
Having Fun in Hi-C Data Analysis

Fengling Chen

Supervisor: Michael Q. Zhang & Yang Chen

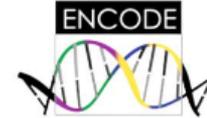
MOE Key Laboratory of Bioinformatics

Bioinformatics Division and Center for Synthetic and Systems Biology, BNRist

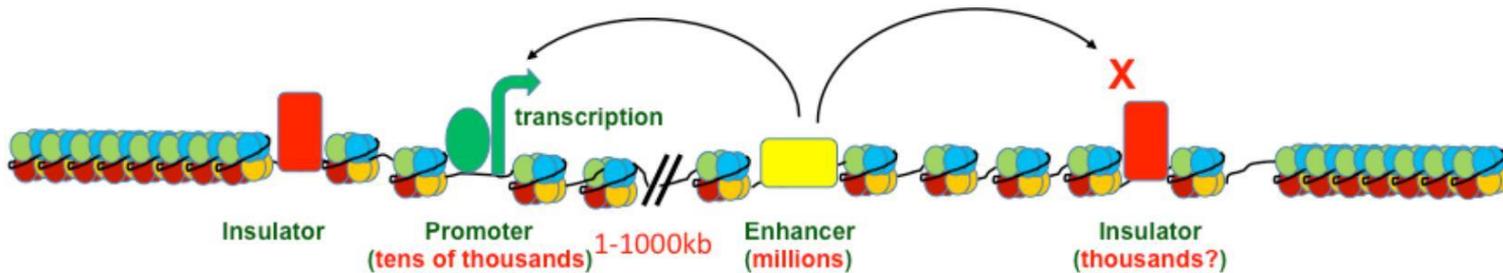
Department of Automation, Tsinghua University

5th January 2019

Background



Human genome: 20,000 genes ; 100,000 promoters ;
500,000 enhancers; 5,000,000 regulatory elements



In development and disease, how these elements interact to regulate gene expression ?

Background

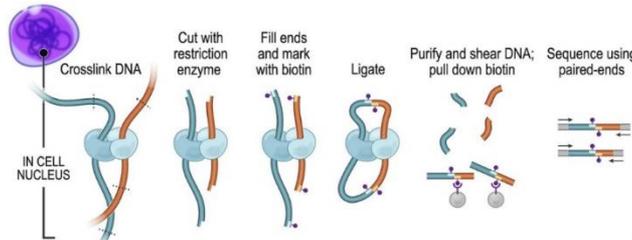


SIMPSONS CONTACT MAP

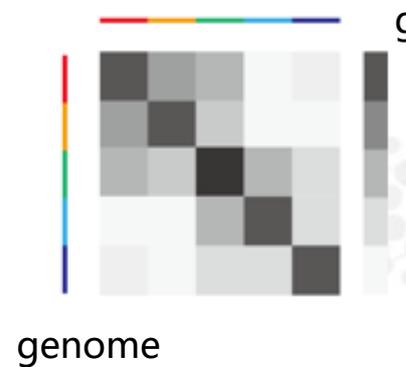
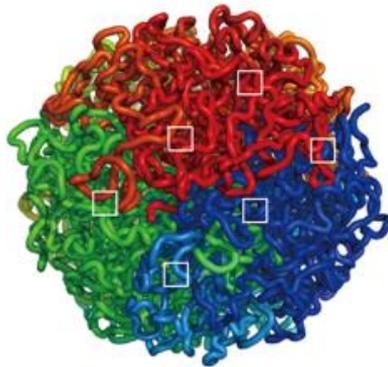
of Pictures Together



	2	0	1	2	1	0	1	0	0
	0	3	2	1	0	0	0	0	0
	1	2	16	6	5	4	11	1	1
	2	1	6	8	6	3	4	0	0
	1	0	5	6	8	4	5	1	0
	0	0	4	3	4	5	5	0	0
	1	0	11	4	5	5	13	1	1
	0	0	1	0	1	0	1	2	1
	0	0	1	0	0	0	1	1	1



Hi-C method



Contact matrix (pairwise)

Background

1. Hierarchical folding of chromatin

Chromatin territory

Compartmentalization of megabase-scale chromatin.

Topologically associating domains;

Chromatin loops

2. Mechanisms to organize chromatin in 3D

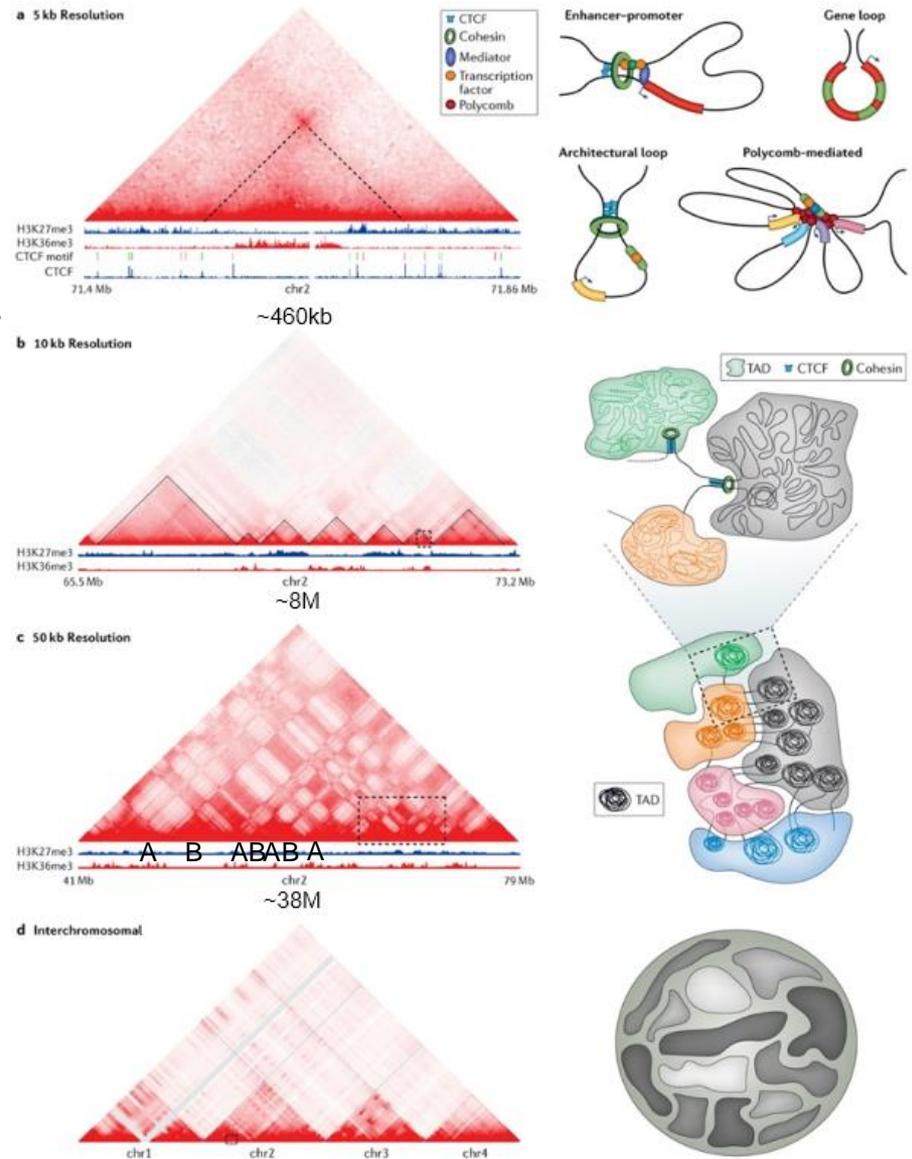
Architectural proteins: CTCF; cohesin; mediator

