Weill Cornell Medical College

Institute for Computational Biomedicine
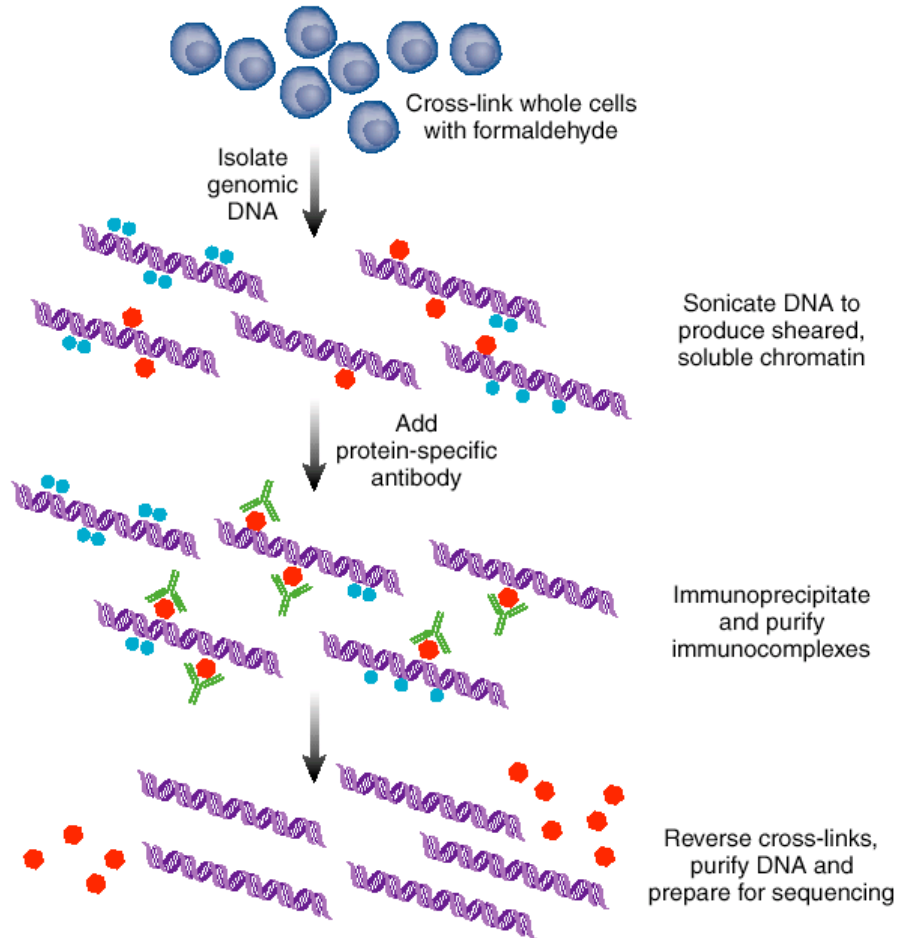
# ChIP-seq and Hi-C

Olivier Elemento, PhD

Laboratory of Cancer Systems Biology

# Plan

1. ChIP-seq
2. A few interesting ChIP-seq papers
3. Quality Control of ChIP-seq data
4. ChIP-seq Peak detection
5. Peak Analysis and Interpretation
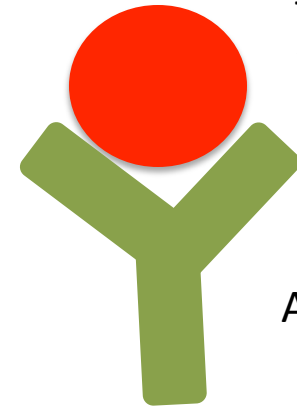6. Mapping chromatin interactions using Hi-C

# ChIP-seq



Cross-link whole cells with formaldehyde

Isolate genomic DNA

Sonicate DNA to produce sheared, soluble chromatin
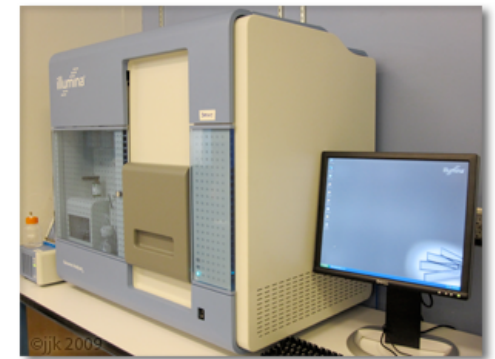
Add protein-specific antibody

Immunoprecipitate and purify immunocomplexes

Reverse cross-links, purify DNA and prepare for sequencing

Katie Ris

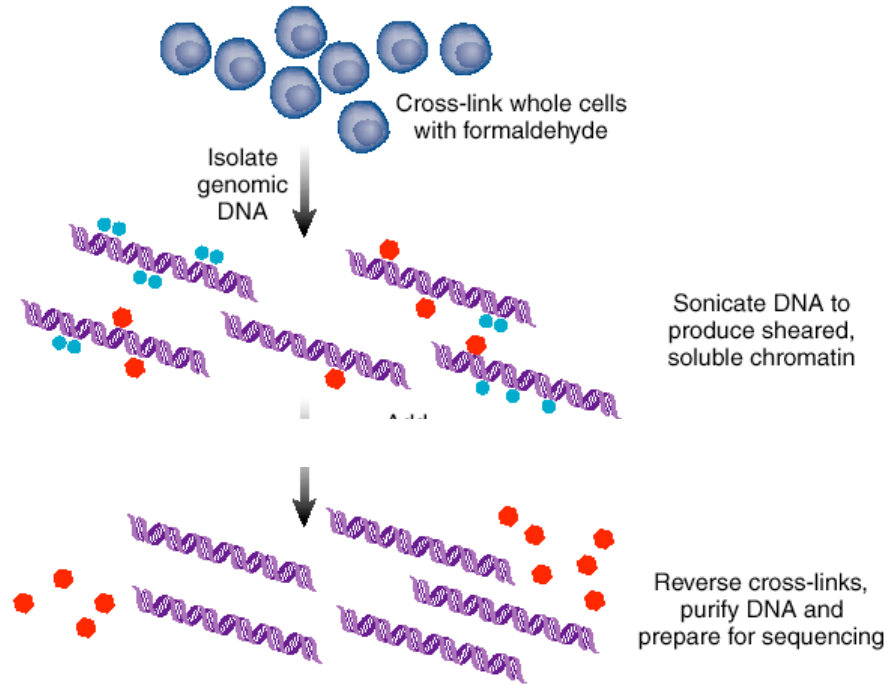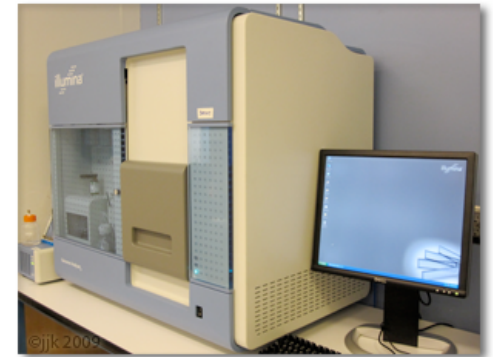Transcription factor of interest (or histone modification)

Antibody

Illumina

# Control: input DNA



Cross-link whole cells with formaldehyde

Isolate genomic DNA

Sonicate DNA to produce sheared, soluble chromatin

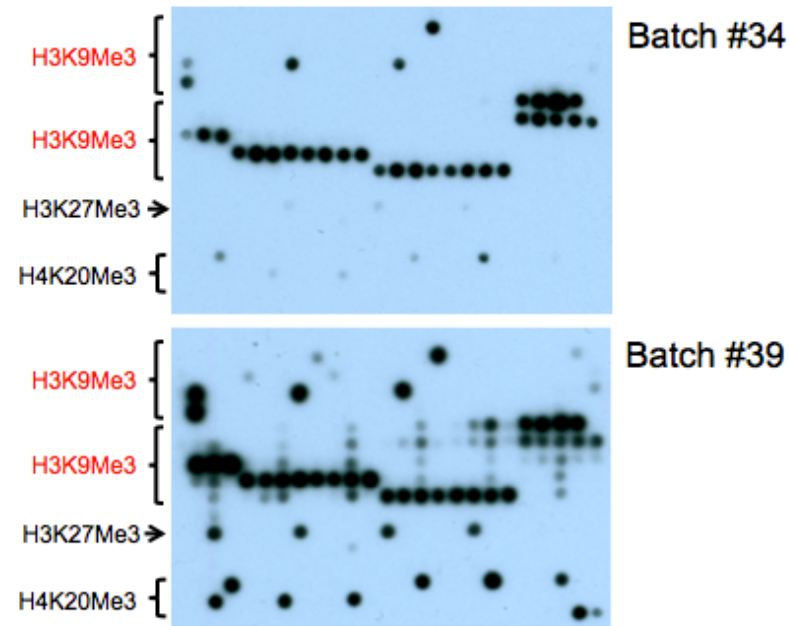Reverse cross-links, purify DNA and prepare for sequencing

Katie Ris

Illumina

Can use IgG as additional control

# ChIP-seq methodology

- Identify ChIP-grade antibody, determine specificity (Western, histone peptide array)

- Optimize conditions using single-locus ChIP-PCR (positive and negative controls)

- Sequence ChIP product using 1 Illumina lane per sample (no TruSeq ChIP-seq), single end

- Sequence input/IgG as control

Abcam H3K9Me3 rabbit polyclonal (ab8898)



Assessing the specificity of a commercial H3K9m3 antibody using histone peptide arrays, K. Bunting & B. Swed, WCMC

# First ChIP-seq paper

## Genome-Wide Mapping of in Vivo Protein-DNA Interactions

David S. Johnson,[1]* Ali Mortazavi,[2]* Richard M. Myers,[1]† Barbara Wold[2,3]†

In vivo protein-DNA interactions connect each transcription factor with its direct targets to form a gene network scaffold. To map these protein-DNA interactions comprehensively across entire mammalian genomes, we developed a large-scale chromatin immunoprecipitation assay (ChIPSeq) based on direct ultrahigh-throughput DNA sequencing. This sequence census method was then used to map in vivo binding of the neuron-restrictive silencer factor (NRSF; also known as REST, for repressor element–1 silencing transcription factor) to 1946 locations in the human genome. The data display sharp resolution of binding position [±50 base pairs (bp)], which facilitated our finding motifs and allowed us to identify noncanonical NRSF-binding motifs. These ChIPSeq data also have high sensitivity and specificity [ROC (receiver operator characteristic) area $\geq$ 0.96] and statistical confidence ($P < 10^{-4}$), properties that were important for inferring new candidate interactions. These include key transcription factors in the gene network that regulates pancreatic islet cell development.

Although much is known about transcription factor binding and action at specific genes, far less is known about the composition and function of entire factor-DNA chromosome can be detected by chromatin immunoprecipitation (ChIP) (*1*). In ChIP experiments, an immune reagent specific for a DNA binding factor is used to enrich target DNA putational discovery of binding motifs feasible, this dictates the quality of regulatory site annotation relative to other gene anatomy landmarks, such as transcription start sites, enhancers, introns and exons, and conserved noncoding features (*2*). Finally, if high-quality protein-DNA interactome measurements can be performed routinely and at reasonable cost, it will open the way to detailed studies of interactome dynamics in response to specific signaling stimuli or genetic mutations. To address these issues, we turned to ultrahigh-throughput DNA sequencing to gain sampling power and applied size selection on immuno-enriched DNA to enhance positional resolution.

The ChIPSeq assay shown here differs from other large-scale ChIP methods such as ChIPArray, also called ChIPchip (*1*); ChIPSAGE (SACO) (*3*); or ChIPPet (*4*) in design, data produced, and cost. The design is simple (Fig. 1A) and, unlike SACO or ChIPPet, it involves no plasmid library construction. Unlike microarray assays, the vast majority of single-copy sites in the genome is accessible for ChIPSeq assay (*5*), rather than a subset selected to be array features.

