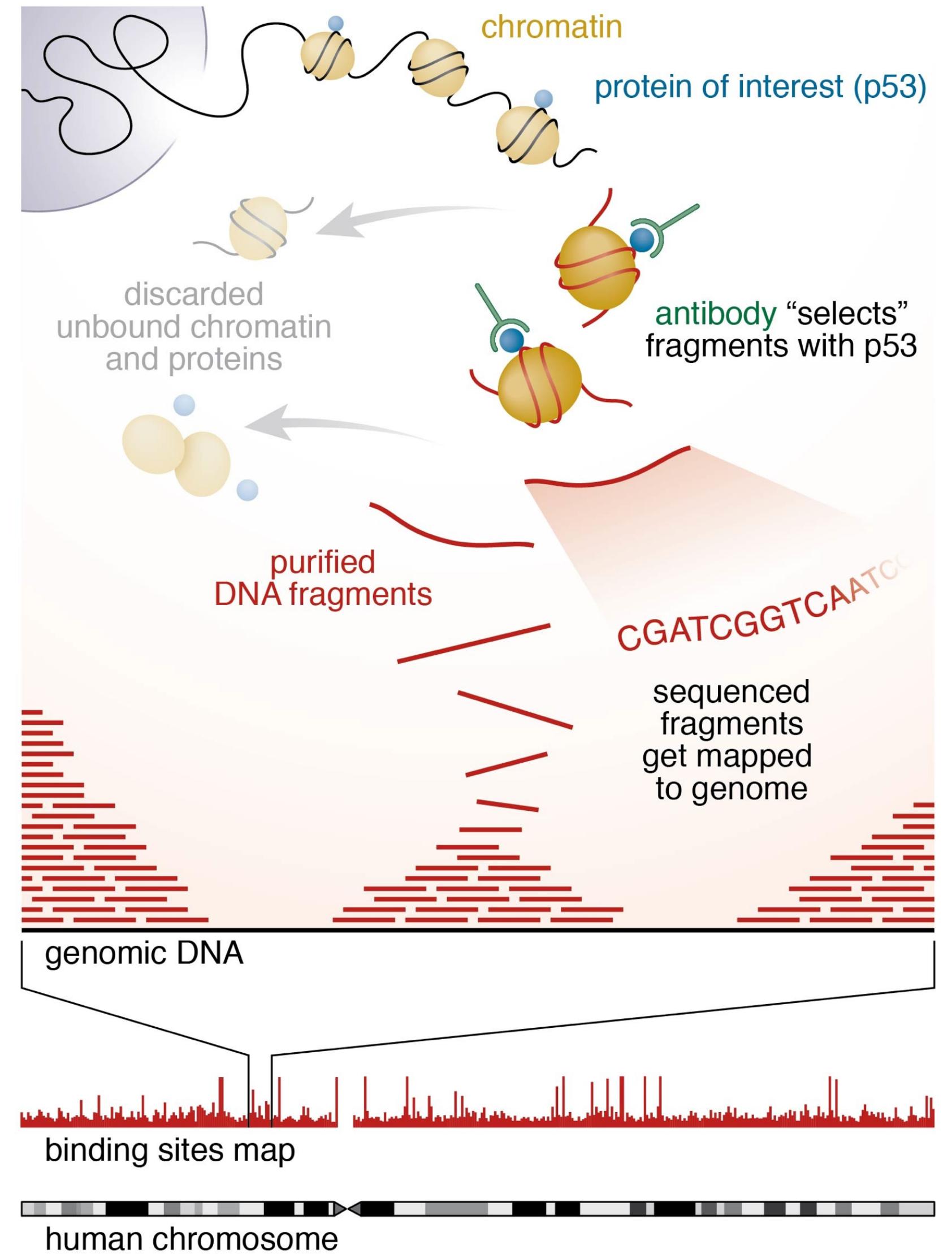
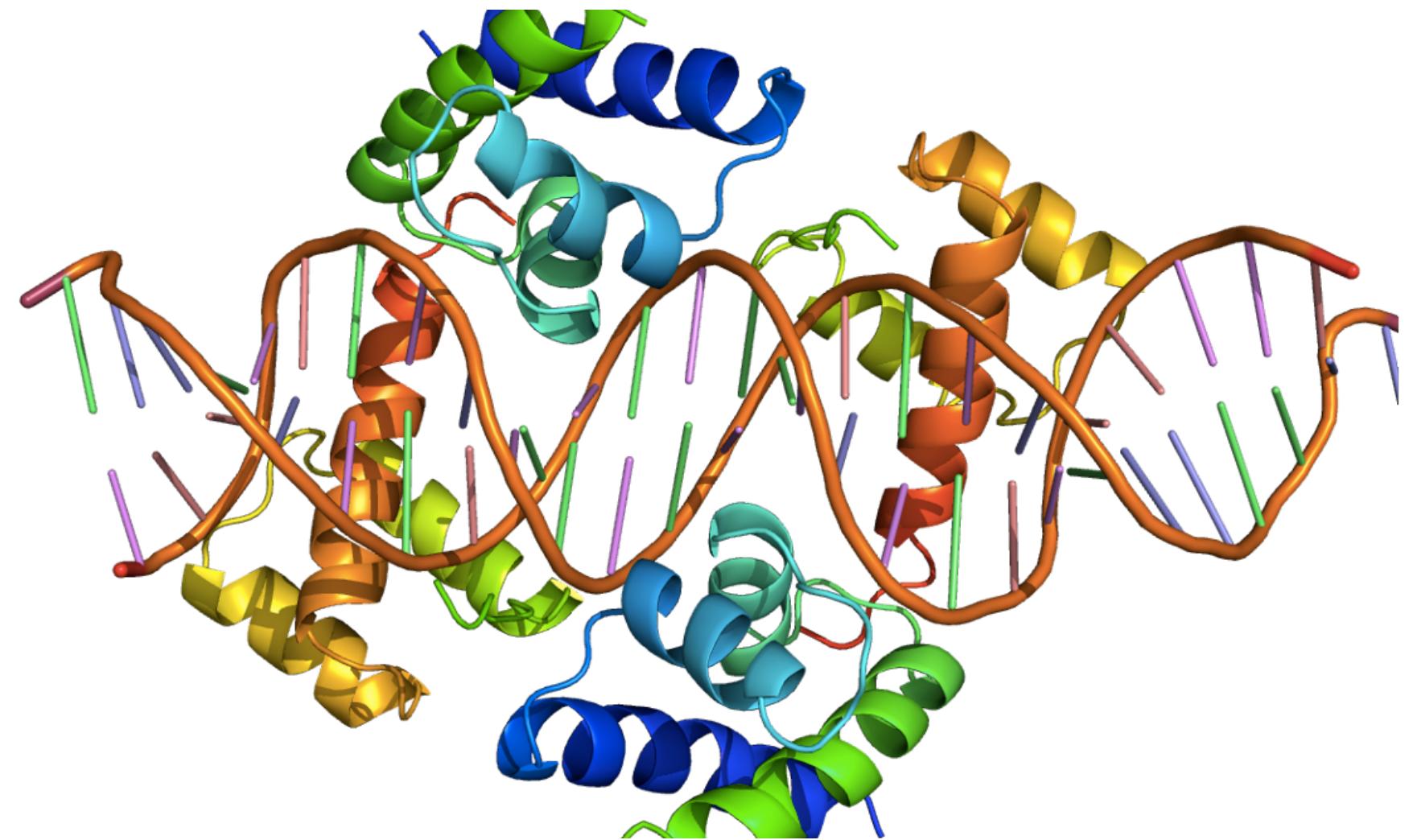


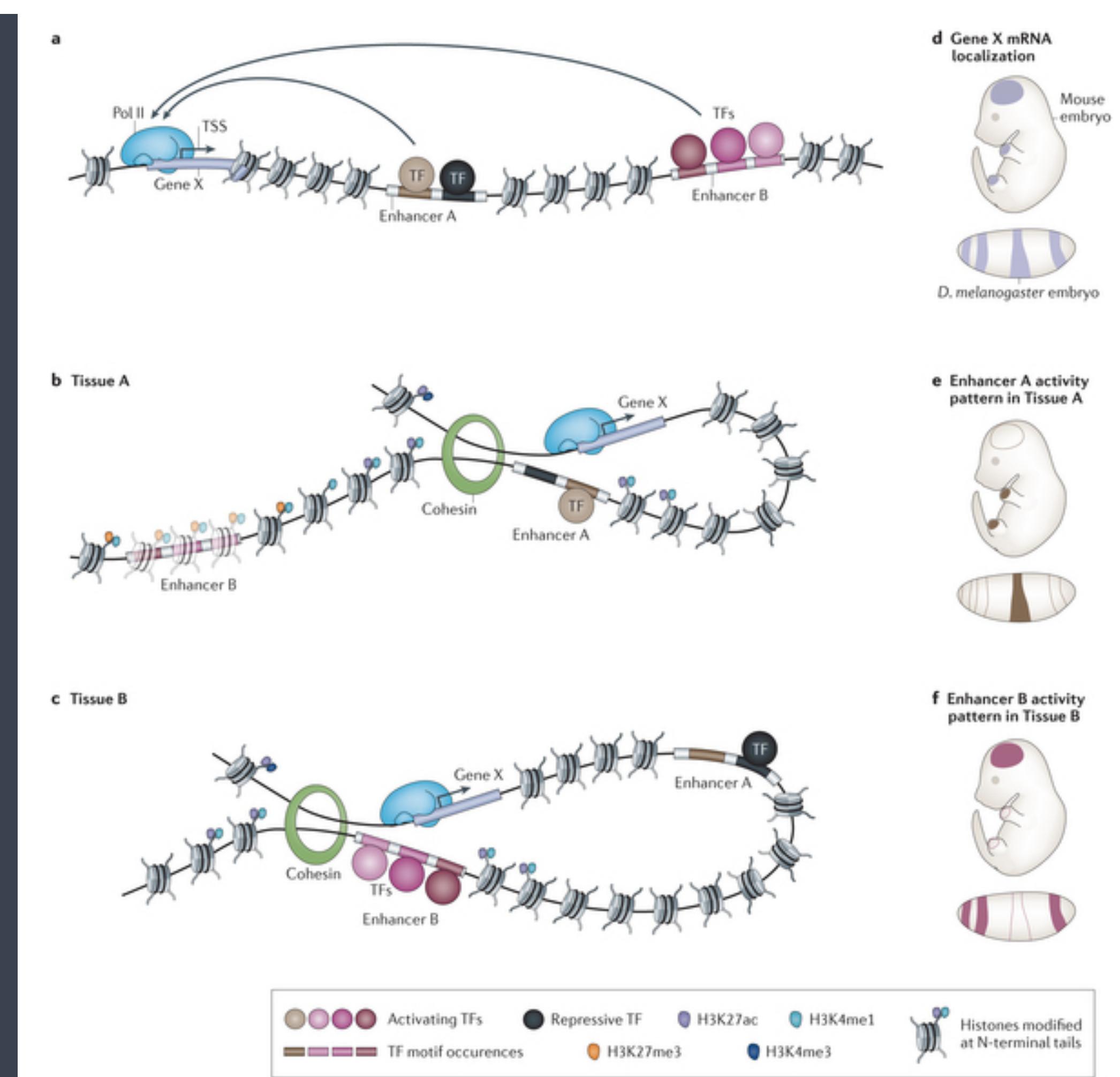
# Introduction to ChIP-seq

HSPH Bioinformatics Core

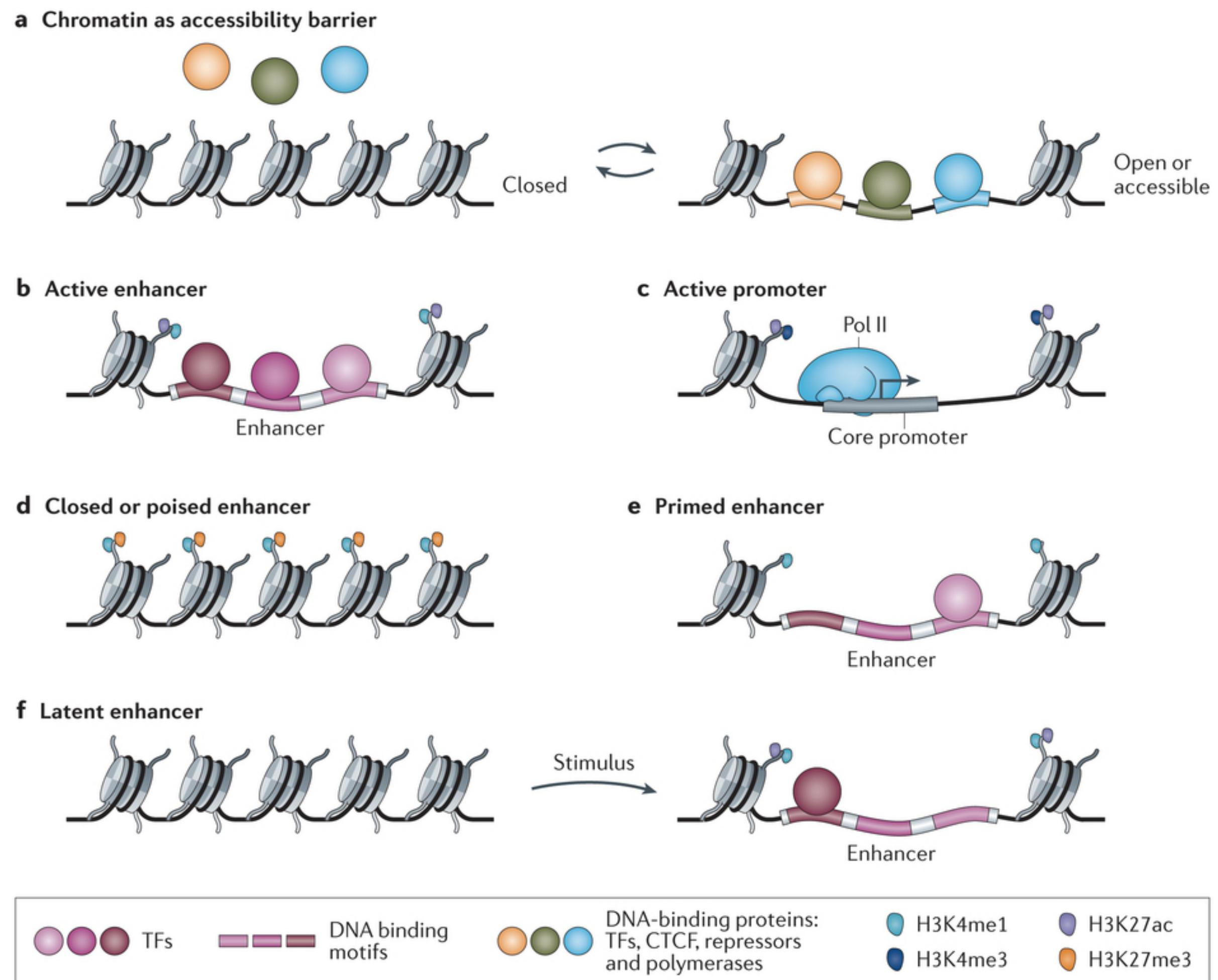


# Complexity in transcriptional regulation

Diverse mechanisms to ensure that genes are expressed at the right time, in appropriate tissues and under specific conditions



Shlyueva, et al (2014). Transcriptional enhancers: from properties to genome-wide predictions.