Non-coding RNAs

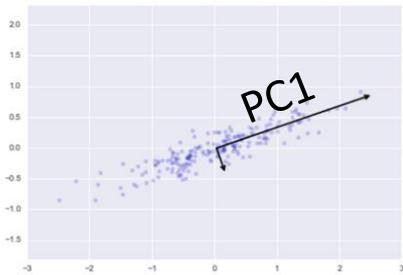
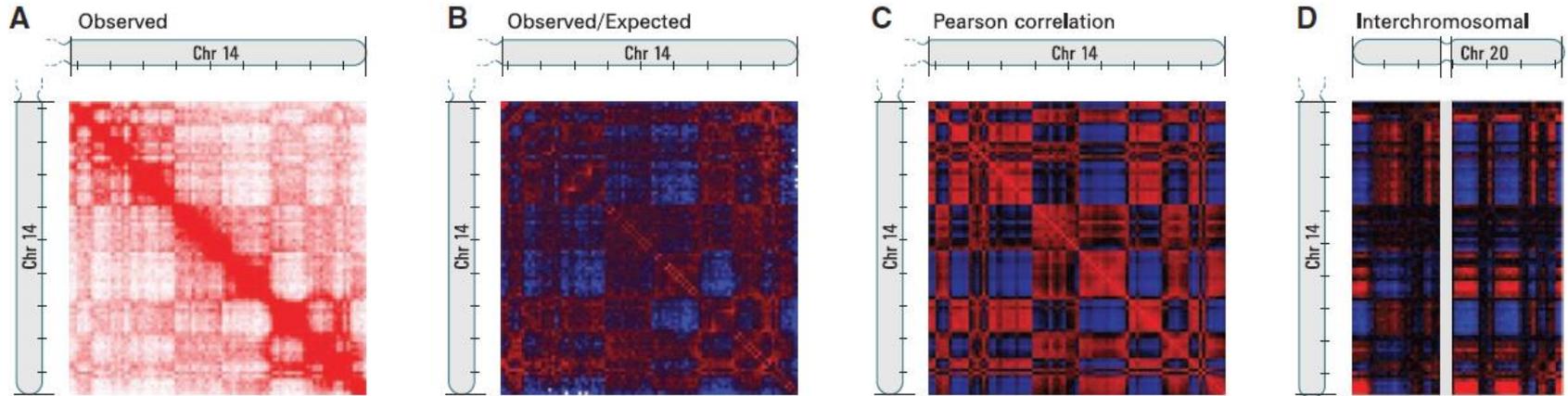
Histone modification

Phase separation

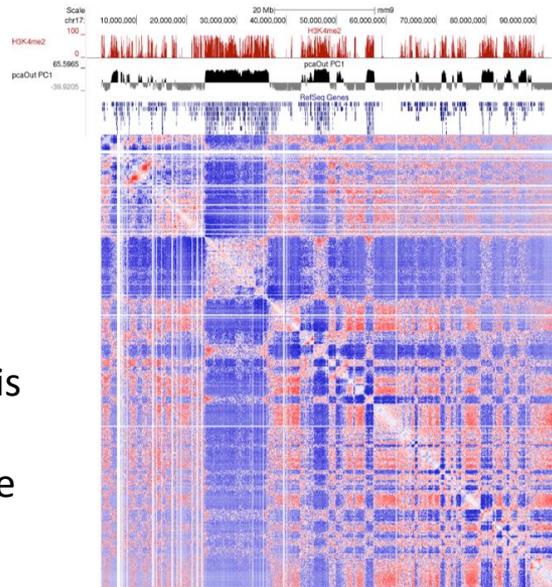
How to extract information on Hi-C maps ?



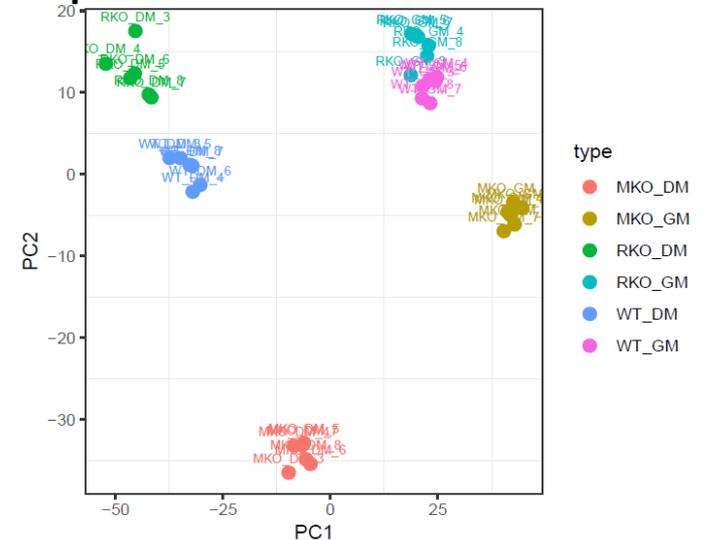
1. Compartment



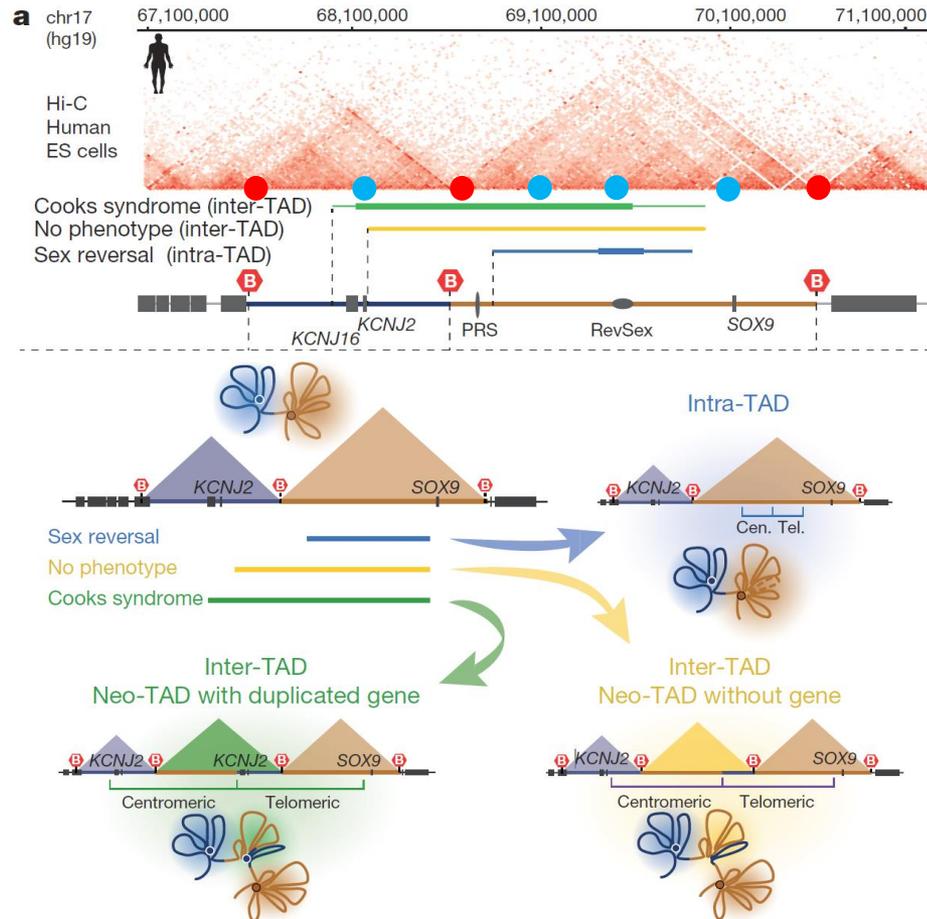
The basic idea behind **PCA** is to redefine the coordinate system such the data can be "described" with as few dimensions as possible.