2007

# Epigenetic modifications at enhancer regions

nature

## LETTERS

## Histone modifications at human enhancers reflect global cell-type-specific gene expression

Nathaniel D. Heintzman[1,2]*, Gary C. Hon[1,3]*, R. David Hawkins[1]*, Pouya Kheradpour[5], Alexander Stark[5,6], Lindsey F. Harp[1], Zhen Ye[1], Leonard K. Lee[1], Rhona K. Stuart[1], Christina W. Ching[1], Keith A. Ching[1], Jessica E. Antosiewicz-Bourget[7], Hui Liu[8], Xinmin Zhang[8], Roland D. Green[8], Victor V. Lobanenkov[9], Ron Stewart[7], James A. Thomson[7,10], Gregory E. Crawford[11], Manolis Kellis[5,6] & Bing Ren[1,4]

The human body is composed of diverse cell types with distinct functions. Although it is known that lineage specification depends on cell-specific gene expression, which in turn is driven by promoters, enhancers, insulators and other cis-regulatory DNA sequences for each gene[1–3], the relative roles of these regulatory elements in this process are not clear. We have previously developed a chromatin-immunoprecipitation-based microarray method (ChIP-chip) to locate promoters, enhancers and insulators in the human genome[4–6]. Here we use the same approach to

Next, we identified putative insulators in the ENCODE regions for these cell types based on CTCF binding, because mammalian insulators are generally understood to require CTCF to block promoter–enhancer interactions[3]. We observed nearly identical CTCF occupancy (Supplementary Table 1 and Supplementary Fig. 1e) and highly correlated CTCF enrichment patterns across all five cell types (Supplementary Fig. 1b), providing experimental support for the mostly cell-type-invariant function of CTCF as suggested by DNase hypersensitivity mapping results[8].
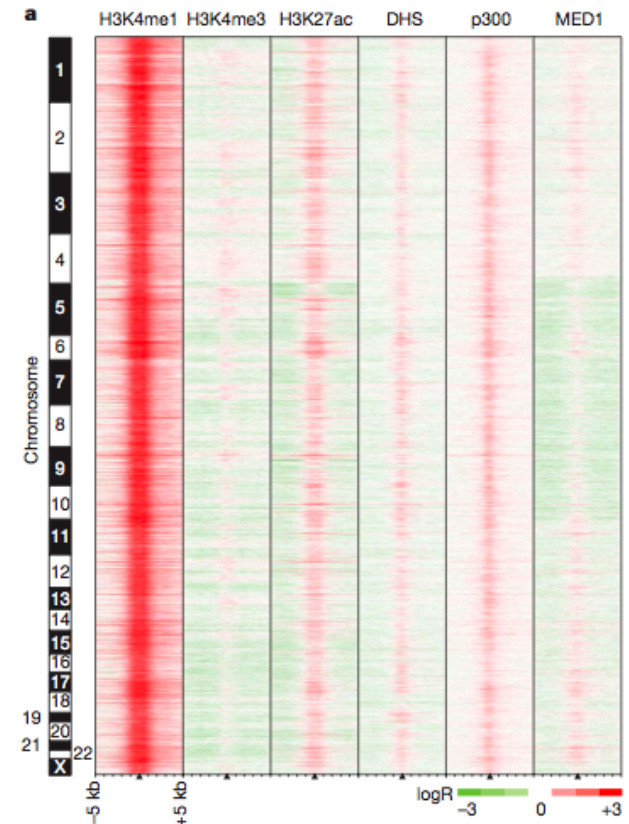
**Figure 2 | Genome-wide enhancer predictions in human cells. a,** We predict 36,589 enhancers in HeLa cells on the basis of chromatin signatures for H3K4me1 and H3K4me3 as determined by ChIP-chip using genome-wide tiling microarrays and condensed enhancer microarrays (see Supplementary Information). Enhancer predictions are located at the centre of 10-kb

# Chromatin states

# LETTER

# Differential oestrogen receptor binding is associated with clinical outcome in breast cancer

Caryn S. Ross-Innes[1], Rory Stark[1], Andrew E. Teschendorff[2], Kelly A. Holmes[1], H. Raza Ali[1,8], Mark J. Dunning[1], Gordon D. Brown[1], Ondrej Gojis[3,4,5], Ian O. Ellis[6], Andrew R. Green[6], Simak Ali[3], Suet-Feung Chin[1], Carlo Palmieri[3], Carlos Caldas[1,7,8,9] & Jason S. Carroll[1,7]
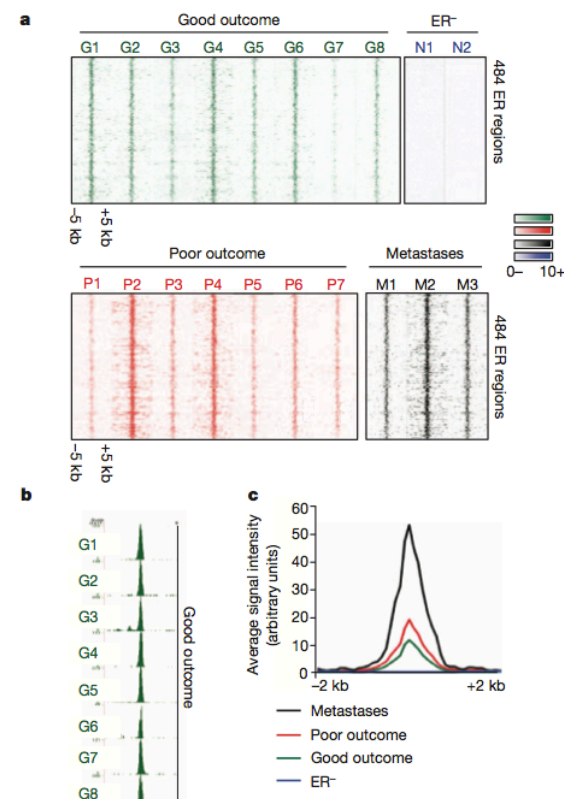
Oestrogen receptor-α (ER) is the defining and driving transcription factor in the majority of breast cancers and its target genes dictate cell growth and endocrine response, yet genomic understanding of ER function has been restricted to model systems[1–3]. Here we map genome-wide ER-binding events, by chromatin immunoprecipitation followed by high-throughput sequencing (ChIP-seq), in primary breast cancers from patients with different clinical outcomes and in distant ER-positive metastases. We find that drug-resistant cancers still recruit ER to the chromatin, but that ER binding is a dynamic process, with the acquisition of unique ER-binding regions in tumours from patients that are likely to relapse. The acquired ER regulatory regions associated with poor clinical outcome observed in primary tumours reveal gene signatures that predict clinical outcome in ER-positive disease exclusively. We find that the differential ER-binding programme observed in tumours from patients with poor outcome is not due to the selection of a rare subpopulation of cells, but is due to the FOXA1-mediated reprogramming of ER binding on a rapid timescale. The parallel redistribution of ER and FOXA1 binding events

breast cancer (Supplementary F[...] metastatic samples from women[...] metastatic locations and samp[...] plementary Fig. 1. As a control,[...] that were ER⁻ (ER-α negative),[...] ER-β.
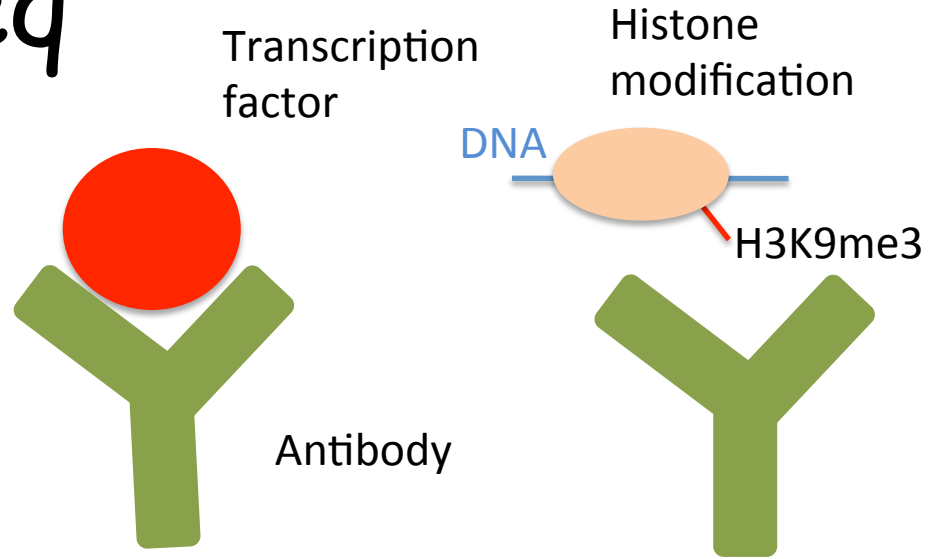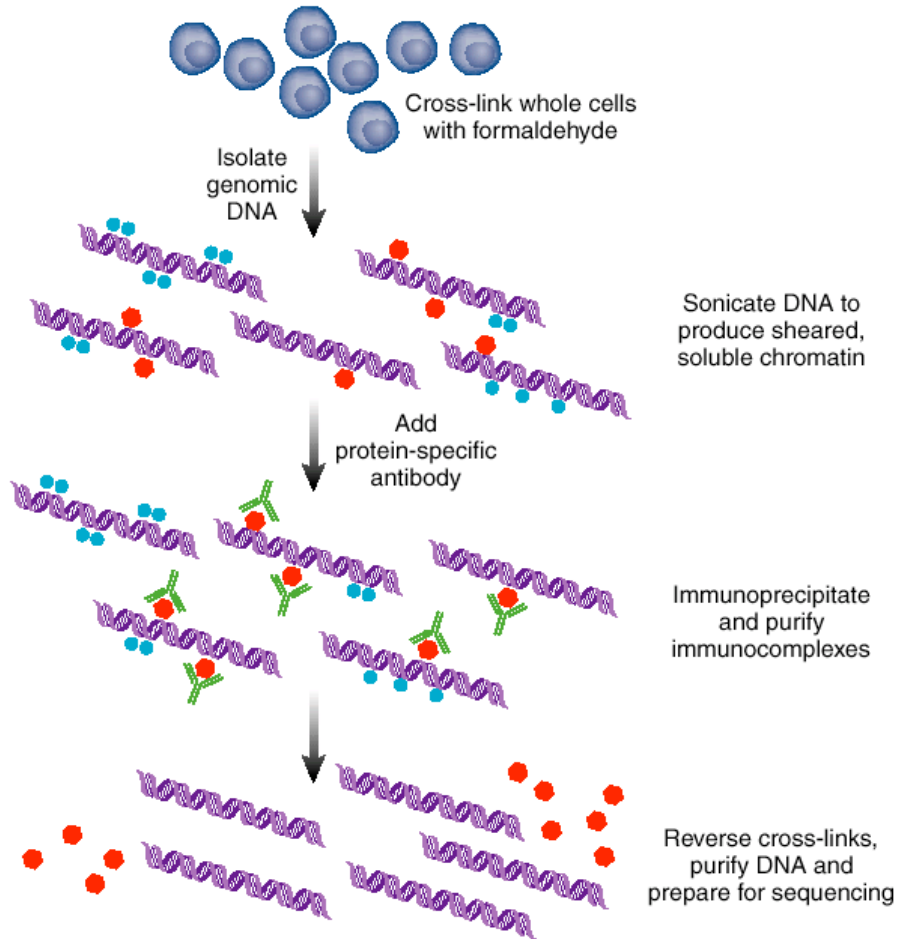
ER ChIP-seq was conducte[...] using two different algorithms,[...] ebi.ac.uk/~swilder/SWEMBL/)[...] number of sequencing reads an[...] is shown in Supplementary Fig.[...] tumours, but total peak intensit[...] events differed. Three tumour[...] ChIP-seq was conducted on the[...] concordance when comparing[...] ($R^2 = 0.954$) suggesting that tu[...] tially influence the ER-binding[...] plementary Fig. 3).
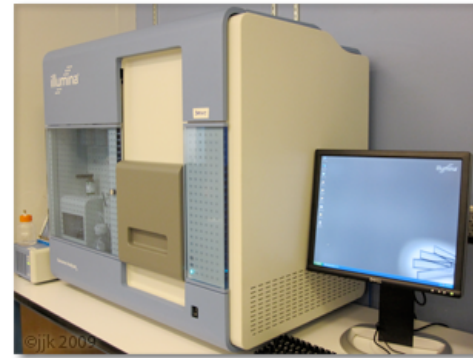
We initially assessed whethe[...]

