture media supplemented with formaldehyde added to

a final concentration of 1%. After 10 minute incubation at room temperature, fixation was terminated by adding 2.5M glycine stock to a final concentration of 125 mM glycine. Cross-linking termination was carried out for 5 minutes at room temperature followed by 15 minutes at 4°C. Cells were harvested with silicone scraper and pelleted, washed once with PBS, snap-frozen and stored at -80°C until processing.

Pellets were resuspended in lysis buffer consisting of 10 mM Tris-HCl (pH 8.0), 10 mM NaCl, 0.2% Igepal CA630 and 1x protease inhibitor (Sigma) in sterile water and incubated on ice for 30 minutes. Cells were lysed with a dounce homogenizer and washed with NEB2 buffer. SDS was added to a final concentration of 0.1% and chromatin was solubilized by incubating at 65°C for 10 minutes. Triton X-100 was added to quench the SDS, and HindIII digestion was performed overnight at 37°C. The next day, the HindIII was inactivated and ligation was performed under dilute conditions at 16°C for 2 hours using T4 DNA ligase (Invitrogen) in ligation buffer consisting of 1% Triton X-100, 0.1mg/mL BSA, 1mM ATP, 50mM Tris-HCl, 50mM NaCl, 10mM MgCl₂ and 1mM DTT. After ligation, cross-links were reversed via incubation with 63.5µg/mL Proteinase K (Invitrogen) for 4 hours at 65°C, at which point the Proteinase K concentration was doubled and the solution was incubated overnight at 65°C. The 3C template DNA was then purified via a phenol extraction and a subsequent phenol-chloroform extraction before precipitation in ethanol. The resulting DNA pellet was resuspended in TE buffer consisting of 10 mM Tris-HCl (pH 8.0) and 1 mM EDTA (pH 8.0), and again purified by a series of phenol-chloroform extractions and precipitated in ethanol. The resulting DNA pellet was resuspended in TE buffer and treated with 100 ug/mL RNase A for 3 hours at 37°C. To confirm successful 3C template generation, the presence of expected ligation products was assayed via conventional PCR.

5C primer design

5C primers were designed at HindIII restriction sites using the my5Csuite primer design tools ([Lajoie et al., 2009](#)), as described in detail in ([Phillips-Cremins et al., 2013](#)).

5C library generation and sequencing

5C libraries were generated as described previously ([Bau et al., 2011](#); [Dostie and Dekker, 2007](#); [Dostie et al., 2006](#); [Phillips-Cremins et al., 2013](#); [van Berkum and Dekker, 2009](#)). 600 ng of each 3C template was mixed with final concentration 1 fmol of each 5C primer in 1x NEB4 buffer. Solution was incubated at 55°C for 16 hr to anneal primers to 3C templates. 5C primers annealed to 3C ligation junctions were ligated via the addition of 1x Taq ligase buffer containing 10 U Taq DNA ligase. Solution was mixed by pipetting and incubated for 1 hour at 55°C. Ligated 5C primers were then selectively amplified via the addition of universal forward (T7) and reverse (T3) primers, which anneal to the complementary universal primer tails of the 5C primers. 5C libraries (400 ng per library) were prepared for sequencing using the NEBNext Ultra DNA Library Prep Kit (NEB # E7370S) and NEBNext Multiplex Oligos for Illumina (NEB # E7335S). After ligation of adapters following manufacturer's protocol, nuclease-free water was added to bring the reaction volume to 100 uL. Fragments of size ~ 220 bp (100 bp 5C product + 120 bp Illumina adapters) were preferentially selected using AgenCourt Ampure XP beads (Beckman Coulter A63881), by first adding 70 uL beads and retaining the supernatant, then adding 25 uL beads, removing the supernatant, and washing and eluting sample from the beads following the manufacturer's protocol. Following adapter ligation and size selection, the libraries with Illumina adapters were amplified with 10 cycles of PCR. The size distribution of the purified libraries were assessed on the Agilent BioAnalyzer using the DNA 1000 kit (Agilent 5067-1505). The resulting 5C libraries were pooled and sequenced with 37-cycles per paired-end on the Illumina NextSeq500.

iPS cell transgene integration detection by 5C primers

This iPS clone was generated via integration of transgenic Oct4, Klf4, and c-Myc genes ([Eminli et al., 2008](#)). Hochedlinger and colleagues demonstrated that this iPS clone exhibits transgene-independent

self-renewal potential, which would exclude that these cells still depended on transgenic OKM expression. We note that our 5C approach does not exclude detection of the exogenous *Oct4* and *Klf4* genes (which were likely virally integrated at sites distal to our 5C regions) with 5C primers that directly bind to the Oct4/Klf4 coding sequence. However, short-range, cis interactions represent the majority of the 5C signal, and we do not analyze trans interactions in this study. Thus, we would expect the transgenes to contribute relatively little to the interaction counts between these genes and other sites within our designed primer set.

RNA-seq library preparation

900,000 cells of each cell type were lysed with Trizol (Life Technologies 15596-026) and snap frozen. Total RNA was extracted and purified using the miRvana miRNA Isolation Kit (Ambion AM 1561) and samples were eluted into 100 uL nuclease free water. All RNA samples had an RNA Integrity Number >9 as assessed by Agilent BioAnalyzer. 50 uL of each RNA sample was treated with 1 uL rDNase I (Ambion 1906) to remove residual genomic DNA. 350 ng DNase-treated total RNA was prepared for sequencing using the Illumina TruSeq Stranded Total RNA Library Prep kit with RiboZero (Illumina RS-122-2202) following the supplier's protocol. cDNA libraries with Illumina adapters were amplified with 15 cycles of PCR. Libraries were purified using AgenCourt Ampure XP beads (Beckman Coulter A63881) with two rounds of 1:1 bead:sample selection. The size distributions of the purified cDNA libraries were assessed on the Agilent BioAnalyzer using the DNA 1000 kit (Agilent 5067-1505). Libraries were pooled and sequenced with 75-cycles per paired-end on the Illumina NextSeq500.

RNAseq data processing

RNAseq reads were aligned to the mouse genome (build mm9) using the Tophat (Tophat v2.1.0) alignment tool ([Trapnell et al., 2009](#)) with the parameters: -r 100 --no-coverage-search --library-type fr-

firststrand and UCSC gene annotations (**Table S3**). Gene level read counts were computed using the htseq-count tool (<http://www-huber.embl.de/users/anders/HTSeq/doc/count.html>) with parameters: -m union --stranded=reverse and UCSC gene annotations. For analyses of all 10 samples (ES_Rep1, ES_Rep2, pNPC_Rep1, pNPC_Rep2, iPS_Rep1, iPS_Rep2, ES2i_Rep1, ES2i_Rep2, iPS2i_Rep1, iPS2i_Rep2), genes with more than three counts in at least five libraries were retained, resulting in a total of 11,767 genes analyzed. To account for library-specific differences in sequencing depth, log2-transformed libraries were normalized by read depth of the 75%tile gene. Libraries were assessed for the absence of batch effects before proceeding to downstream biological analyses (**Figure S4**).

CTCF binding detection by CHIP-qPCR

Approximately 20 million cells were fixed in serum-free culture media supplemented with formaldehyde added to a final concentration of 1%. After 10 minute incubation at room temperature, fixation was terminated by adding 2.5M glycine stock to a final concentration of 125 mM glycine. Cross-linking termination was carried out for 5 minutes at room temperature followed by 15 minutes at 4°C. Cells were harvested with silicone scraper and pelleted, washed once with PBS, snap-frozen and stored at -80°C until processing.