PCA of PC1 values in muscle development



2. CDB

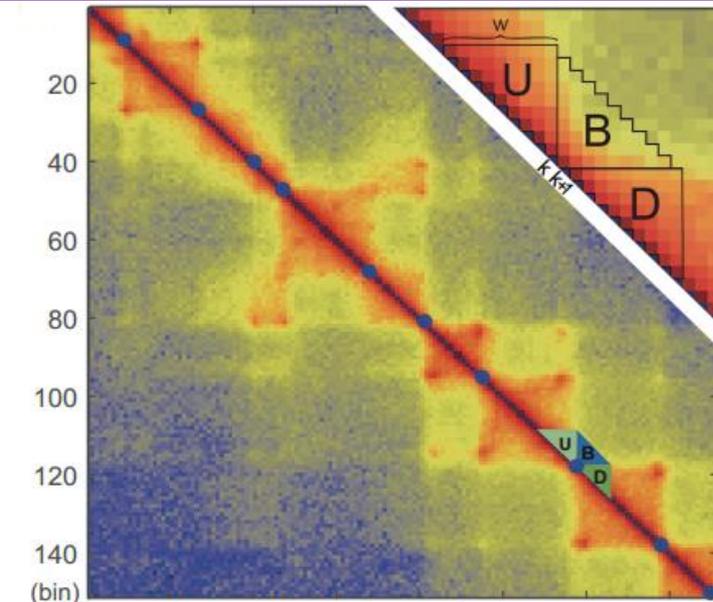


contact domain boundaries (CDBs) includes TAD boundaries and sub-TAD boundaries.

TADs: invariant/conserved; absolute insulated

sub-TADs : varied ; cell-type specific gene regulation; relatively insulated

HiCDB Method



- ✓ Calculate relative insulation (RI) under different window size

$$RI(w, s) = \frac{U(w, s) + D(w, s) - B(w, s)}{U(w, s) + D(w, s) + B(w, s)}$$

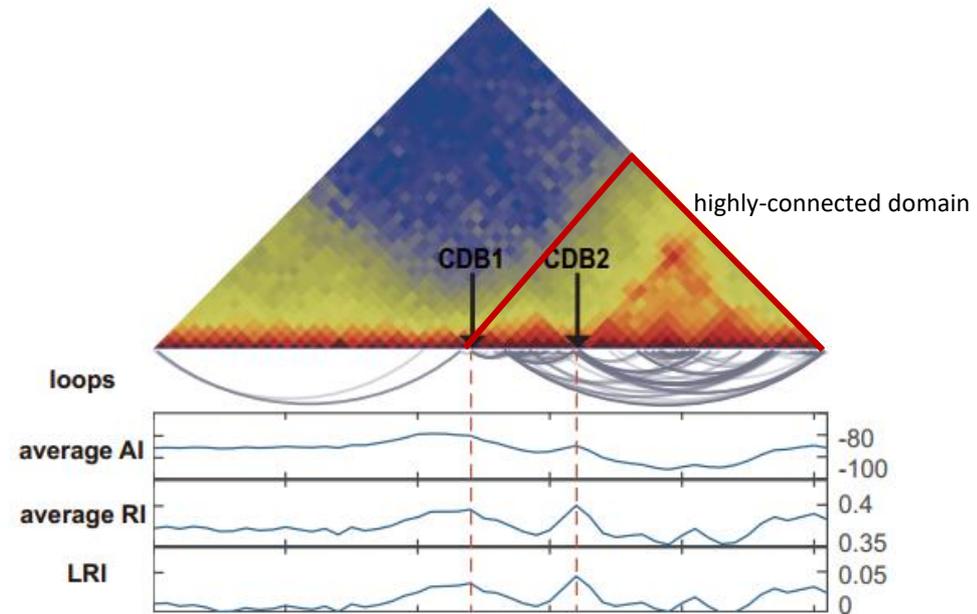
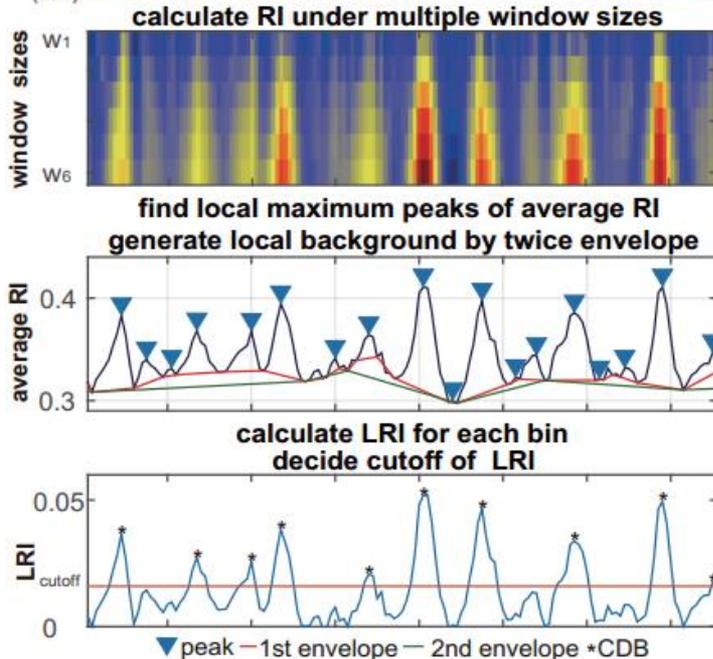
- ✓ Peak detection on average RI