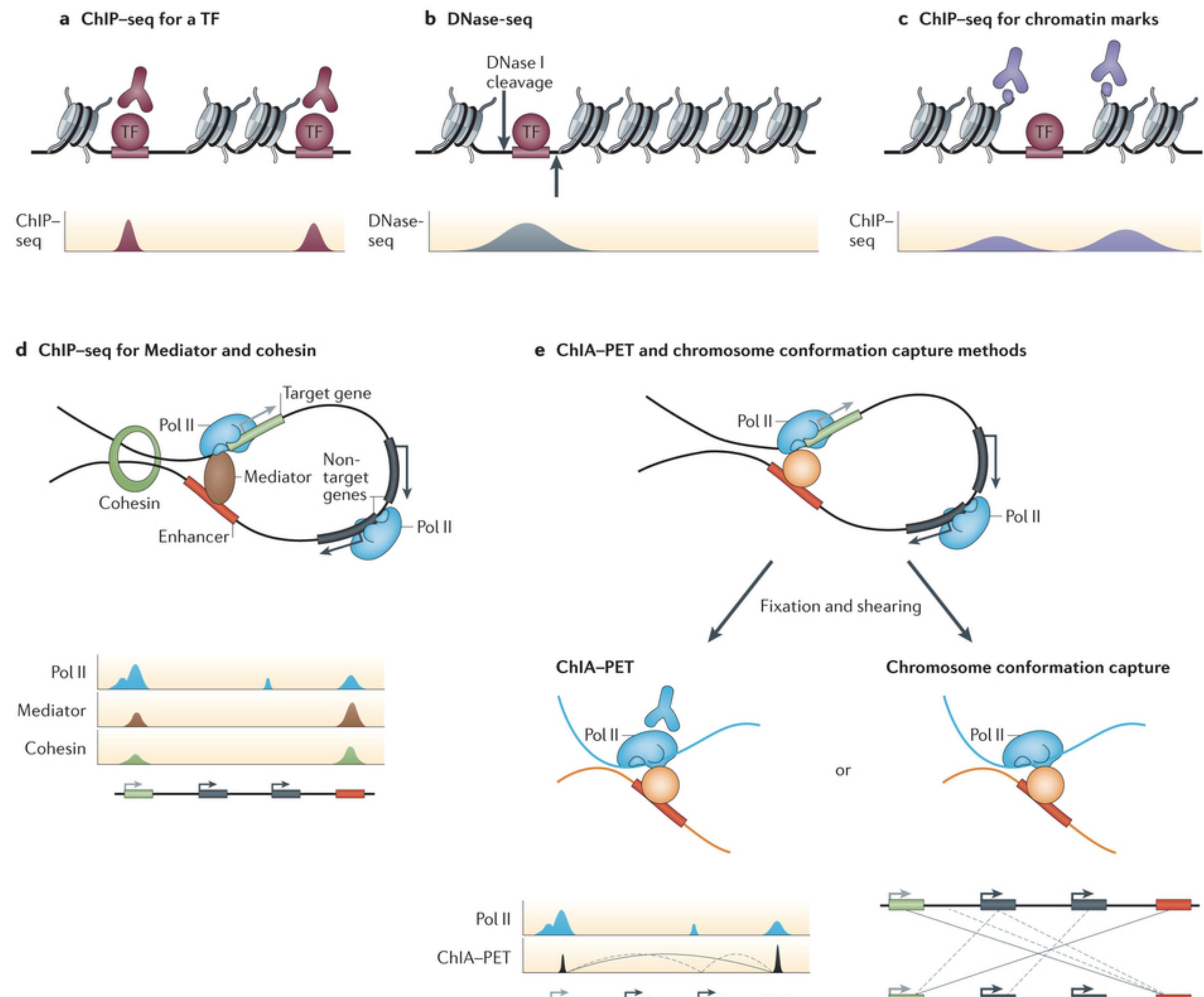


Shlyueva, et al (2014)

Nature Reviews | Genetics

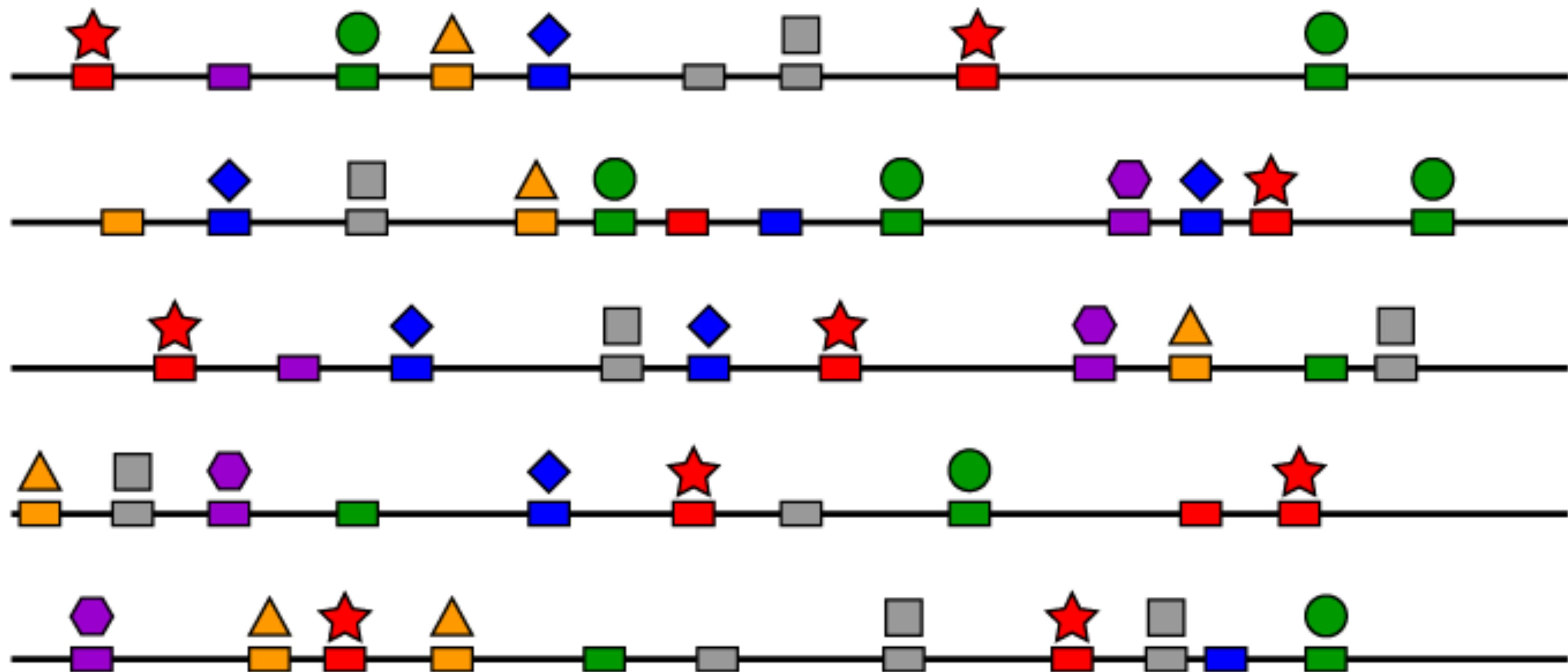
Chromatin structure determines if a gene is expressed or not

# Genomic methods for detecting regulatory elements

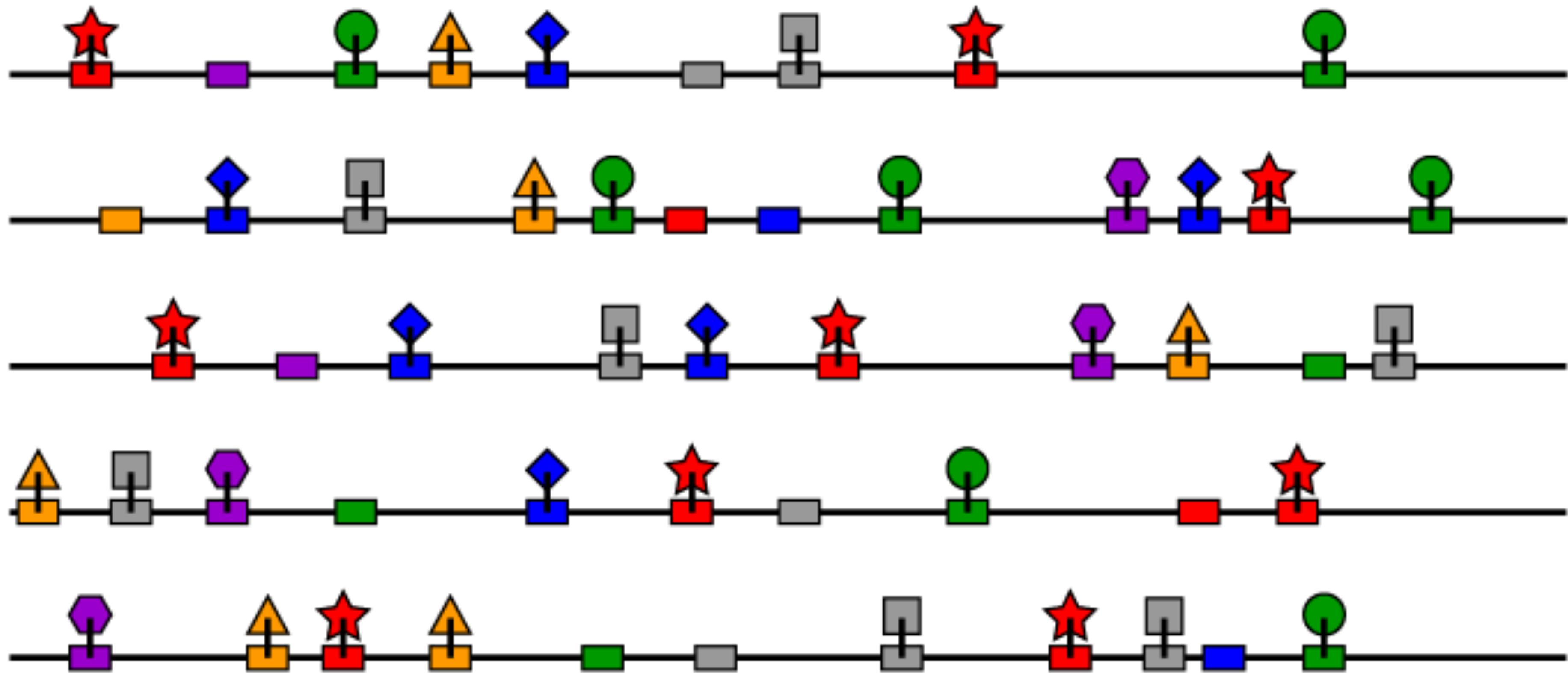


Shlyueva, et al (2014)