Cell pellets were thawed for 10 min on ice before use. Nuclei were isolated by resuspending each pellet in 1 mL Cell Lysis Buffer (10 mM Tris pH 8.0, 10 mM NaCl, 0.2% NP-40/Igepal, Protease Inhibitor, PMSF), incubating on ice for 10 min, and spinning to pellet. Nuclei were resuspended in 500 μ L Nuclear Lysis Buffer (50 mM Tris pH 8.0, 10 mM EDTA, 1% SDS, Protease Inhibitor, PMSF) and incubated on ice for 20 min. After bringing the samples up to volume by the addition of 300 μ L IP Dilution Buffer (20 mM Tris pH 8.0, 2 mM EDTA, 150 mM NaCl, 1% Triton X-100, 0.01% SDS, Protease Inhibitor, PMSF), samples were sonicated for 45 minutes using an Epishear sonicator set at 100% amplitude, with cycles of 30 seconds on and 30 seconds off. The resulting sheared chromatin was spun down, and the

supernatant was transferred to a preclearing solution of 3.7 mL IP Dilution Buffer, 0.5 mL Nuclear Lysis Buffer, 175 uL of Agarose Protein A/G beads, and 50 ug Rabbit IgG, and rotated at 4°C. 35 uL Protein A/G agarose beads were pre-bound with 10 uL anti-CTCF antibody (Millipore #07-729) and incubated for 2 hours during the pre-clear stage. After a two hour pre-clear incubation, the beads were pelleted, and 4.5 mL supernatant was removed. 200 uL was reserved for input control, while the remaining supernatant was transferred to agarose beads pre-bound with antibody and rotated overnight at 4°C. Bound bead complexes were washed once with 1 mL IP Wash Buffer 1 (20 mM Tris pH 8.0, 2 mM EDTA, 50 mM NaCl, 1% Triton X-100, 0.1% SDS), twice with 1 mL High-Salt Buffer (20 mM Tris pH 8.0, 2 mM EDTA, 500 mM NaCl, 1% Triton X-100, 0.01% SDS), once with IP Wash Buffer 2 (10 mM Tris pH 8.0, 1 mM EDTA, 0.25 M LiCl, 1% NP-40/Igepal, 1% Na-deoxycholate), and finally once with 1x TE. Complexes were eluted by twice resuspending bound beads in 110 uL Elution Buffer (100 mM NaHCO₃, 1% SDS), pelleting the beads after each elution and transferring 100 uL supernatant to a new tube. Finally, 12 uL of 5M NaCl and 20 ug RNase A were added to both 200 uL IP and input samples and incubated at 65 degrees for 1 hour, followed by the addition of 60 ug of Proteinase K and overnight incubation at 65 degrees. DNA was isolated via phenol-chloroform extraction and ethanol precipitation, and concentration was quantified using Qubit fluorometer.

ChIP libraries were prepared from 3 ng of IP and input DNA using the NEBNext Ultra Library Prep Kit (NEB #E7370) following the manufacturers protocol for preparation of ChIP libraries. After adapter ligation, no size selection step was performed, and ligated samples were enriched through 18 PCR cycles using NEBNext Multiplex Oligos for Illumina (NEB #E7335). Libraries were eluted in 30 uL 0.1x TE, and a fragment size distribution between 250 and 1200 bp including sequencing adapters was confirmed using a High-Sensitivity assay on a Agilent Bioanalyzer.

Primers were designed to query specific CTCF binding sites:

Figure Panel	Forward Primer	Reverse Primer	Genomic Coordinates
--------------	----------------	----------------	---------------------

5G (NPC-iPS)	TGTGGTCCTTTGTCCTTCCTG	TGTCACGCATCCTGAATCTTC	Chr3:35002112-35002461
5G (ES only)	AACTCACTAAGTGGCCCGAAG	ACCCAGCTCCACGAAAATG	Chr3:34658834-34659306
6H	GTGTACAAGCACGCACGTATG	AAAGGGAGGTGCTCAATGGTC	Chr4:54936308-54936574
S7G	TAACCCTCACTGCTTGCCTAG	TGTGTCCTTAGCAGACGTGTC	Chr16:90635525-90635762

Quantitative PCR was performed by loading 1 ng of each sample library into each 20 uL reaction, including 10 uL Power SYBR Green PCR Master Mix (Applied Biosystems # 4367659), and corresponding primers (200 nM final concentration). Reactions were loaded onto an Applied Biosystems StepOnePlus in three replicates and assayed using standard qPCR cycling conditions (95°C for 10 min, followed by 40 cycles of 95°C for 15 sec and 65°C for 1 min). The CT threshold was set at 1900 so as to fall in the middle of the exponential phase for all primers and to capture the CT value for all samples. To facilitate comparison among the five cellular conditions, relative enrichment in CTCF ChIP signal was assessed by normalizing data by a reference control primer representing a constitutively bound CTCF site.

5C data processing pipeline

Paired-end read mapping and counting

5C data were generated with paired-end sequencing (37 bp paired-end reads) on the Illumina NextSeq 500 instrument. The two ends of paired-end (PE) reads were aligned independently to a pseudo-genome consisting of all 5C primers using Bowtie with default parameters (<http://bowtie-bio.sourceforge.net/index.shtml>) (Langmead, 2010). Only reads with one unique alignment were considered for downstream analyses. Interactions were counted when both paired-end reads could be uniquely mapped to the 5C primer pseudo-genome. Only interactions between forward-reverse primer pairs were tallied as true counts (**Table S1**).

Low count primer removal

Primers with fewer than 100 total reads across all possible cis primer ligation partners were excluded from further analysis. Removed primers are listed below:

#track	Start	Stop	Primer ID
chr3	87677389	87683794	5C_326_Nestin_FOR_117:0
chr3	88032708	88035039	5C_326_Nestin_FOR_192:0
chr3	88124897	88125644	5C_326_Nestin_FOR_214:0
chr3	88283586	88286361	5C_326_Nestin_FOR_248:0
chr16	91242594	91247280	5C_325_Olig1- Olig2_FOR_193:0
chr17	35285175	35292115	5C_327_Oct4_FOR_191:0
chr17	36018525	36020858	5C_327_Oct4_FOR_378:0
chr17	36023358	36024542	5C_327_Oct4_FOR_380:0
chr17	36393683	36395722	5C_327_Oct4_FOR_472:0
chr3	34546431	34549386	5C_329_Sox2_REV_154:0

Raw contact matrix visualization

First we designated the restriction fragments to which 5C primers were designed as “queried restriction fragments”. Raw contact matrices were generated for each region by placing the number of counts read for the interaction of the i th queried restriction fragment in the region with the j th queried restriction fragment in the region in the ij th entry of the contact matrix. This created a square, symmetric matrix of contacts with dimensions equal to the number of primers in the region. Because interactions between fragments whose corresponding primers are oriented in the same direction cannot be detected with our 5C primer design, not every entry of this matrix corresponds to a detectable fragment-fragment interaction.

Because approximately half of the entries in this contact matrix represent undetectable fragment-fragment interactions, we visualized raw contact matrices at the fragment level by arranging the forward primers on the x-axis and the reverse primers on the y-axis, in order of primer number, which corresponded directly with the sorted order of genomic coordinates (heatmaps in **Fig. S1A**). Thus, the ij th cell of the resulting heatmap represents the number of counts for the interaction of the

fragment queried by the j th forward primer with that queried by the i th reverse primer. This heatmap, used only for initial visualization, is therefore asymmetric and not necessarily square.