$$\overline{RI}(s) = \frac{1}{w_n - w_1} \sum_{w=w_1}^{w_n} RI(w, s)$$

- ✓ Remove background

$$LRI(s) = \overline{RI}(s) - lower_envelope(lower_envelope(\overline{RI}(s)))$$

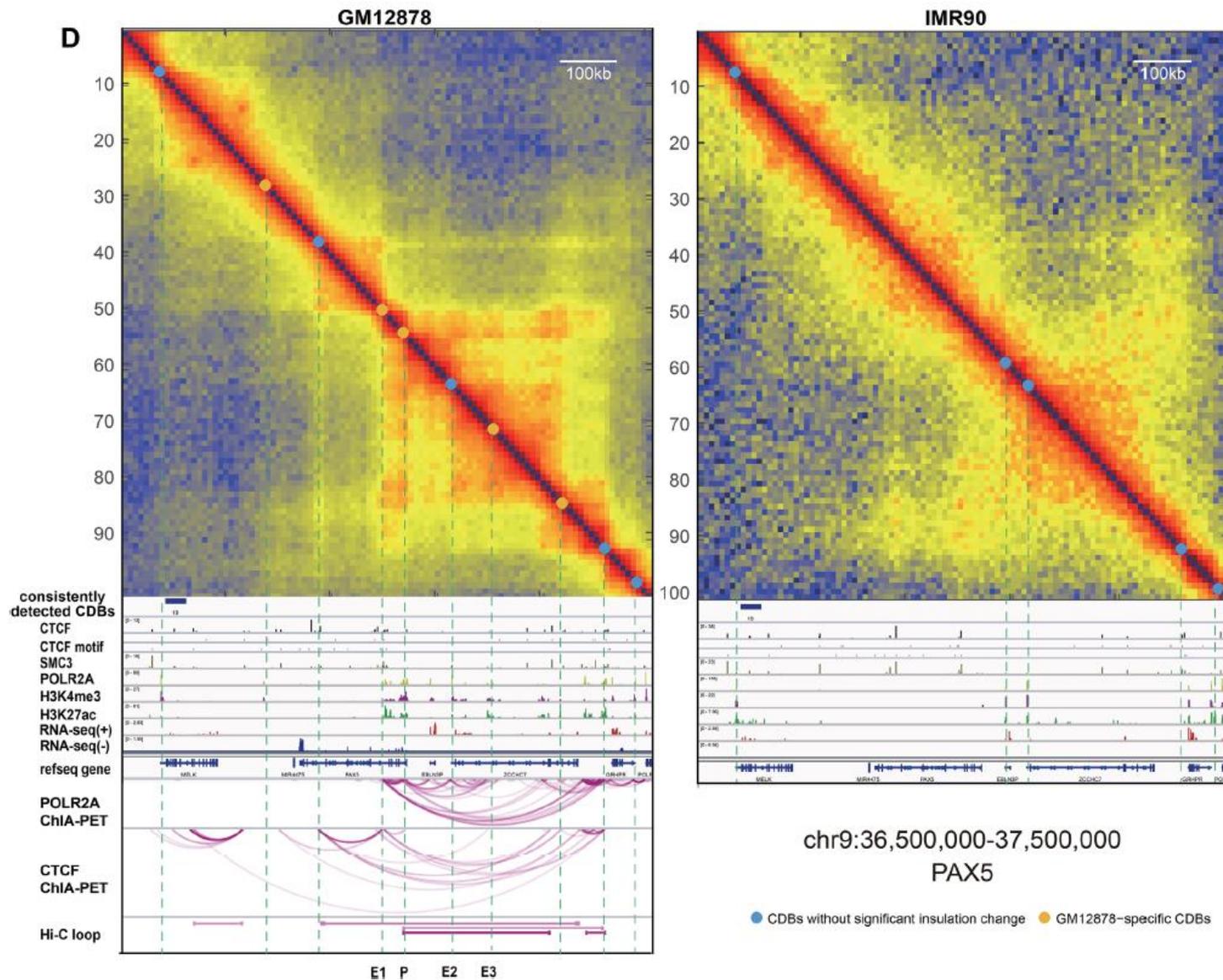
- ✓ Cut-off chosen (GSEA-like method)



aRI or LRI can be compared genome-wide

aRI or LRI help to find CDBs under highly-connected domains.