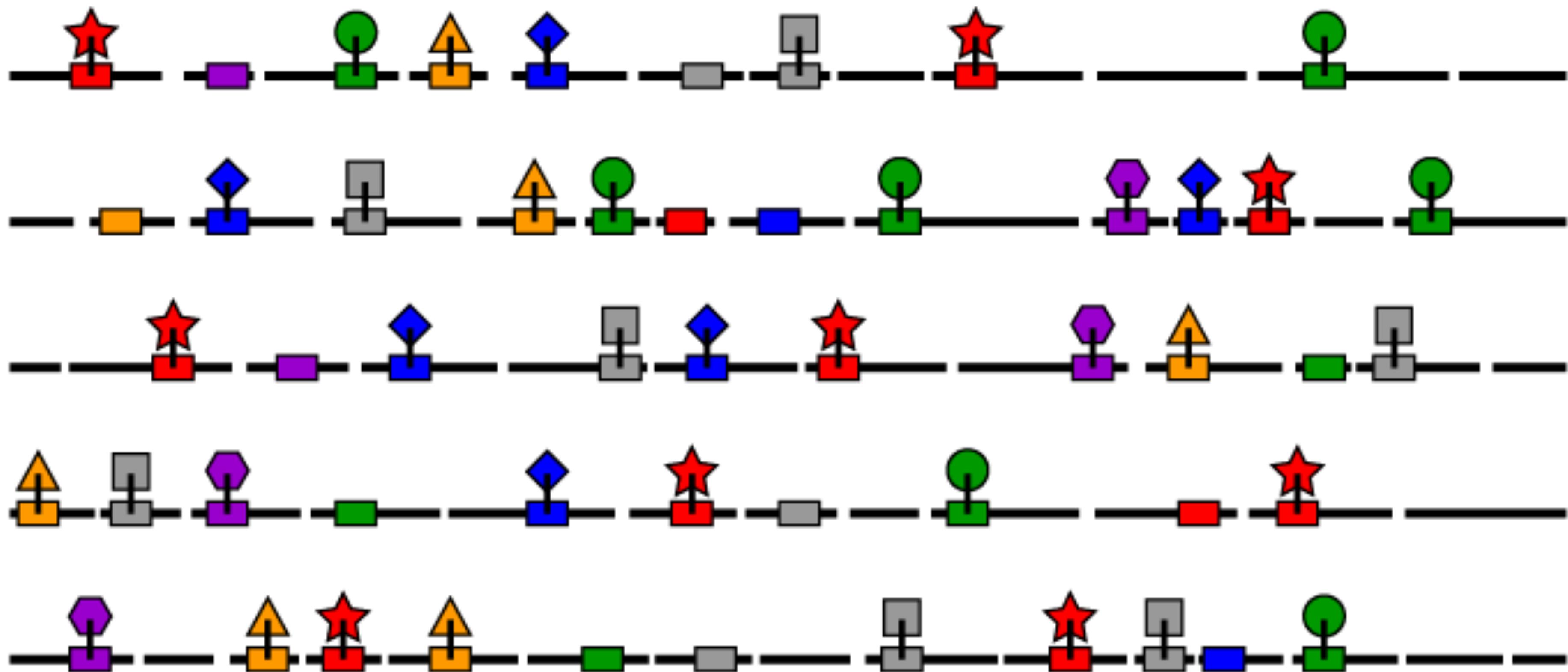
Nature Reviews | Genetics



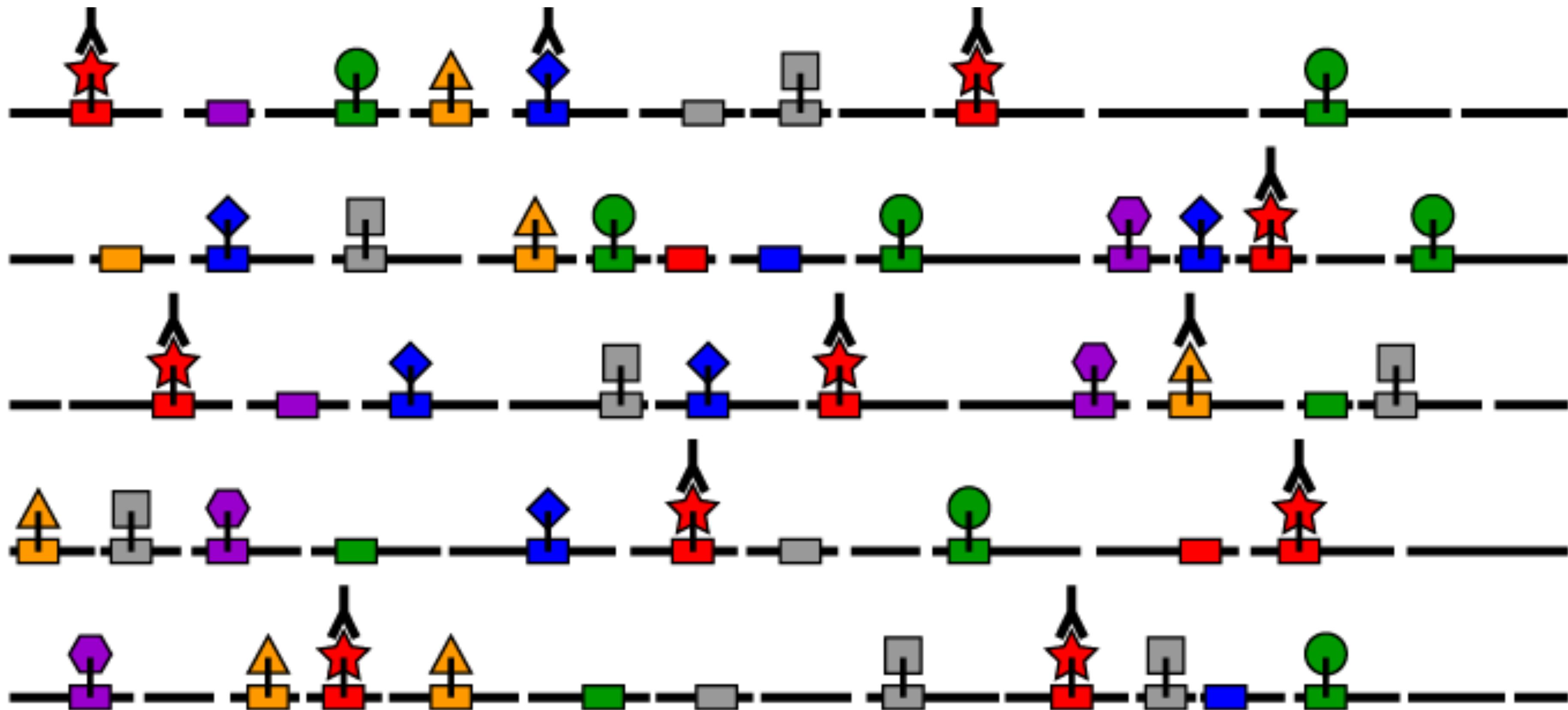
ChIP-Seq: Library preparation



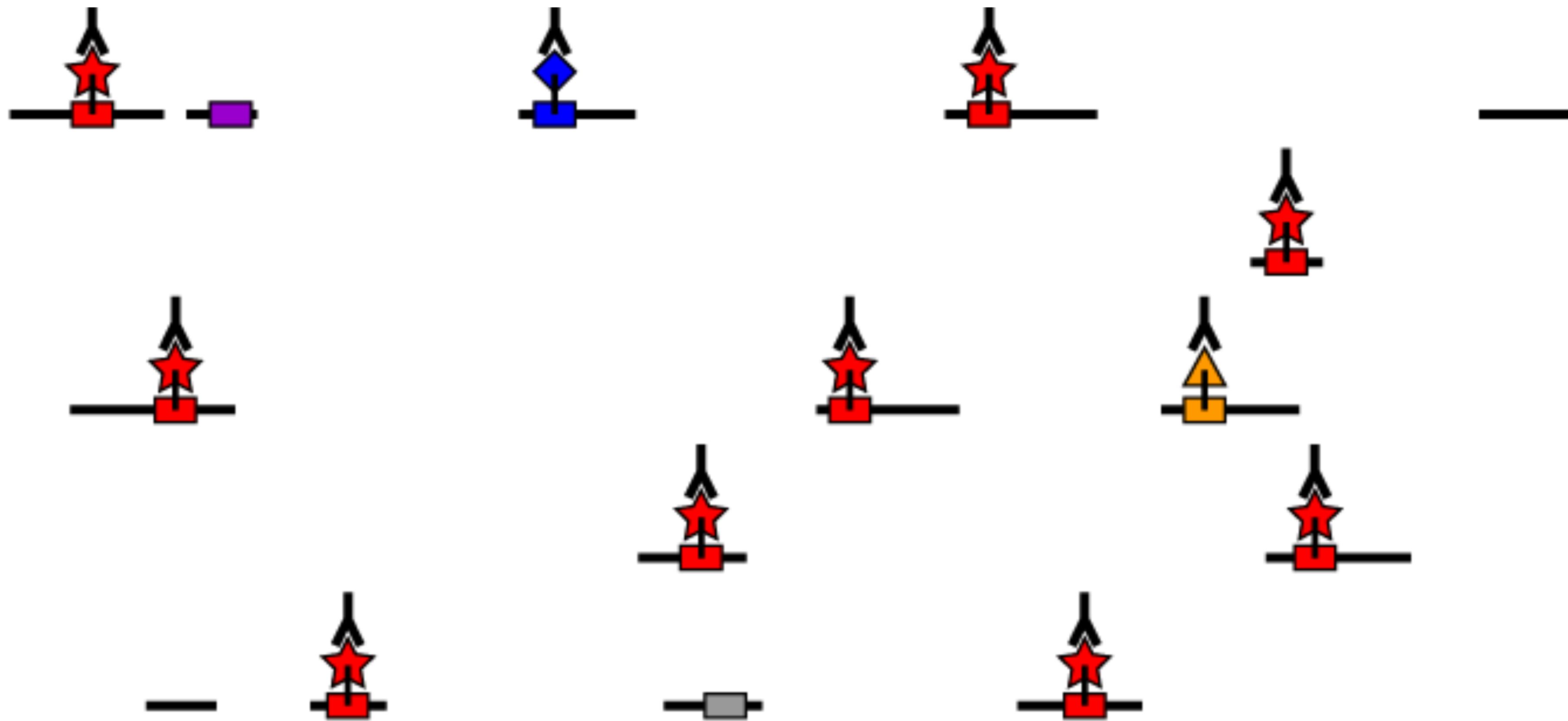
Crosslink protein to the DNA



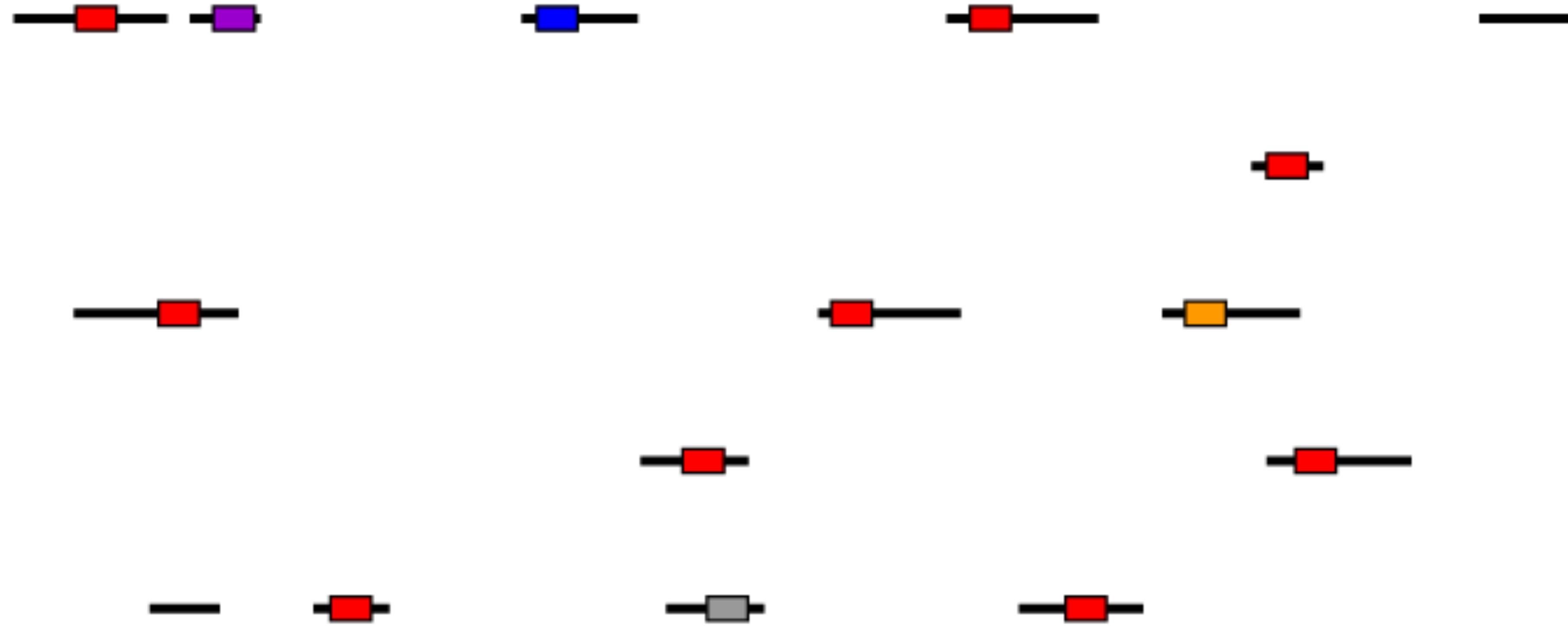
Fragment



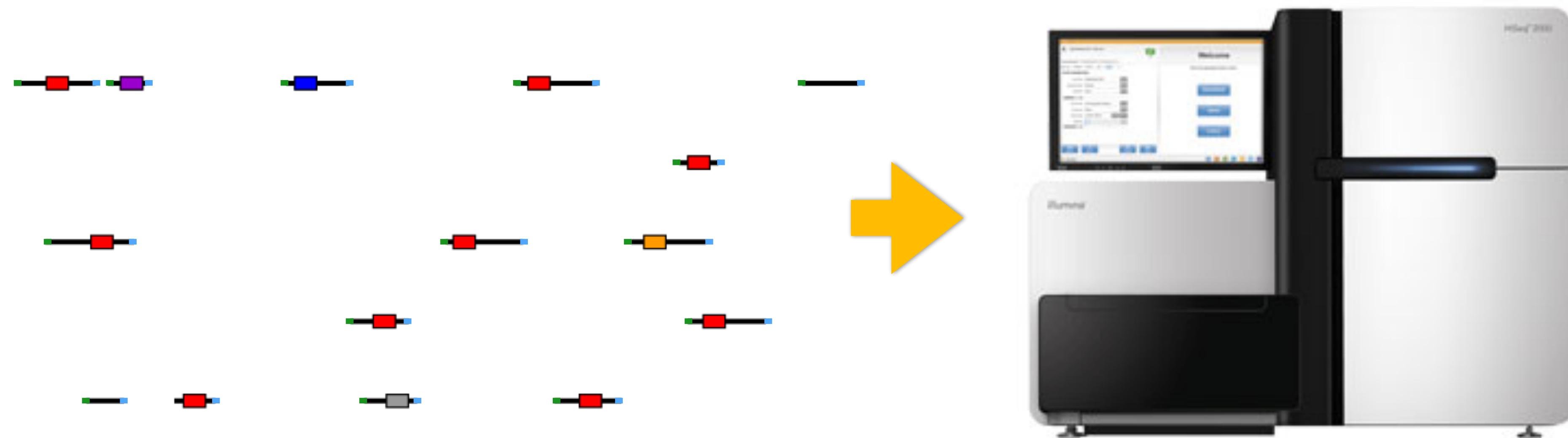
TF-specific antibody



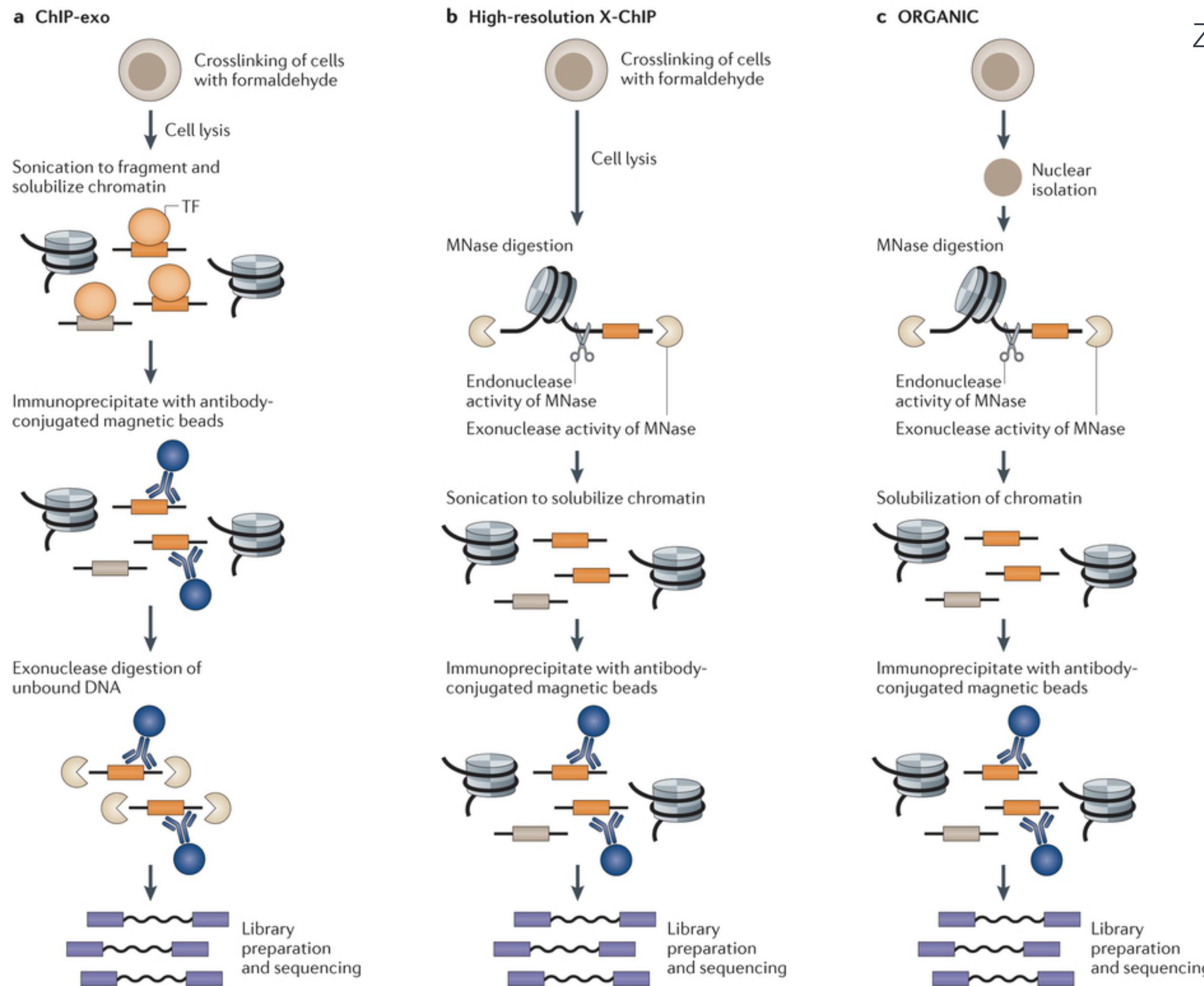
Immunoprecipitate



Reverse crosslink and purify DNA



Ligate adaptors and sequence



# High resolution variations of ChIP-seq

# ChIP Output

- ▶ DNA sample *enriched* for fragments associated with event
- ▶ Only a fraction correspond to actual signal
- ▶ Depends on number of active signal sites, number of starting genomes, and efficiency f IP

# ChIP Output

- ▶ DNA sample *enriched* for fragments associated with event
- ▶ Only a fraction correspond to actual signal
- ▶ Depends on number of active signal sites, number of starting genomes, and efficiency f IP

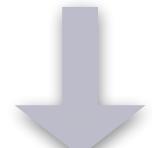
Proper controls are essential...

# Why are controls necessary?

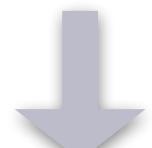