Quantile normalization

It is essential to account for technical variation among 5C replicates - in particular, batch effects for experiments processed or sequenced on different days - before comparing dynamic architecture between biological conditions. Indeed, we have found that two important factors driving experimental variability between biological replicates are (i) library complexity and (ii) sequencing depth differences between each batch of processed samples. We have found that a simple normalization factor is insufficient to remove bias due to sequencing depth because the differences in read counts between replicates tend to compound in a nonlinear manner based on the underlying complexity of the library.

Quantile normalization is a rank-based approach that has successfully been used to normalize microarray (Bolstad et al., 2003), RNAseq (Bullard et al., 2010) and Hi-C (Dixon et al., 2015) data prior to downstream modeling. Here we also find that quantile normalization is effective at placing different 5C libraries on the same distributional scale (compare distance dependence and histograms in **Fig. S1A-B**) while preserving biologically significant architectural features (compare heatmaps in **Fig. S1A-B**). We have noticed that quantile normalization is particularly effective on 5C datasets because the strongest signal in the raw data is the distance-dependence background, providing a smooth, ubiquitous rank-order gradient for the comparison of contacts across replicates and conditions. Indeed, we found that our analysis was largely insensitive to the exact placement of the quantile normalization step relative to the other steps. For example, we moved the quantile normalization step to the end of our 5C analysis pipeline (**Fig. S2A+B,E-G**) and found that all views of the data show striking similarity to the corresponding stages of our implemented data processing pipeline (**Fig. S1A-F**).

Primer correction

Consistent with our findings in (Phillips-Cremins et al., 2013), we noticed the presence of primer-specific bias in our 5C data. For example, we observed strongly underenriched or overenriched stripes in our raw heatmaps – indicating that entire rows/columns can have increased or decreased counts (heatmaps in **Fig. S1A**). Consistent with this observation, the cis interactions for each primer show up to an ~8500-fold variation in mean interaction frequency, suggesting the presence of artifacts independent from the biology that influence the 5C signal (boxplots in **Fig. S1A**). To account for primer-specific artifacts, we applied our previously developed primer correction method that uses stochastic gradient descent to compute primer-effect normalization factors (Phillips-Cremins et al., 2013). After the primer correction step, we observed a marked attenuation of primer-specific artifacts (heatmaps and boxplots, **Fig. S1C**).

Low count fragment-fragment pair removal

Fragment-fragment pairs with primer-corrected counts below 10 in any replicate were flagged as low outliers with essentially unreliable values and were excluded from further analysis.

Contact matrix binning

We next generated a binned contact frequency matrix by binning each of our queried regions at regular 4 kb intervals (approximately equal to the average cut frequency of our chosen restriction enzyme, HindIII). To assign a value to each element of the binned contact probability matrix, we computed an arithmetic mean of logged counts using a square, 20 kb smoothing window as:

$$b_{i,j} = \frac{\sum_{k,l \in \mathcal{M}} \mathbf{1}(|m_k - M_i| \leq 10 \text{ kb}, |m_l - M_j| \leq 10 \text{ kb}) \log_2(n_{k,l} + 1)}{\sum_{k,l \in \mathcal{M}} \mathbf{1}(|m_k - M_i| \leq 10 \text{ kb}, |m_l - M_j| \leq 10 \text{ kb}) \mathbf{1}(d_k \neq d_l)}$$

where $b_{i,j}$ is the value assigned to the ij th entry of the binned contact matrix for the region and represents the contact frequency of the i th and j th bins in the region, m_k represents the midpoint of the k th primer in the region, M_i represents the midpoint of the i th bin in the region, and $n_{k,l}$ represents the number of counts for the interaction of the k th queried fragment in the region with the l th queried fragment in the region after primer normalization. $\mathbf{1}(d_k \neq d_l)$ represents an indicator function that checks whether the k th and l th primer in the region have the same directionality. This ensures that the average is computed only over the possible primer-primer interactions.

If more than 80% of all the fragment-fragment pairs in a bin-bin pair's smoothing window had values that were zero, impossible, or had been previously removed as low outliers, that bin-bin pair was determined to be located in a low-confidence region and was excluded from further analysis. The bin-bin pair removal condition can be represented as:

$$\frac{\sum_{i,j \ni |m_i - M_k| \leq 10 \text{ kb}, |m_j - M_l| \leq 10 \text{ kb}} \mathbf{1}(n_{i,j} > 0)}{\sum_{i,j \ni |m_i - M_k| \leq 10 \text{ kb}, |m_j - M_l| \leq 10 \text{ kb}} \mathbf{1}} < 20\% \Rightarrow b_{k,l} \text{ excluded from further analysis}$$

We selected the 20 kb smoothing window size and the 4 kb matrix resolution through a process of (1) iteratively testing window sizes and matrix resolutions, (2) visually inspecting the resultant heatmaps and (3) qualitatively comparing heatmaps to classic epigenetic marks. Our final strategy optimally accounted for sampling noise in 5C data while retaining what we term a pseudo-fragment (~12 kb) resolution (discussed in detail below). We chose to assign values to the entries of the binned contact matrix using an average rather than a sum because HindIII has been previously shown to exhibit highly variable restriction site density across the genome. To attenuate the spatial noise present in our fragment-level data, our binning strategy effectively averages counts across a 20 kb window (compare

heatmaps in **Fig. S1C+D** and **Fig. S2B+E**). This reduction of spatial noise is concurrent with a tightening of the distribution of counts across this step (compare histograms in **Fig. S1C+D**).

Pseudo-fragment level 5C mapping resolution

Many definitions of 3C/4C/5C/Hi-C resolution have been reported. Therefore, it is important to clarify our definition of resolution and our strategy for matrix binning. In a recent publication, the so-called “mapping resolution” of a Hi-C contact density map was defined as the smallest locus size such that 80% of the loci have at least 1000 contacts ([Rao et al., 2014](#)). Importantly, Rao et al. reported the numbers in this definition as the finest scale at which they could reliably discern and distinguish architectural features in a Hi-C heatmap. By contrast to the “mapping resolution” metric, Rao et al. also define an alternative “matrix resolution” metric which is simply the bin size selected by the investigator when constructing a contact density matrix. In our lowest read depth replicate, iPS+2i Rep 1, 97% of the queried fragments have more than 1000 contacts. Thus, if we define our loci as the individual restriction fragments queried by the assay, all our datasets have a mapping resolution equal to the fragment size (~4 kb). We find a 4 kb bin size as the finest scale at which we can discern architectural features in our 5C contact density matrix. On the basis of a strictly “matrix resolution” definition, the resolution of our 5C data would be 4 kb. However, because we use a square 20 kb smoothing function (discussed below), there are hypothetical situations in which we cannot resolve two perfectly punctate features that are within 20 kb of each other. Thus, our “mapping” resolution falls in the range of 4-20 kb.

The design and orientation of 5C primers is another critical factor unique to 5C that must be considered in calculating resolution. Importantly, the true alternating 5C primer design used here and in ([Phillips-Cremins et al., 2013](#)) only queries a subset of possible fragment-fragment interactions. Specifically, forward and reverse primers were tiled in a true alternating manner across our genomic regions. Only forward-reverse (F-R) and reverse-forward (R-F) ligation products can be detected with the

ligation-mediated amplification approach. Thus, although we can distinguish most interactions at a ~4 kb resolution, our more generalized resolution due to the alternating primer design is at the level of F-R-F or R-F-R fragment sequences (~12 kb; also the midpoint between our 4-20 kb mapping resolution).