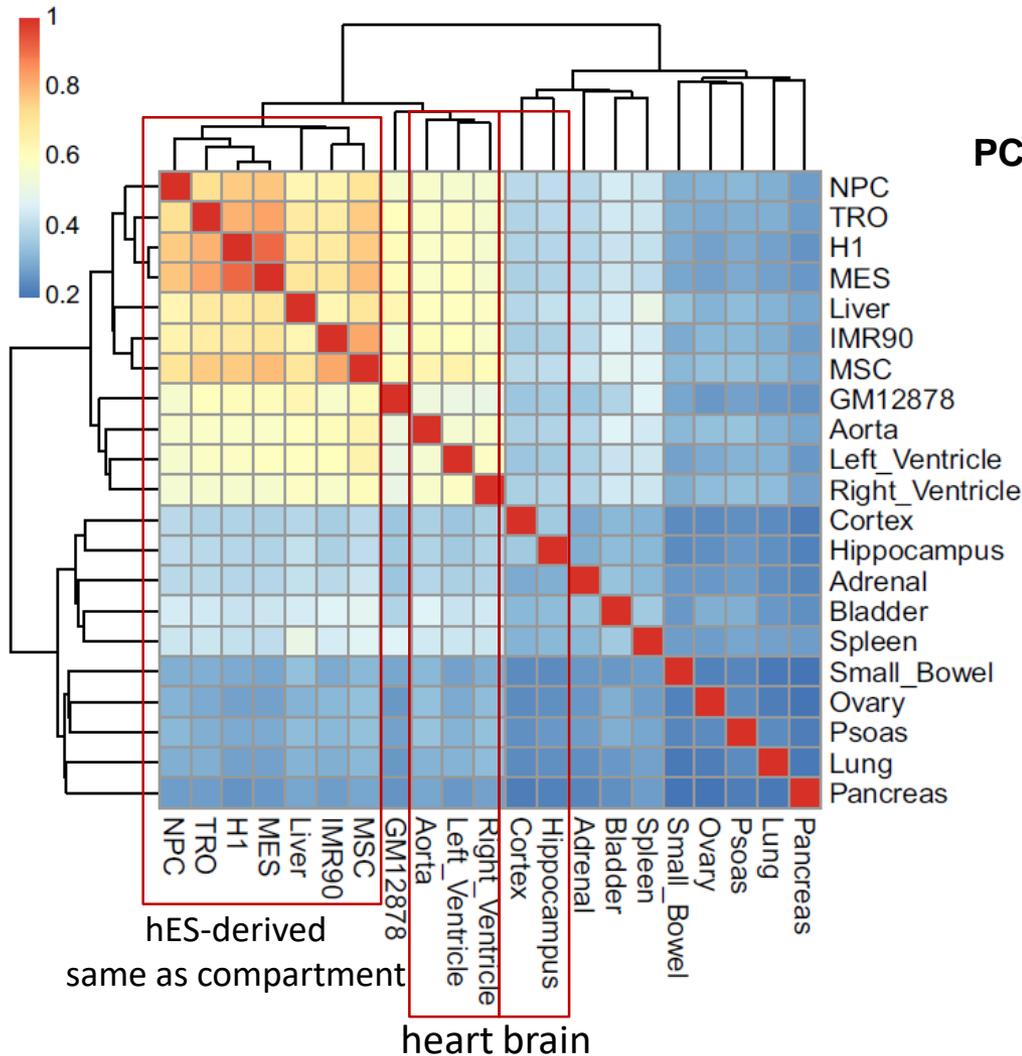
Differential CDB example



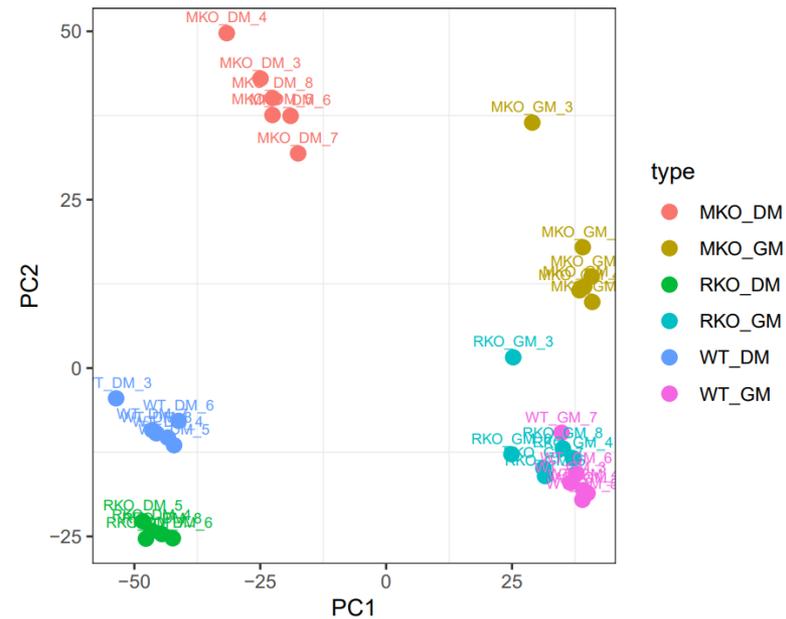
Cell-type-specific CDBs correlate with cell-type-specific histone modification and gene activation.

Clustering based on aRI values

Spearman correlation of CDB aRI values



PCA of CDB aRI values in muscle development



aRI score is more sensitive to sequencing depth than compartment.

3. Loop Calling

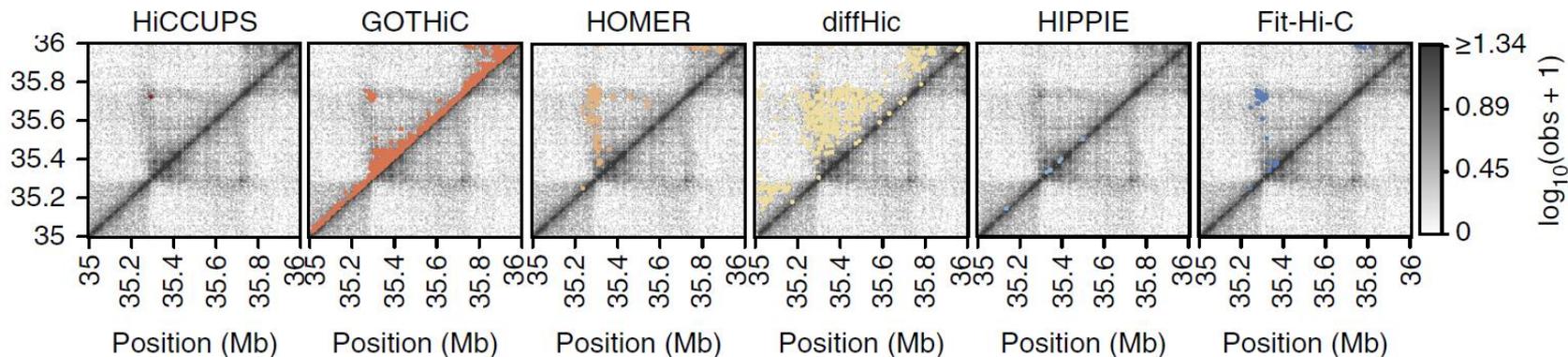
Table 3 Methods of identifying significant interactions

表 3 显著交互作用鉴别方法

算法或软件	模型	针对性	特点
HIPPIE ^[42]	负二项分布	显著交互作用	在酶切片段分辨率下，鉴别近距离交互作用
Fit-HiC ^[80]	二项分布	显著交互作用	对交互作用与距离的关系两次拟合
GOTHIC ^[81]	二项分布	显著交互作用	利用覆盖率估计期望的交互作用
HiC-DC ^[82]	零截尾负二项分布	显著交互作用	模型考虑了 Hi-C 数据的多零和高散度的特性
HOMER ^[43]	二项分布	显著或差异交互作用	多功能集成软件
HMRFBayes ^[83] & FastHiC ^[84]	负二项分布、隐 Markov 随机场	显著交互作用	考虑相邻交互作用的影响
HiCCUPS ^[31]	泊松过程	显著交互作用	去除 TAD 结构的影响
PSYCHIC ^[85]	对数正态分布	显著交互作用	基于 TAD 结构构建背景模型
CHiCAGO ^[86]	负二项方分布、泊松分布	显著交互作用	针对 Capture Hi-C 实验数据
Dynamic Interactions ^[13]	二项分布	差异交互作用	利用生物学重复
HiBrowse ^[87] & DiffHiC ^[88]	负二项分布	差异交互作用	借助 edgeR, 利用生物学重复
FIND ^[89]	空间泊松过程	差异交互作用	考虑相邻交互作用的关联性

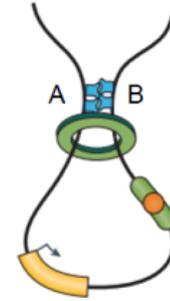
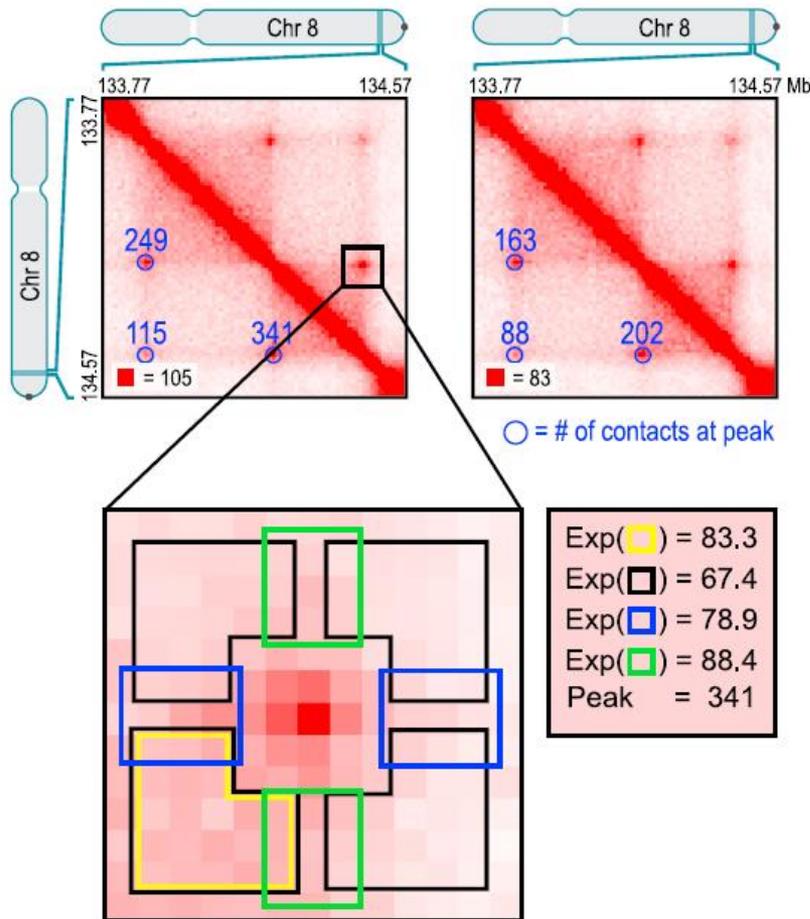
local background model