- Open chromatin regions are fragmented more easily than closed regions
- Repetitive sequences might seem to be enriched (inaccurately assessed repeat copy numbers in the assembled genome)
- Uneven distribution of sequence tags across the genome
- A ChIP-Seq peak should be compared with the same region of the genome in a matched control

Biological samples/Library preparation

Crosslink proteins to DNA



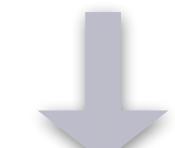
Shear DNA



Immunoprecipitation



Reverse crosslink



Size selection and PCR

Specific antibody (ChIP enrichment)

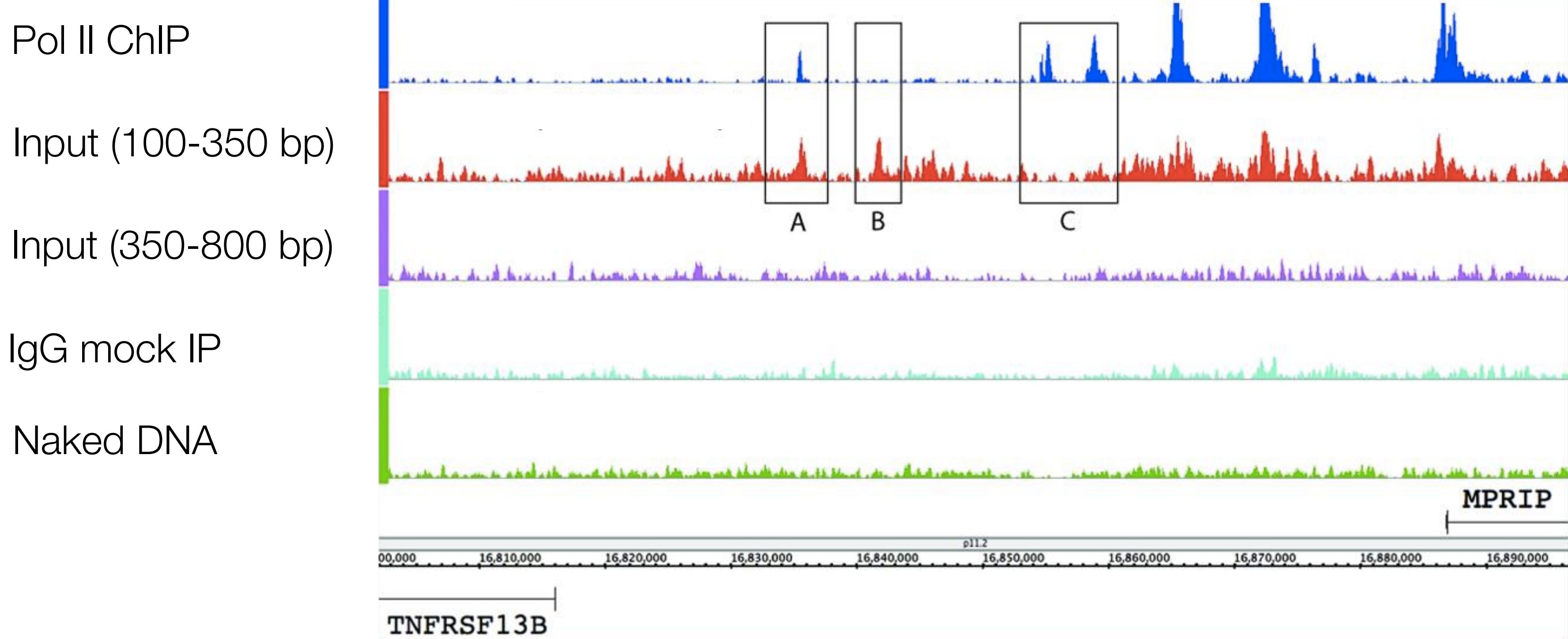


No IP (Input DNA)

No antibody ("mock IP")

Non-specific antibody (IgG "mock IP")

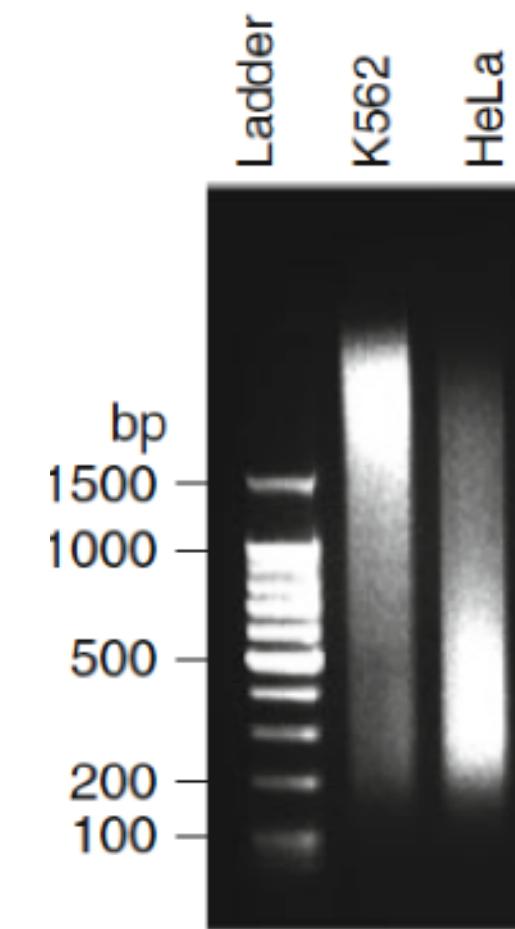
# ChIP-Seq Controls



Map of ChIP-seq versus control signals

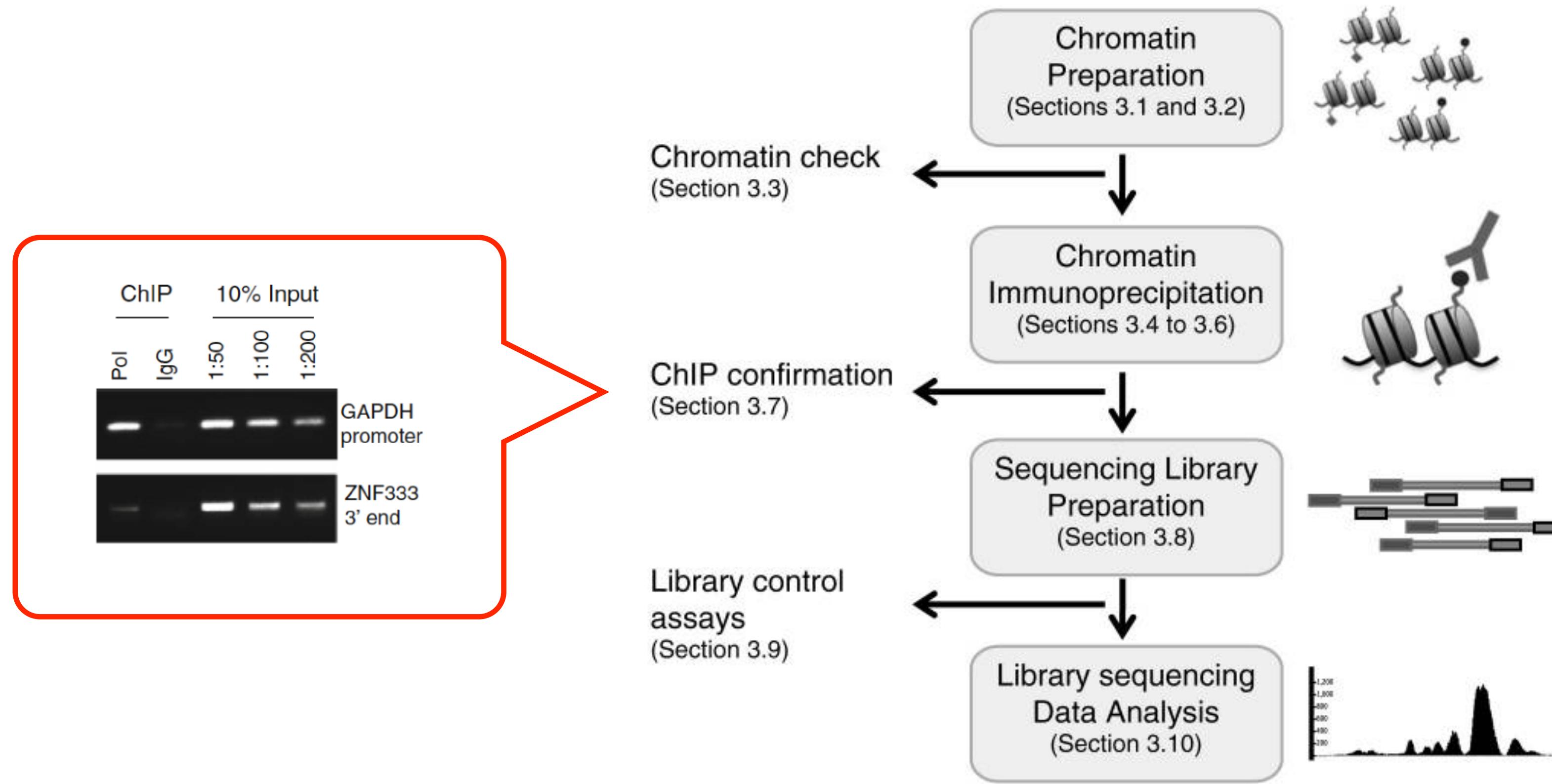
# Parameters for a successful ChIP

- ▶ Antibody!
- ▶ Amount of starting material
- ▶ Chromatin fragmentation
  - ▶ Size matters (not too big and not too small)
  - ▶ Can vary between cell types
  - ▶ Stringency of washes