# ChIP-seq

Transcription factor

Histone modification

DNA

H3K9me3

Antibody



Cross-link whole cells with formaldehyde

Isolate genomic DNA

Sonicate DNA to produce sheared, soluble chromatin

Add protein-specific antibody

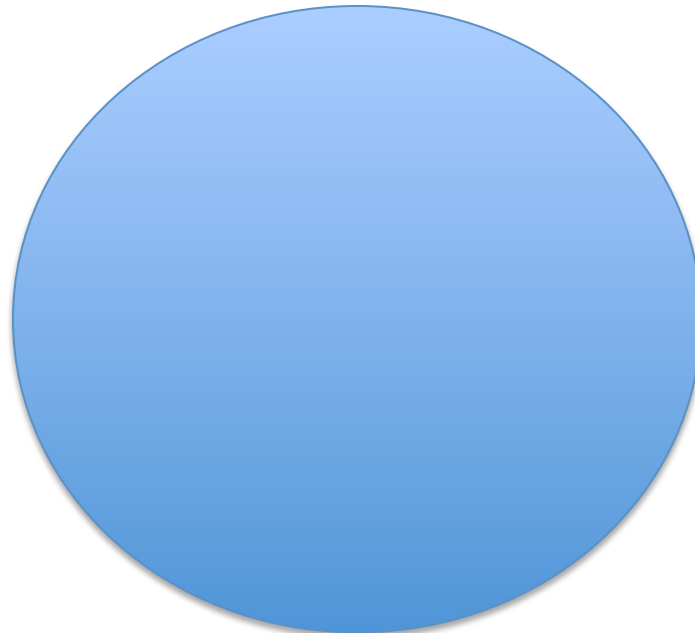Immunoprecipitate and purify immunocomplexes
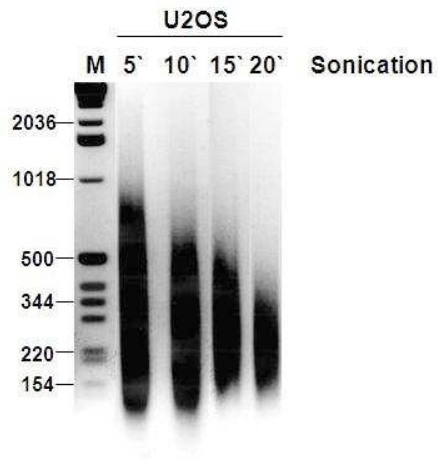
Reverse cross-links, purify DNA and prepare for sequencing

Katie Ris

Illumina

ACCAATAACCGAGGCTCATGCTAAGGCGTTAGCCACAGATG**GAAGTCCGA**CGGCTTGATCCAGAATGGTGTGTGGATTGCCTTGGAACTGATTAGTGAATTC
TGGTTATTGGCTCCGAGTACGATTCCGCAATCGGTGTCTAC**CTTCAGGCT**GCCGAACTAGGTCTTACCACACACCTAACGGAACCTTGACTAATCACTTAAG
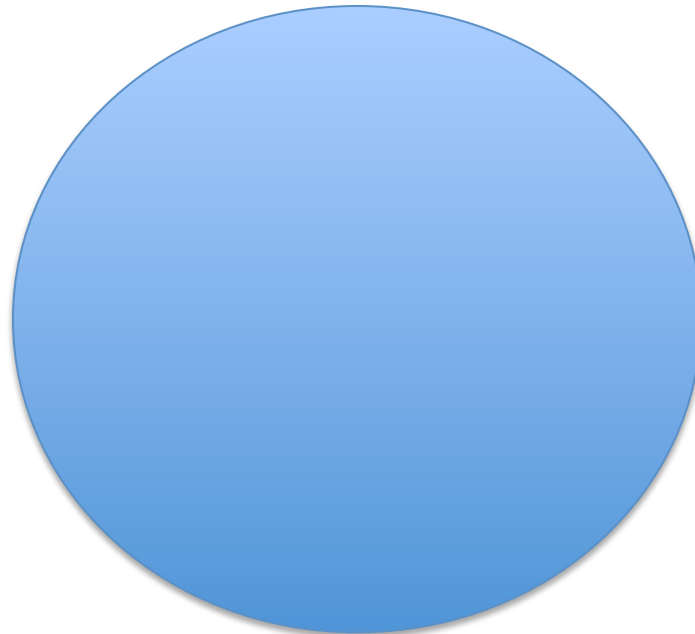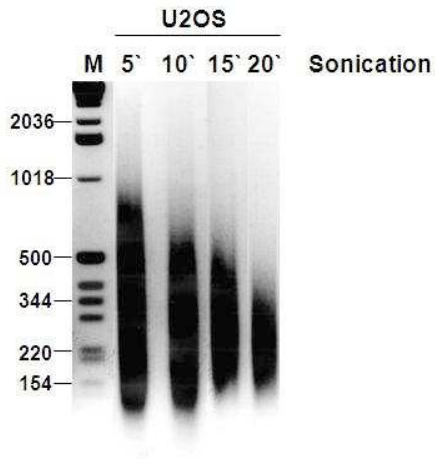
Average length ~ 170bp

40-100bp

←——————————→

**ACCAATAACCGAGGCTCA**TGCTAAGGCGTTAGCCACAGATG**GAAGTCCGA**CGGCTTGATCCAGAATGGTGTGTGGATTGCCTTGGAACTGATTAGTGAATTC
TGGTTATTGGCTCCGAGTACGATTCCGCAATCGGTGTCTAC**CTTCAGGCT**GCCGAACTAGGTCTTACCACACACCTAACGGAACCTTGACTAATCACTTAAG

←————————————————————————————————————————————→

Average length ~ 170bp

U2OS

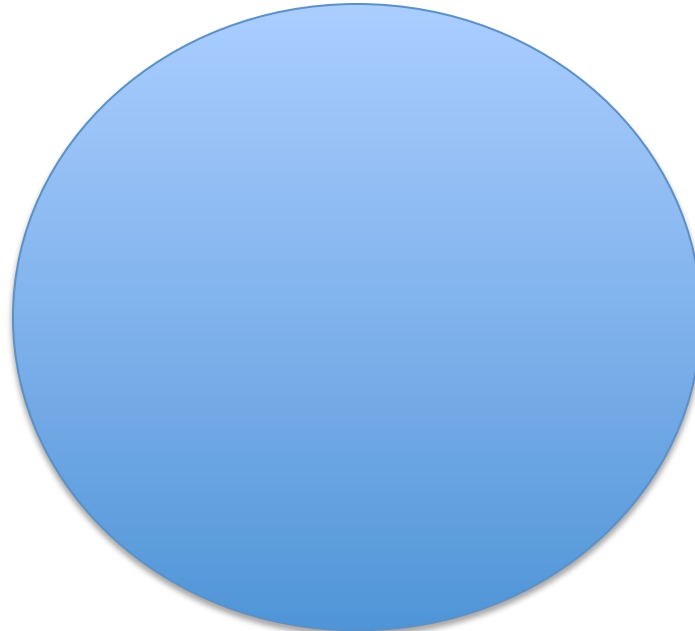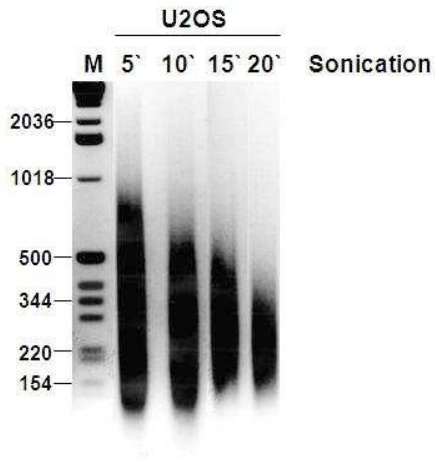M  5`  10` 15` 20`  Sonication

2036

1018

500

344

220

154

40-100bp

ACCAATAACCGAGGCTCATGCTAAGGCGTTAGCCACAGATGGAAGTCCGACGGCTTGATCCAGAATGGTGTGTGGATTGCCTTGGAACTGATTAGTGAATTC
TGGTTATTGGCTCCGAGTACGATTCCGCAATCGGTGTCTACCTTCAGGCTGCCGAACTAGGTCTTACCACACACCTAACGGAACCTTGACTAATCACTTAAG

Average length ~ 170bp

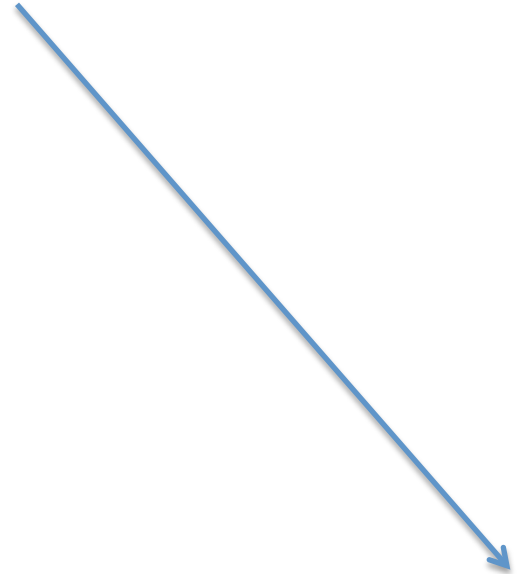# BWA tutorial (for aligning single end reads to genome)

- Get genome, e.g., from UCSC
  - http://hgdownload.cse.ucsc.edu/goldenPath/hg19/bigZips/chromFa.tar.gz

- Combine into 1 file
  - tar zvfx chromFa.tar.gz
  - cat *.fa > wg.fa

- Indexing the genome
  - bwa index -p hg19bwaidx -a bwtsw wg.fa

- Align ChIP reads to reference genome
  - bwa aln -t 4 hg19bwaidx s_3_sequence.txt.gz > s_3_sequence.txt.bwa
- Convert to SAM format
  - bwa samse hg19bwaidx s_3_sequence.txt.bwa s_3_sequence.txt.gz > s_3_sequence.txt.sam

- Align input reads to same reference genome
  - bwa aln -t 4 hg19bwaidx s_4_sequence.txt.gz > s_4_sequence.txt.bwa
- Convert to SAM format
  - bwa samse hg19bwaidx s_4_sequence.txt.bwa s_4_sequence.txt.gz > s_4_sequence.txt.sam

# Reads can map to multiple locations/chromosomes



Read 1

Read 2

Reference Human Genome (hg19)

# Reads map to one strand or the other

Read 1

Read 2

hg18

# SAM format

```
DH1608P1_0130:6:1103:10579:166379#TTAGGC    16    chr1   1249828 37    51M    *    0    0
GGGCGTGACTCTGATCTCAGGCATCGTCTCCGCCGCGCTCCCGGACCCGCG    eb`XXYbZdadee^ceV]X][ccTcc^ebeece
eeeWbeeeeeeeceeaee    XX:Z:NM_017871,32    NM:i:0 MD:Z:51

DH1608P1_0130:6:1102:3415:150915#TTAGGC 16    chr1   1249828 37    51M    *    0    0
GGGCGGGACTCTGATCTCAGGCATCGTCTCCGCCGCGCTCCCGGACCCGCG    BBBBBBBBBBBac]bbbceedaeddeZceeea_ba_\_eee
eeeedaeeee    XX:Z:NM_017871,32    NM:i:1 MD:Z:5T45

DH1608P1_0130:6:1102:13118:62644#TTAGGC 16    chr1   1249828 37    51M    *    0    0
GGGCGTGCCTCGGATCTCAGGCATCGTCTCCGCCGCGCTCCCGGACCCGCG    BBBBBBBBBBBBBBBBBBBBBB`XTbSa`cffegdggeccbe
effdeggggg    XX:Z:NM_017871,32    NM:i:2 MD:Z:7A3T39

DH1608P1_0130:6:1203:3012:157120#TTAGGC 16    chr1   1249826 25    51M    *    0    0
AAGGCCGTGACTCTGATCTCAGCCCTCGTCTCCGCCGCGCTCCCGGACCCG    BBBBBBBB^`QWZZ]UXYSZSTFRU]Z__SO[adcc[acdV
\`Y]YWY][_    XX:Z:NM_017871,34    NM:i:3 MD:Z:4G17G1A26

DH1608P1_0130:6:2206:4445:12756#TTAGGC 16    chr1   1246336 25    1M3487N50M    *    0    0
CCAAAGGGTGTGACTCTGATCTCGGGCATCGTCTCCGCCGCGCTCCCGGAC    BBBBBBBBBBBBBBBBBBBBBBBB`YdddYdc\
cacaNddddcdddaeeee    XX:Z:NM_017871,37    NM:i:3 MD:Z:2C5C14A27

DH1608P1_0130:6:2203:7903:43788#TTAGGC 16    chr1   1246336 37    1M3487N50M    *    0    0
CCCAAGGGCGTGACTCTGATCTCAGGCATCGTCTCCGCCGCGCTCCCGGAC    adbe[fbcbccb_cb^cb^^c^edgeggggggdf
ggefffgfbfggggegeg    XX:Z:NM_017871,37    NM:i:0 MD:Z:51
```

