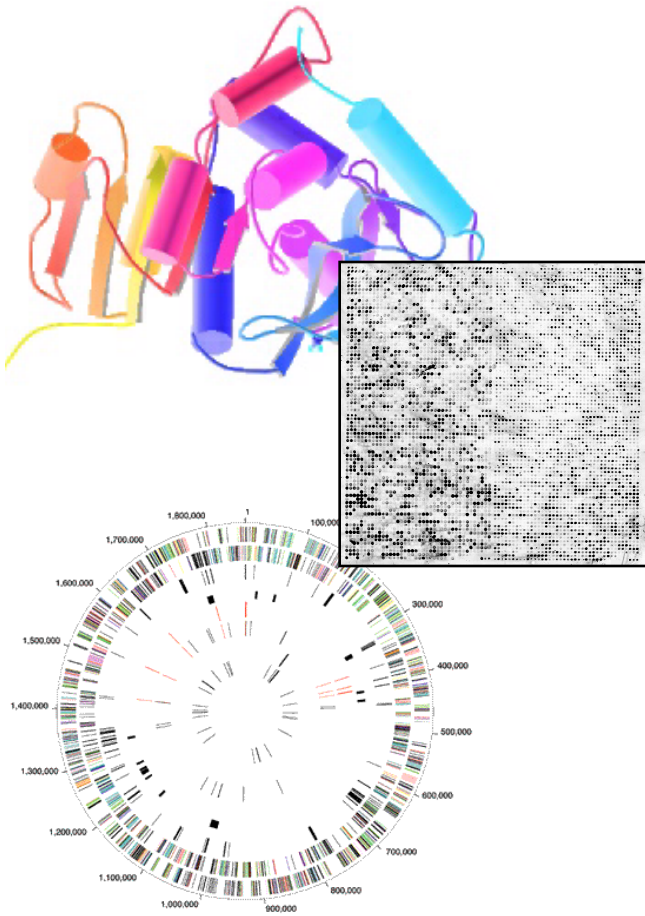


Biomed. Data Science:

# Basic Multi-omic Analyses

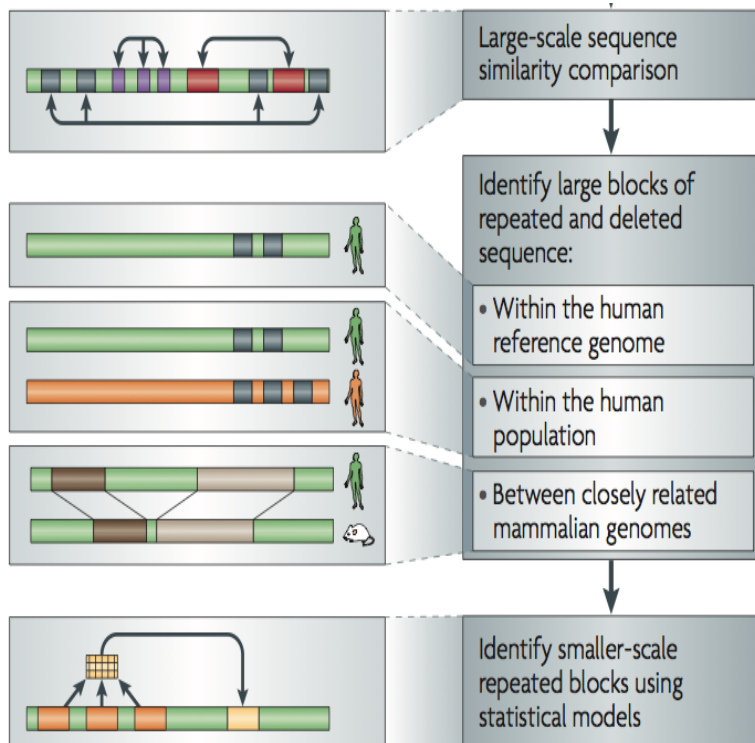


Mark Gerstein, Yale University  
[gersteinlab.org/courses/452](http://gersteinlab.org/courses/452)  
(last edit in spring '20)

# Non-coding Annotations: Overview

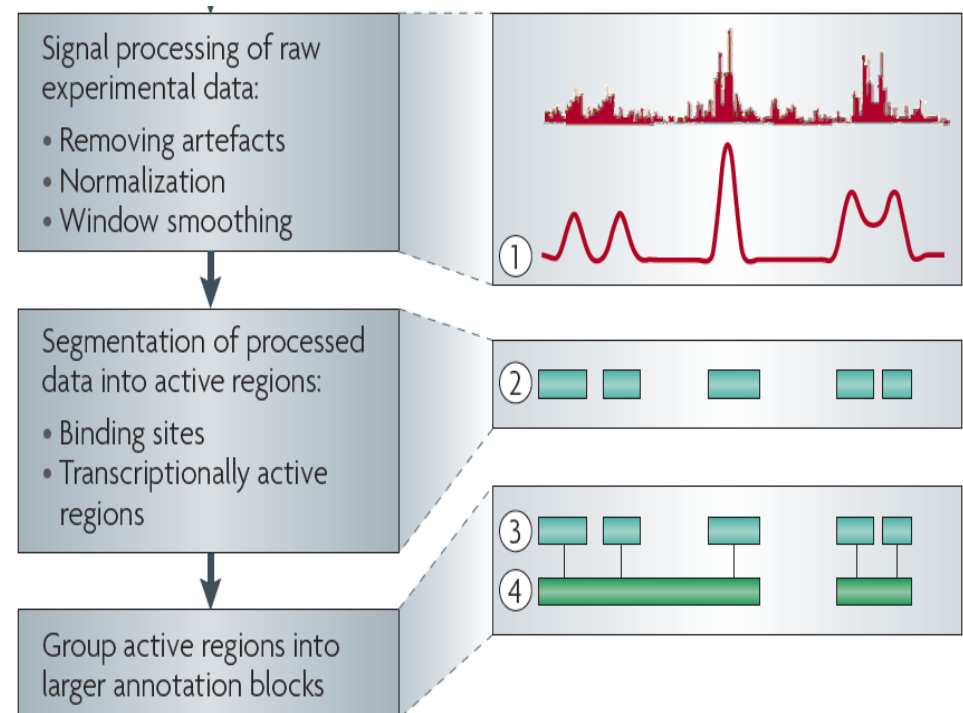
Features are often present on multiple "scale" (eg elements and connected networks)

## Sequence features, incl. Conservation



## Functional Genomics

### Chip-seq (Epigenome & seq. specific TF) and ncRNA & un-annotated transcription



[Nat. Rev. Genet. (2010) 11: 559]

hESC-A

N1-A

N2-A

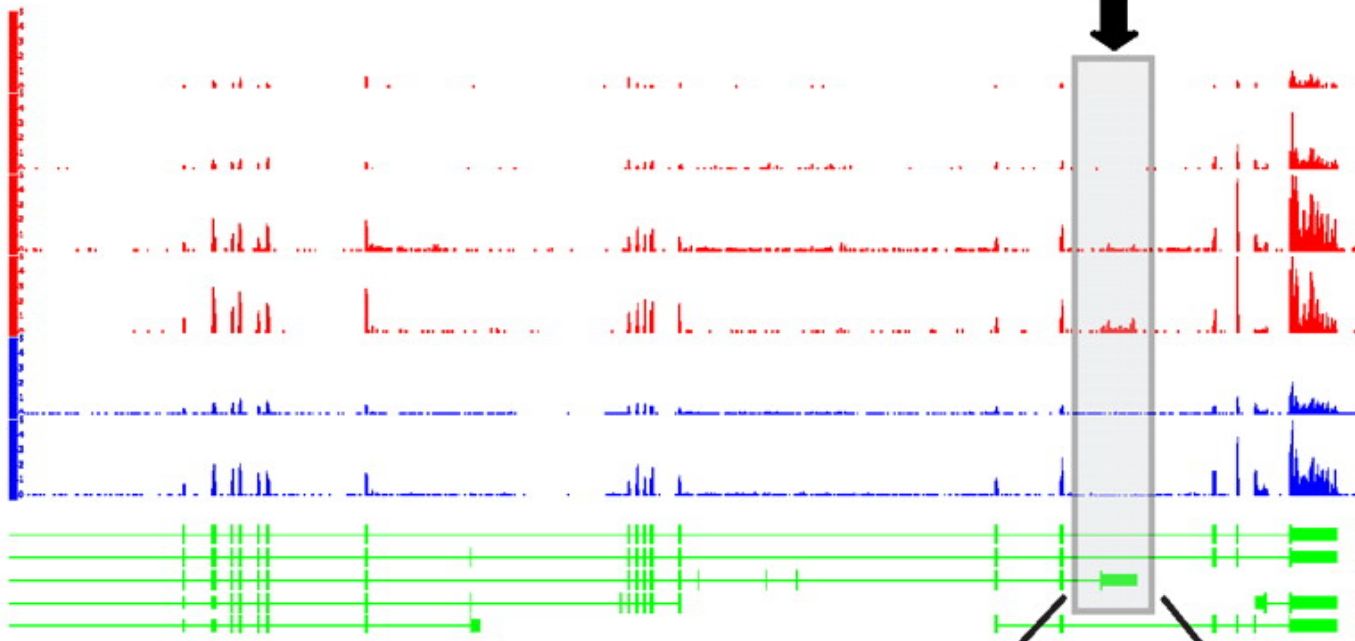
N3-A

hESC-B

N2-B

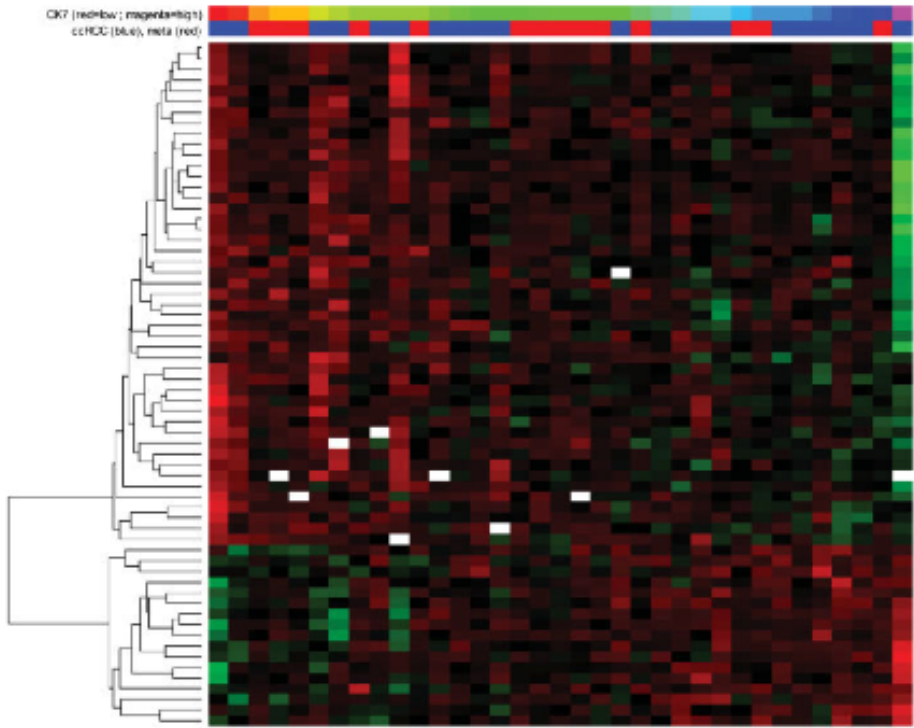
UCSC

Genes



DK7 (red/blue; magenta/white)  
ccRCC (blue), meta (red)