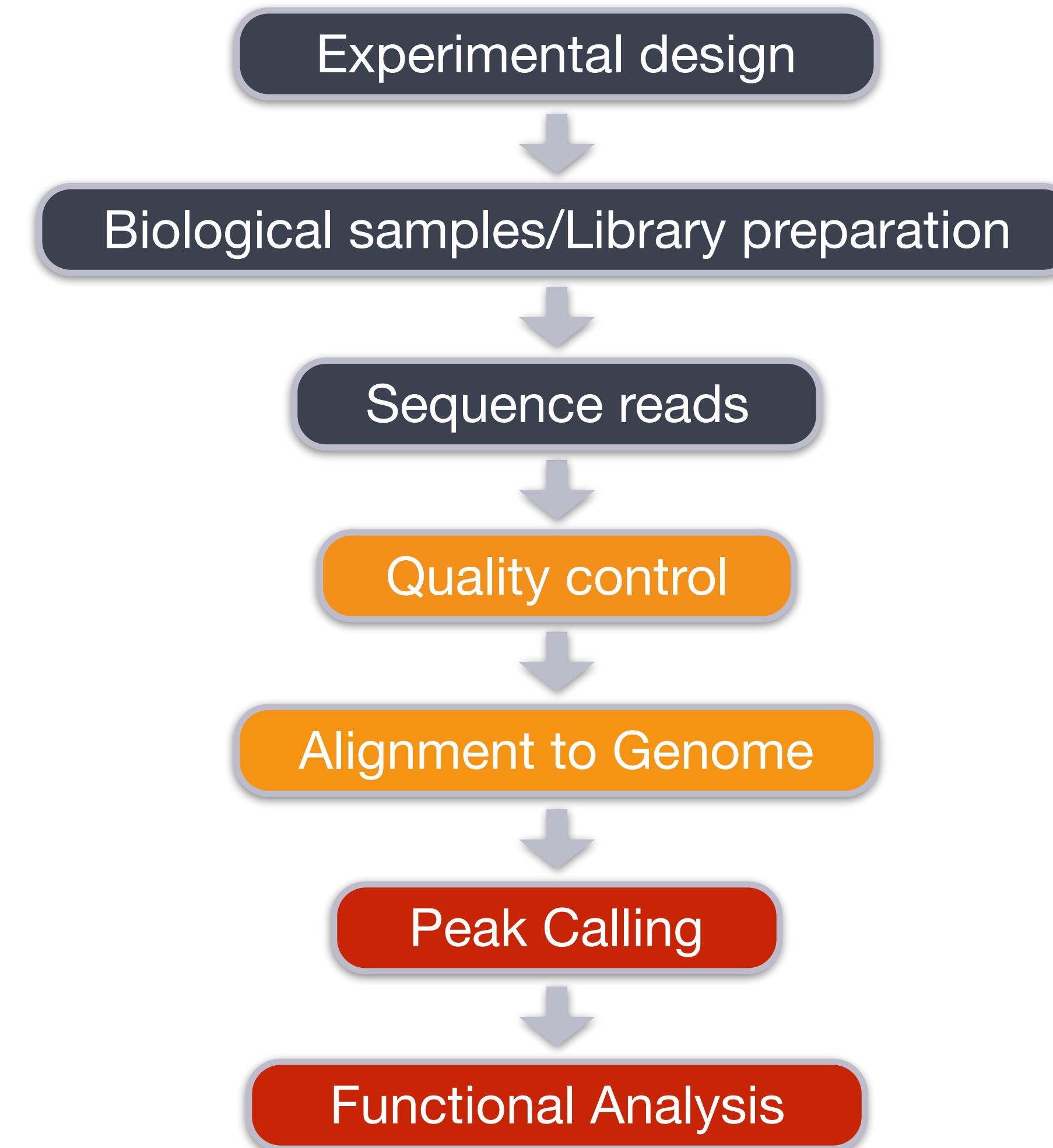
**Fragments too big:**  
Reduced signal to noise ratio  
in ChIP-seq

**Oversonication:**  
Fragmentation biased towards  
promoter regions causes  
ChIP-seq enrichments at  
promoters in both, ChIP AND  
control (input) sample



## Quality check points

O'Geen et al (2011), Methods Mol Biol: <http://www.ncbi.nlm.nih.gov/pmc/articles/PMC4151291/>



ChIP-seq workflow

# Sequencing considerations

- ▶ Read length (25- to 150-bp)
  - > Longer reads and paired-end reads improve mappability (allele-specific chromatin events, investigations of transposable elements)
  - > Balance cost with value of more informative reads
- ▶ Avoiding batches or distributing samples evenly over batches
- ▶ # sequences generated (5-10M minimum; 20-40M as standard)

# Sequencing considerations

- ▶ Biological replicates are essential to understand variation and for differential binding analysis
- ▶ Sequencing multiple biological replicates or sample replicates to a lower sequencing depth is preferable to sequencing a single sample to a greater depth
- ▶ Better to sequence high-quality sample at low depth than low-quality sample to high depth
- ▶ Input controls should be sequenced to equal or greater depths than IP samples

# Sequencing QC and alignment

- ▶ FastQC to assess sequencing quality
- ▶ Alignment using BWA or Bowtie
- ▶ Keep only uniquely mapped reads (i.e. reads aligned to a single genomic location)
- ▶ Evaluate library complexity (uniquely mapped locations / uniquely mapped reads)
  - ▶ Low complexity is characterized by a significant proportion of reads sharing identical start sites. This results in a lot of redundant sequence reads, which just end up in the trash.
  - ▶ A larger number of reads from a sufficiently complex library increases the chances of finding all true binding sites

# Wilbanks & Facciotti, PLoS One (2010)

Program	Reference Version	Graphical user interface?	Window-based scan	Tag clustering	Gaussian kernel density estimator	Strand-specific scoring	Peak height or fold enrichment (FE)	Background subtraction	Compensates for genomic duplications or deletions	False Discovery Rate	Compare to normalized control data (FE)	Compare to statistical model fitted with control data	Statistical model or test
CisGenome	28 1.1	X*	X			X	X		X		X		conditional binomial model
Minimal ChipSeq Peak Finder	16 2.0.1			X		X				X			
E-RANGE	27 3.1			X		X				X	X		chromosome scale Poisson dist.
MACS	13 1.3.5		X			X			X		X		local Poisson dist.
QuEST	14 2.3				X	X			X**		X		chromosome scale Poisson dist.
HPeak	29 1.1		X			X					X		Hidden Markov Model
Sole-Search	23 1	X	X			X		X			X		One sample t-test
PeakSeq	21 1.01			X		X					X		conditional binomial model
SISSRS	32 1.4		X			X				X			
spp package (wtd & mtc)	31 1.7		X			X		X	X'	X			
Generating density profiles				Peak assignment	Adjustments w. control data	Significance relative to control data							

X\* = Windows-only GUI or cross-platform command line interface

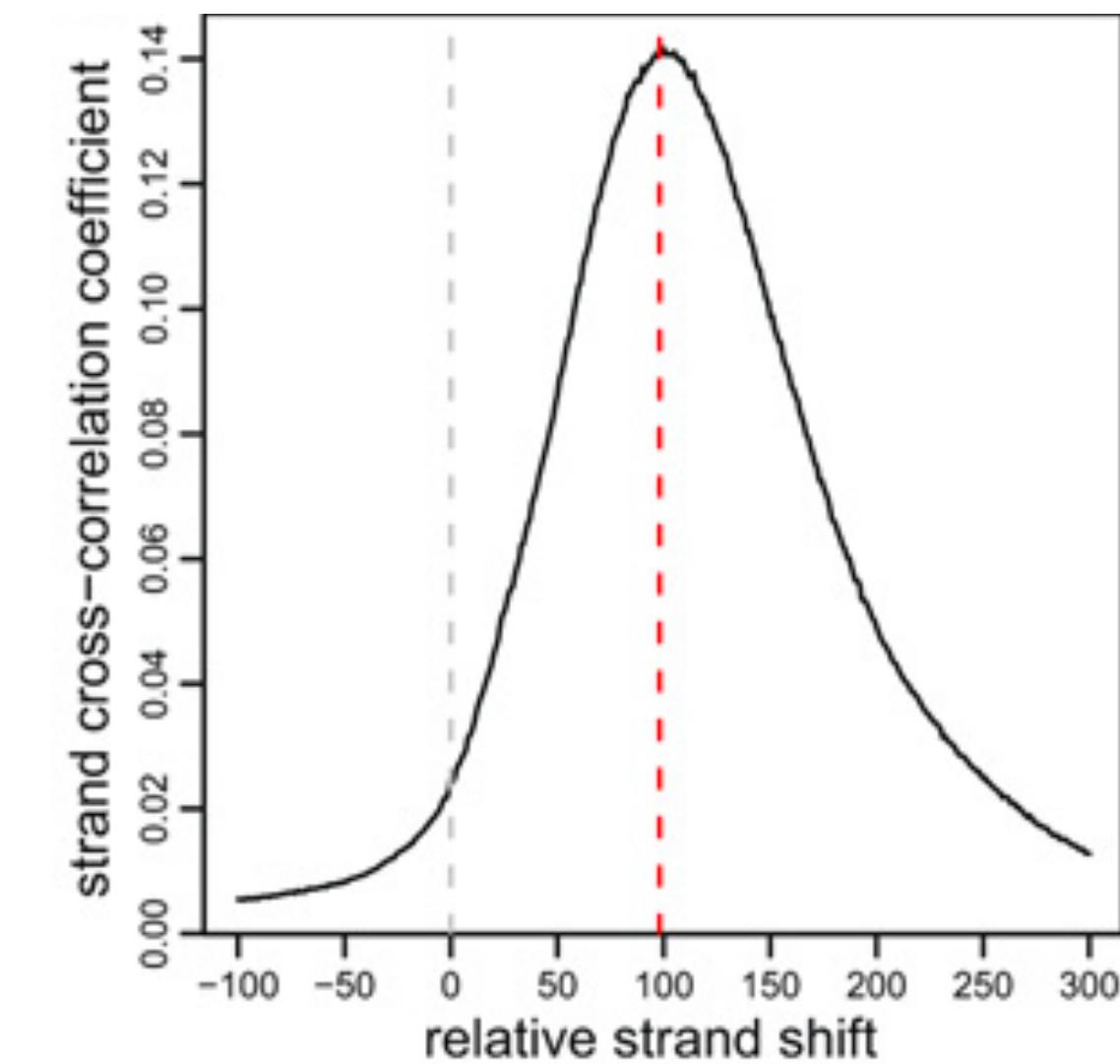
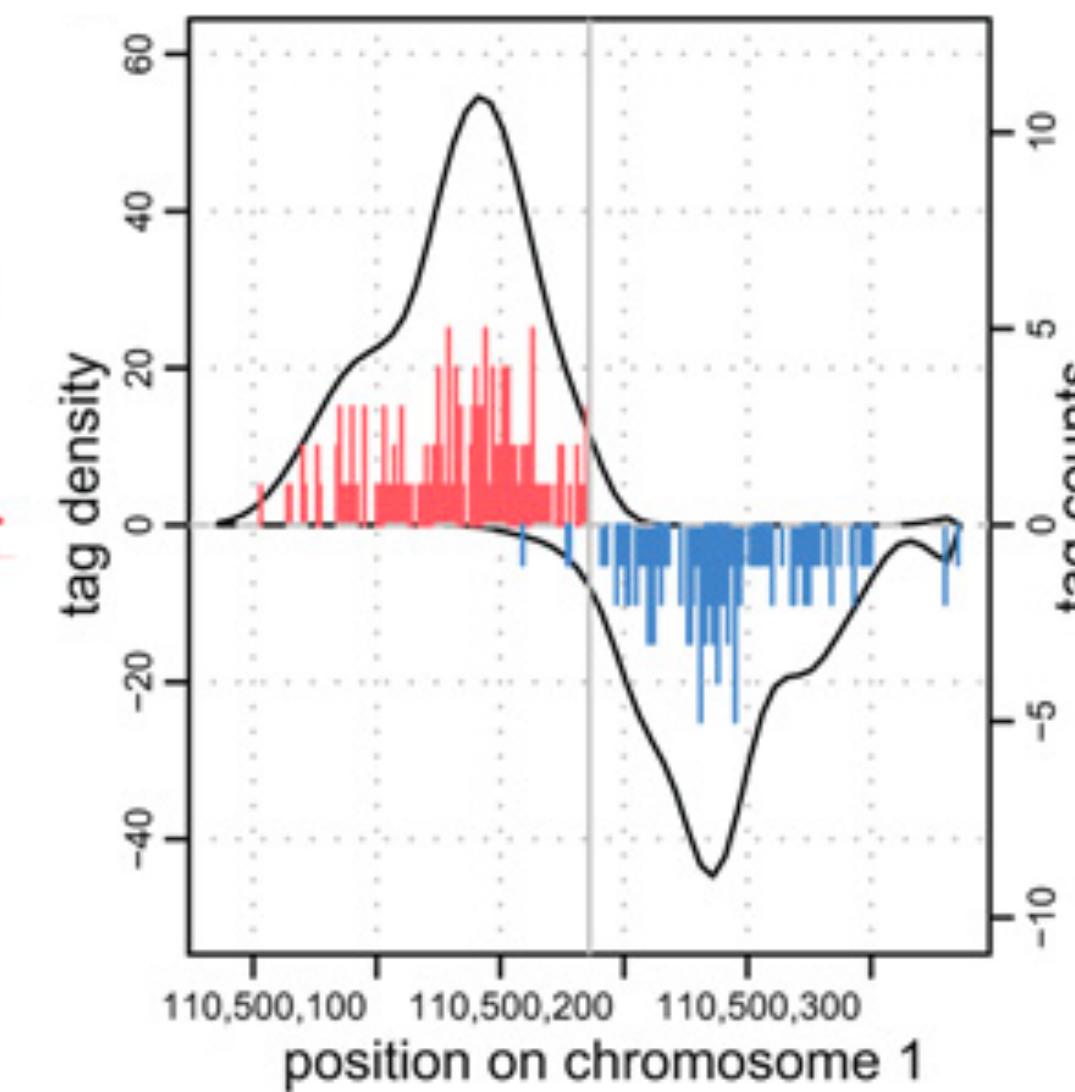
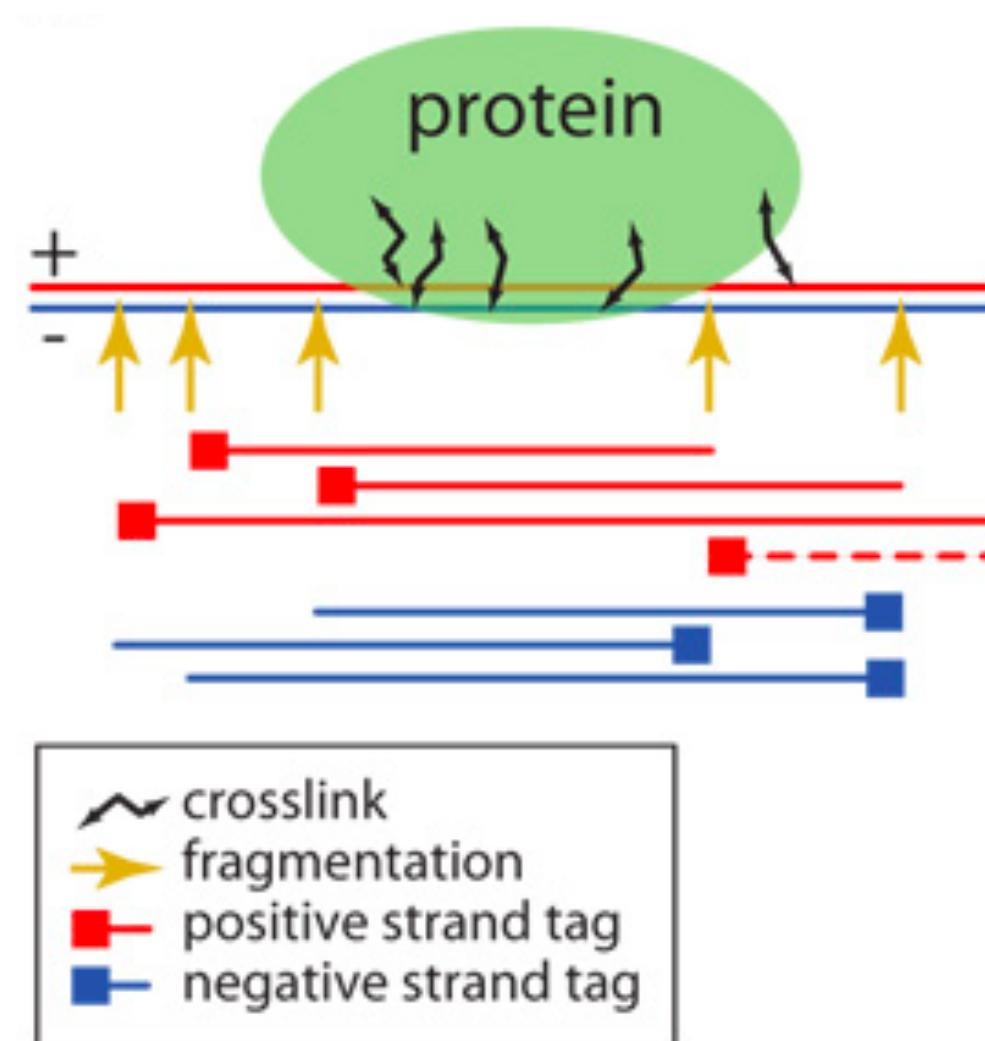
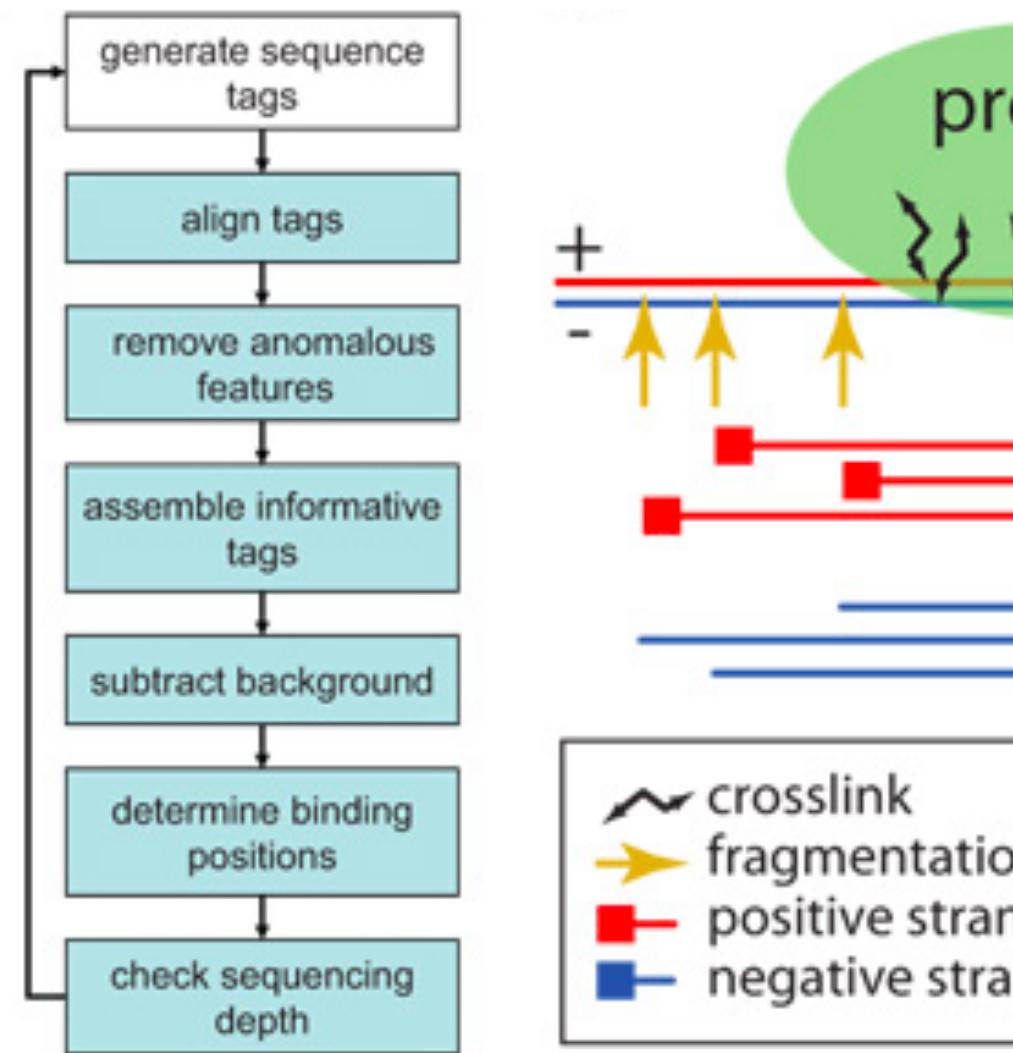
X\*\* = optional if sufficient data is available to split control data