To our knowledge, no Hi-C map has been reported at true single-fragment resolution as even the highest density maps have been binned to 1-5 kb resolution with a 4 bp cutter that cuts approximately every 200-300 bp in the genome. Thus, the highest resolution maps to date still average or sum information from at least 4 (1 kb resolution) but as many as 1000's (1 Mb resolution) of adjacent restriction fragments prior to modeling, parameterization of models, and downstream analyses. The reason for this requisite binning step is that the sampling noise in 5C/Hi-C contact matrices represents a significant barrier in obtaining high-confidence information for the read counts in every bin across the genome. However, a high-confidence understanding of the interaction frequency can be modeled at the expense of losing some resolution by averaging or summing counts from nearby fragment-fragment pairs. Here, we use 5C, which offers key advantages over Hi-C in its ability to obtain high complexity contact density maps with a logistically reasonable sequencing depth. Thus, we have high complexity libraries (i.e. most restriction fragment ligation products have been sampled at an ultra-high count density). For example, in iPS+2i Rep 1, our lowest-mapping replicate, 80% of our originally queried fragments received >5340 counts. Ultimately, to account for spatial noise, we chose a 20 kb windowing function to yield a search space over an approximately 5x5 grid of primer-primer pairs (F-R-F-R-F or R-F-R-F-R). Overall, we propose that our resolution falls between 4 and 20 kb – with approximately a 12 kb resolution due to the true alternating primer design.

Identification of bad primer gaps

Restriction site density varies widely across the genome. Additionally, it is possible that certain primers fail to produce any counts due to technical error. Finally, many restriction fragments did not receive a

primer due to low quality scores, leaving several loci unqueried by the assay. All three factors may affect the distance between one existing "working" primer and the next downstream "working" primer. When this distance is small compared to the smoothing window, the gap will be successfully spanned by multiple unique smoothing windows. When this distance is on a similar scale to the smoothing window, the smoothing window will be too small to reliably smooth across the gap. Within each region, we identified columns of bins that contained no positive counts from any primer ligation. When the length of a run of consecutive missing or zero fragments was greater than half the size of the smoothing window plus one bin, we classified the gap as "unsmoothable." Unsmoothable gaps are marked with dark gray on the heatmaps and excluded from all statistical analyses.

Distance-dependence normalization

To account for the distance-dependence background inherent in 3C-related assays, we computed an empirical expected distance-dependence model (**Fig. S1G**). Within each region and replicate, we first grouped the bin-bin pairs according to their interaction distance d , as measured by the number of bins separating the constituent bins in the bin-bin pair. We then computed the mean of the binned interaction frequencies within each group, as follows:

$$\mu_d = \text{mean}_i [b_{i,i+d}]$$

where μ_d is the mean value at distance d (measured in number of bins of separation), and $(b_{i,i+d})_i$ is the sequence of binned contact frequencies for bin-bin pairs at distance d . Since the number of matrix entries included in each average will decrease with increasing distance d , these mean values are statistically weak predictors at long (> 600-700 kb for a 1 Mb region) distance scales. To account for any noise in our empirical distance-dependence estimations, we lowess-smoothed a subset of the empirical expected values in order to obtain a smooth approximation to the empirical expected values. Due to the

high number of matrix entries at distances ≤ 300 kb, we retained the original mean values at short distance scales (≤ 300 kb for a 1 Mb region).

We next used our empirical expected model to normalize the binned contact matrices by computing a fold-enrichment of counts relative to the expected (**Figs. S1E, S2G**). Since the values in our binned contact matrices were already log-transformed, we directly computed a log-scale fold-enrichment as:

$$f_{i,j} = b_{i,j} - \mu_{|i-j|}$$

where $f_{i,j}$, the ij th entry of the distance-normalized contact matrix, represents the log-scale fold-enrichment of interactions between the i th and j th bins in the region, $b_{i,j}$ is the ij th element of the binned interaction matrix, and $\mu_{|i-j|}$ represents the distance-dependence normalization factor appropriate for a bin-bin pair at distance $d = |i - j|$ within the region under consideration (described above). Distance dependence-normalized counts show no discernable relationship with interaction distance compared to data at earlier stages of the analysis (histograms in **Figs. S1E, S2G**).

Noteworthy, the Klf4 region spans two distinct sub-TADs with markedly different interaction frequencies. We divided Klf4 into two separate sub-regions and created independent expected models for sub-region_1 (single block: chr4:54,899,978-55,371,978 x chr4:54,899,978-55,371,978) and sub-region_2 (the union of three blocks: chr4:54,899,978-55,371,978 x chr4:55,371,978-55,887,978, chr4:55,371,978-55,887,978 x chr4:55,371,978-55,887,978 and chr4:55,371,978-55,887,978 x chr4:54,899,978-55,371,978). Spearman's rank-order correlation coefficients for distance-corrected interaction frequencies of ES, NPC, and iPS replicates can be found in **Table S2**.

Probabilistic model fitting and distance-corrected interaction scores

We modeled our distance-corrected interaction frequency values as a continuous random variable using a logistic distribution parameterized independently for each region and replicate (**Fig. S3A**). We fit the

logistic distribution by computing region-specific and replicate-specific location (l) and scale (s) parameters with maximum likelihood estimation through the R `fitdistr()` function. We computed right-tail p-values for every entry of distance-normalized contact matrices via the R `plogis()` algorithm, the `lower.tail=FALSE` argument and the below logistic cumulative distribution function:

$$p_{i,j} = 1 - \frac{1}{1 + e^{-(f_{i,j}-l)/s}}$$

where $p_{i,j}$ represents the right-tailed p-value for the relative interaction frequency found in the ij th entry of the distance-normalized contact matrix.

Prior to downstream thresholding/classification of significant 3-D interactions, p-values were transformed into distance-corrected interaction scores with:

$$IS_{i,j} = -10 \times \log_2(p_{i,j})$$

Our computed distance-corrected interaction score offers a specific metric for identification/detection of significant 3-D interactions that are visually evident but difficult to disentangle from the underlying noise in the raw data (illustrated in heatmaps **Fig. S1F**). The highest (red/black) bins in ES and NPC heatmaps show strong cell type-specific correlation with known cell type-specific chromatin marks (heatmaps in **Fig. S1F**) while exhibiting strong attenuation of primer effects, absence of distance-dependence background signal and minimal distribution skewing due to technical differences in library complexity (boxplots and histograms in **Fig. S1F**).

GC content bias investigation

We assessed the degree of GC content bias in our original data and the degree to which our primer correction step attenuated the bias. First, we grouped restriction fragments into strata according to the GC content of the genome-binding portion of each 5C primer (i.e. the full 5C primer sequence minus the

universal T7/T3 tail). We computed the sums of cis interactions for all primers in each strata and plotted each data point as an enrichment over the average cis interaction sum across all primers (**Fig. S1H**). A comparison of G-C content bias for each of the first three stages of our analysis pipeline demonstrated that primers with extreme GC content are relatively depleted for counts in our raw data and that this bias is attenuated after primer correction (**Fig. S1H**). The attenuation in primer bias in extreme GC content strata is consistent with the goal of our primer correction scheme to push all primers towards equal visibility.