Global background model



3. Loop Calling

Local background model



Step1. build background model.

$$\text{donut filter: } E_{i,j}^d = \frac{\sum_{a=i-w}^{i+w} \sum_{b=j-w}^{j+w} M_{a,b}^* - \sum_{a=i-p}^{i+p} \sum_{b=j-p}^{j+p} M_{a,b}^* - \sum_{a=i-w}^{i+w} M_{a,j}^* - \sum_{a=i+p+1}^{i+w} M_{a,j}^* - \sum_{b=j-w}^{j-p-1} M_{i,b}^* - \sum_{b=j+p+1}^{j+w} M_{i,b}^*}{\sum_{a=i-w}^{i+w} \sum_{b=j-w}^{j+w} E_{a,b}^* - \sum_{a=i-p}^{i+p} \sum_{b=j-p}^{j+p} E_{a,b}^* - \sum_{a=i-w}^{i+w} E_{a,j}^* - \sum_{a=i+p+1}^{i+w} E_{a,j}^* - \sum_{b=j-w}^{j-p-1} E_{i,b}^* - \sum_{b=j+p+1}^{j+w} E_{i,b}^*} \times E_{i,j}^*$$

$$\text{lower left filter: } E_{i,j}^{ll*} = \frac{\sum_{a=i+1}^{i+w} \sum_{b=j-w}^{j-1} M_{a,b}^* - \sum_{a=i+1}^{i+p} \sum_{b=j-p}^{j-1} M_{a,b}^*}{\sum_{a=i+1}^{i+w} \sum_{b=j-w}^{j-1} E_{a,b}^* - \sum_{a=i+1}^{i+p} \sum_{b=j-p}^{j-1} E_{a,b}^*} \times E_{i,j}^*$$

$$\text{vertical filter: } E_{i,j}^{v*} = \frac{\sum_{a=i-w}^{i-p-1} \sum_{b=j-1}^{j+1} M_{a,b}^* - \sum_{a=i+p+1}^{i+w} \sum_{b=j-1}^{j+1} M_{a,b}^*}{\sum_{a=i-w}^{i-p-1} \sum_{b=j-1}^{j+1} E_{a,b}^* - \sum_{a=i+p+1}^{i+w} \sum_{b=j-1}^{j+1} E_{a,b}^*} \times E_{i,j}^*$$

$$\text{horizontal filter: } E_{i,j}^{h*} = \frac{\sum_{b=j-w}^{j-p-1} \sum_{a=i-1}^{i+1} M_{a,b}^* - \sum_{b=j+p+1}^{j+w} \sum_{a=i-1}^{i+1} M_{a,b}^*}{\sum_{b=j-w}^{j-p-1} \sum_{a=i-1}^{i+1} E_{a,b}^* - \sum_{b=j+p+1}^{j+w} \sum_{a=i-1}^{i+1} E_{a,b}^*} \times E_{i,j}^*$$

Step2. calculate p-value ,fold and FDR. (FDR is distance related)

Step3. calculate enriched pixels.

Step4. merge pixels nearby.

a Poisson process with parameter $\lambda = E_local \times C_i \times C_j$.

4. Differential Loop Detection

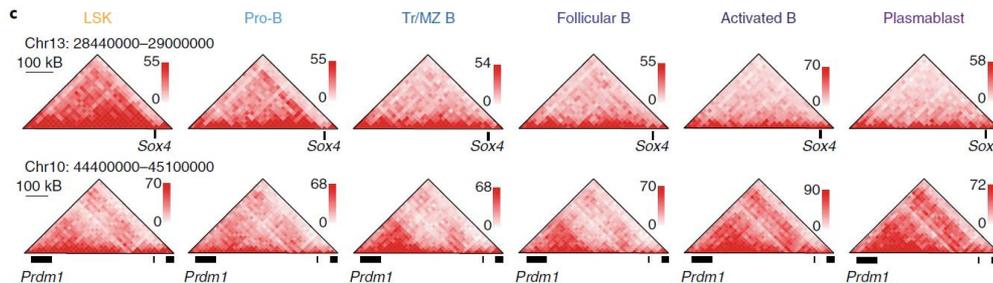
Differential interaction: test for each pixel (diffHiC, HiCcompare);
test for each pixel considering neighborhood (FIND)
Differential loop: test for each loop considering neighborhood



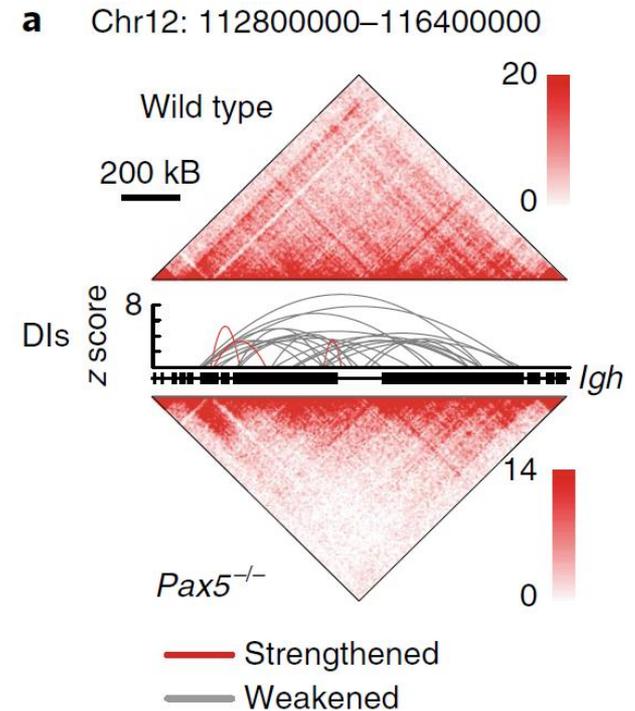
Transcription-factor-mediated supervision of global genome architecture maintains B cell identity

Timothy M. Johanson^{1,2,4}, Aaron T. L. Lun^{1,2,4}, Hannah D. Coughlan^{1,2,4}, Tania Tan^{1,2}, Gordon K. Smyth^{1,3}, Stephen L. Nutt^{1,2,5*} and Rhys S. Allan^{1,2,5*}

Recent studies have elucidated cell-lineage-specific three-dimensional genome organization; however, how such specific architecture is established or maintained is unclear. We hypothesized that lineage-defining transcription factors maintain cell identity via global control of genome organization. These factors bind many genomic sites outside of the genes that they directly regulate and thus are potentially implicated in three-dimensional genome organization. Using chromosome-conformation-capture techniques, we show that the transcription factor Paired box 5 (Pax5) is critical for the establishment and maintenance of the global lineage-specific architecture of B cells. Pax5 was found to supervise genome architecture throughout B cell differentiation, until the plasmablast stage, in which Pax5 is naturally silenced and B cell-specific genome structure is lost. Crucially, Pax5 did not rely on ongoing transcription to organize the genome. These results implicate sequence-specific DNA-binding proteins in global genome organization to establish and maintain lineage fidelity.