X' = method excludes putative duplicated regions, no treatment of deletions

# Peak Calling

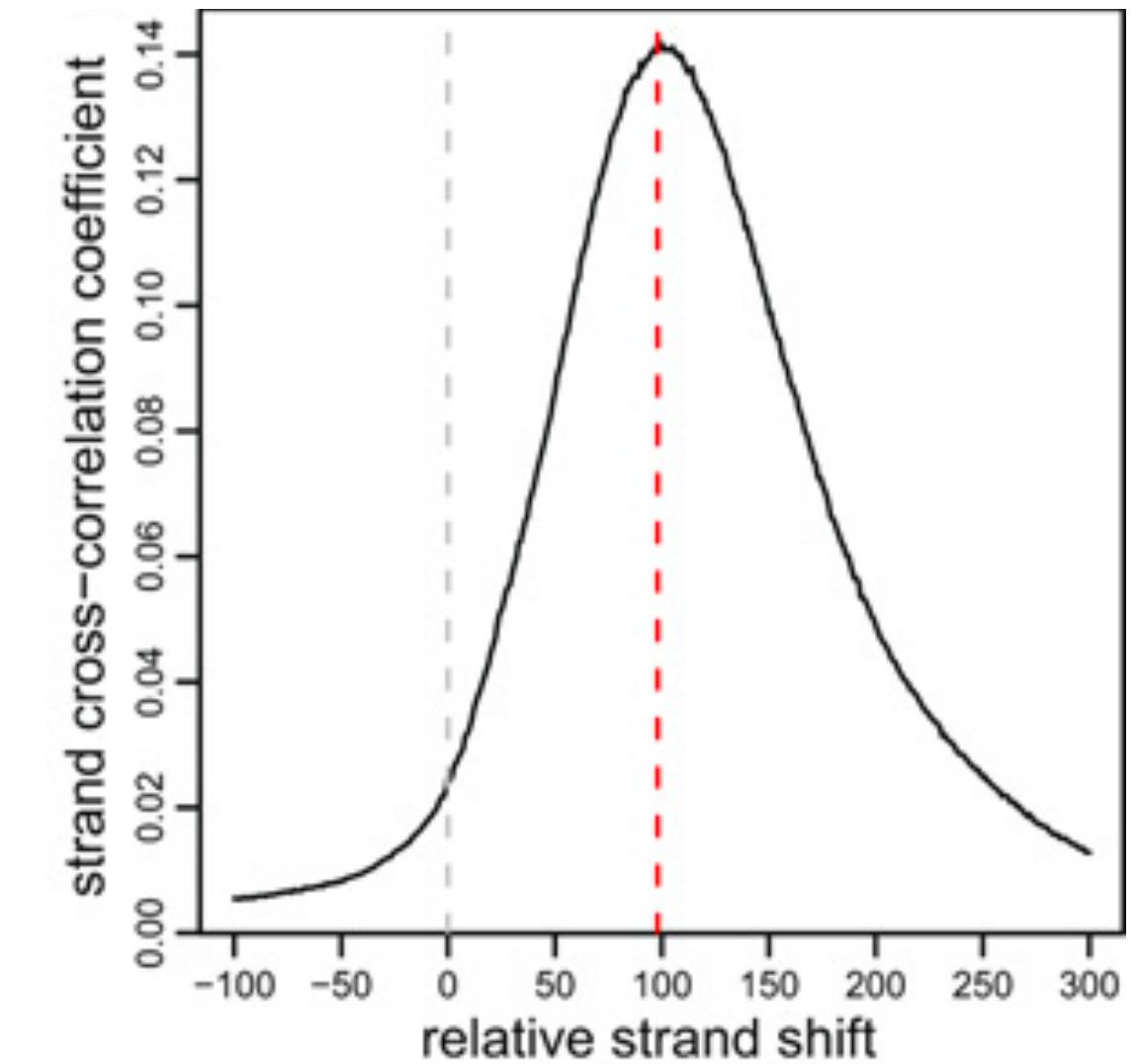
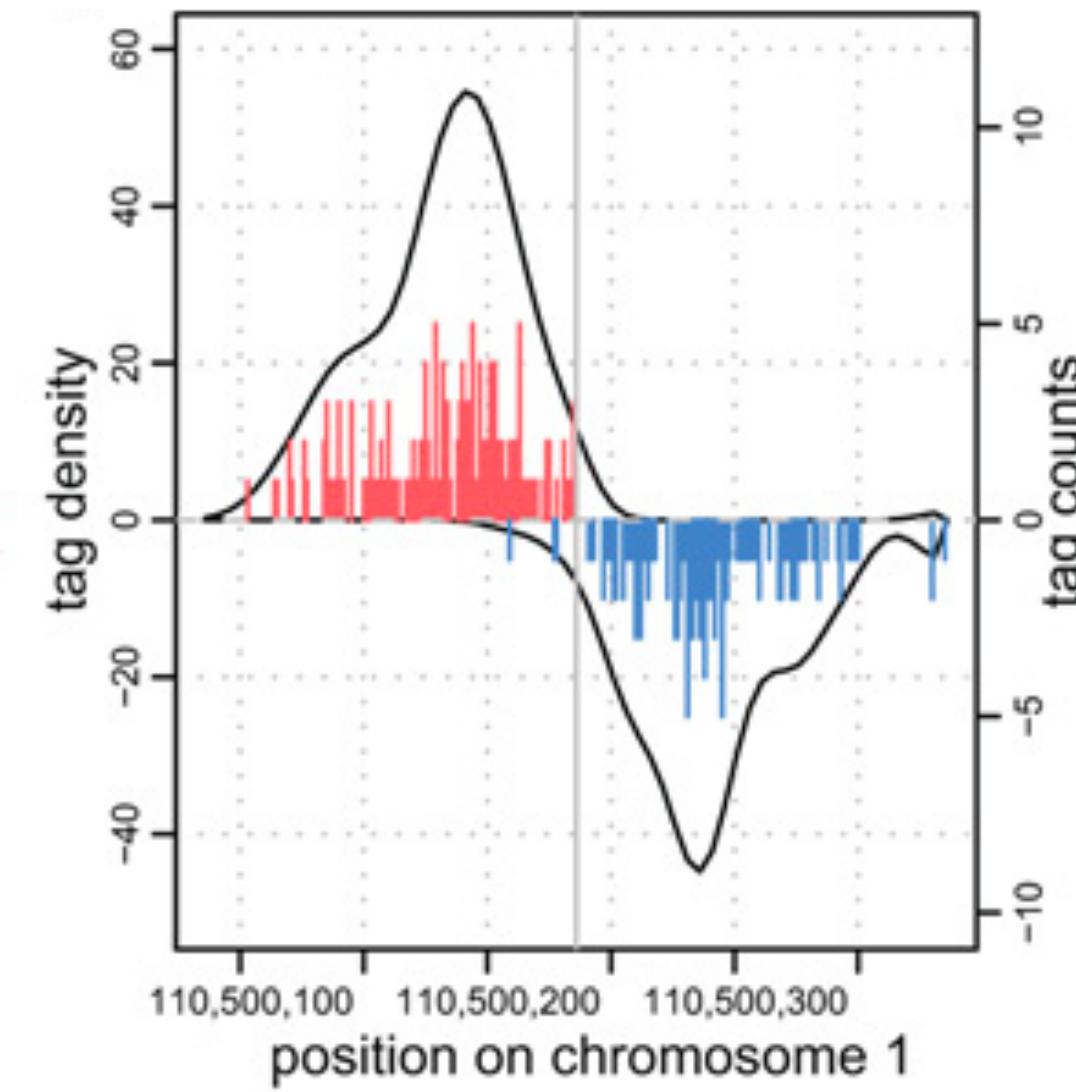
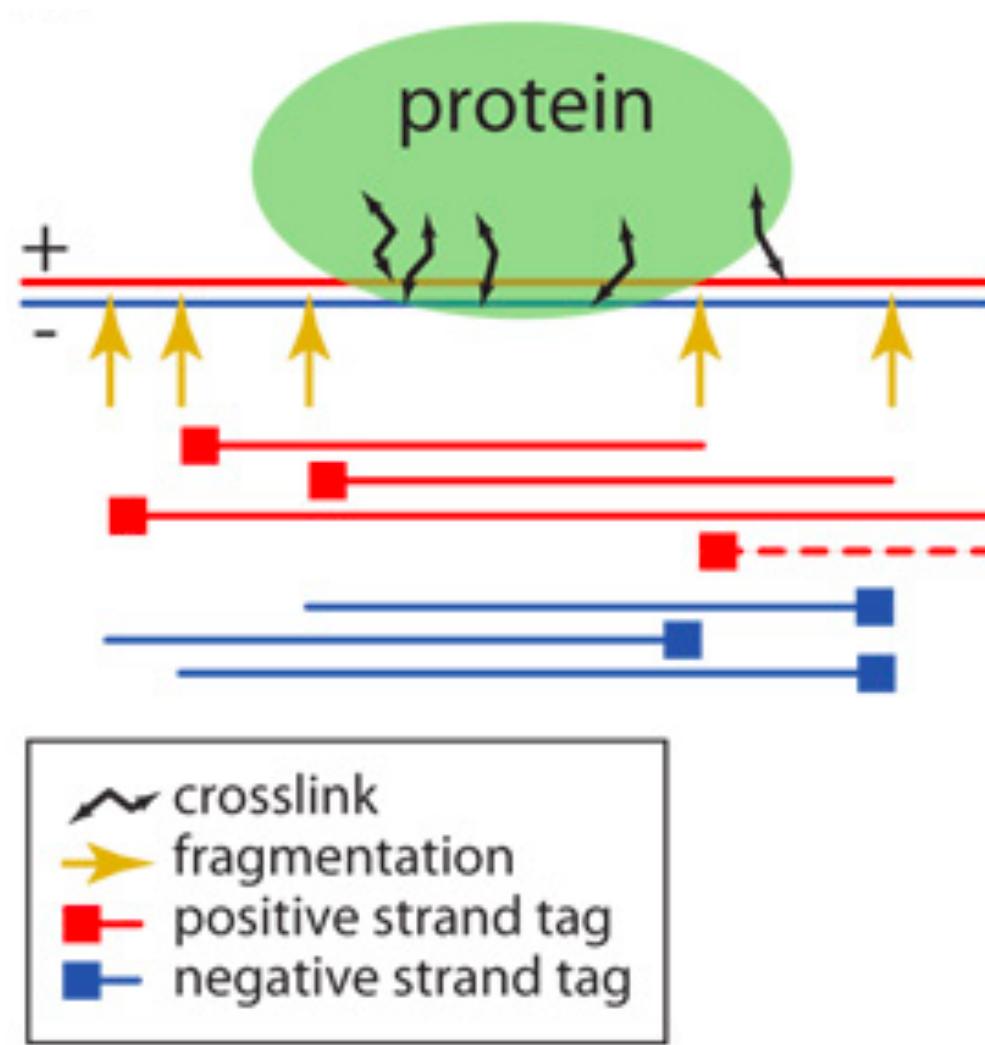
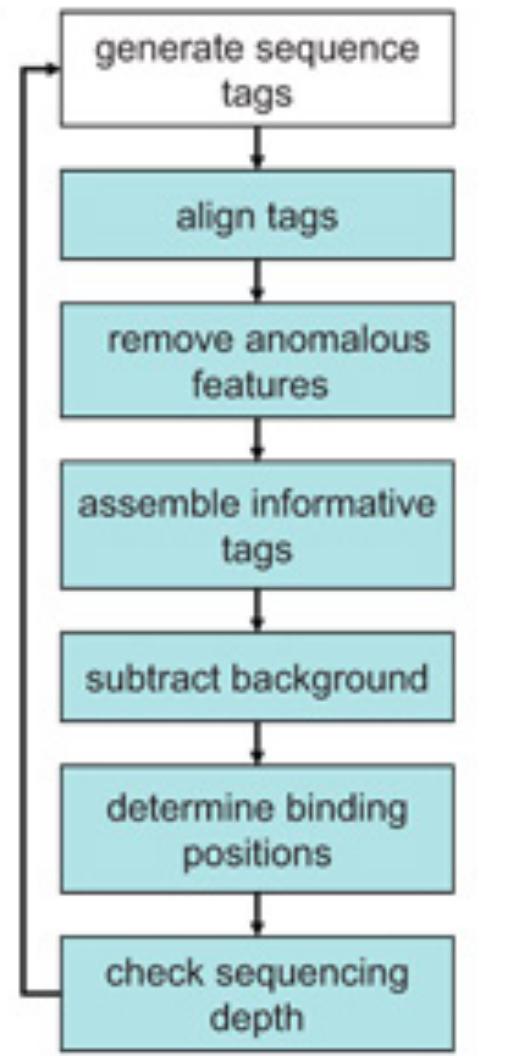
# How to choose one?