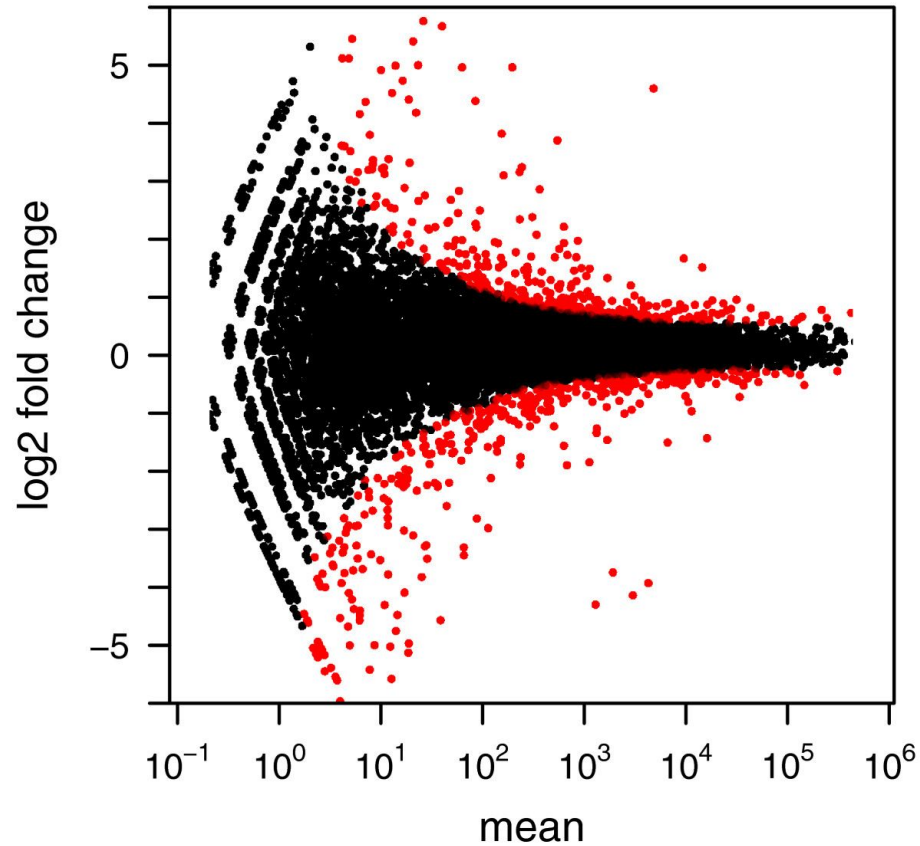
CAM1



**Information from  
RNA-seq:  
Avg. signal at exons &  
TARs (RPKMs)**

[PNAS 4:107: 5254 ; IJC 123:569]

# Differential expression analysis



*Genome Biology, 2010 11:R106*



# Differential expression analysis: Count-based

1. **DESeq** -- based on negative binomial distribution
2. **edgeR** -- use an overdispersed Poisson model
3. **baySeq** -- use an empirical Bayes approach
4. **TSPM** -- use a two-stage poisson model

Anders and Huber *Genome Biology* 2010, 11:R106  
<http://genomebiology.com/2010/11/10/R106>



METHOD

Open Access

Differential expression analysis for sequence count data

Simon Anders\*, Wolfgang Huber

BIOINFORMATICS APPLICATIONS NOTE

Vol. 26 no. 1 2010, pages 139–140  
doi:10.1093/bioinformatics/btp616

Gene expression

**edgeR: a Bioconductor package for differential expression analysis of digital gene expression data**

Mark D. Robinson<sup>1,2,\*</sup>, Davis J. McCarthy<sup>2,†</sup> and Gordon K. Smyth<sup>2</sup>

<sup>1</sup>Cancer Program, Garvan Institute of Medical Research, 384 Victoria Street, Darlinghurst, NSW 2010 and  
<sup>2</sup>Bioinformatics Division, The Walter and Eliza Hall Institute of Medical Research, 1G Royal Parade, Parkville, Victoria 3052, Australia

Hardcastle and Kelly *BMC Bioinformatics* 2010, 11:422  
<http://www.biomedcentral.com/1471-2105/11/422>



RESEARCH ARTICLE

Open Access

baySeq: Empirical Bayesian methods for identifying differential expression in sequence count data

Thomas J Hardcastle\*, Krystyna A Kelly

*Statistical Applications in Genetics and Molecular Biology*

Volume 10, Issue 1

2011

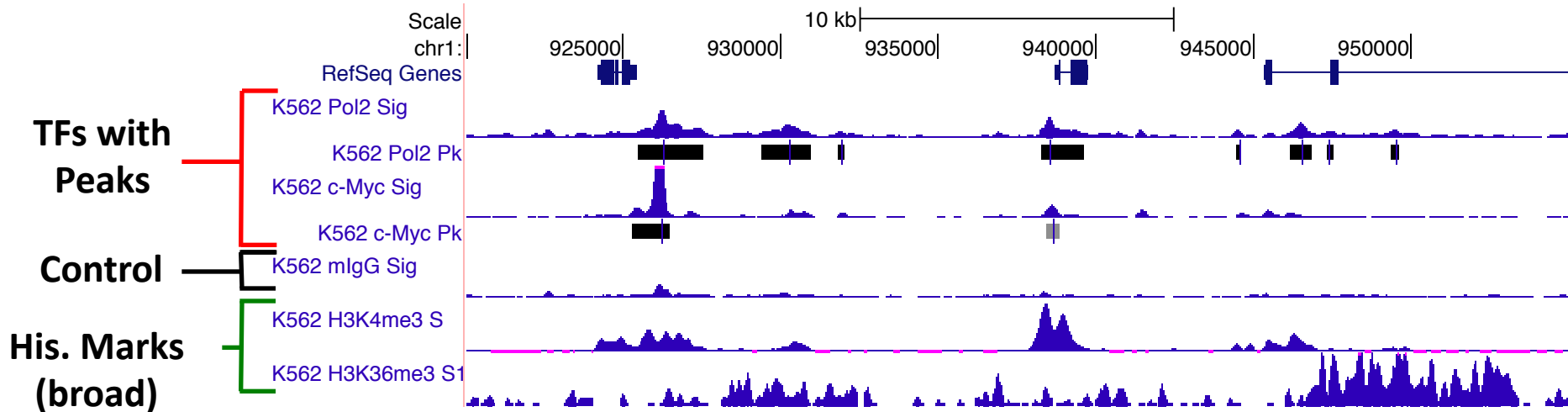
Article 26

A Two-Stage Poisson Model for Testing RNA-Seq Data

Paul L. Auer, Fred Hutchinson Cancer Research Center  
Rebecca W. Doerge, Purdue University

# chip-seq

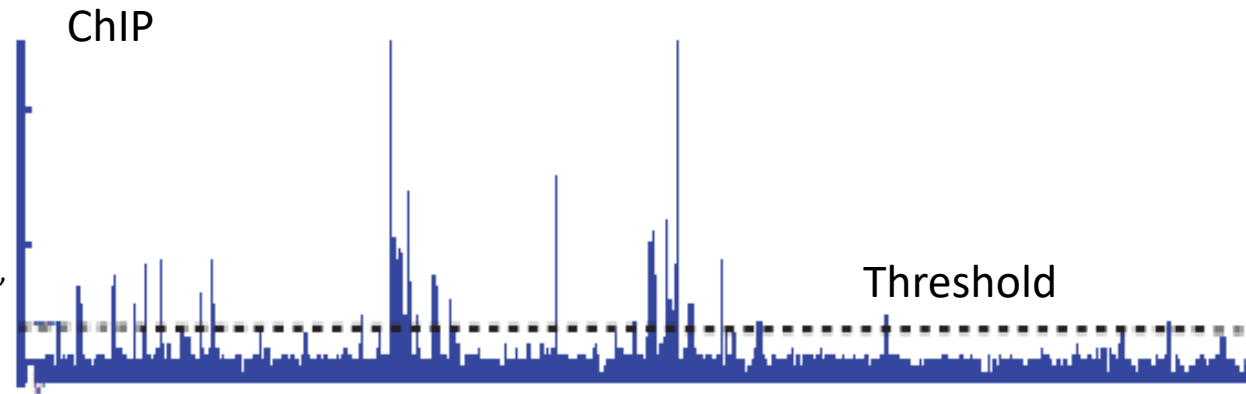
# Information from Chip-seq



[*Science* 330: 1775  
+ ENCODE Data  
Sources  
TFs & Control: Yale  
HMs: UW & Broad ]

# Summarizing the Signal: "Traditional" ChipSeq Peak Calling

- Generate & threshold the signal profile to identify candidate target regions
  - Simulation (PeakSeq),
  - Local window based Poisson (MACS),
  - Fold change statistics (SPP)



Potential Targets



- Score against the control



Significantly Enriched targets



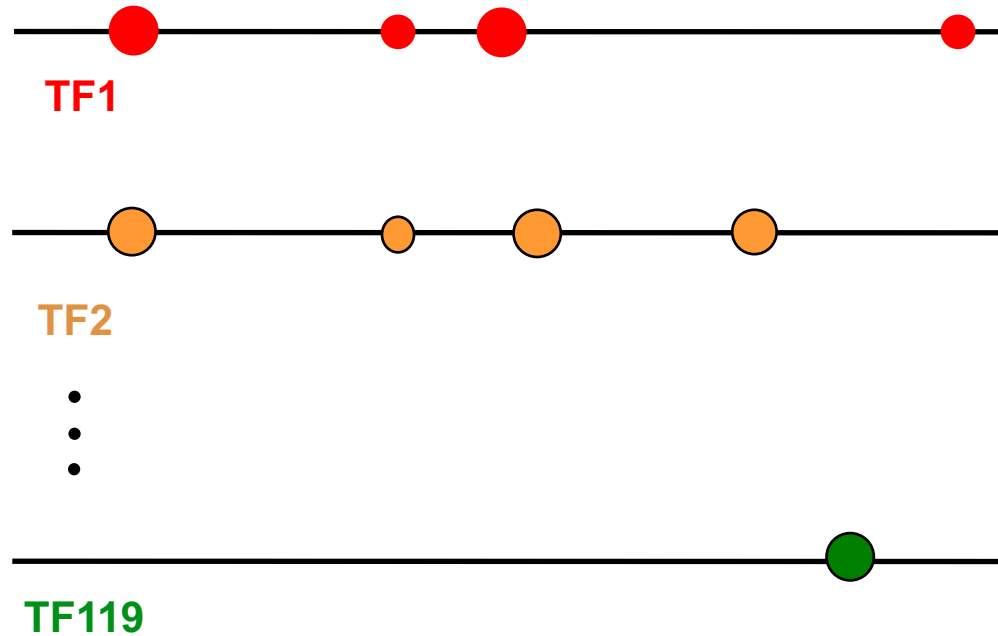
# Data Flow: Chip-seq expts. to co-associating peaks