To further investigate the GC bias relationships in our data, we stratified our primer-primer pairs into a 2-D grid of strata depending on the GC content of the upstream and downstream primer comprising the forward-reverse primer pair. We then visualized the enrichment of counts within each stratum, computed as described by Ren and colleagues ([Jin et al., 2013](#)) as:

$$E_{a,b} = \frac{\sum_{i,j \ni l_a < g_i \leq u_a, l_b < g_j \leq u_b, i > j} c_{i,j}}{\sum_{i,j \ni l_a < g_i \leq u_a, l_b < g_j \leq u_b, i > j} \mu}$$

where $E_{a,b}$ is the enrichment value for the abth stratum in the grid, l_a and u_a are the lower and upper GC content limits, respectively, of the ath stratum, l_b and u_b are the lower and upper GC content limits, respectively, of the bth stratum, g_i is the GC content of the ith primer, $c_{i,j}$ is the number of counts for the interaction of the ith primer with the jth primer, and μ is the mean number of counts across all primer-primer pairs.

We generated GC strata heatmaps for raw and primer corrected data (**Fig. S1I**). Although the strata with the most extreme GC contents show less bias after normalization, there was still a noticeable enrichment of counts centered on the 50-60% to 50-60% pairwise GC content range. This result is consistent with previous observations by Ren and colleagues suggesting that there might be a

biologically significant enrichment for 3-D interactions between genomic elements with high GC content levels for cis contacts at distance scales < 2 Mb (Jin et al., 2013).

Comparison of 5C analysis pipeline to alternative approaches

We compared the results from our current 5C data analysis steps to the results of the corresponding steps in our previously published 5C analysis pipeline (**Fig. S2A-D**). In our previous approach, data were not quantile normalized, the distance-dependence background was modeled parametrically with a Weibull distribution, no binning was performed and pvalues were computed via modeling single fragment resolution data with a compound normal-lognormal distribution (Phillips-Cremins et al., 2013).

First, we corrected for primer effects by employing the same primer normalization strategy in our current and original analysis pipelines. The primer correction step attenuated under/over-enriched stripes in the heatmaps, pushing all rows/columns toward equal visibility, independent of whether or not the data were quantile normalized (compare boxplots and heatmaps in **Figs. S1C and S2B**). Second, our 2016 empirical, region-specific distance-dependence models show improved ability to correct for the short-range distance-dependence relationship than our previous 2013 parametric distance-dependence model (compare heatmaps and distance-dependence curves in **Figs. S1E and S2C**). Third, our 2015 binning approach at ~12 kb 'pseudo-fragment resolution' (discussed above) offers key improvements in highlighting the true looping signal vs. noise when compared to our 2013 ~4 kb 'single fragment resolution' maps (compare heatmaps in **Figs. S1D-F and S2C-D**). Finally, our 2016 approach to model distance-corrected interaction frequencies as a continuous random variable with the logistic distribution results in the clear illumination of underlying looping patterns in distance-corrected interaction score heatmaps. Our previous approach modeling single fragment resolution data with a compound normal-lognormal distribution did allow for the identification of a few of the strongest structural features that change dynamically between cell types. However, distance-corrected interaction

score maps from the 2013 pipeline exhibited a much greater degree of spatial noise that obscured many important 3-D interactions (compare heatmaps in **Figs. S1F and S2D**). Finally, we moved the order of our current pipeline steps - conducting quantile normalization after binning, performing the binning step on unlogged data and logging only for visualization – and the resultant heatmaps showed similar results to our current pipeline steps, suggesting that the biological conclusions are robust to the order at which we conduct our pre-processing steps (**Figs. S2E-G**).

Overall, our 5C methods were chosen because they yield highly sensitive and quantitative identification/detection of significant 3-D interactions while exhibiting strong attenuation of primer effects, absence of distance-dependence background signal and minimal distribution skewing due to technical differences in library complexity (**Fig. S1F**).

Principal component analysis

Principal component analysis was performed to scatter the six experimental replicates according to their distance-corrected interaction frequencies at each bin-bin pair. The R `prcomp()` function with active center and scale parameters was used to compute the principal components for our six conditions. We plotted the projection of our six conditions onto the first two principle components as a scatterplot.

Classification of cell type-specific 3-D interactions

To classify cell type-specific 3-D interactions, we generated scatterplots of distance-corrected interaction scores for pairwise combinations of ES cells, NPCs and iPS cells (**Fig. 3A-F**). Specifically, for every 4 kb bin, the minimum distance-corrected interaction score between the two replicates for each cell type was plotted to ensure both replicates must fall above any threshold to be considered for classification. Distance-corrected interaction scores ≤ 3.219 in ES cells, NPCs and iPS cells were classified

as “background” interactions. Interactions for which all cell types had a distance-corrected interaction score ≤ 30 were not considered in the parsing of any 3-D interaction class.

For each pairwise comparison, distance-corrected interaction scores were classified as: (i) ‘present in both cell types’, (ii) ‘present in cell type 1’, (iii) ‘present in cell type 2’, (iv) ‘unable to be differentially assigned with confidence’, or (v) a ‘background’ interaction (i.e. low interaction score) in both cell types (**Fig. 3**). Pairwise interaction classifications were then combined to determine differential interactions among all three cell types.

Reproducible distance-corrected interaction scores $\geq 53.219^*$ in cell type 1 *and* cell type 2 were considered ‘present in both cell types’. Similarly, if the difference between the minimum interaction scores of both cell types did not exceed 14, the interaction was also classified as ‘present in both cell types’. Interactions with differences between the distance-corrected interaction scores of the two cell types greater than 14 that also had interaction scores ≥ 43.219 but < 53.219 in all cell types were removed from consideration because of uncertainty whether to classify them as constitutive or cell-type specific. The remaining interactions (i.e. at least one cell type interaction score > 30 , at least one cell type interaction score < 43.219 , and the difference between the minimum replicates of the cell types > 14) were classified as ‘present in cell type 1’ if the interaction score in ‘cell type 1’ was greater and ‘present in cell type 2’ if the interaction score in ‘cell type 2’ was greater.

Pairwise classifications were combined to construct the 3-D interaction categories between the three cell types. Interactions that were considered ‘present in both cell types’ in all pairwise comparisons were parsed into the “constitutive” (grey class) 3-D interaction category. Interactions that were classified as ‘present in both ES and iPS cells’ but were found to be ES- and iPS-specific when comparing these cell types to NPCs were parsed into the “ES-iPS” (purple class) 3-D interaction category. Interactions that were classified as ‘present in ES cells’ when thresholded against both iPS and NPC distance-corrected interaction scores were parsed into the “ES-only” (red class) 3-D interaction

category. Similarly, interactions classified as ‘present in both iPS cells and NPCs’ but were found to be iPS- and NPC-specific in comparison with ES cells were parsed into the “NPC-iPS” (blue class) 3-D interaction category. Similarly, ‘present in both ES cells and NPCs’ interactions were parsed into the “ES-NPC” (yellow class) 3-D interaction category if the interactions were not present when compared to iPS cells. Finally, interactions classified as ‘present in iPS cells’ when thresholded against both ES cells and NPCs were parsed into the “iPS-only” (orange class) 3-D interaction category, and interactions classified as ‘present in NPCs’ when thresholded against both ES and iPS cells were parsed into the “NPC-only” (green class) 3-D interaction category. We subsequently removed any interaction that was classified but had a size less than 20 kb. Additionally, we removed interactions with size greater than 400 kb if they did not form an adjacency cluster (See “Interaction Adjacency Clustering” below) of at least 5 pixels. The bin numbers of the interactions whose interaction scores are presented in barplots in **Figs. 4, S5, 5, 6, S7** can be found in **Table S7**.