CIGAR string, eg 5M3487N46M = 5bp-long block, 3487 insert, 46bp-long block

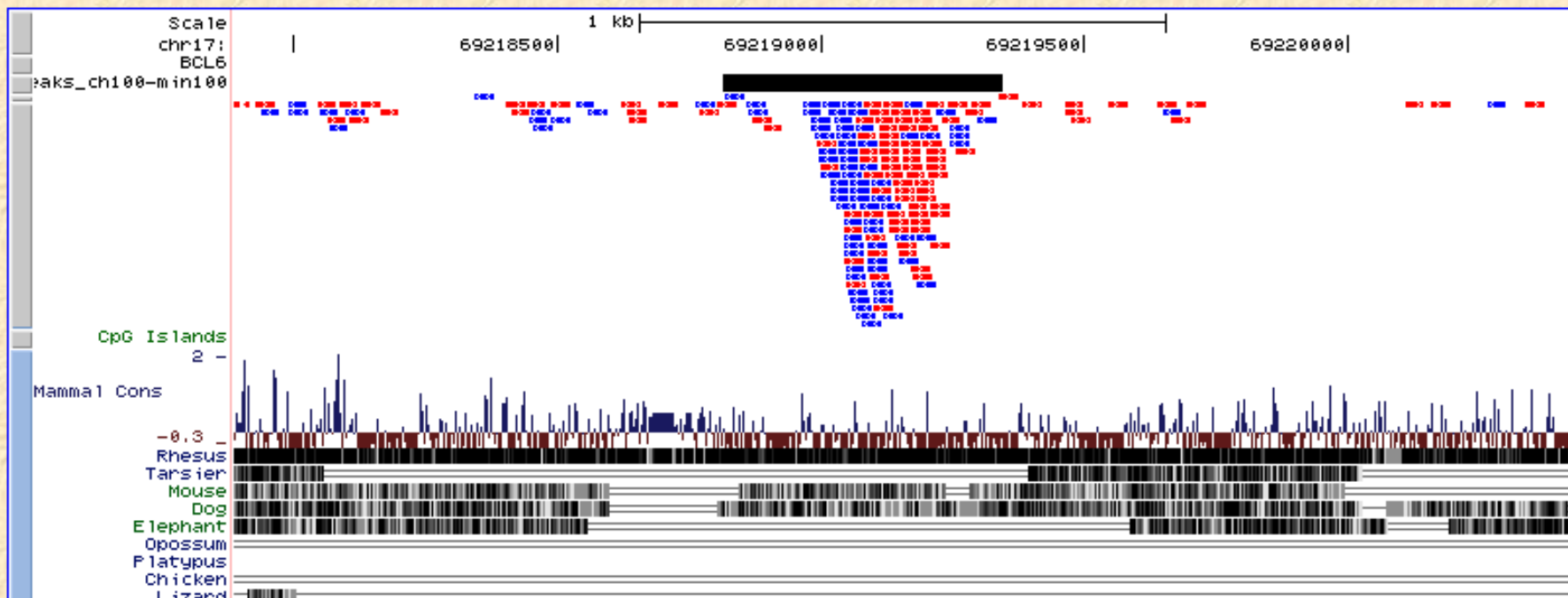MD tag, e.g, MD:Z:4T46 = 5 matches, 1 mismatch (T in read), 46 matches

XT tag, e.g. XT:A:U = unique mapper; XT:A:R = more than 1 high-scoring matches

91930500          91931000          91931500          91932000          91932500

User Supplied Track

# ChIP-seq peak calling

# A nice peak

# Not all peaks are nice

# MACS

Method

# Model-based Analysis of ChIP-Seq (MACS)

Yong Zhang[¤*], Tao Liu[¤*], Clifford A Meyer[*], Jérôme Eeckhoute[†],
David S Johnson[‡], Bradley E Bernstein[§¶], Chad Nussbaum[¶],
Richard M Myers[¥], Myles Brown[†], Wei Li[#] and X Shirley Liu[*]

Addresses: [*]Department of Biostatistics and Computational Biology, Dana-Farber Cancer Institute and Harvard School of Public Health, 44 Binney Street, Boston, MA 02115, USA. [†]Division of Molecular and Cellular Oncology, Department of Medical Oncology, Dana-Farber Cancer Institute and Department of Medicine, Brigham and Women's Hospital and Harvard Medical School, 44 Binney Street, Boston, MA 02115, USA. [‡]Gene Security Network, Inc., 2686 Middlefield Road, Redwood City, CA 94063, USA. [§]Molecular Pathology Unit and Center for Cancer Research, Massachusetts General Hospital and Department of Pathology, Harvard Medical School, 13th Street, Charlestown, MA 02129, USA. [¶]Broad Institute of Harvard and MIT, 7 Cambridge Center, Cambridge, MA, 02142, USA. [¥]Department of Genetics, Stanford University Medical Center, Stanford, CA 94305, USA. [#]Division of Biostatistics, Dan L Duncan Cancer Center, Department of Molecular and Cellular Biology, Baylor College of Medicine, One Baylor Plaza, Houston, TX 77030, USA.

¤ These authors contributed equally to this work.

Correspondence: Wei Li. Email: wl1@bcm.edu. X Shirley Liu. Email: xsliu@jimmy.harvard.edu

# MACS



Estimate d based on high quality peaks

The Poisson distribution

$$P(X \geq x) = 1 - \sum_{0}^{x-1} \frac{\lambda^x e^{-\lambda}}{x!}$$

$\lambda$=expected # of reads within an interval of 2*d* bp

$$\lambda_{\text{local}} = \max(\lambda_{BG}, [\lambda_{1k},] \lambda_{5k}, \lambda_{10k})$$

# in R P(X>=5|λ=0.001) is
1-sum(dpois(0:4, 0.001))

# BayesPeak

Research article

## BayesPeak: Bayesian analysis of ChIP-seq data

Christiana Spyrou*[1,3], Rory Stark[3], Andy G Lynch[4] and Simon Tavaré[2,4]

Address: [1]Statistical Laboratory, Centre for Mathematical Sciences, Wilberforce Road, Cambridge, UK, [2]DAMTP, Centre for Mathematical Sciences, Wilberforce Road, Cambridge, UK, [3]Cancer Research UK, Cambridge Research Institute, Li Ka Shing Centre, Robinson Way, Cambridge UK and [4]Department of Oncology, University of Cambridge, Li Ka Shing Centre, Robinson Way, Cambridge, UK

Email: Christiana Spyrou* - C.Spyrou@statslab.cam.ac.uk; Rory Stark - Rory.Stark@cancer.org.uk; Andy G Lynch - Andy.Lynch@cancer.org.uk; Simon Tavaré - st321@cam.ac.uk

* Corresponding author

## Abstract

**Background:** High-throughput sequencing technology has become popular and widely used to study protein and DNA interactions. Chromatin immunoprecipitation, followed by sequencing of the resulting samples, produces large amounts of data that can be used to map genomic features such as transcription factor binding sites and histone modifications.

# BayesPeak (Bayesian Hidden Markov Models)

Observed variable

$$Z_t = \begin{cases} 0 & \text{if} \quad (S_t, S_{t+1}) = (0,0) \\ 1 & \text{if} \quad (S_t, S_{t+1}) = (0,1) \\ 2 & \text{if} \quad (S_t, S_{t+1}) = (1,0) \\ 3 & \text{if} \quad (S_t, S_{t+1}) = (1,1) \end{cases}$$

The emission distributions of the model are

$$Y_t^+, Y_{t+1}^- \mid Z_t = 0 \sim \quad \text{Poisson}(\lambda_0 \gamma^{w_t})$$
$$Y_t^+, Y_{t+1}^- \mid Z_t = 1,2,3 \sim \quad \text{Poisson}((\lambda_0 + \lambda_1)\gamma^{w_t})$$
$$\lambda_0 \sim \quad \Gamma(\alpha_0, \beta_0)$$
$$\lambda_1 \sim \quad \Gamma(\alpha_1, \beta_1)$$

Parameters estimated using Bayesian treatment

Hidden states

# Other peak finders

**Table 1: Comparison of different peak-calling algorithms**

| Method | A | B | C | D | E | F | G |
|---|---|---|---|---|---|---|---|
| CSPF | control or IP only | read length no orientation | merge strands no shift | N | simple height criteria | ROC curve (empirically) | both |
| XSET | IP only | fragment length orientation | merge strands no shift | Y | simple height criteria | FDR estimate using Poisson distribution | both |
| Mikkelsen et al. | IP only | no orientation | no merge no shift | Y | p-values from permutations | no official FDR | both |
| MACS | control or IP only | fragment length orientation no duplicated reads | shift reads merge strands | N | Poisson p-values | FDR estimate by peaks in control:IP | both |
| QuEST | control | orientation | shift reads merge strands | N | kernel density estimation | FDR estimate by permutations of the control | better for TF |
| FindPeaks | IP only | fragment length orientation | no merge no shift | N | simple height criteria | FDR estimate by permutations of the IP | both |
| SISSR | control or IP only | fragment length orientation | no merge no shift | N | compares reads on different strands | FDR estimate by peaks in background:IP | better for TF |
| Kharchenko et al. | control | orientation | no merge no shift | N | Poisson distribution | FDR estimate by permutations of the control | better for TF |
| PeakSeq | control | fragment length orientation | merge strands | Y | sample normalisation Binomial distribution | FDR estimate, q-values (BH correction) | both |
| BayesPeak | control or IP only | fragment length orientation | no merge no shift | N | negative binomial distribution, Bayesian posterior probabilities | posterior enrichment probabilities | both |