**119 TFs** from 458 ChIP-Seq experiments (2 Tb tot.)

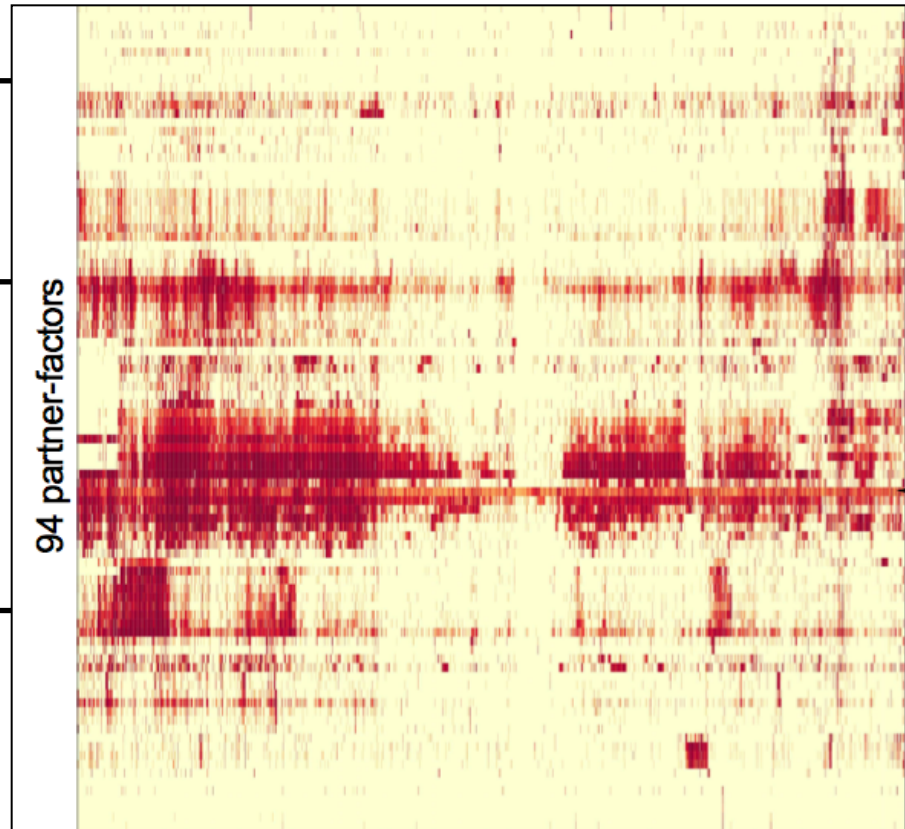
↓  
Signal Tracks



↓  
**7M Peaks** from Uniform Peak Calling



- Mostly in Tier 1 cell lines
  - K562, GM12878, H1h-ESC...
- Matching RNA-Seq data in all cell-lines
- SPP & PeakSeq
- thresholding w. IDR (replicas)



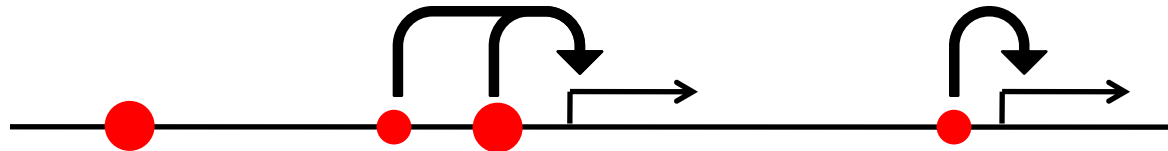


# Data Flow: peaks to proximal & distal networks

Peak Calling

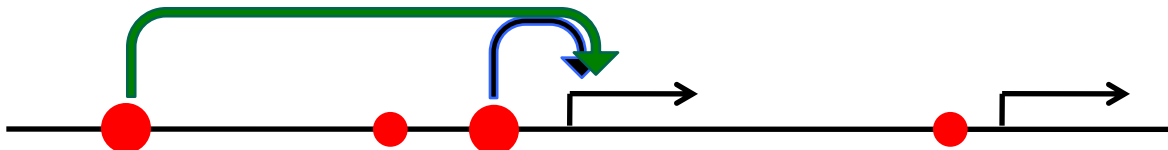


Assigning TF binding sites to targets



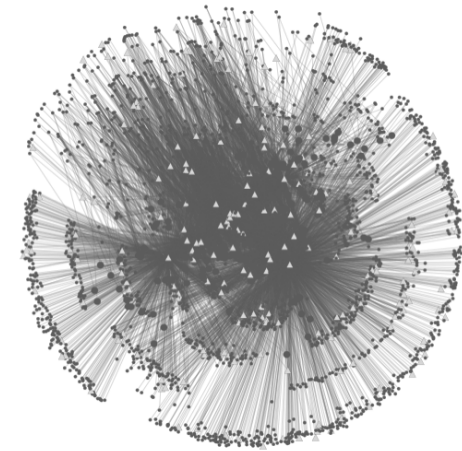
Filtering high confidence edges & distal regulation

Based on stat. model combining signal strength & location relative to typical binding

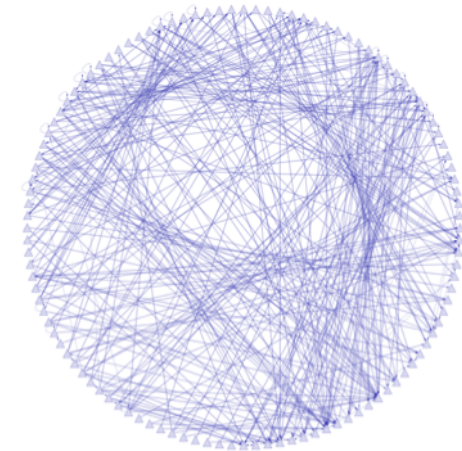


~500K Edges

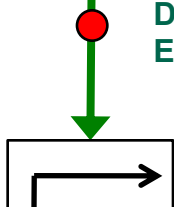
[ Cheng et al., *Bioinfo.* ('11);  
Gerstein et al. *Nature* (in press, '12) ;  
Yip et al., *GenomeBiology* (in press, '12)]



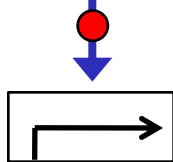
~26K Edges



Potential Distal Edge

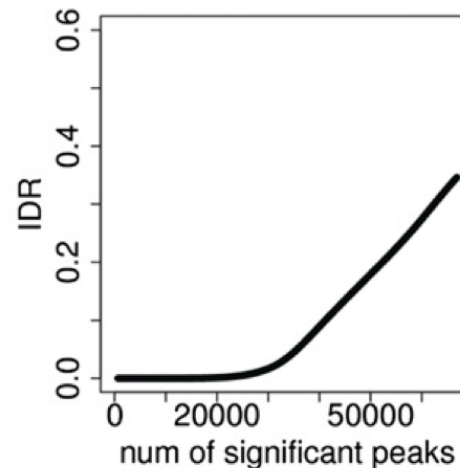
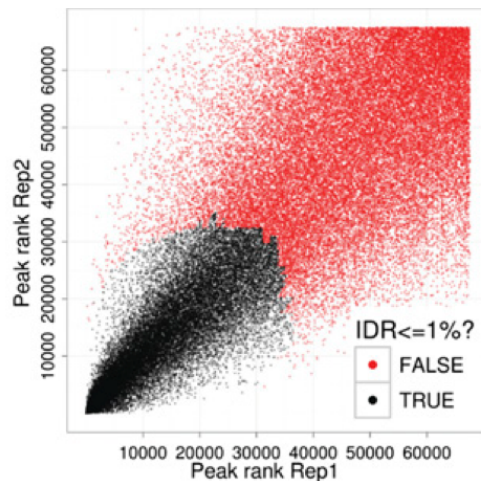


Strong Proximal Edge

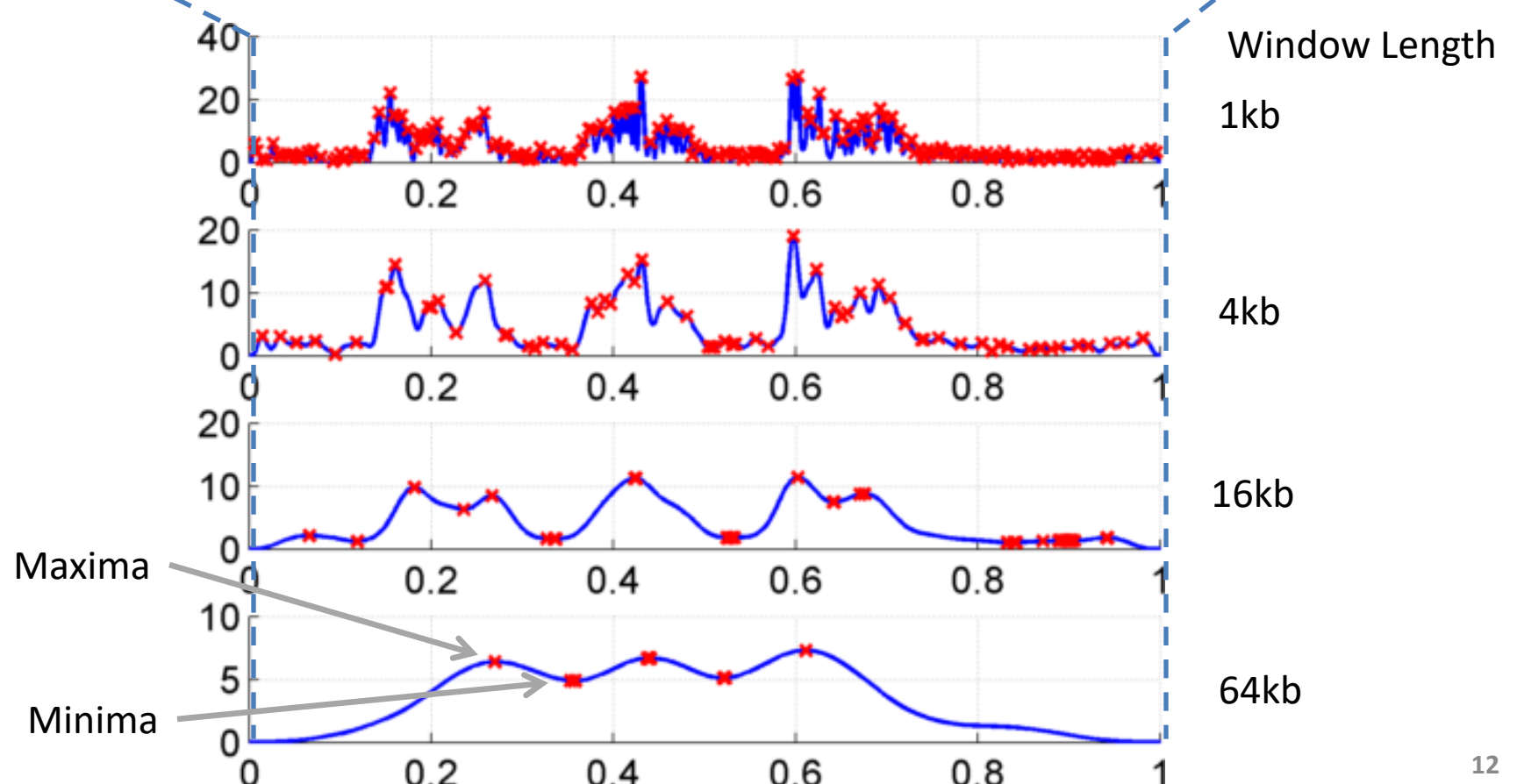
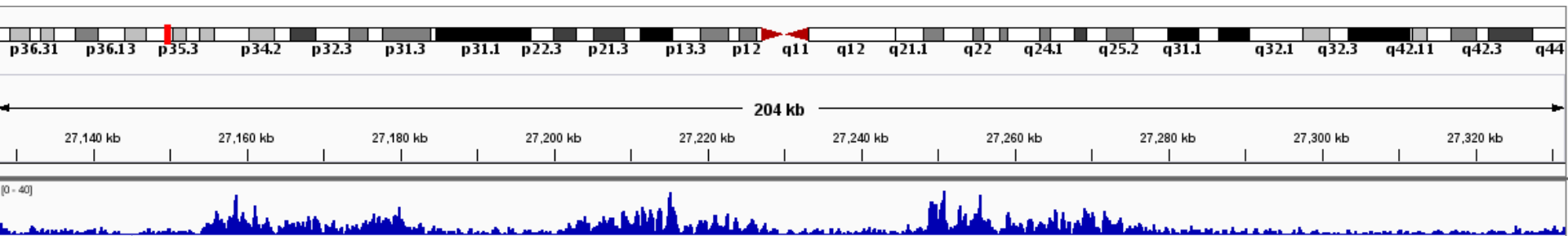