50kb diffHiC



4. Differential Loop Detection

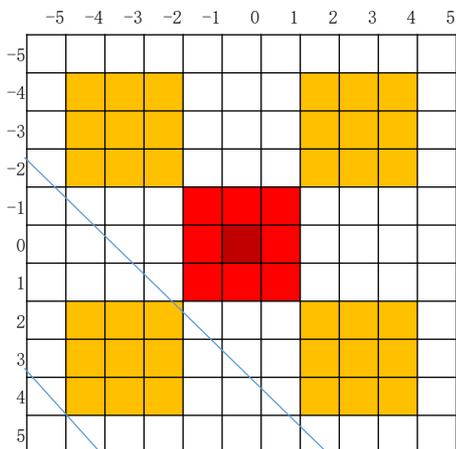
DEseq2

Read count matrix K_{ij} (gene i , sample j)

Design matrix X_{jr} (sample j , condition r)

- $K_{ij} \sim NB(\mu_{ij}, \alpha_{ij})$
 - NB generalizes Poisson, mean: μ variance: $\mu + \alpha\mu^2$
 - α : dispersion parameter
- $\mu_{ij} = s_j q_{ij}$
 - $s_j = \text{median}_{i:K_i^R \neq 0} \frac{K_{ij}}{K_i^R}, K_i^R = (\sum_{j=1}^m K_{ij})^{\frac{1}{m}}$
- $\log(q_{ij}) = \sum_r X_{jr} \beta_{ir}$
 - β : log fold change (LFC)

Null hypothesis: Beta=0 (gene is not differential as gene does not change when condition changes)

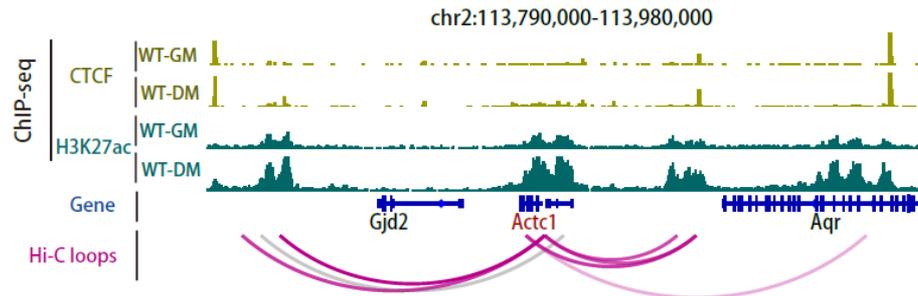


$$\frac{u_{ij} \times norm1 \times norm2}{norm_expected} = q_{ij}$$

$$S_{ij} = \frac{norm_expected}{norm1 \times norm2}$$

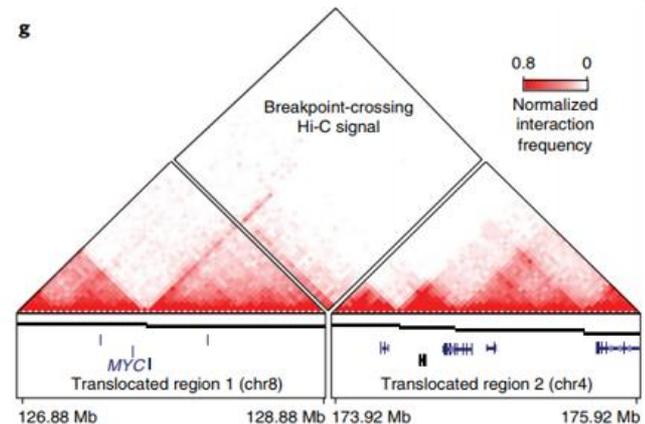
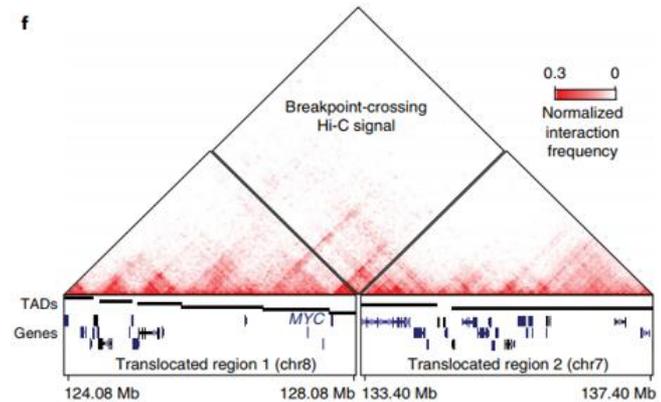
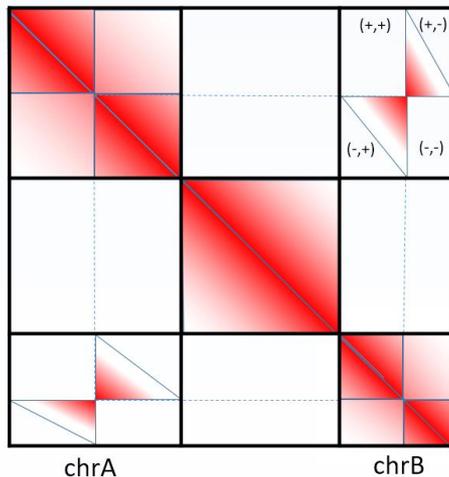
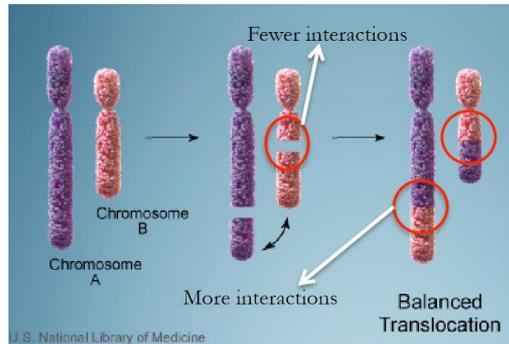
local	sample	rep
D,P	GM,DM	1,2,3 (GM) 1,2,3 (DM)

design formula \sim rep + sample + rep:sample + local + sample:local



5. Structural Variation

- ✓ 20-100M reads (1-5X coverage);
- ✓ complex / large-scale structural variation
- ✓ breakpoint in repeat region



Features local maximal peak in distal regions → NMS

loss of interaction in breakpoint(possible)

right angle with direction

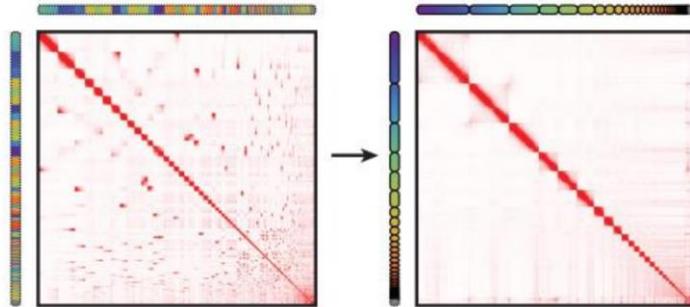
→ gradient/template matching

In development

6. Genome Assembly



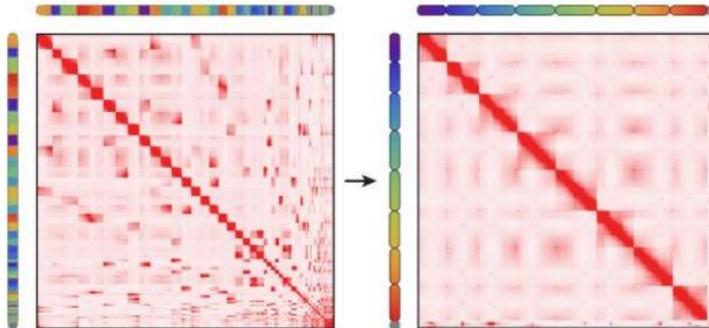
Band-tailed pigeon (*Patagioenas fasciata*)