The methods shown are compared with respect to the following features:
A. whether they require a control sample (control) or whether they only use the ChIP sample (IP only)
B. whether they take into account the (average) length of the reads/fragments and their orientation
C. whether they take into account the different DNA strands or if they merge the reads, and whether the reads are shifted towards the 3' end
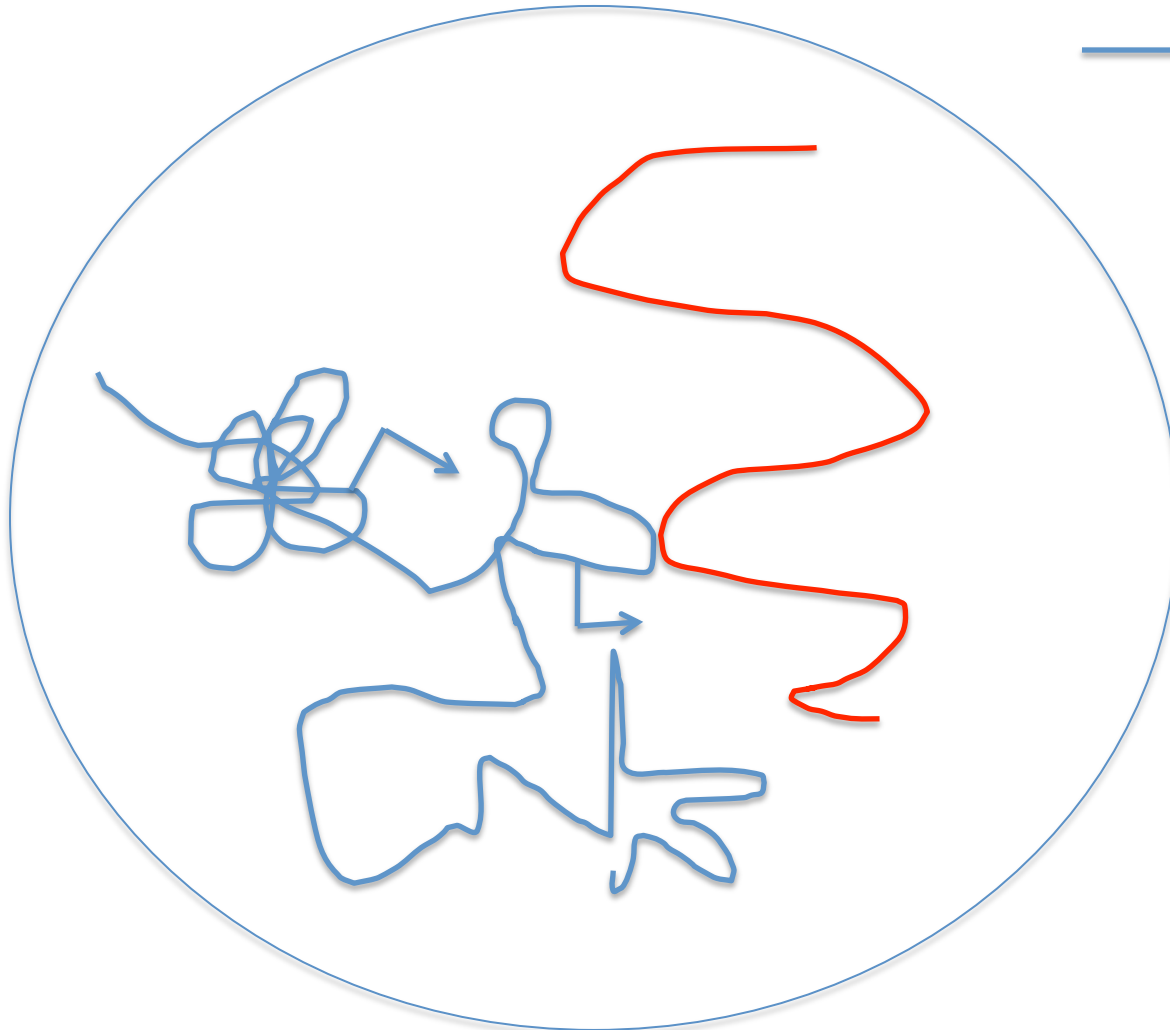D. whether an externally estimated mappability file is used
E. how the scores, on which the classifications are based, are estimated
F. whether/how any FDR or sensitivity/specificity estimates are calculated
G. whether or not the method is applicable to both transcription factor (TF) and histone mark data.

# Mapping chromatin interactions using the Hi-C technique
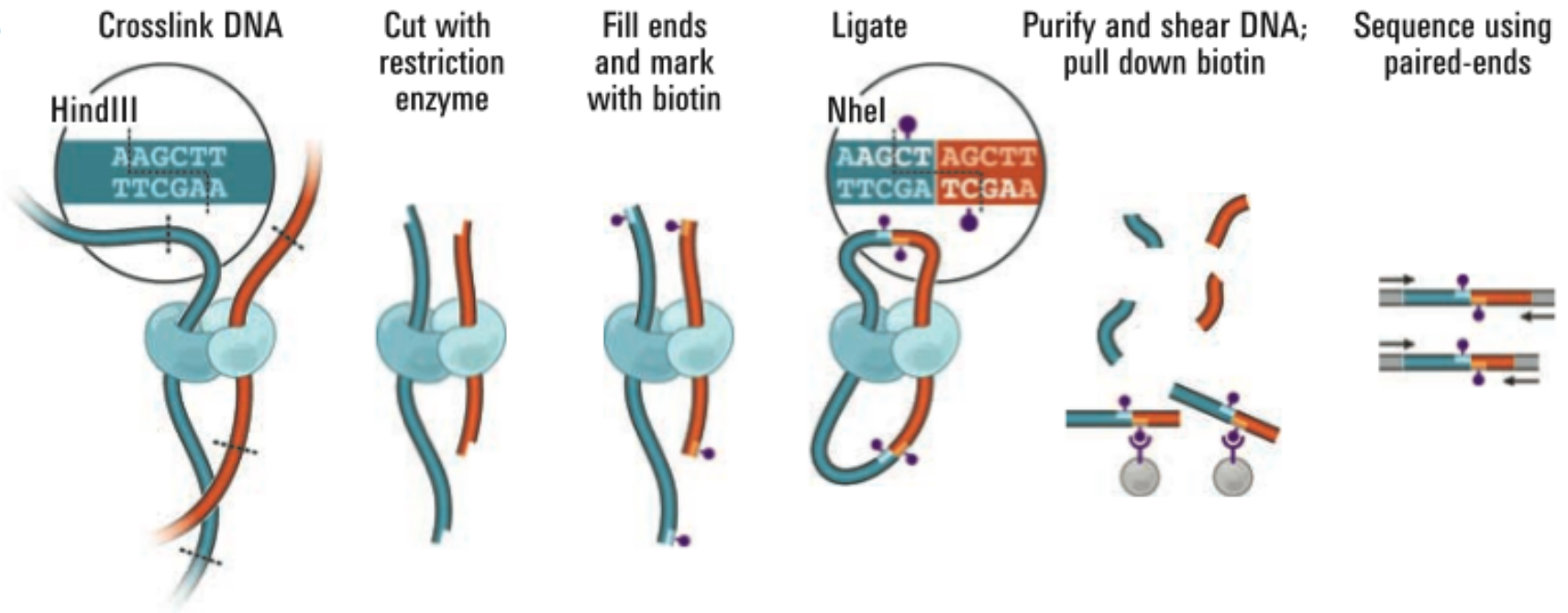
# The human genome is not linear



Terminology:
Intra-chromosomal interaction
Inter-chromosomal interaction
Interaction hub

How dynamic is the three dimensional architecture of the human genome ?

How does 3D localization impact gene expression ?

Is the 3D genome different in normal cells and cancer cells ?

# Genome-scale mapping of chromatin interactions using HiC



Lieberman-Aiden et al, 2009

Can a single (oncogenic) transcription factor induce global changes in chromatin structure ?

# RESEARCH ARTICLE

## Recurrent Fusion of *TMPRSS2* and ETS Transcription Factor Genes in Prostate Cancer

Scott A. Tomlins,[1] Daniel R. Rhodes,[1,2] Sven Perner,[7,9]
Saravana M. Dhanasekaran,[1] Rohit Mehra,[1] Xiao-Wei Sun,[7]
Sooryanarayana Varambally,[1,6] Xuhong Cao,[1] Joelle Tchinda,[7]
Rainer Kuefer,[10] Charles Lee,[7] James E. Montie,[3,5,6]
Rajal B. Shah,[1,3,5,6] Kenneth J. Pienta,[3,4,5,6] Mark A. Rubin,[7,8]
Arul M. Chinnaiyan[1,2,3,5,6]*

Recurrent chromosomal rearrangements have not been well characterized in common carcinomas. We used a bioinformatics approach to discover candidate oncogenic chromosomal aberrations on the basis of outlier gene expression. Two ETS transcription factors, *ERG* and *ETV1*, were identified as outliers in prostate cancer. We identified recurrent gene fusions of the 5' untranslated region of *TMPRSS2* to *ERG* or *ETV1* in prostate cancer tissues with outlier expression. By using fluorescence in situ hybridization, we dem-

(6). This karyotypic complexity is thought to reflect secondary genomic alterations acquired during tumor progression.
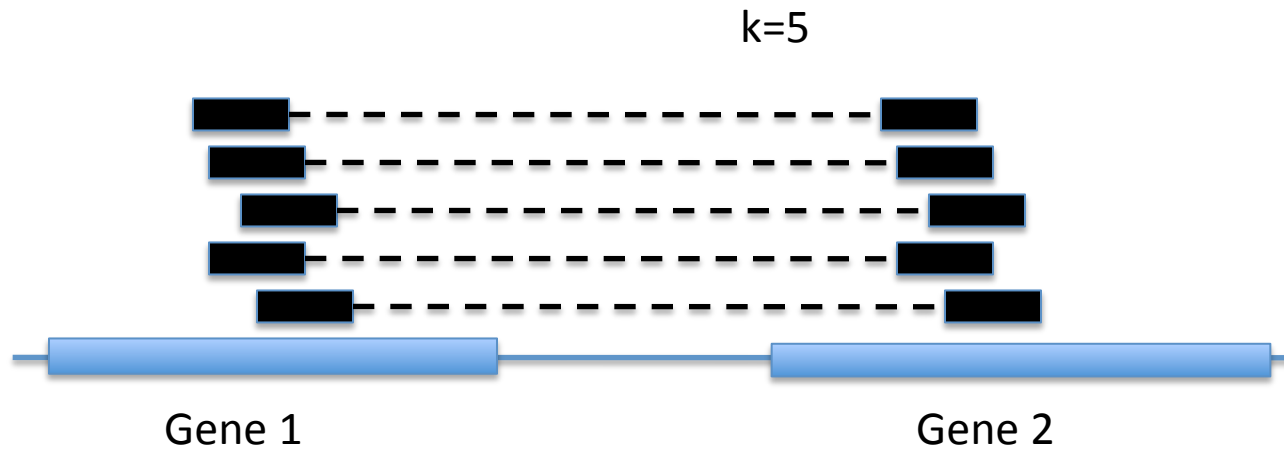
We hypothesized that rearrangements and high-level copy number changes that result in marked overexpression of an oncogene should be evident in DNA microarray data but not necessarily by traditional analytical approaches. In the majority of cancer types, heterogeneous patterns of oncogene activation have been observed; thus, traditional analytical methods that search for common activation of genes across a class of cancer samples (e.g., $t$ test or signal-to-noise ratio) will fail to find such oncogene expression profiles. Instead, a method that searches for marked overexpression in a subset of cases is needed. Toward this end, we developed a method termed cancer outlier profile analysis (COPA). COPA seeks to accentuate and identify outlier profiles by applying a simple numerical transformation based on the median and median absolute deviation of a gene expression profile (7) (fig. S1A).
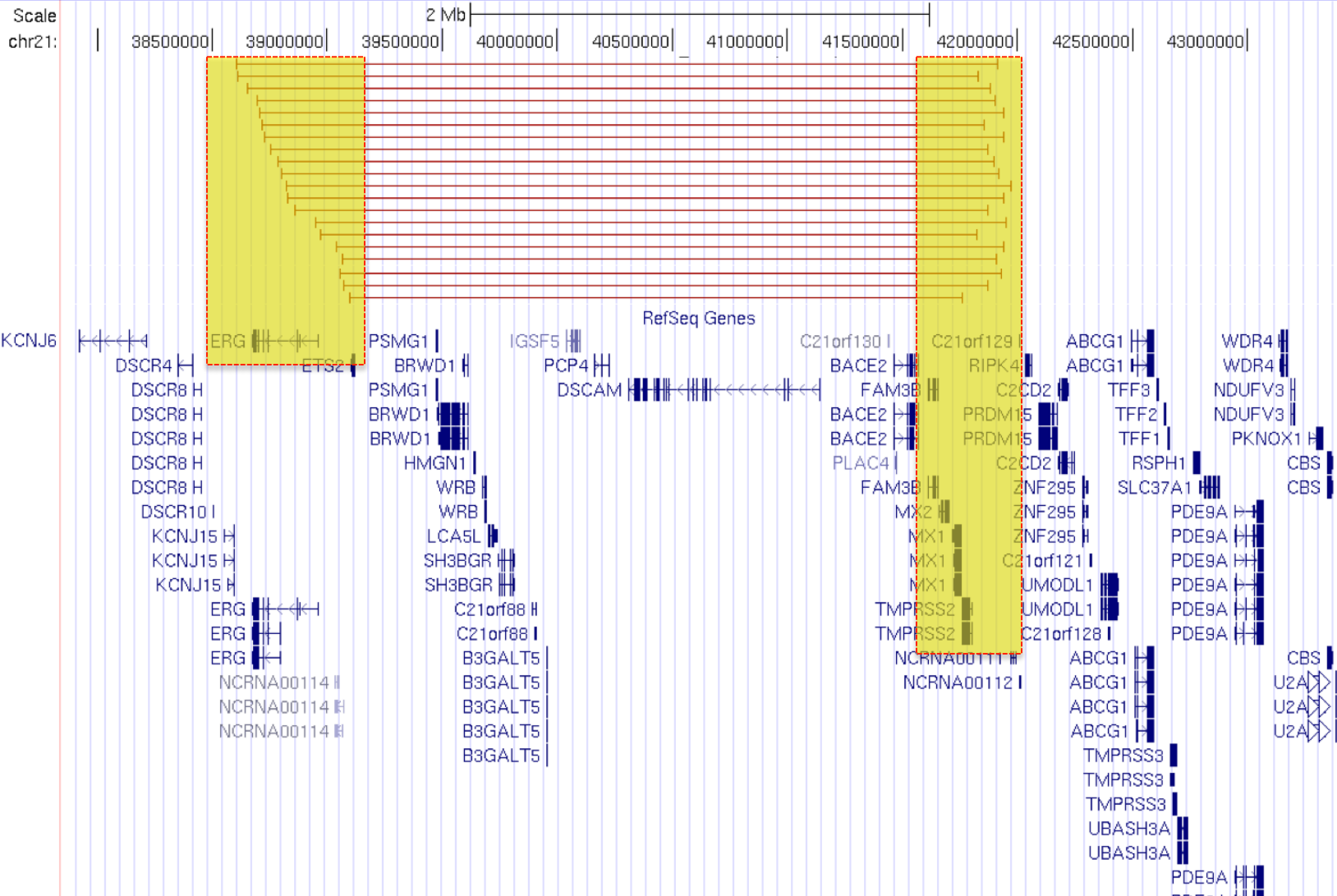
ERG is fused to TMPRSS2 and over-expressed in > 50% of prostate tumors