*Note on thresholds: $53.219 = -10 * \log_2(0.025)$; $43.219 = -10 * \log_2(0.05)$; $30 = -10 * \log_2(0.125)$; $3.219 = -10 * \log_2(0.8)$, thus interaction scores of 53.219, 43.219, 30, and 3.219 correspond to interaction p-values of 0.025, 0.05, 0.125, and 0.8, respectively.

Empirical false discovery rate calculation

Justification of strategy

To compute an empirical false discovery rate (eFDR) for our interaction score thresholds, we employed a strategy in which we simulated 5C experiments consisting of three identical cellular conditions with two replicates each. The motivation/rationale for this strategy was that we wanted to determine how many 3-D interactions would be called by our thresholding/classification scheme (**Figs. 3, S3**) when comparing three cellular states (n=2 biological replicates each) that have been simulated to contain equivalent 3-D architecture. For example, we simulated ES1_Rep1, ES1_Rep2, ES2_Rep1, ES2_Rep2, ES3_Rep1, and

ES3_Rep2, where all six replicates were generated from the same model (modeled based on our experimental ES data, discussed below). After the creation of the simulated replicates, ES1, ES2, and ES3 were treated as the distinct conditions for categorization purposes. By quantifying the number of interactions that we would expect by chance to pass our thresholds (discussed above), we can compute an eFDR for each 3-D interaction class identified when comparing ES vs. NPC vs. iPS cells.

Model generation – mean parameter estimation

First, we generated simulations of 5C data. To generate each of the simulations, we created three independent models, each of which was based on one of three cell type subsets (ES, NPC, iPS) of our experimental data. For each of these three models, we first computed a mean parameter by calculating the mean distance-corrected interaction frequency for that bin-bin pair among the two experimental replicates for the cell type the model was based on. We represent this mathematically as:

$$\mu_{c,s,i,j} = \frac{\sum_{r=1}^2 f_{c,r,s,i,j}}{2}$$

where $\mu_{c,s,i,j}$ is the mean distance-corrected interaction frequency for the ij th bin-bin pair of the s th region in the model for cell type c and $f_{c,r,s,i,j}$ is the distance-corrected interaction frequency for the ij th bin-bin pair of the s th region in the experimental data for replicate r in cell type c .

Model generation – estimating the mean-variance relationship

Second, to obtain reasonable estimates for variance, we estimated a region-specific mean-variance relationship by performing a linear regression on the scatterplot of mean versus sample standard deviation of the distance-corrected interaction frequency for each bin-bin pair in each region among the two experimental replicates for the cell type being considered. This linear regression allowed us to compute a predicted standard deviation given a mean as:

$$\hat{\sigma}_{c,s,i,j} = m_{c,s}\mu_{c,s,i,j} + b_{c,s}$$

where $\hat{\sigma}_{c,s,i,j}$ is the predicted standard deviation of distance-corrected interaction frequency for the ij th bin-bin pair of the s th region in the model for cell type c , $\mu_{c,s,i,j}$ is the mean distance-corrected interaction frequency for the ij th bin-bin pair of the s th region in the model for cell type c , and $m_{c,s}$ and $b_{c,s}$ are the slope and y-intercept parameters obtained from the linear regression of mean versus standard deviation for the s th region in the experimental data from cell type c .

Model generation – variance parameter estimation

Third, we used the mean-variance relationship to estimate the standard deviation parameter. We set the simulation standard deviation at each bin-bin pair to a linear combination of the observed standard deviation for that bin-bin pair in the experimental data for that cell type and our predicted standard deviation at that bin-bin pair as follows:

$$\sigma_{c,s,i,j} = \alpha \hat{\sigma}_{c,s,i,j} + \beta \sqrt{\frac{1}{2} \sum_{r=1}^2 (f_{c,r,s,i,j} - \mu_{c,s,i,j})^2}$$

where $\sigma_{c,s,i,j}$ is the final standard deviation parameter for ij th bin-bin pair of the s th region in the model for cell type c , $\sqrt{\frac{1}{2} \sum_{r=1}^2 (f_{c,r,s,i,j} - \mu_{c,s,i,j})^2}$ is the sample standard deviation of the distance-corrected interaction frequencies of the ij th bin-bin pair of the s th region in the experimental data from cell type c (r indexes the replicates), and α and β are constants chosen to ensure that the noise in the data generated by the model closely approximates the noise in the actual experimental data.

Simulations

Fourth, after computing the model parameters $\mu_{c,s,i,j}$ and $\sigma_{c,s,i,j}$, we generated simulated 5C experiments by drawing simulated distance-corrected interaction frequencies from a normal distribution with mean, variance parameters as follows:

$$F_{c,s,i,j} \sim N(\mu_{c,s,i,j}, \sigma_{c,s,i,j})$$

where $F_{c,s,i,j}$ is a random variable representing the simulated distance-corrected interaction frequency for the ij th bin-bin pair of the s th region for a simulation of cell type c and $\mu_{c,s,i,j}$ and $\sigma_{c,s,i,j}$ are the mean distance-corrected interaction frequency and the final standard deviation parameter, respectively, for the ij th bin-bin pair of the s th region in the model for cell type c . We chose a normal distribution in accordance with our assumption that the replicate-to-replicate noise for repeated measurement of the same exact bin-bin interaction would be normally distributed.

Monte Carlo, pvalue calculation, classification

Fifth, we used the above approach to generate six simulated 5C experiments from the same model, and then applied our logistic fits and our thresholding/classification scheme (described above) to each of the simulations. As in our real 5C data, we modeled the distribution of simulated distance-corrected interaction frequencies with a logistic distribution parameterized independently for each region. Logistic fits were used to assign p-values to every bin-bin pair in the simulation. P-values were converted to interaction scores as described above. The six independently constructed simulations were grouped into three equivalent categories containing two replicates each and subjected to the same thresholding/classification scheme as our experimental data. The number of simulated bin-bin pairs that were categorized into each of our 3-D interaction classes was recorded. This process was repeated 1000 times for each of our three cell types, and the numbers of simulated bin-bin pairs falling into each category were averaged across the 1000 trials and across the three cell types. We confirmed that our simulations fairly recapitulated the noise seen in the experimental data by comparing Spearman's and Pearson's correlation coefficients as well as histograms and empirical cumulative distribution functions for our simulations to those we observed in our experimental data.

Computing the false discovery rates for each 3-D interaction class

Finally, we computed false discovery rates. Because the six simulated experiments represent simulated biological replicates, any bin-bin pair that was categorized into any category other than constitutive or background represents a false positive. Therefore, we estimated the false positive rate (FPR) for our thresholds for each of the other categories as the number of simulated bin-bin pairs falling into that category divided by the total number of bin-bin pairs in the simulation. Mathematically, this is represented as:

$$\text{FPR}_t^{\text{sim}} = \frac{\bar{n}_t^{\text{sim}}}{N}$$

where $\text{FPR}_t^{\text{sim}}$ is the simulation false positive rate for category t , \bar{n}_t^{sim} is the average number of bin-bin pairs categorized into category t across all simulations, and N is the total number of bin-bin pairs in each simulation. We then assumed that the FPR for our simulation was a good estimate for the FPR in the categorization of our real experimental data.

$$\text{FPR}_t^{\text{sim}} \approx \text{FPR}_t^{\text{exp}}$$

where $\text{FPR}_t^{\text{sim}}$ is the simulation false positive rate for category t and $\text{FPR}_t^{\text{exp}}$ is the experimental false positive rate for category t . Our real experimental data and our simulations had the same number of bins and therefore the same number of bin-bin pairs to be categorized. Therefore, we estimated that for each category other than background and constitutive, the number of false positives observed in our simulations was equal to the number of false positives in our experimental data.

$$\text{FPR}_t^{\text{sim}} \approx \text{FPR}_t^{\text{exp}} \Rightarrow \bar{n}_t^{\text{sim}} \approx \text{FP}_t^{\text{exp}}$$

where \bar{n}_t^{sim} is the average number of bin-bin pairs categorized into category t across all simulations and FP_t^{exp} is the experimental number of false positives in category t .