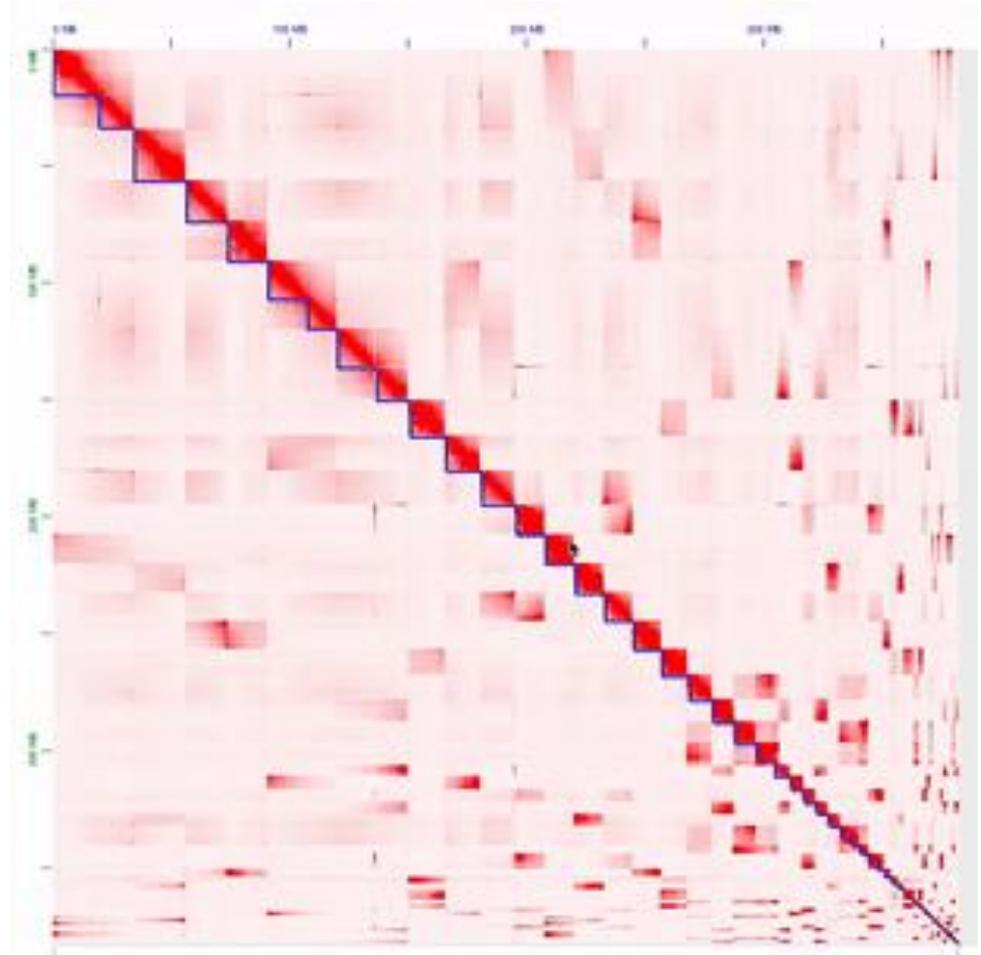
Draft assembly:
Murray et al. *Science* (Nov 2017) (Illumina DNA-Seq + Dovetail)

R108 v. 1.0

R108 v. HIC



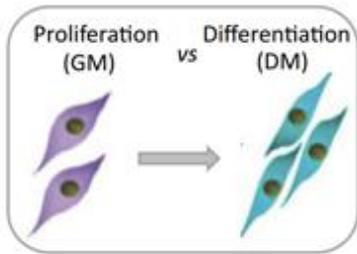
Draft assembly:
Moll et al. *BMC Genomics* (Aug 2017) (PacBio DNA-Seq + Bionano + Dovetail)



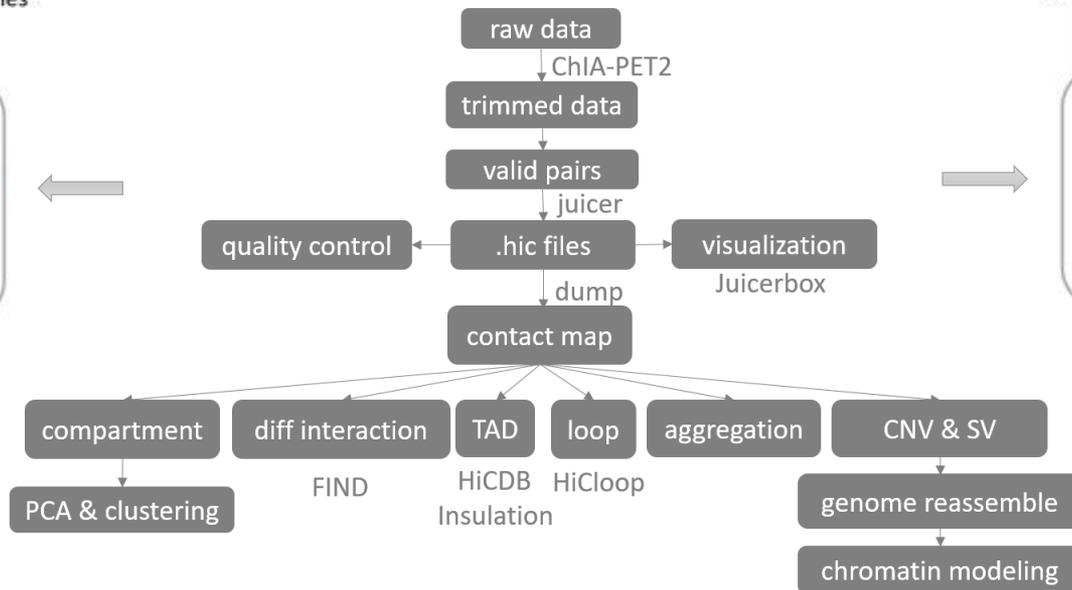
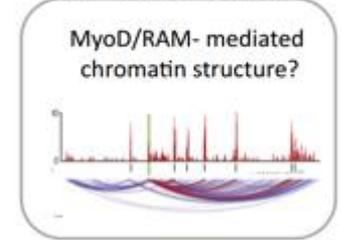
Take home message

Practice is the sole criterion for testing truth.

- Muscle stem cell-derived cell lines
C2C12 cells



- Sorted muscle stem cells from KO mice
MyoD^{-/-} muscle stem cells
RAM^{-/-} muscle stem cells



efficient; QC; visualization; multiple feature detection

Have fun in exploring biological data !

Michael Q. Zhang Lab @ THU



Thank you!

Fengling Chen 陈凤玲: Ph.D. candidate

Email: cf15@mails.tsinghua.edu.cn

[fchen@github](https://github.com/fchen)