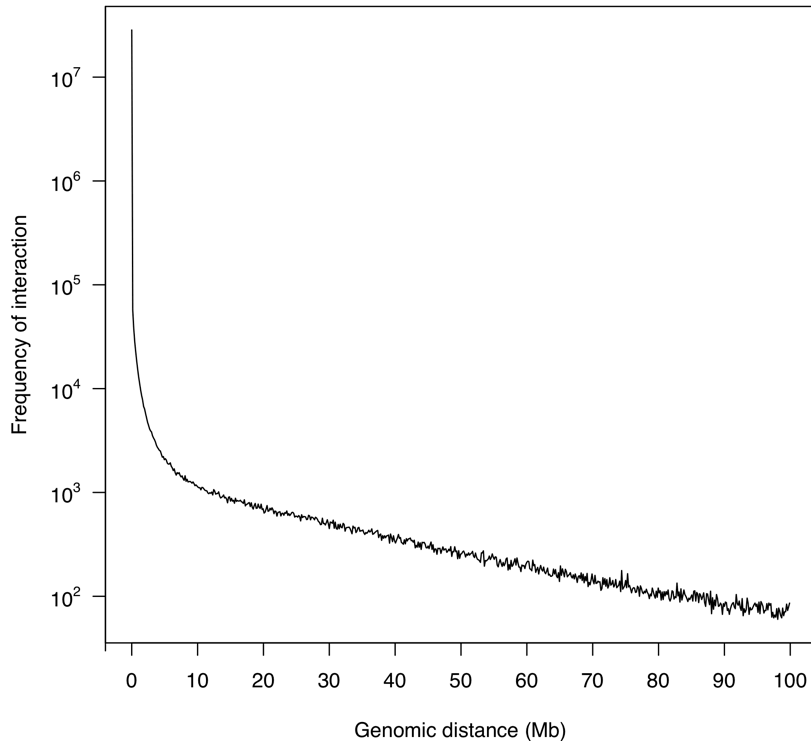# Prostate epithelial lesion cell line (RWPE1)

# Quantifying proximity using HiC reads



k=5

Gene 1

Gene 2

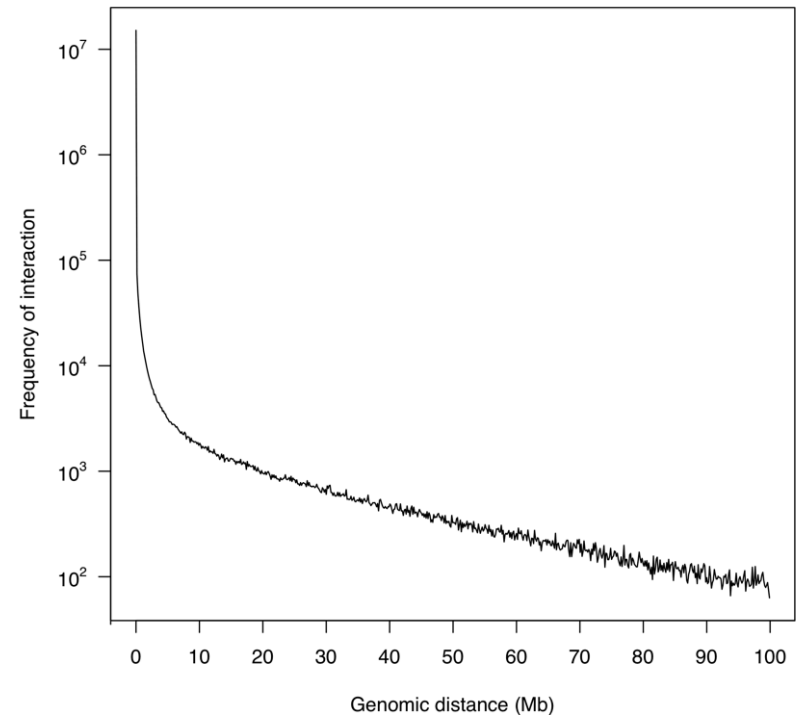# Interaction frequency decreases with genomic distance



ERG

GFP

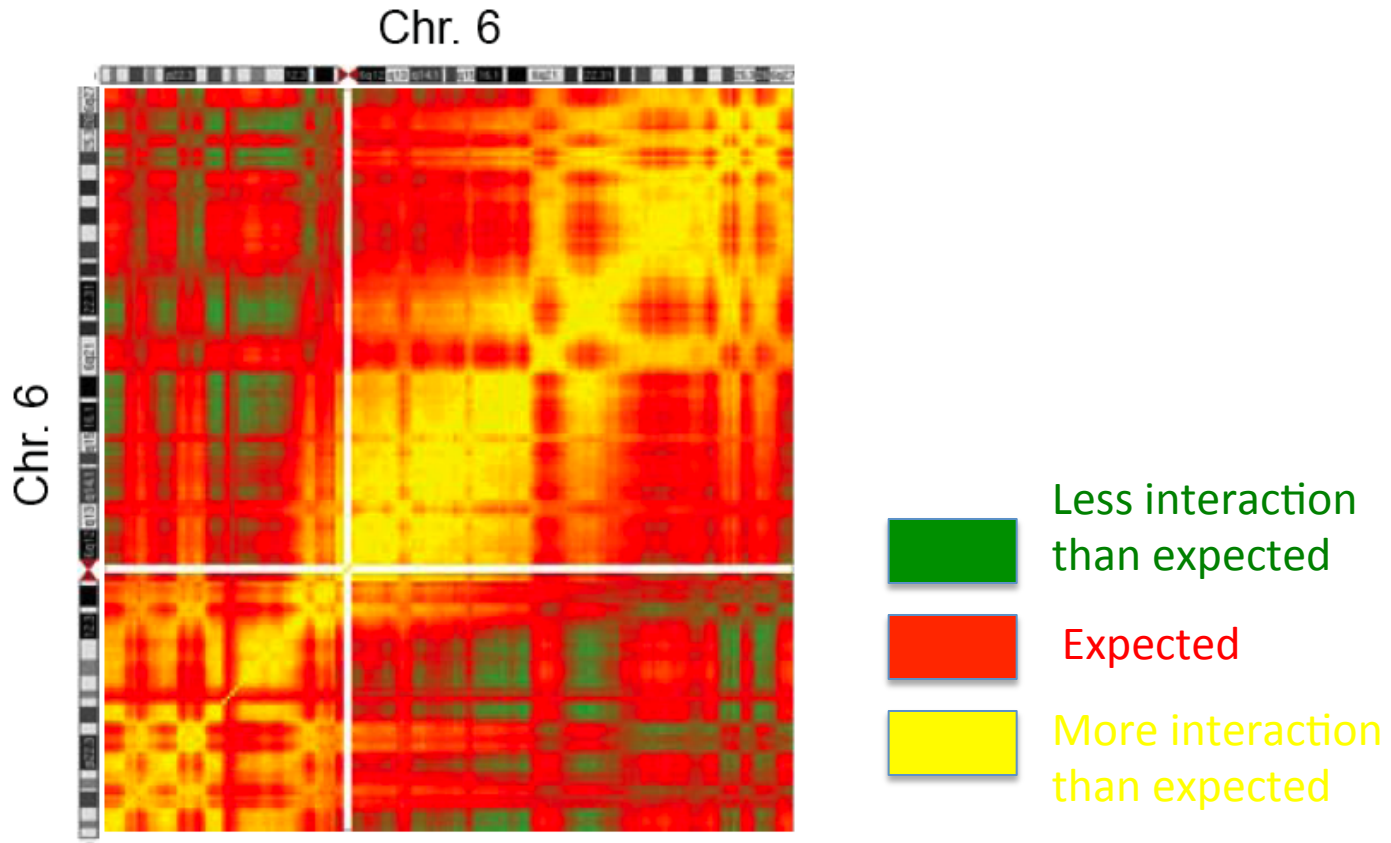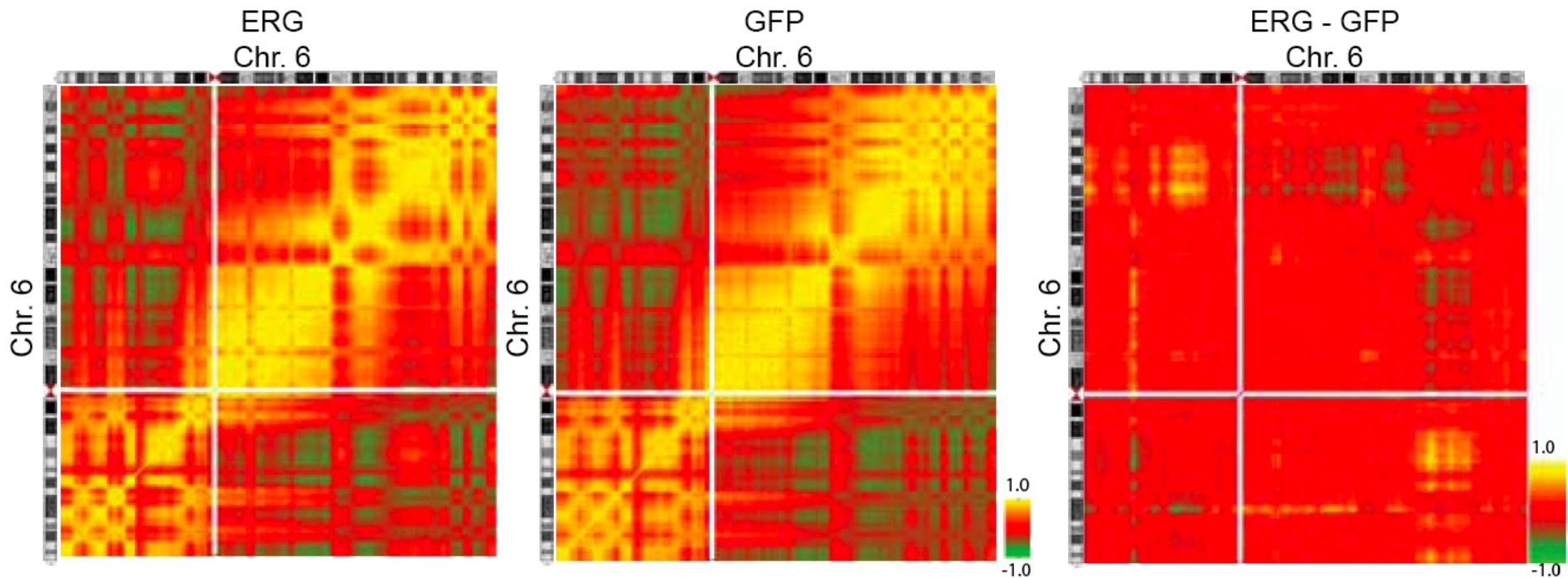$$P_{\text{interaction}}(s) \sim s^{-1}$$
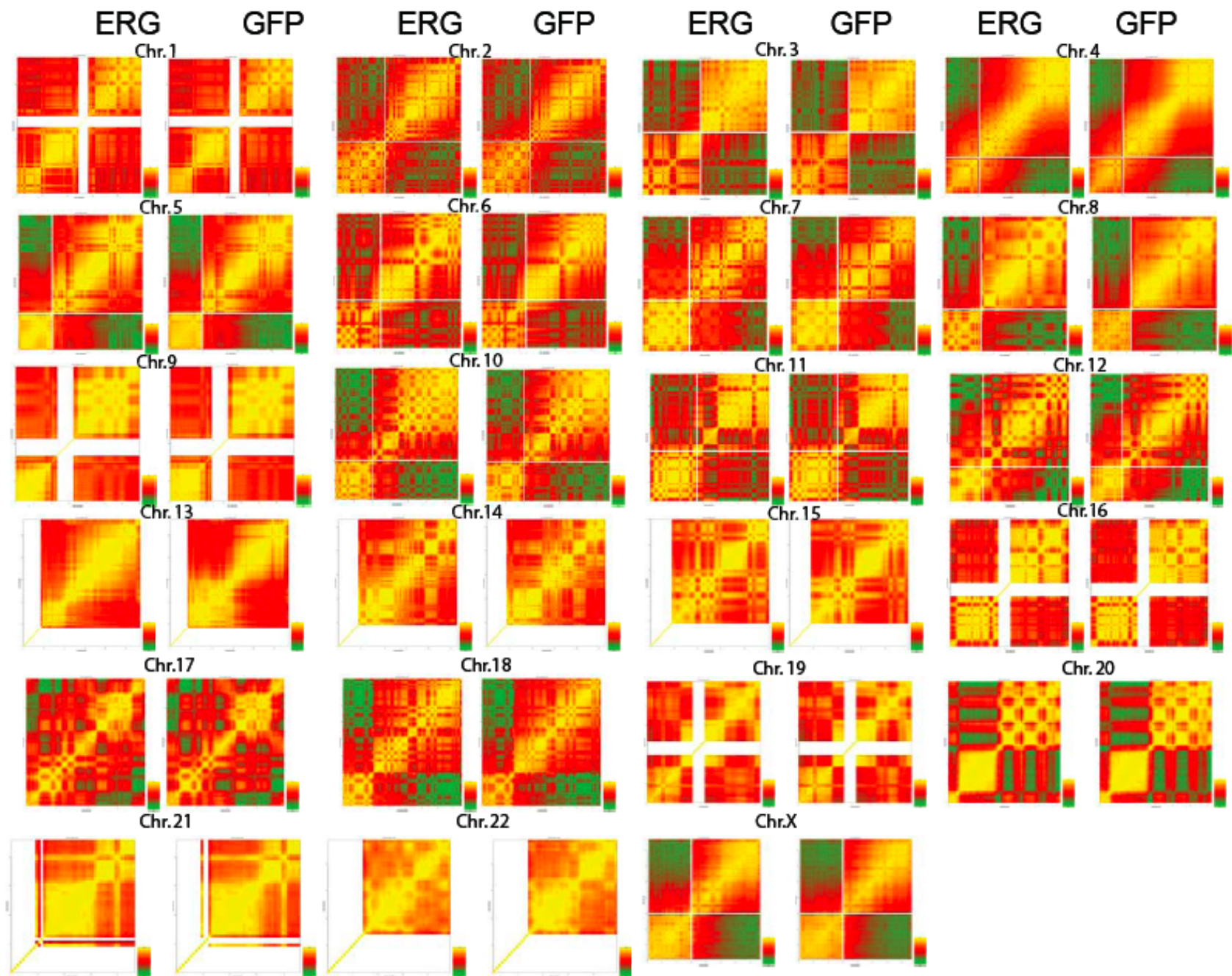
(500Kb-1Mb range)