# The irreproducible discovery rate (IDR)

- Unified approach to measure the reproducibility of findings identified from replicate high-throughput experiments.
- Idea: call peaks with low cutoff and classify peaks as reproducible or not (bivariate rank distributions) based on overlap of ranked peaks (consistency)

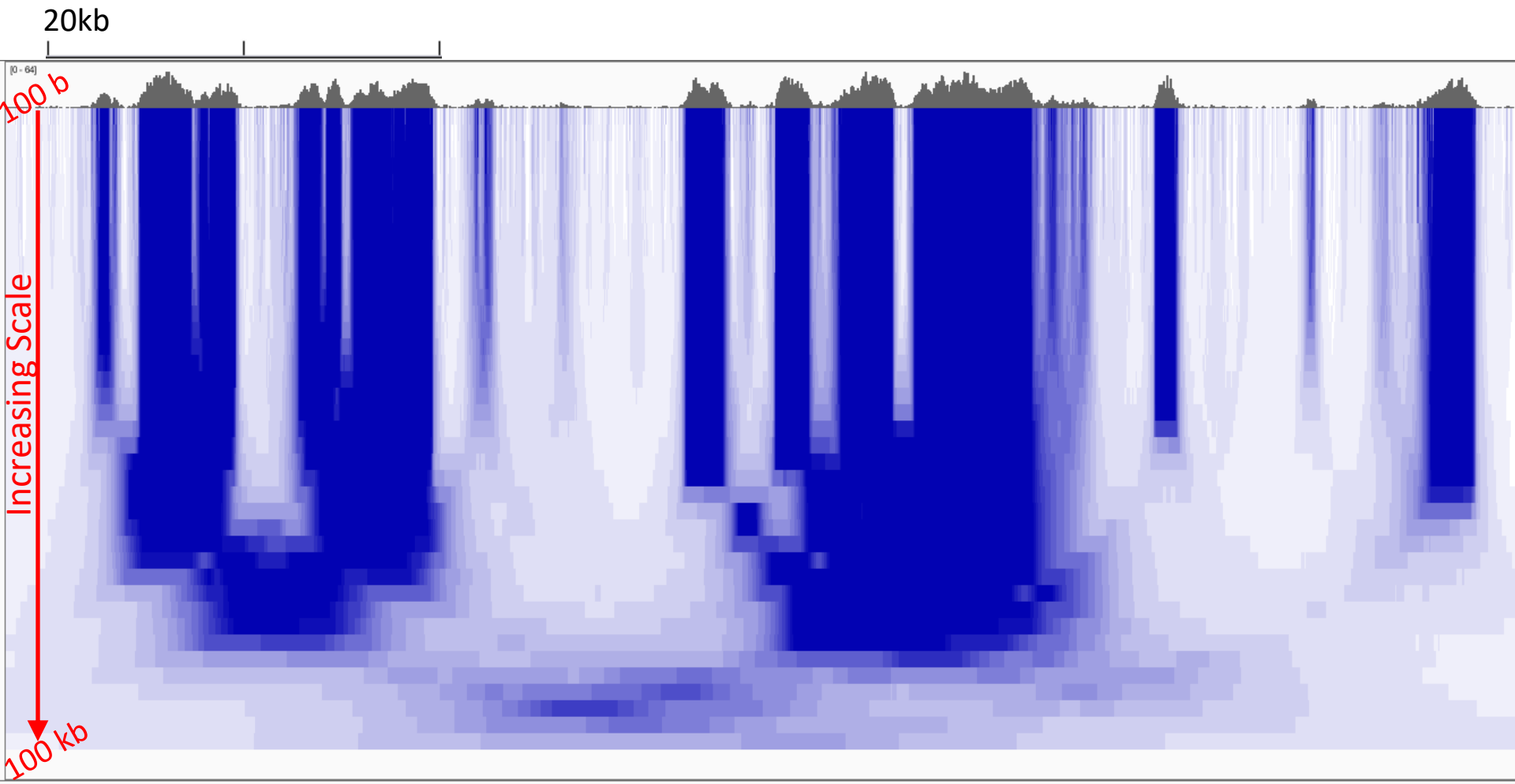


# Multiscale Analysis, Minima/Maxima based Coarse Segmentation

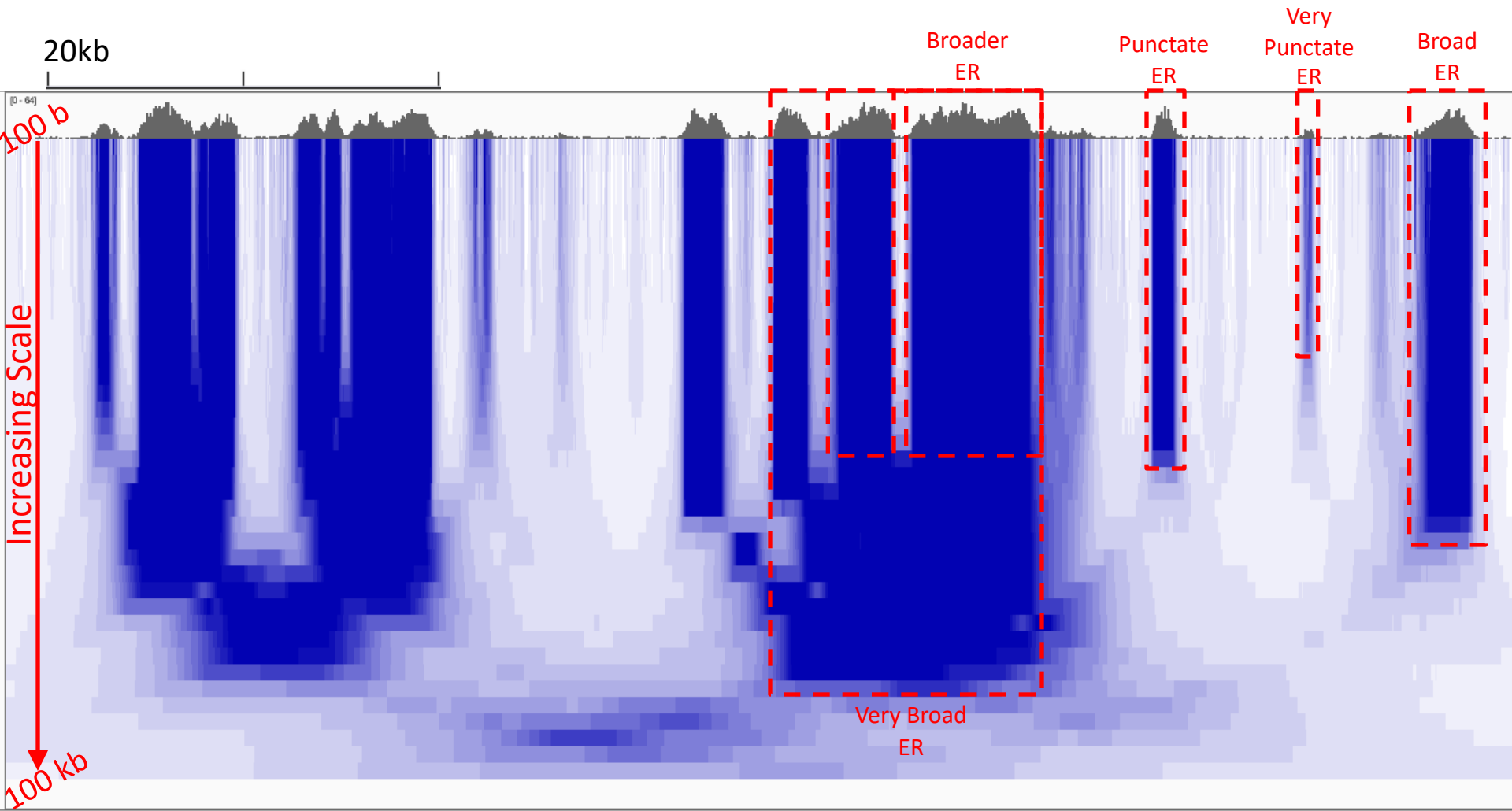


Harmanci et al, Genome Biology 2014, MUSIC.gersteinlab.org

# Multiscale Decomposition



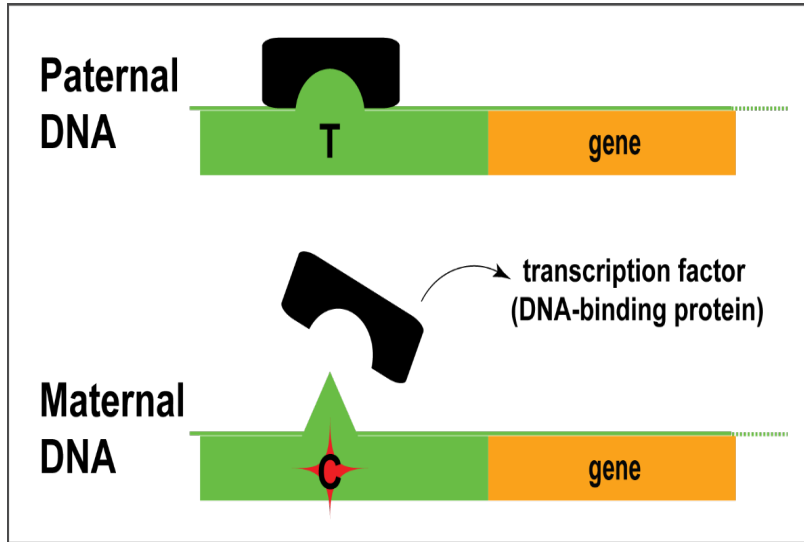
# Multiscale Decomposition



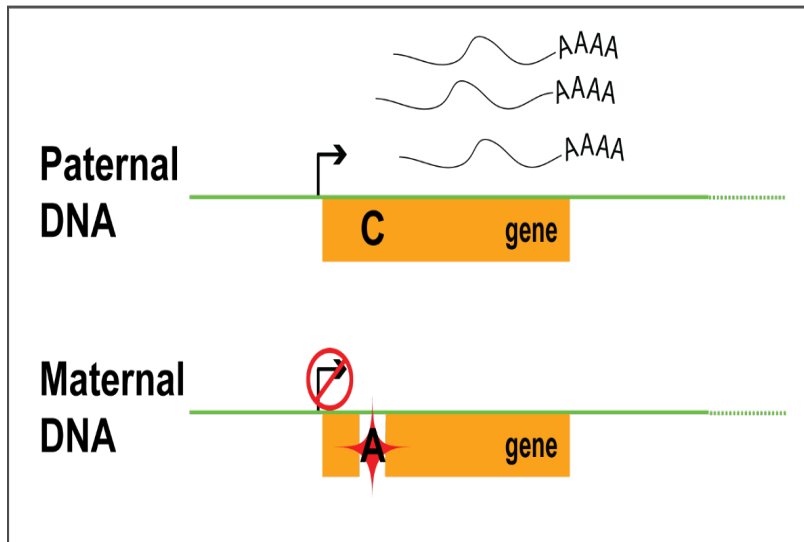


# ASB/ASE & eQTL

# Allele-specific binding and expression



Genomic variants affecting allele-specific behavior e.g. allele-specific binding (ASB)



e.g. allele-specific expression (ASE)