We then estimated the false discovery rate (FDR) in our experimental data by dividing this estimated number of false positives by the total number of bin-bin pairs declared significant in the experimental data. Mathematically, this is represented as:

$$\text{FDR}_t^{\text{exp}} = \frac{\text{FP}_t^{\text{exp}}}{n_t^{\text{exp}}} \approx \frac{\bar{n}_t^{\text{sim}}}{n_t^{\text{exp}}}$$

where n_t^{exp} is the number of bin-bin pairs categorized into category t in the experimental data. Because a different number of bin-bin pairs were declared significant in different categories, we computed different FDRs for different categories (**Fig. 3H-I**).

6 sample vs 10 sample 5C data processing

5C data was processed either in a 6 sample batch, which includes only ES, NPC, and iPS replicates, or a 10 sample batch, which includes all 2i replicates in addition to the core 6 samples. Cell-type specific 3D interactions were classified using the ‘6-sample’ group of ES, NPC, and iPS replicates. In instances where heatmaps are displayed for only these three cell types (i.e. Fig. 4, S5B, S6), we use ‘6-sample’ normalized data, whereas when data is displayed for all 5 cell types (i.e. Fig. 5, S5F, 6, S7), we present ‘10-sample’ normalized data.

Interaction adjacency clustering

Spatially adjacent interactions of the same classification were iteratively grouped into clusters in order to quantify the number of interaction clusters present in our data. For a given classified pixel, we queried if that pixel was adjacent to an already identified cluster – if adjacent, the pixel was appended to that cluster - if not adjacent, the pixel was assigned its own cluster. Clusters of the same classification that were directly adjacent to themselves at the end of the iterative process were merged.

ChIPseq peakcalling

A summary of all ChIP-seq data sets re-analyzed in this study is provided in **Table S4**. Data was downloaded from GEO (<http://www.ncbi.nlm.nih.gov/geo/>). Sequences were aligned to NCBI Build 37 (UCSC mm9) using default parameters (-v1 -m1) in Bowtie. Only sequences that mapped uniquely to the genome were used for further analysis. Model-based Analysis for ChIP Sequencing (MACS) was used for peak calling (<http://liulab.dfci.harvard.edu/MACS/00README.html>). For CTCF ChIPseq, default parameters were used with a p-value cutoff of $p < 1 \times 10^{-8}$. For histone modification ChIPseq (e.g. H3K4me1, H3K27ac, H3K4me3), we skipped the model-building step by calling the parameter --no model with a p-value cutoff of either $p < 1 \times 10^{-8}$, $p < 1 \times 10^{-6}$ or $p < 1 \times 10^{-4}$.

Parsing ES-specific and NPC-specific genes

Normalized RNA-seq counts were parsed by fold change between ES cells and NPCs into ES-specific and NPC-specific gene expression categories. Genes that were at least two-fold upregulated in ES cells compared to NPCs were classified as ES-specific, whereas genes that were at least two-fold upregulated in NPCs compared to ES cells were classified as NPC-specific. ES-specific genes were further refined by required overlap with high-confidence H3K27ac signal (peaks called at $p < 1 \times 10^{-6}$) in ES cells (found in **Table S5**). NPC-specific genes were further refined by required overlap with high-confidence H3K27ac signal (peaks called at $p < 1 \times 10^{-4}$) in NPCs (found in **Table S6**). Inactive genes were parsed by identifying those genes falling within queried 5C regions that did not exhibit H3K27ac signal (peaks called at $p < 1 \times 10^{-2}$) in either ES cells or NPCs.

Parsing ES-specific and NPC-specific enhancers

H3K27ac peaks (ES, $p < 1 \times 10^{-6}$; NPC, $p < 1 \times 10^{-4}$) were merged if they fell within 500 bp end-to-end distance of each other. NPC H3K27ac was peak-called at a lower threshold than the ES H3K27ac after

visual observation that there appeared to be a smaller dynamic range of the NPC H3K27ac ChIPseq data between the active and inactive state. ES-specific enhancers were defined by overlap between merged H3K27ac peaks and H3K4me1 peaks ($p < 1 \times 10^{-4}$) in ES cells and the absence H3K27ac in NPCs (defined by subtraction of low-confidence NPC-binding sites for H3K27ac ($p < 1 \times 10^{-2}$)). NPC-specific enhancers were defined by overlap between merged H3K27ac peaks and H3K4me1 peaks ($p < 1 \times 10^{-4}$) in NPCs and the absence H3K27ac in ES cells (defined by subtraction of low-confidence ES-binding sites for H3K27ac ($p < 1 \times 10^{-2}$)). To ensure subtraction of all potential genes, it was required that parsed ES-specific and NPC-specific enhancers did not fall within 2 kb of a transcription start site. A summary of all ChIP-seq datasets utilized can be found in **Table S4**.

Parsing ES-specific and NPC-specific CTCF sites

ES-specific CTCF was defined by the presence of high-confidence binding sites ($p < 1 \times 10^{-8}$) in ES cells and the absence of CTCF in NPCs (defined by subtraction of low-confidence NPC-binding sites for CTCF ($p < 1 \times 10^{-2}$)). NPC-specific CTCF was defined by the presence of high-confidence binding sites ($p < 1 \times 10^{-8}$) in NPCs and the absence of CTCF in ES cells (defined by subtraction of low-confidence ES-binding sites for CTCF ($p < 1 \times 10^{-2}$)). Constitutive CTCF was defined by the presence of high-confidence binding sites ($p < 1 \times 10^{-8}$) in both cell types. A summary of all ChIP-seq datasets utilized can be found in **Table S4**.

Computing enrichments

Annotation intersections

For each bin in each of our 5C regions, we identified the genomic elements that overlapped that bin, or the neighboring 2 bins on either side (matching our 20 kb window, see *Contact matrix binning* above); the bin was then considered to 'contain' those genomic elements. Next, to interrogate pairwise connections between distinct genomic elements, we found all the bin-bin pairs whose upstream bin

contained the first type of genomic element and whose downstream bin contained the second type of genomic element, or the reverse. For each of these bin-bin pairs, we checked which interaction classification category, if any, they fell into. We recorded the total number of intersections of this interaction class for every pair of types of genomic elements being considered and for every category in our interaction categorization scheme. By considering pairs of genomic elements in this way, we attempted to identify instances of one type of genomic element interacting with another type of genomic element. In our analysis, we included pairs of the same type of genomic elements (e.g., ES-specific genes interacting to ES-specific genes). We also created an artificial type of genomic element (referred to as “wildcard” element) that was present in every bin of every 5C region. Including this “wildcard” genomic element allowed us to query interactions that involved one specified type of genomic element interacting with any other location, irrespective of what genomic elements were present on the other side (see **Fig. 6D**).