$$P_{\text{interaction}}(s) \sim s^{-1}$$

(500Kb-1Mb range)

# Contact maps show domains of interactions



RWPE1-ERG

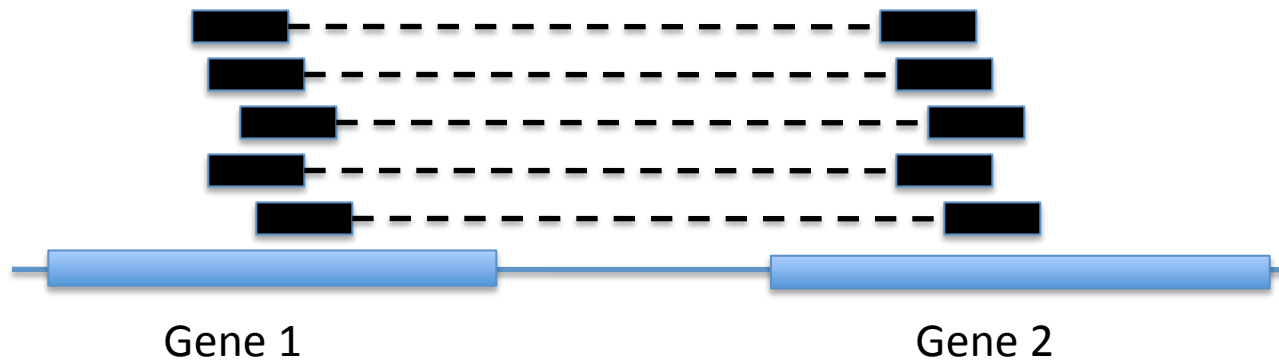# There are differences in interactions between ERG and GFP cells