- ▶ Widely used
- ▶ Actively maintained and updated
- ▶ Default settings are a good start but know your parameters for your peak caller
- ▶ Be critical! Visually inspect your data (IGV)



*Kharchenko, Nature Biotech, 2008*

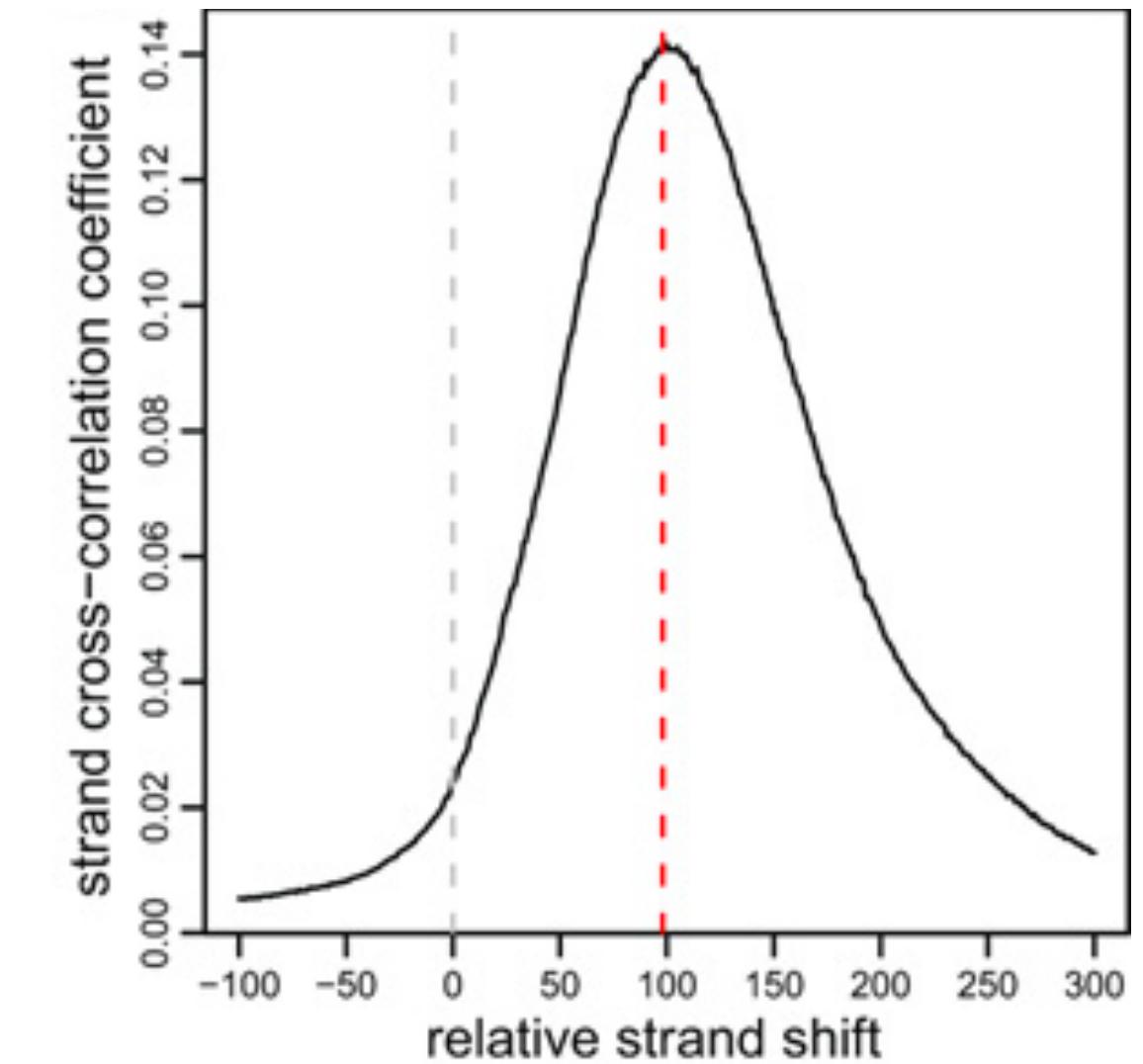
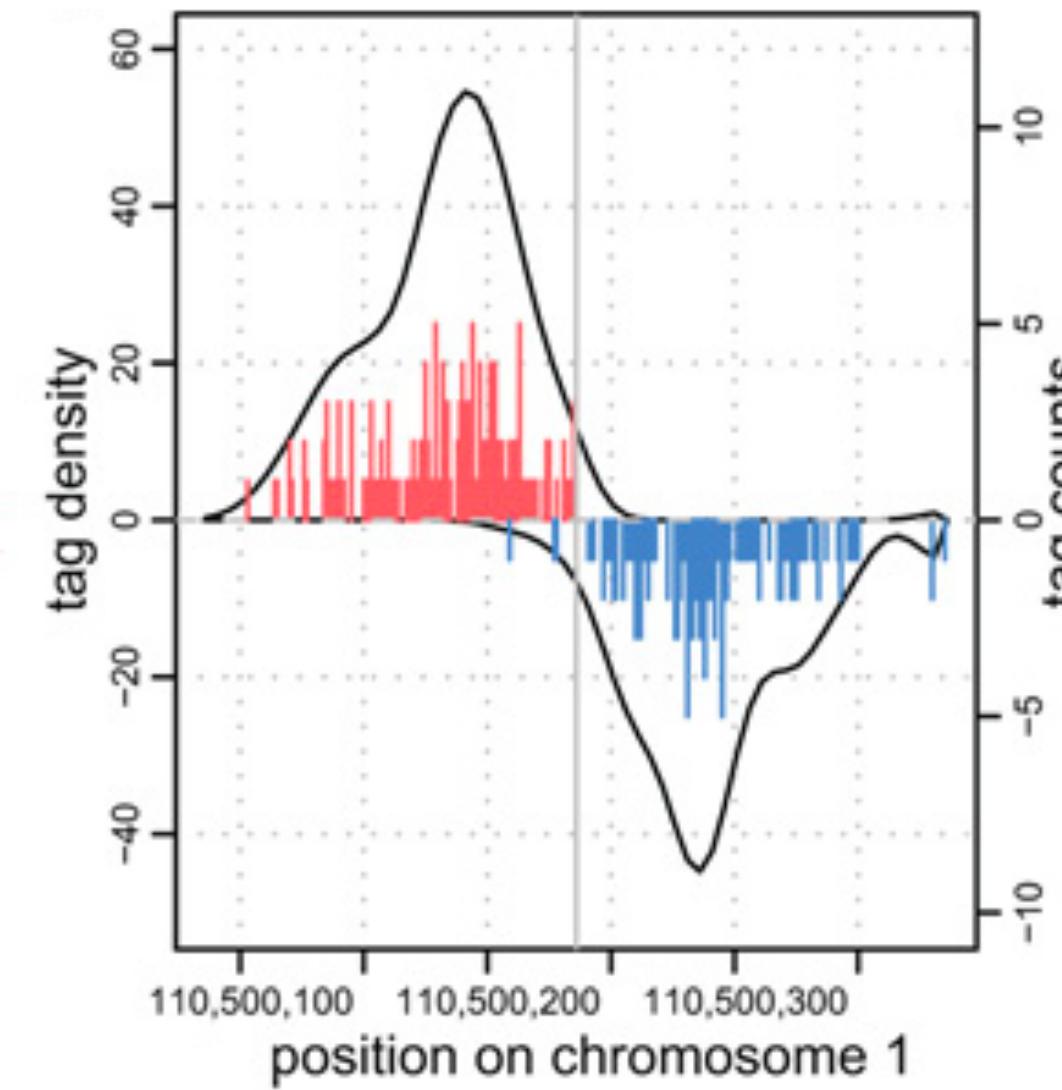
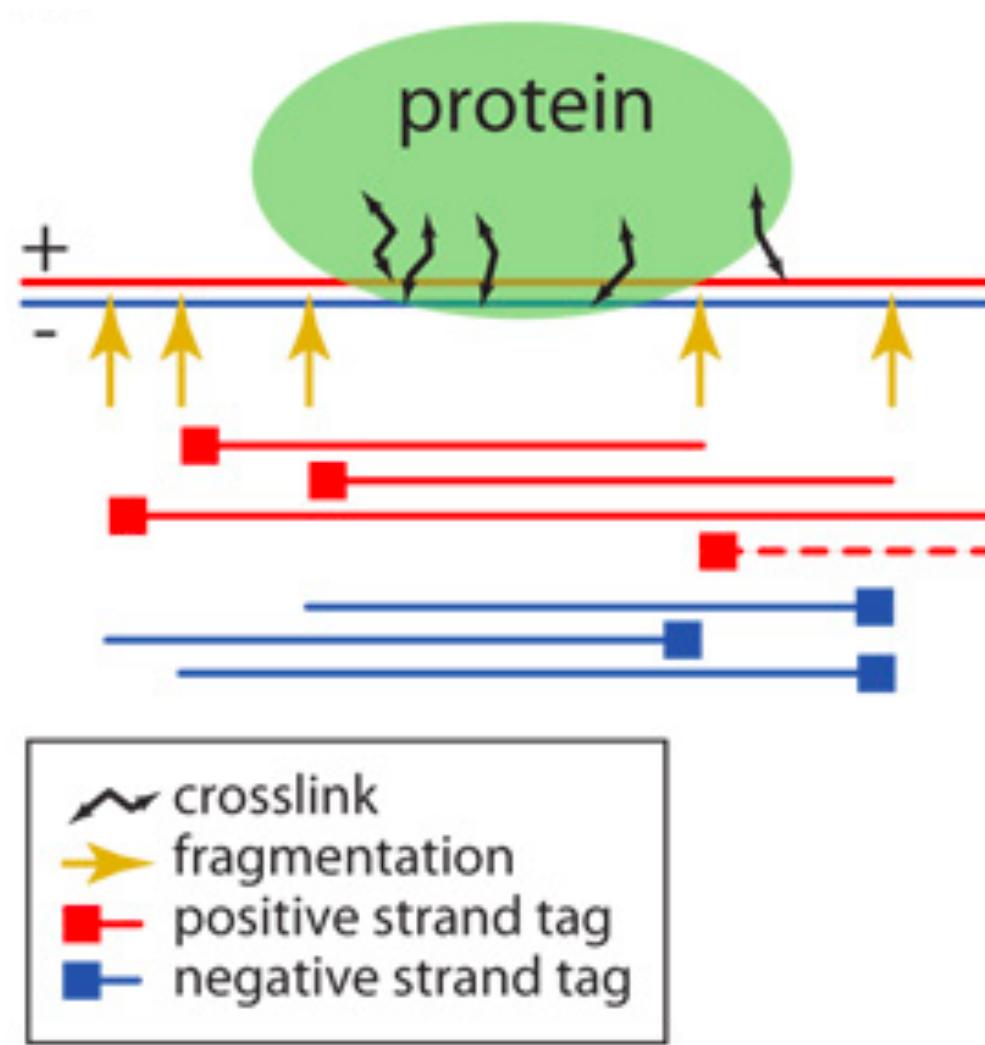
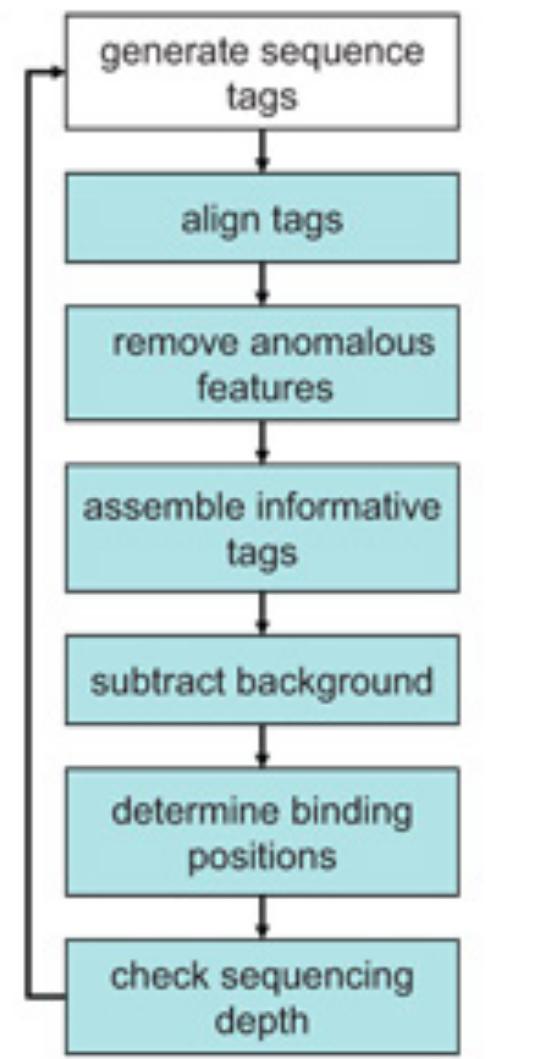
# The SPP peak calling pipeline