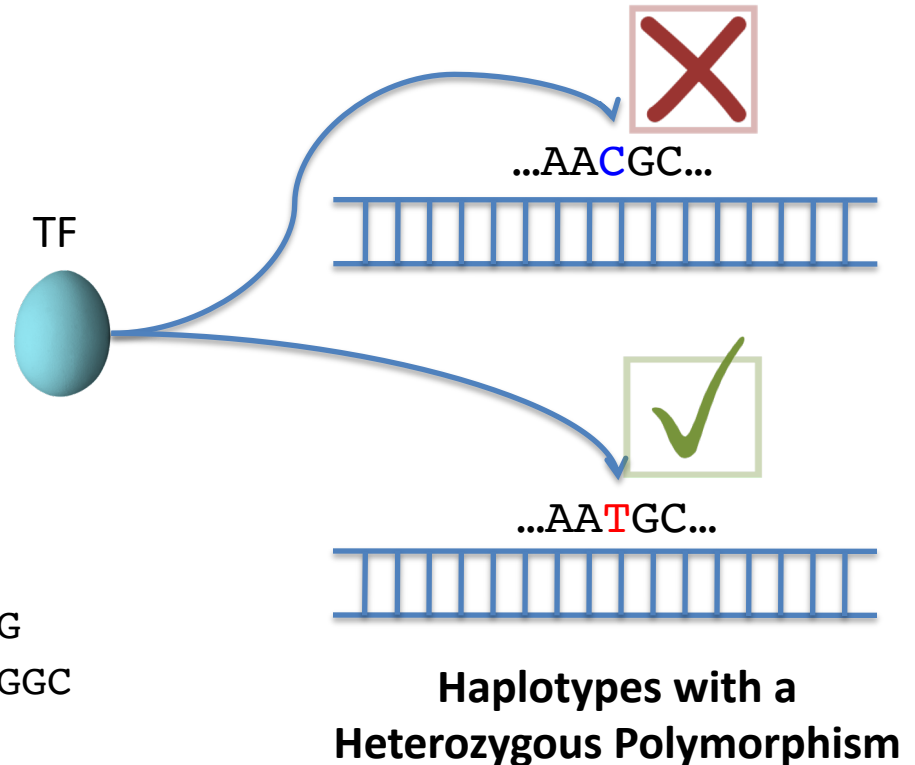
# Inferring Allele Specific Binding/Expression using Sequence Reads

## RNA/ChIP-Seq Reads

ACTTTGATAGCGTCAAT**T**G  
CTTTGATAGCGTCAAT**T**GC  
CTTTGATAGCGTCAAC**C**GC  
TTGACAGCGTCAAT**T**GCAC  
TGATAGCGTCAAT**T**GCACG  
ATAGCGTCAAT**T**GCACGTC  
TAGCGTCAAT**T**GCACGTCG  
CGTCAAC**C**GCACGTCGGGA  
GTCAAT**T**GCACGTCGAGAG  
CAAT**T**GCACGTCGGGAGTT  
AA**T**GCACGTCGGGAGTTG  
T**T**GCACGTTGGGAGTTGGC

10 x **T**

2 x **C**



Interplay of the annotation and individual sequence variants

# Many Technical Issues in Determining ASE/ASB: Reference Bias (naïve alignment against reference)

## ASE/ASB Example:

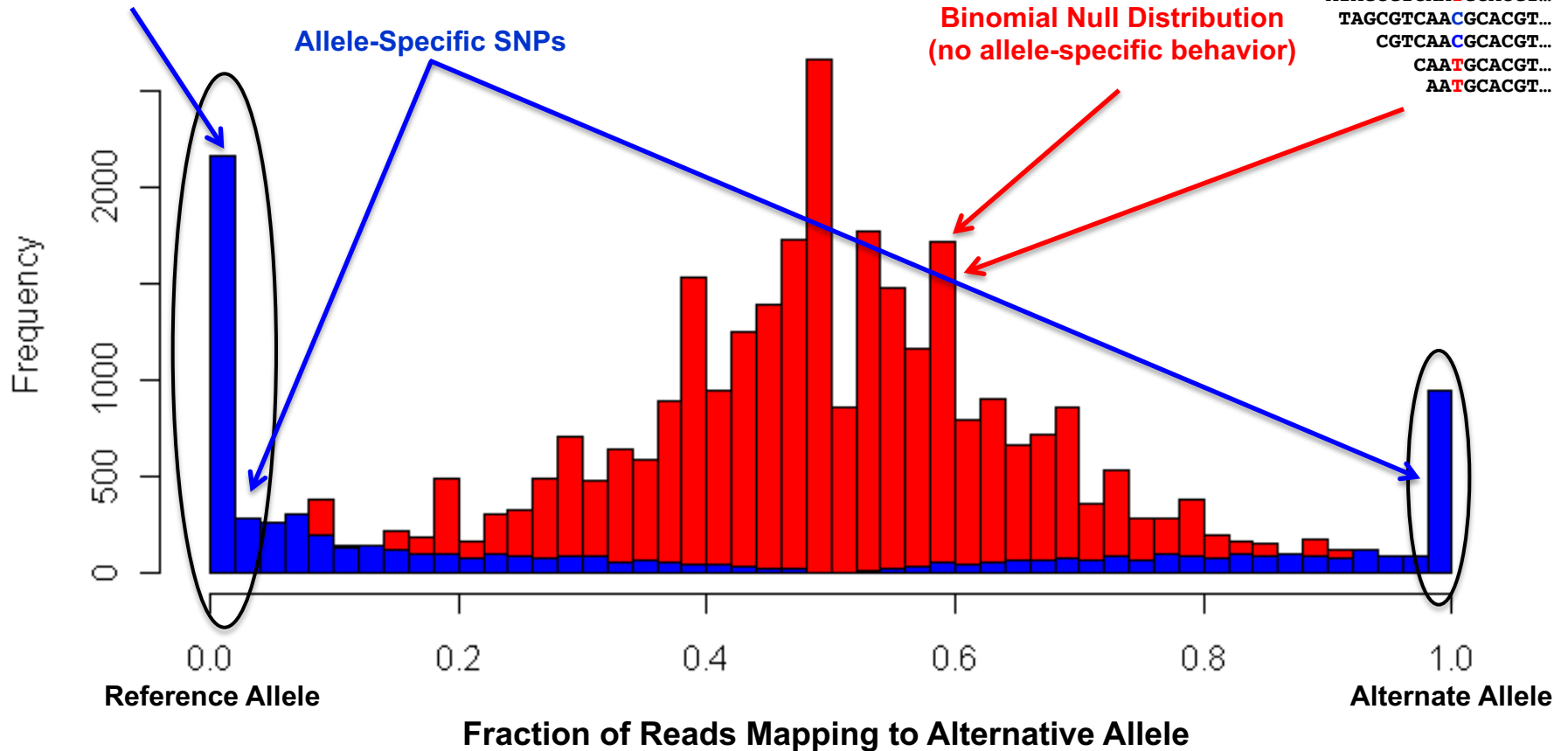
```

...GTCAATGCAC
...GTCAATGCACG
...GTCAATGCACGTC
...GTCAATGCACGTCG
...GTCAACGCACGTCGGGA
GTCAATGCACGTCGAGAG
CAATGCACGTCGGGAGTT
AATGCACGTCGGGAGTTG
    
```

## Null Example:

```

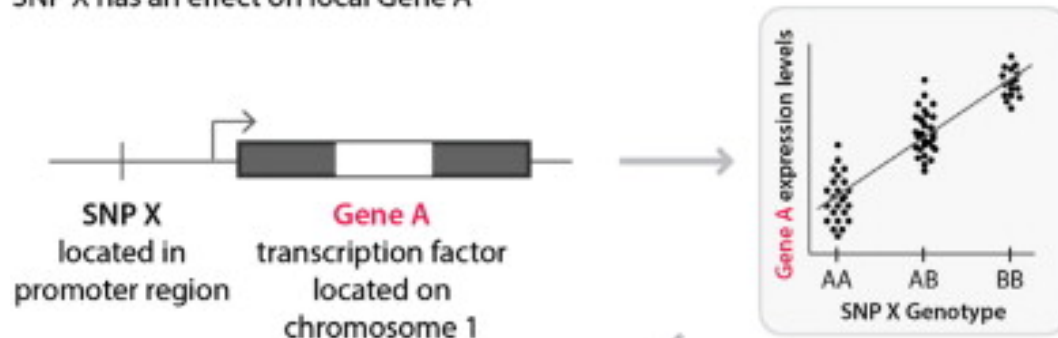
ACTTTGATAGCGTCAATG
CTTTGATAGCGTCAACGC
TTGACAGCGTCAATGCAC
ATAGCGTCAATGCACGT...
TAGCGTCAACGCACGT...
CGTCAACGCACGT...
CAATGCACGT...
AATGCACGT...
    
```