Computing percentage incidence, fold-enrichment above background, and p-values

Next, we divided the interaction counts for each pair of genomic element classes in each interaction category by the total number of interactions in that category to obtain the percentage of interactions in that category that involved an interaction between the two types of genomic elements in the pair. We then computed a fold-enrichment for each interaction type’s percentage above the background interaction type’s percentage. Finally, we computed p-values for the enrichment by applying Fisher’s exact test to the contingency table below:

Number of interactions in the selected category involving the two selected annotations	Number of interactions in the background category involving the two selected annotations	Number of interactions in either the selected or the background category involving the two selected annotations
--	--	---

Number of interactions in the selected category not involving the two selected annotations	Number of interactions in the background category not involving the two selected annotations	Number of interactions in the selected or the background category not involving the two selected annotations
Total number of interactions in the selected category	Total number of interactions in the background category	

We used the p-value for the particular tail of the distribution that matched the direction of the enrichment (i.e., the right-tail p-value if the interaction was enriched over background, and the left-tail p-value if the interaction was depleted below background, generally equivalent to the lesser of the two p-values). P-values were computed using the `scipy.stats.fisher_exact` function from the `scipy` Python computational library.

Visualizing enrichments

These enrichment quantification strategies were employed to investigate the intra-regional interactions of a selected annotation on either side of the interaction (via our “wildcard” annotation), and interactions between one selected annotation and another selected annotation falling within each interaction classification. Enrichments were visualized as either bar plots (showing the percentages of interactions between a pair of annotations falling into each of the interaction categories with the height of the different bars) or heat maps (with the color representing the log base 2 fold-enrichment of a certain interaction category above background for the percentage of interactions between a pair of annotations and the text showing the upper bound for the p-value for that enrichment).

Computing connectivity

To compute the ‘connectivity’ metric for each genomic annotation (**Fig. 7**), we first summed the number of significant interactions present in a given cell type that contained that annotation on at least one side

of the interaction. A 'connectivity' value was computed by dividing the total number of interactions made by each annotation by the total number of interactions called significant in that cell type. For example, for the "ES enhancers in ES cells" category, we counted the number significant interactions that intersected an ES enhancer and were categorized as either ES only, ES-iPS, ES-NPC, or constitutive (the four interaction classes present in ES cells); this sum was then divided by the total number of interactions categorized as ES only, ES-iPS, ES-NPC, or constitutive.

Supplemental References

- Bau, D., Sanyal, A., Lajoie, B.R., Capriotti, E., Byron, M., Lawrence, J.B., Dekker, J., and Marti-Renom, M.A. (2011). The three-dimensional folding of the alpha-globin gene domain reveals formation of chromatin globules. *Nat Struct Mol Biol* *18*, 107-114.
- Bolstad, B.M., Irizarry, R.A., Astrand, M., and Speed, T.P. (2003). A comparison of normalization methods for high density oligonucleotide array data based on variance and bias. *Bioinformatics* *19*, 185-193.
- Bullard, J.H., Purdom, E., Hansen, K.D., and Dudoit, S. (2010). Evaluation of statistical methods for normalization and differential expression in mRNA-Seq experiments. *BMC Bioinformatics* *11*, 94.
- Creyghton, M.P., Cheng, A.W., Welstead, G.G., Kooistra, T., Carey, B.W., Steine, E.J., Hanna, J., Lodato, M.A., Frampton, G.M., Sharp, P.A., *et al.* (2010). Histone H3K27ac separates active from poised enhancers and predicts developmental state. *Proc Natl Acad Sci U S A* *107*, 21931-21936.
- Dekker, J., Rippe, K., Dekker, M., and Kleckner, N. (2002). Capturing chromosome conformation. *Science* *295*, 1306-1311.
- 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.
- Dostie, J., and Dekker, J. (2007). Mapping networks of physical interactions between genomic elements using 5C technology. *Nat Protoc* *2*, 988-1002.
- Dostie, J., Richmond, T.A., Arnaout, R.A., Selzer, R.R., Lee, W.L., Honan, T.A., Rubio, E.D., Krumm, A., Lamb, J., Nusbaum, C., *et al.* (2006). Chromosome Conformation Capture Carbon Copy (5C): a massively parallel solution for mapping interactions between genomic elements. *Genome Res* *16*, 1299-1309.
- Eminli, S., Utikal, J., Arnold, K., Jaenisch, R., and Hochedlinger, K. (2008). Reprogramming of neural progenitor cells into induced pluripotent stem cells in the absence of exogenous Sox2 expression. *Stem Cells* *26*, 2467-2474.
- Gheldof, N., Smith, E.M., Tabuchi, T.M., Koch, C.M., Dunham, I., Stamatoyannopoulos, J.A., and Dekker, J. (2010). Cell-type-specific long-range looping interactions identify distant regulatory elements of the CFTR gene. *Nucleic Acids Res* *38*, 4325-4336.
- Jin, F., Li, Y., Dixon, J.R., Selvaraj, S., Ye, Z., Lee, A.Y., Yen, C.A., Schmitt, A.D., Espinoza, C.A., and Ren, B. (2013). A high-resolution map of the three-dimensional chromatin interactome in human cells. *Nature* *503*, 290-294.
- Lajoie, B.R., van Berkum, N.L., Sanyal, A., and Dekker, J. (2009). My5C: web tools for chromosome conformation capture studies. *Nat Methods* *6*, 690-691.

Langmead, B. (2010). Aligning short sequencing reads with Bowtie. *Curr Protoc Bioinformatics Chapter 11*, Unit 11 17.

Meissner, A., Mikkelsen, T.S., Gu, H., Wernig, M., Hanna, J., Sivachenko, A., Zhang, X., Bernstein, B.E., Nusbaum, C., Jaffe, D.B., *et al.* (2008). Genome-scale DNA methylation maps of pluripotent and differentiated cells. *Nature* *454*, 766-770.

Mikkelsen, T.S., Ku, M., Jaffe, D.B., Issac, B., Lieberman, E., Giannoukos, G., Alvarez, P., Brockman, W., Kim, T.K., Koche, R.P., *et al.* (2007). Genome-wide maps of chromatin state in pluripotent and lineage-committed cells. *Nature* *448*, 553-560.

Phillips-Cremins, J.E., Sauria, M.E., Sanyal, A., Gerasimova, T.I., Lajoie, B.R., Bell, J.S., Ong, C.T., Hookway, T.A., Guo, C., Sun, Y., *et al.* (2013). Architectural protein subclasses shape 3D organization of genomes during lineage commitment. *Cell* *153*, 1281-1295.

Rais, Y., Zviran, A., Geula, S., Gafni, O., Chomsky, E., Viukov, S., Mansour, A.A., Caspi, I., Krupalnik, V., Zerbib, M., *et al.* (2013). Deterministic direct reprogramming of somatic cells to pluripotency. *Nature* *502*, 65-70.

Rao, S.S., 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.

Stadler, M.B., Murr, R., Burger, L., Ivanek, R., Lienert, F., Scholer, A., van Nimwegen, E., Wirbelauer, C., Oakeley, E.J., Gaidatzis, D., *et al.* (2011). DNA-binding factors shape the mouse methylome at distal regulatory regions. *Nature* *480*, 490-495.

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

van Berkum, N.L., and Dekker, J. (2009). Determining spatial chromatin organization of large genomic regions using 5C technology. *Methods Mol Biol* *567*, 189-213.