ERG  GFP  ERG  GFP  ERG  GFP  ERG  GFP

Chr. 1  Chr. 2  Chr. 3  Chr. 4

Chr. 5  Chr. 6  Chr. 7  Chr. 8

Chr. 9  Chr. 10  Chr. 11  Chr. 12

Chr. 13  Chr. 14  Chr. 15  Chr. 16

Chr. 17  Chr. 18  Chr. 19  Chr. 20

Chr. 21  Chr. 22  Chr. X

# Fisher exact tests to quantify difference in interactions between GFP and ERG
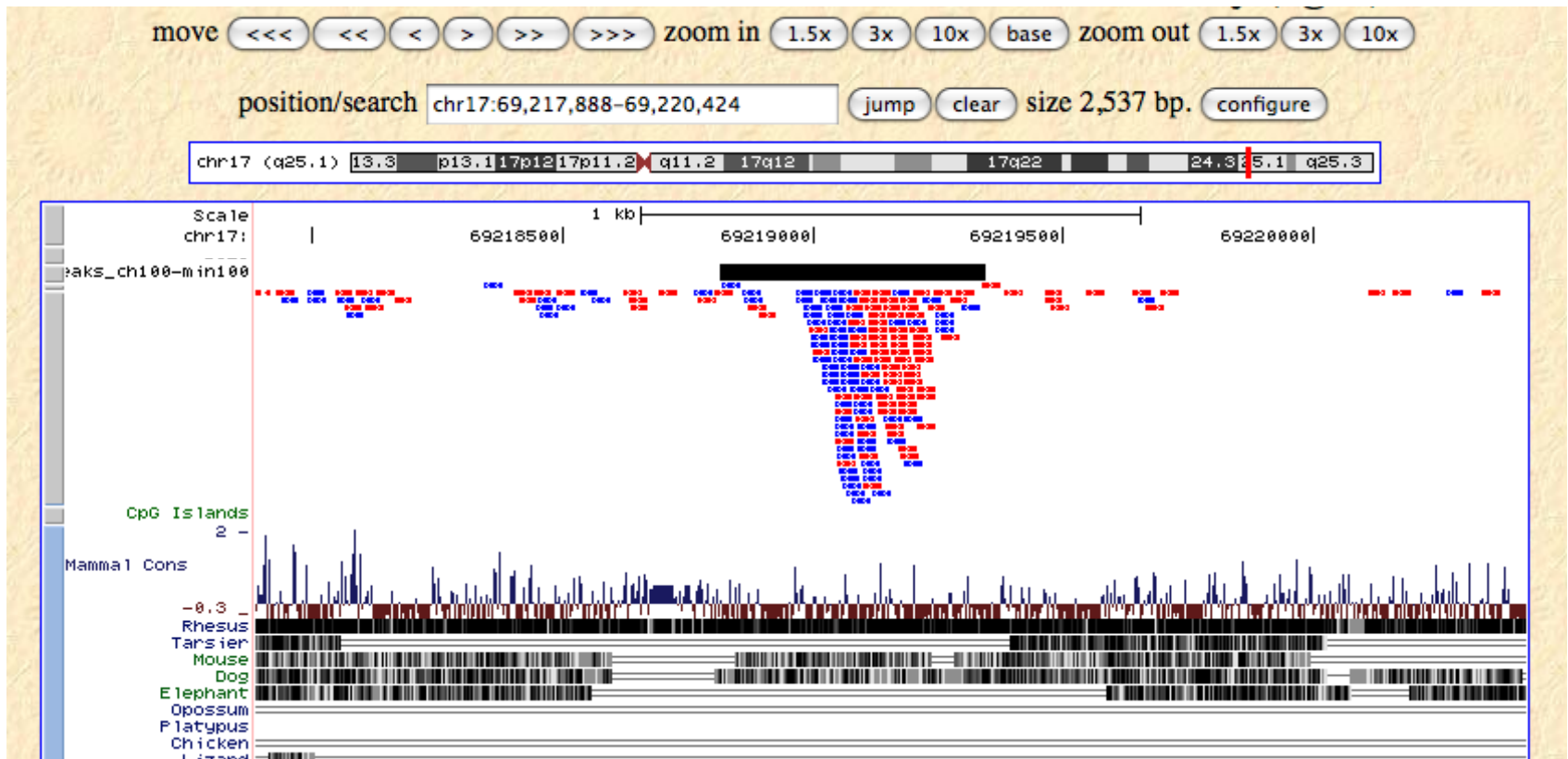


$$P(X = k) = \frac{\binom{m}{k}\binom{N-m}{n-k}}{\binom{N}{n}}.$$

# A graphical representation of all **gains** or **losses** of interaction in ERG cells vs GFP

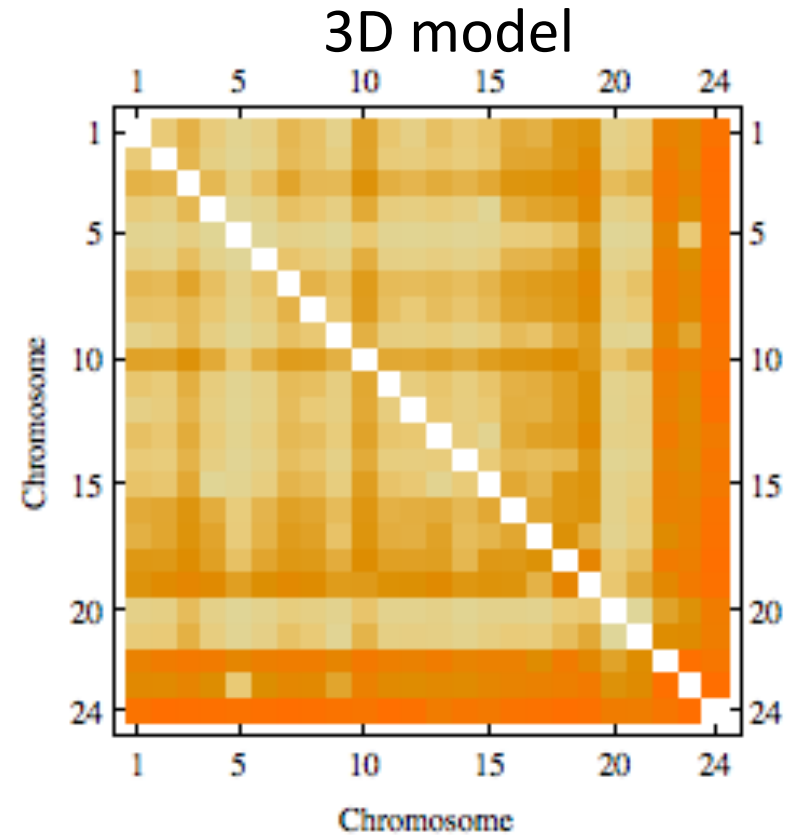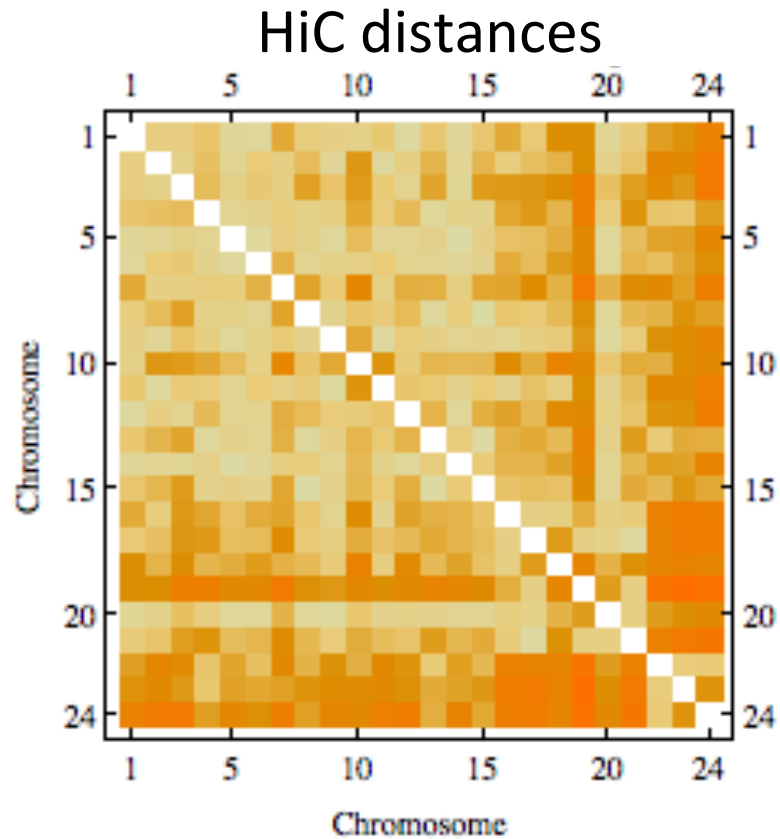# ChIP-seq: ERG binds to >6,000 location in RWPE1-ERG cells

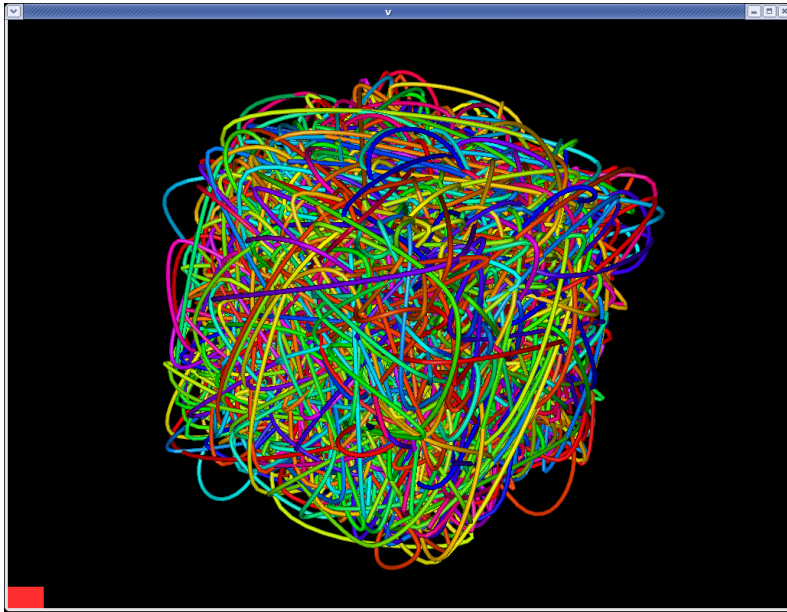# Loci/genes that **lose** or **gain** many interactions are more likely to be bound by ERG



Top N hubs of differential interactions

# 3D genome models

- Bin the genome into 1Mb blocks

- D = 1 / k  where k = number of reads "connecting" two blocks

- Start with random 3D topologies, use gradient descent algorithm to make 3D models that best recapitulate the HiC distances
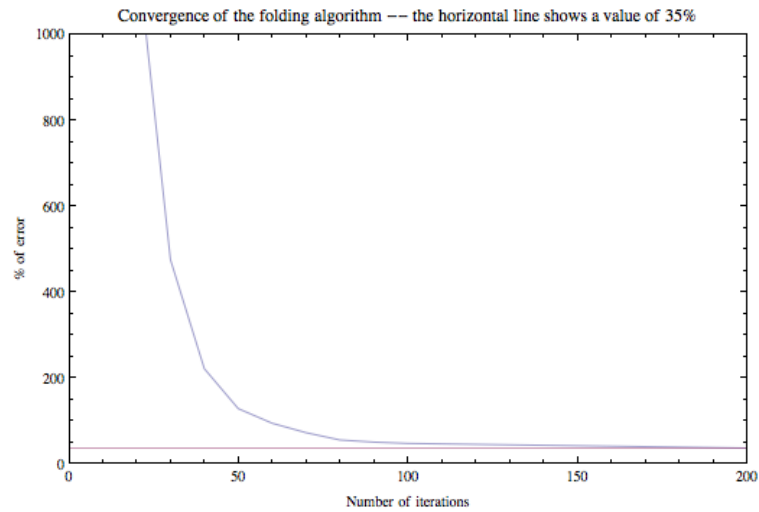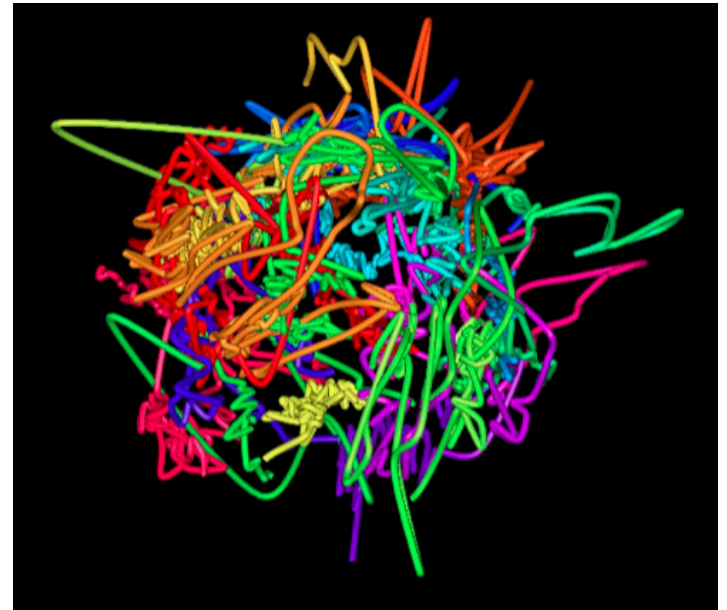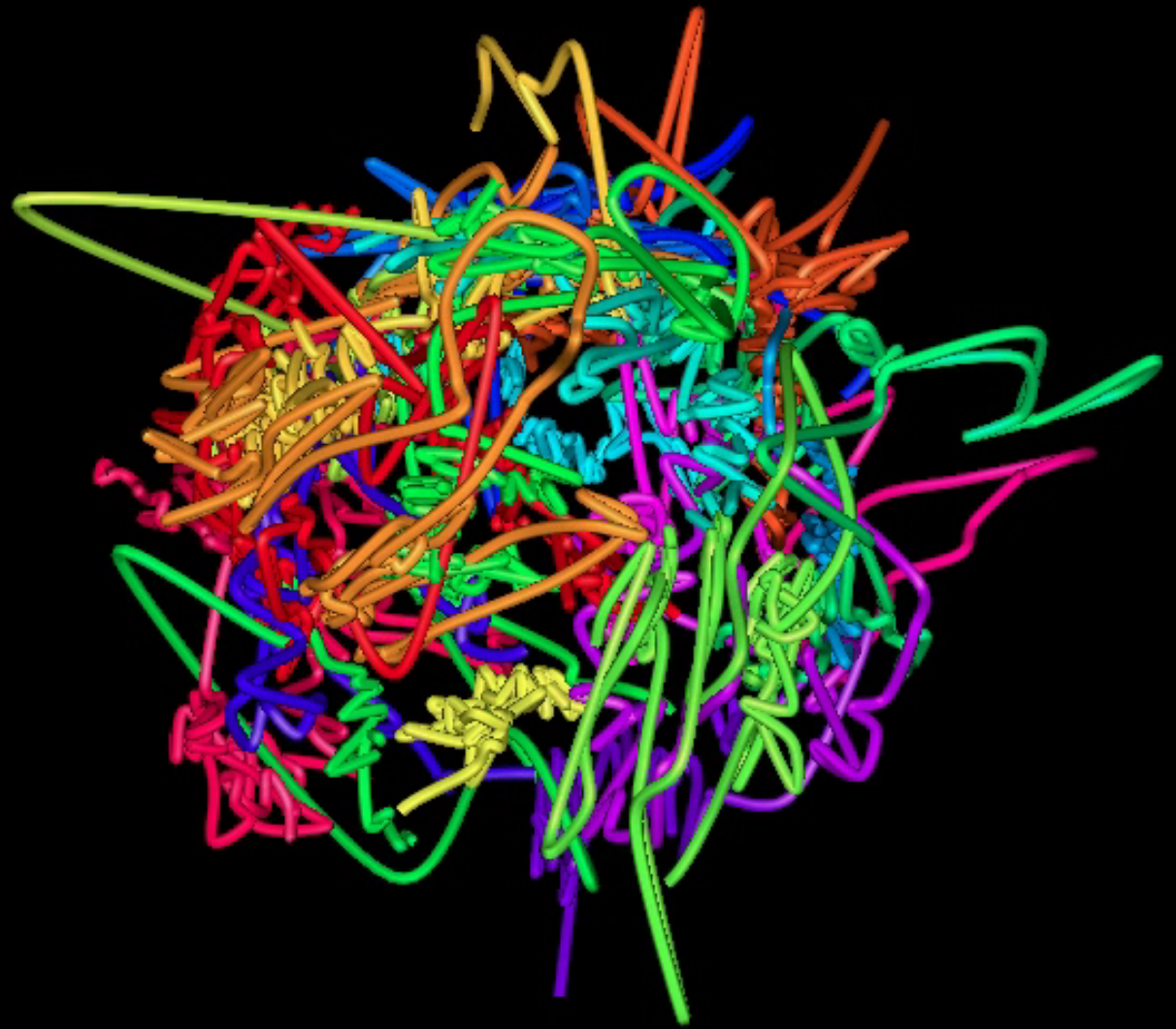
# 3D genome models



HiC distances

3D model

Initial random model

Final model



Convergence of the folding algorithm −− the horizontal line shows a value of 35%

# How does ERG over-expression change chromosome conformation ?