# Expression quantitative trait

## Cis-eQTL

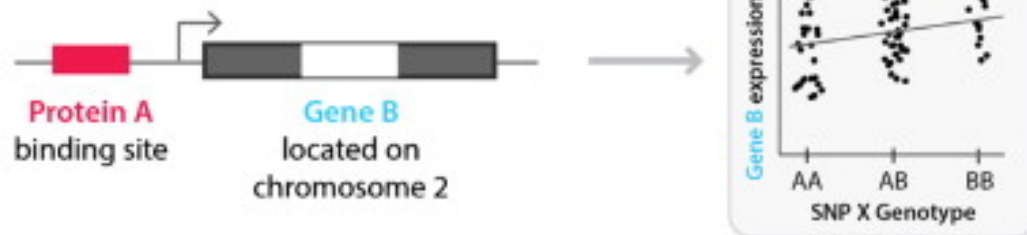
SNP X has an effect on local Gene A



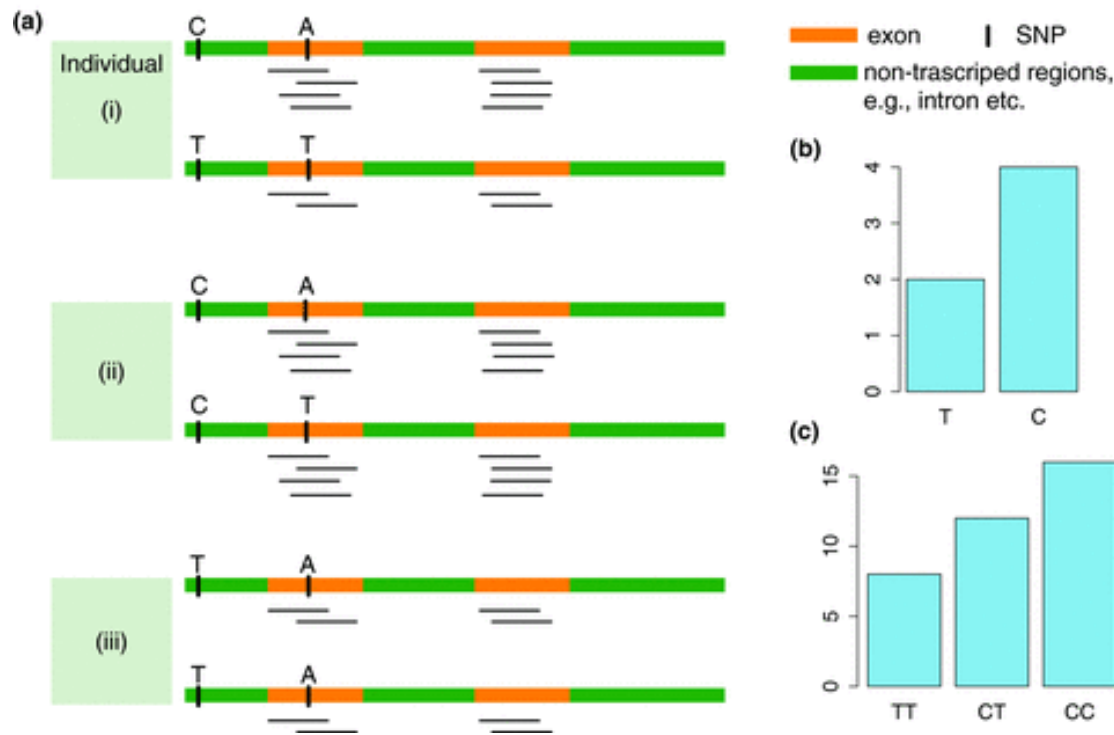
Altered **Protein A** levels, effect on the binding to the transcription factor binding sites of downstream genes

## Trans-eQTL

SNP X has an effect on distant Gene B through an intermediary factor (such as a transcription factor)



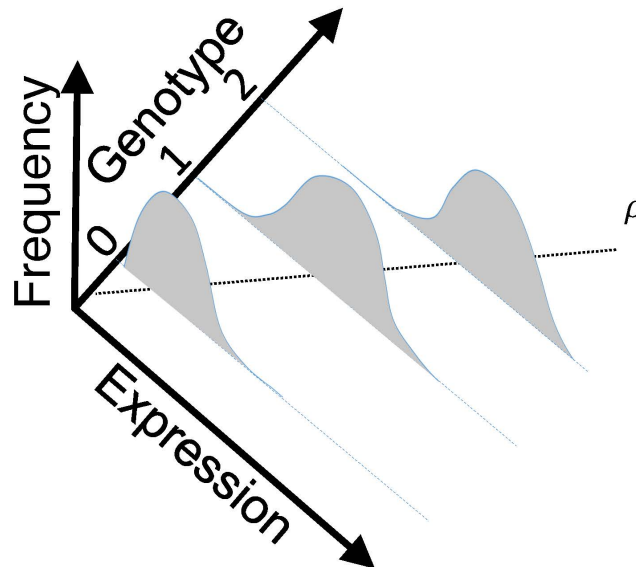




# eQTL Mapping Using RNA-Seq Data

- eQTLs are genomic loci that contribute to variation in mRNA expression levels
- eQTLs provide insights on transcription regulation, and the molecular basis of phenotypic outcomes
- eQTL mapping can be done with RNA-Seq data

[*Biometrics* 68(1) 1–11]



# Hi-C

# 3D organization of genome



"We finished the genome map, now we can't figure out how to fold it."

image credit: Iyer et al. BMC Biophysics 2011, cartoonist John Chase

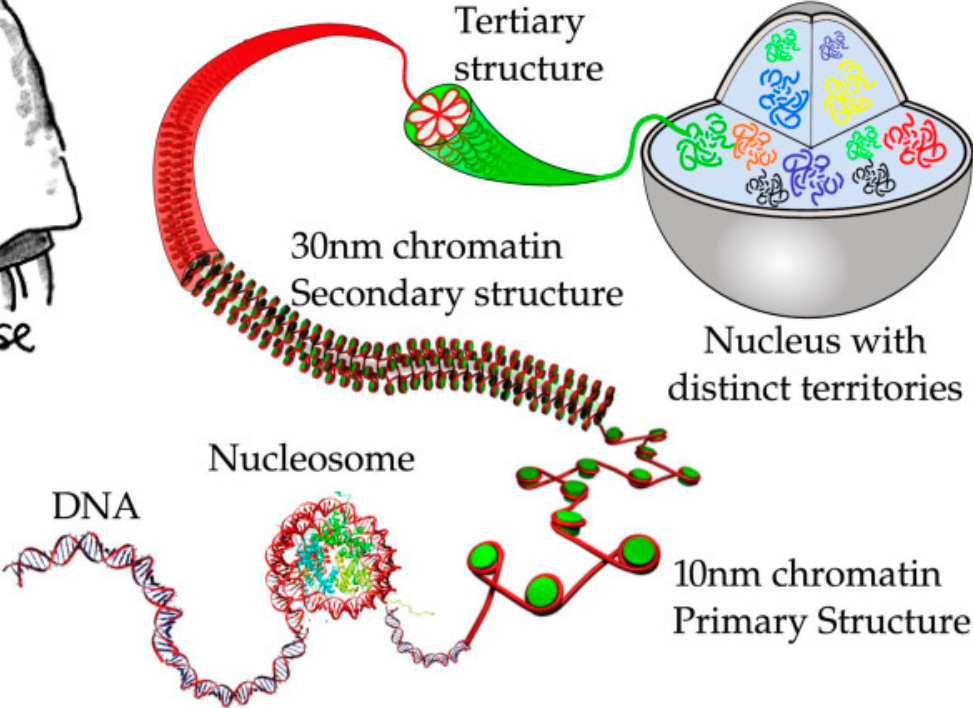
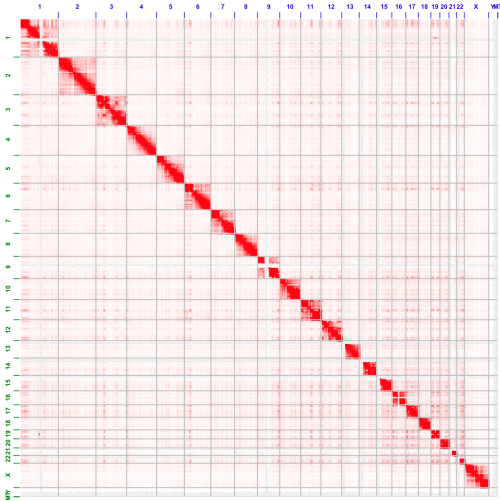
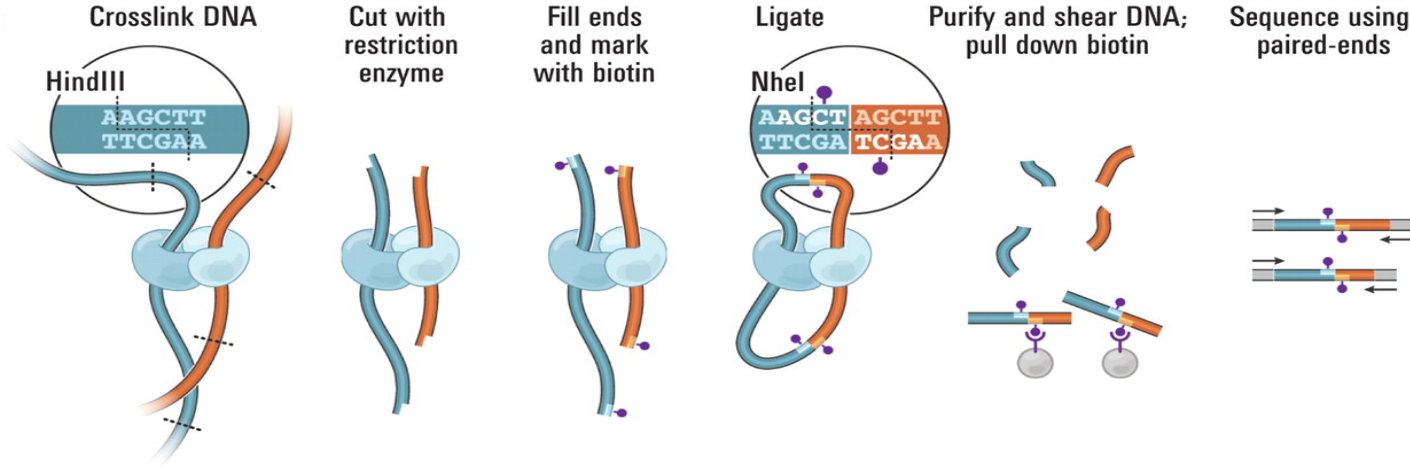


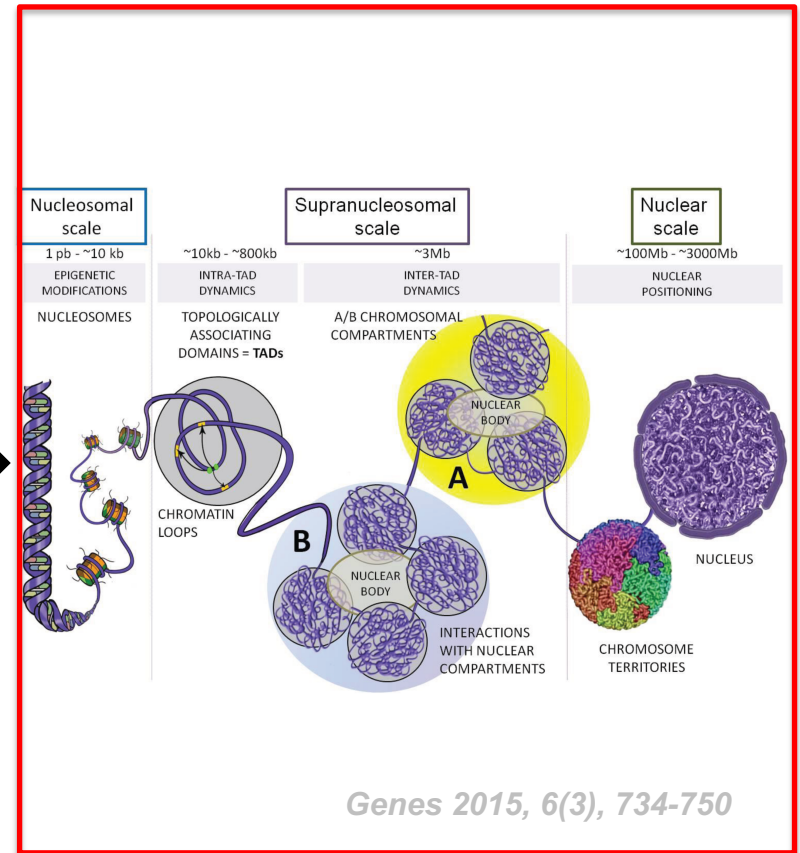
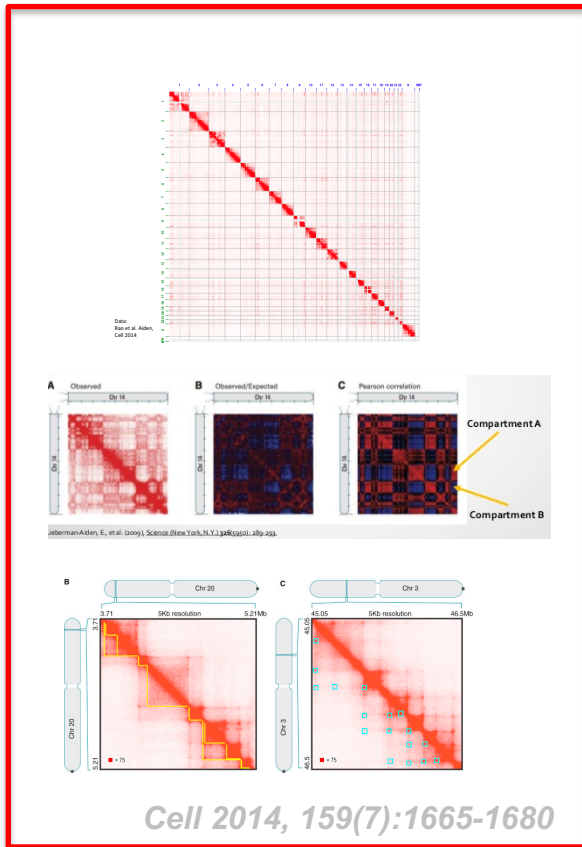
image credit: Iyer et al. BMC Biophysics 2011

# Hi-C contact map



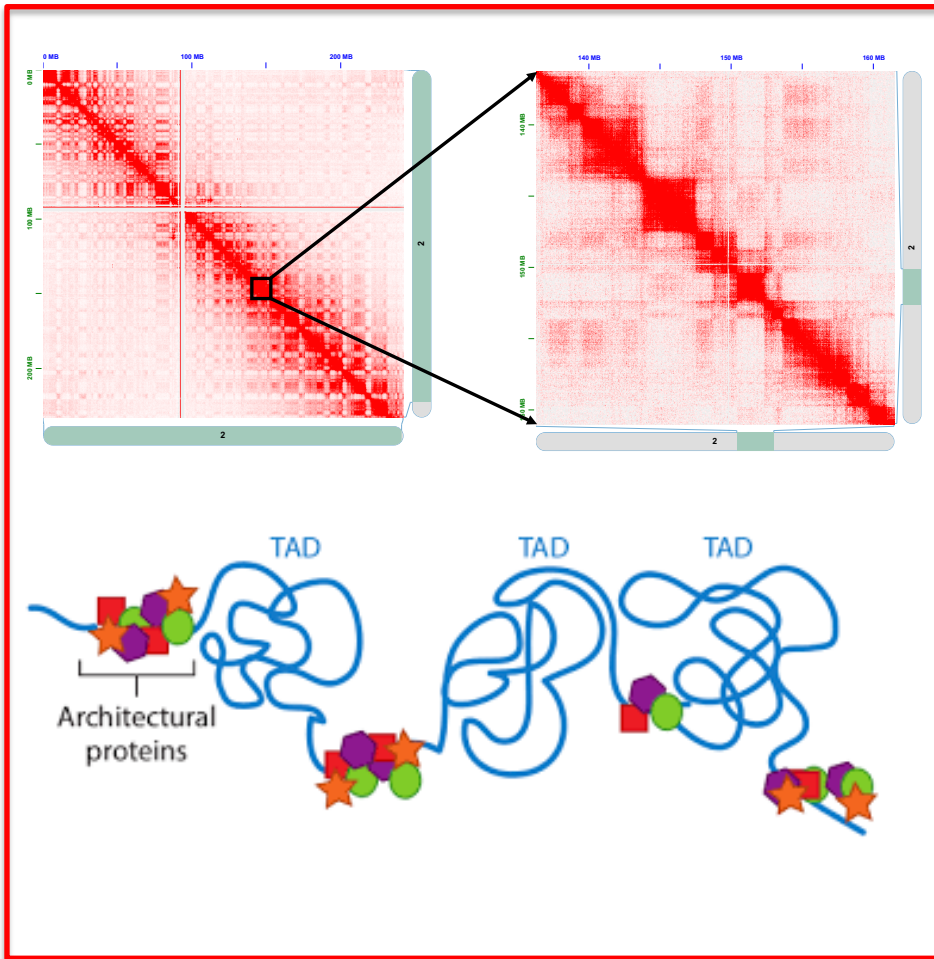
Science 2009, 5950: 289-293

# Hi-C contact map and Genome architecture

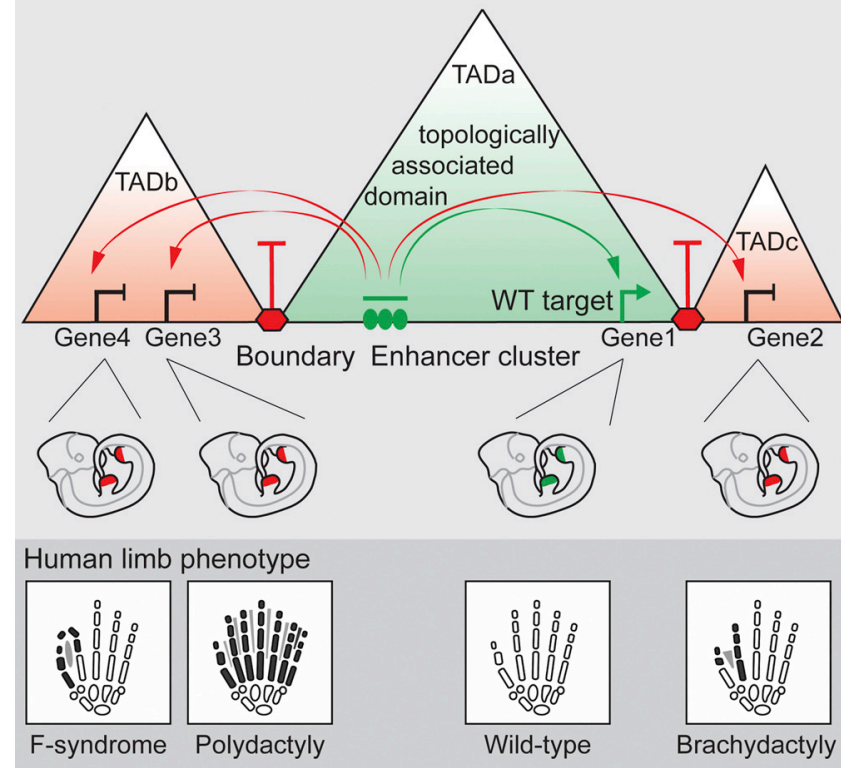




# Topologically Associating Domain



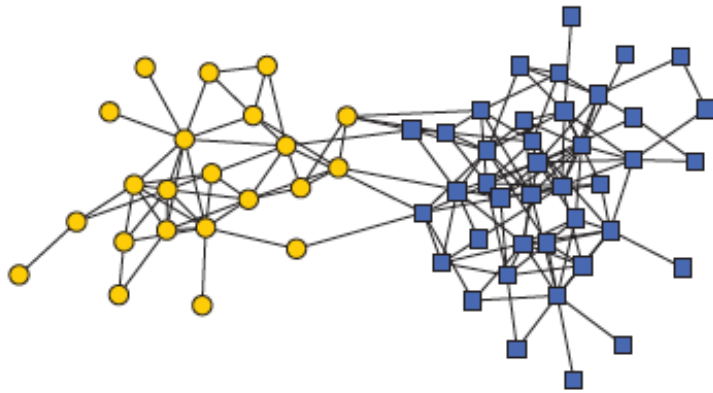
Structural variations affecting TAD boundaries



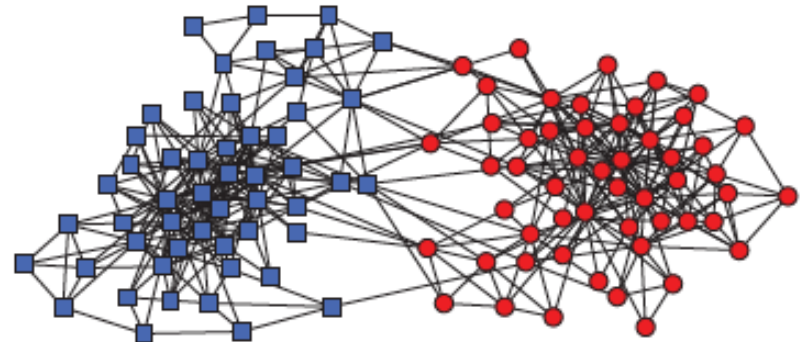
*Cell 2015, 161:1012-1025*

# Modularity

# Network modularity



Dolphin social network



Political books

Newman Phys. Rev. E 2013

adjacency matrix

$$Q = \frac{1}{2m} \sum_{i,j} \left( W_{ij} - \frac{k_i k_j}{2m} \right) \delta_{\sigma_i \sigma_j}$$

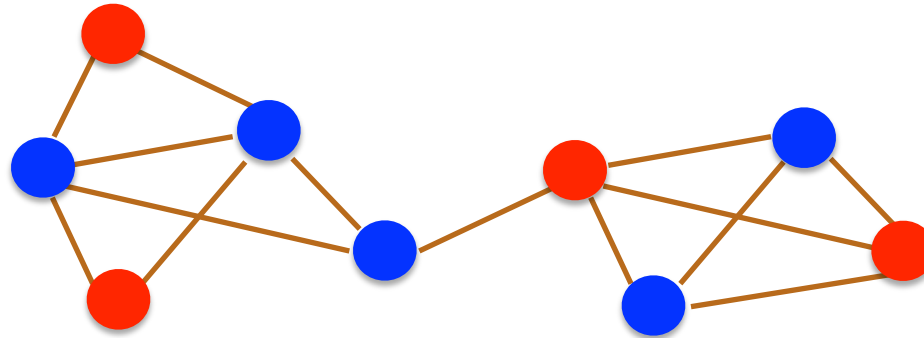
number of edges

degree of node  $i$

whether or not  $i, j$  are in the same module

expected number of edges between  $i$  and  $j$

# Network modularity

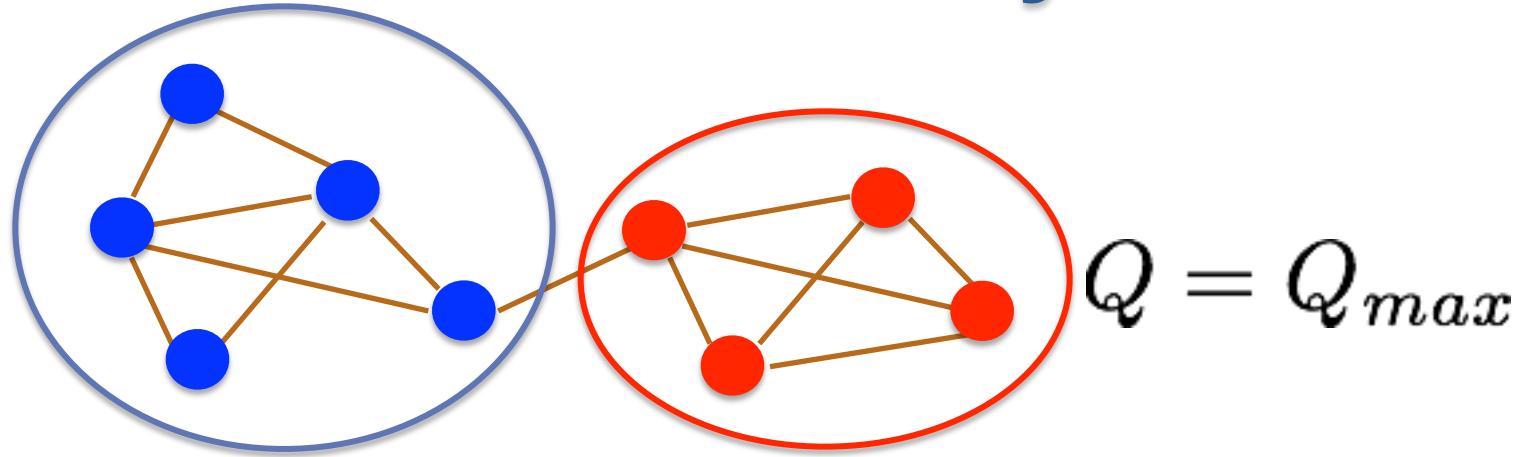


$$Q \approx 0$$

$$Q = \frac{1}{2m} \sum_{i,j} \left( W_{ij} - \frac{k_i k_j}{2m} \right) \delta_{\sigma_i \sigma_j}$$

adjacency matrix  $W_{ij}$   
 degree of node  $i$   $k_i$   
 number of edges  $2m$   
 expected number of edges between  $i$  and  $j$   $\frac{k_i k_j}{2m}$   
 whether or not  $i, j$  are in the same module  $\delta_{\sigma_i \sigma_j}$

# Network modularity



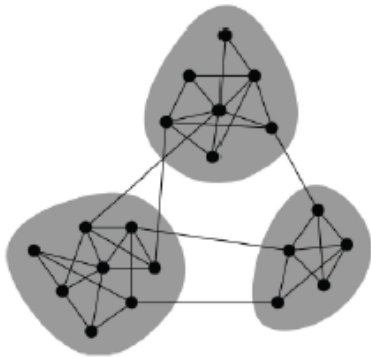
Optimization  
problem  
for sim.  
annealing

$$Q = \frac{1}{2m} \sum_{i,j} \left( W_{ij} - \frac{k_i k_j}{2m} \right) \delta_{\sigma_i \sigma_j}$$

adjacency matrix  $W_{ij}$   
 degree of node  $i$   $k_i$   
 whether or not  $i, j$  are in the same module  $\delta_{\sigma_i \sigma_j}$   
 number of edges  $2m$   
 expected number of edges between  $i$  and  $j$   $\frac{k_i k_j}{2m}$

# TAD Finding

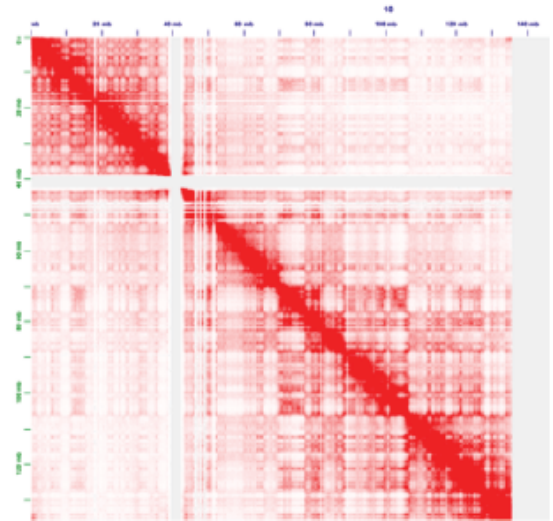
# Identifying TADs in multiple resolutions



Modularity maximization

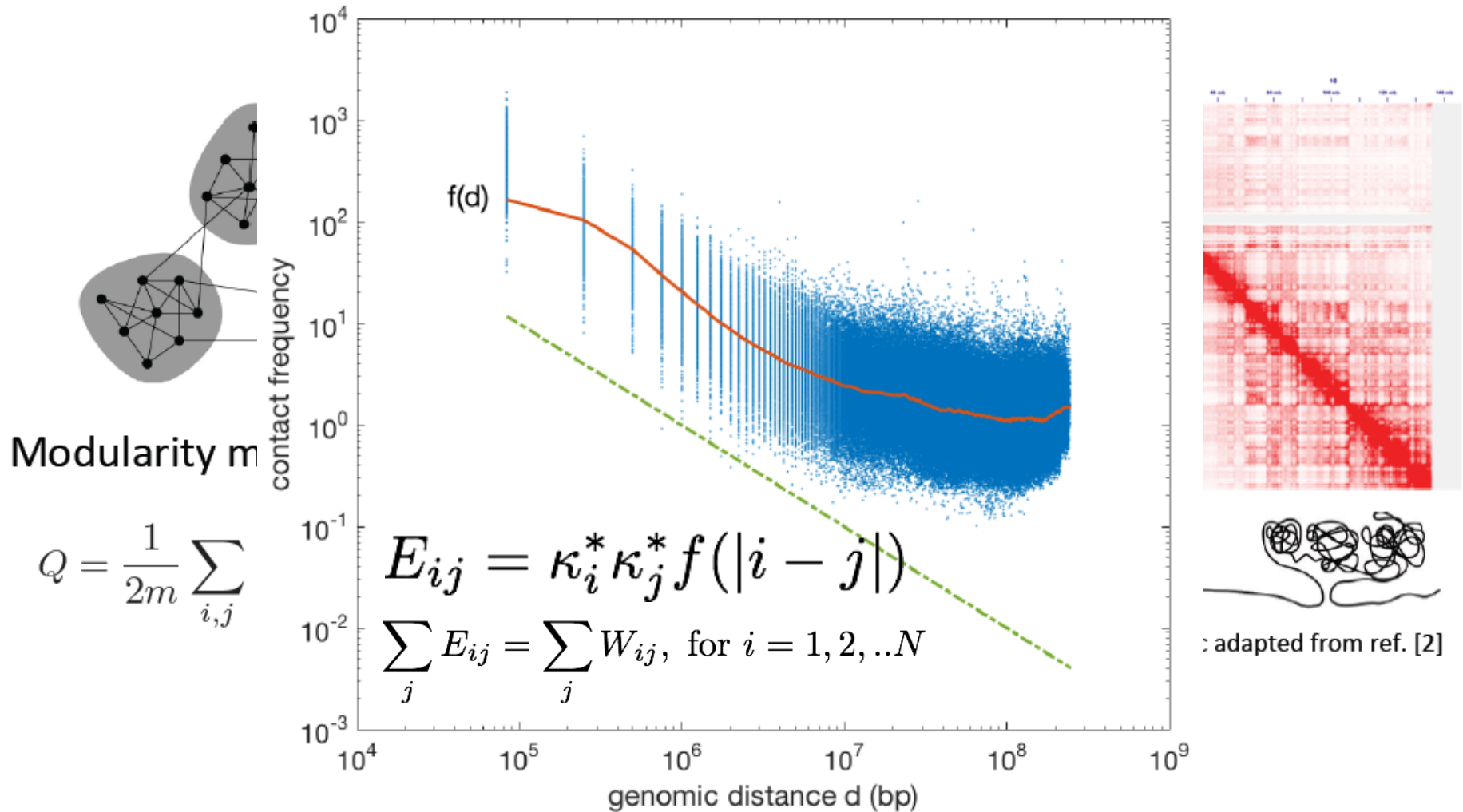
$$Q = \frac{1}{2m} \sum_{i,j} \left( W_{ij} - \frac{k_i k_j}{2m} \right) \delta_{\sigma_i \sigma_j}$$

network	contact map
node	chromosome bin
edge	Hi-C contact
# of connections	coverage
module	domain



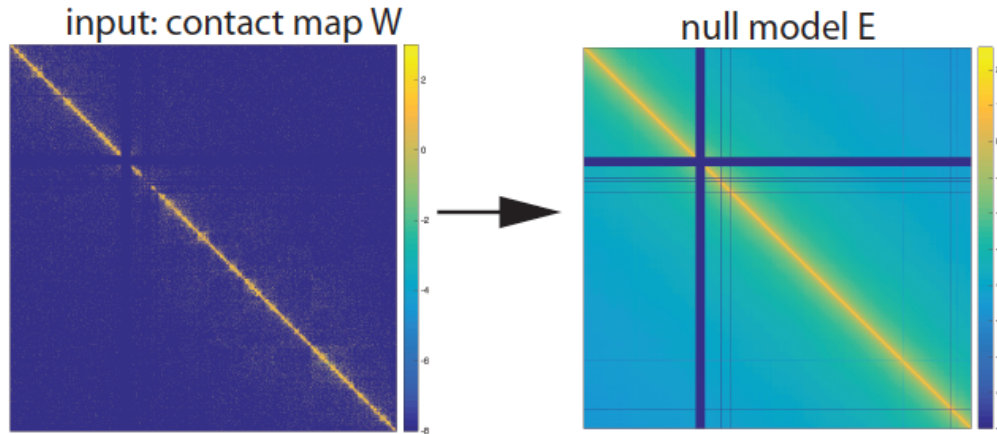
schematic adapted from ref. [2]

# Identifying TADs in multiple resolutions





## Identifying TADs in multiple resolutions



$$E_{ij} = \kappa_i^* \kappa_j^* f(|i - j|)$$

Numerically solve for  $\kappa_i^*$  in equations

$$\sum_j E_{ij} = \sum_j W_{ij}, \text{ for } i = 1, 2, \dots, N$$

Choose a particular resolution  $\gamma$   
Optimize  $Q$  over all possible partitions

$$Q = \frac{1}{2N} \sum_{ij} (W_{ij} - \gamma E_{ij}) \delta_{\sigma_i \sigma_j} \quad \gamma: \text{resolution parameter}$$

Multiple runs to define boundary scores  
for all pairs of adjacent bins

consensus boundaries based on  
the boundary scores

consensus TADs output