*Kharchenko, Nature Biotech, 2008*

- ChIP-seq fragments are sequenced from the 5' end only

# The SPP peak calling pipeline



*Kharchenko, Nature Biotech, 2008*

- ▶ ChIP-seq fragments are sequenced from the 5' end only
- ▶ Alignment generates bimodal pattern which is used to estimate the relative strand shift

# The SPP peak calling pipeline

# Strand cross correlation profile

# Strand cross correlation profile

- ▶ Compute the number of read starts at each position on the + strand and separately on the - strand for each chromosome

# Strand cross correlation profile

- ▶ Compute the number of read starts at each position on the + strand and separately on the - strand for each chromosome
- ▶ Shift these vectors wrt each other and compute the correlation for each shift

# Strand cross correlation profile

- ▶ Compute the number of read starts at each position on the + strand and separately on the - strand for each chromosome
- ▶ Shift these vectors wrt each other and compute the correlation for each shift

$$X(\delta) = \sum_{c \in C} \frac{N_c}{N} P\left[n_c^+(x + \delta/2), n_c^-(x - \delta/2)\right]$$

# Strand cross correlation profile

- ▶ Compute the number of read starts at each position on the + strand and separately on the - strand for each chromosome
- ▶ Shift these vectors wrt each other and compute the correlation for each shift

$$X(\delta) = \sum_{c \in C} \frac{N_c}{N} P\left[n_c^+(x + \delta/2), n_c^-(x - \delta/2)\right]$$

# Strand cross correlation profile

- ▶ Compute the number of read starts at each position on the + strand and separately on the - strand for each chromosome
- ▶ Shift these vectors wrt each other and compute the correlation for each shift

$$X(\delta) = \sum_{c \in C} \frac{N_c}{N} P\left[n_c^+(x + \delta/2), n_c^-(x - \delta/2)\right]$$

where  $\delta$  = strand shift

$n_c^s$  = number of tags whose 5' end maps to the position  $x$  on the strand  $s$

$P[a, b]$  = Pearson linear correlation coefficient

$N_c$  = number of tags mapped to a chromosome  $c$

$N$  = total number of tags

# Strand cross correlation profile

- ▶ Compute the number of read starts at each position on the + strand and separately on the - strand for each chromosome
- ▶ Shift these vectors wrt each other and compute the correlation for each shift

$$X(\delta) = \sum_{c \in C} \frac{N_c}{N} P\left[n_c^+(x + \delta/2), n_c^-(x - \delta/2)\right]$$

where  $\delta$  = strand shift

$n_c^s$  = number of tags whose 5' end maps to the position  $x$  on the strand  $s$

$P[a, b]$  = Pearson linear correlation coefficient

$N_c$  = number of tags mapped to a chromosome  $c$

$N$  = total number of tags

- ▶ Plot a cross-correlation profile as the cross-correlation values on the y-axis and the shift that you used to compute the correlation on the x-axis

# Strand cross correlation profile

- ▶ Compute the number of read starts at each position on the + strand and separately on the - strand for each chromosome
- ▶ Shift these vectors wrt each other and compute the correlation for each shift

$$X(\delta) = \sum_{c \in C} \frac{N_c}{N} P\left[n_c^+(x + \delta/2), n_c^-(x - \delta/2)\right]$$

where  $\delta$  = strand shift

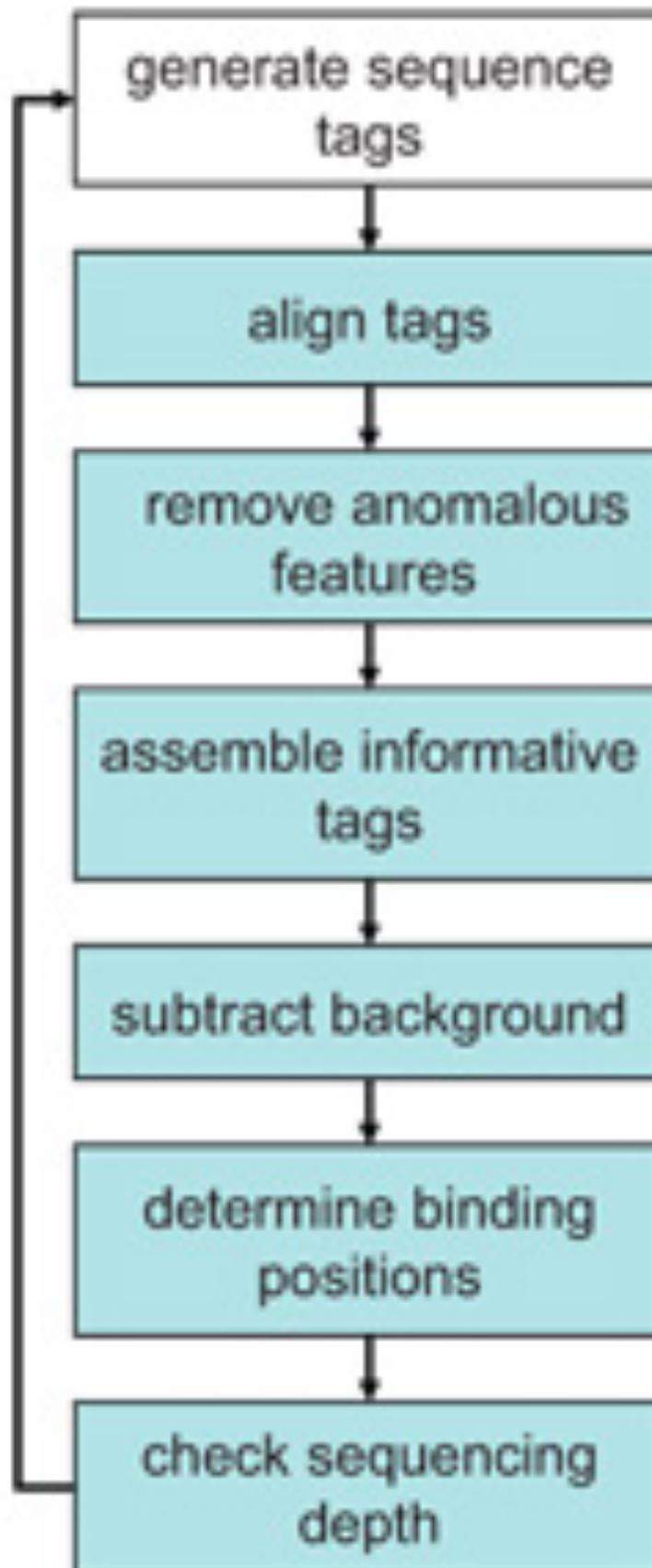
$n_c^s$  = number of tags whose 5' end maps to the position  $x$  on the strand  $s$

$P[a, b]$  = Pearson linear correlation coefficient

$N_c$  = number of tags mapped to a chromosome  $c$

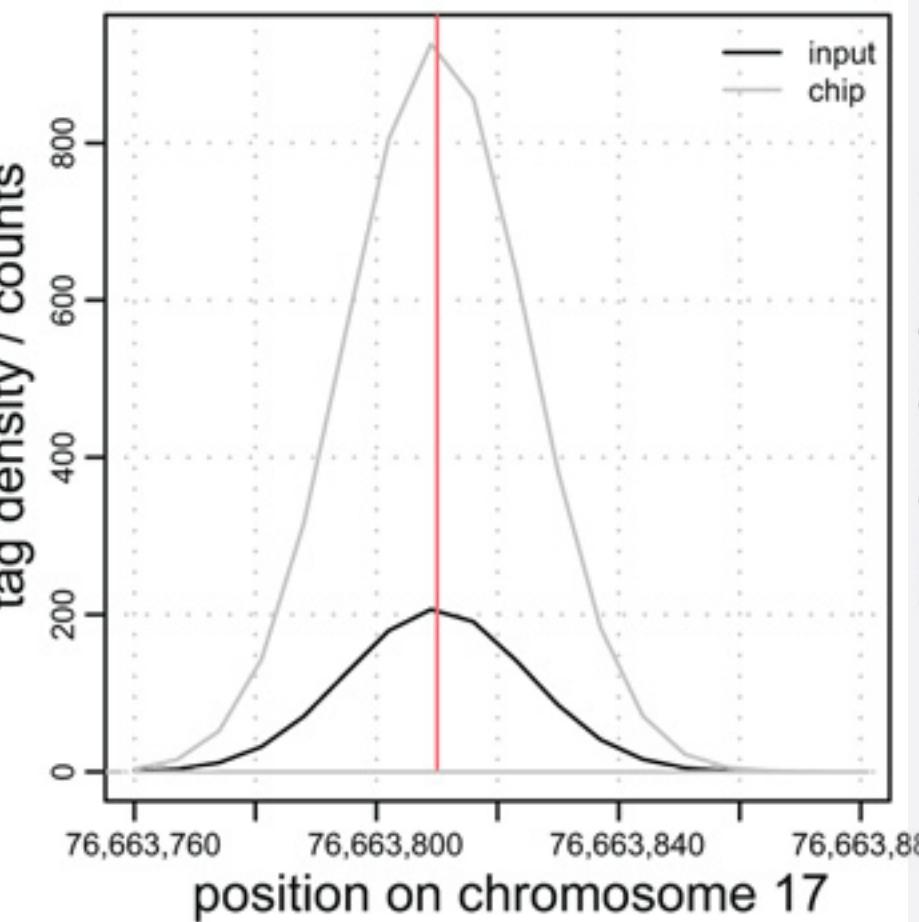
$N$  = total number of tags

- ▶ Plot a cross-correlation profile as the cross-correlation values on the y-axis and the shift that you used to compute the correlation on the x-axis
- ▶ Due to the 'shift' phenomenon of reads on the + and - strand around true binding sites, one would get a peak in the cross-correlation profile at the predominant fragment length



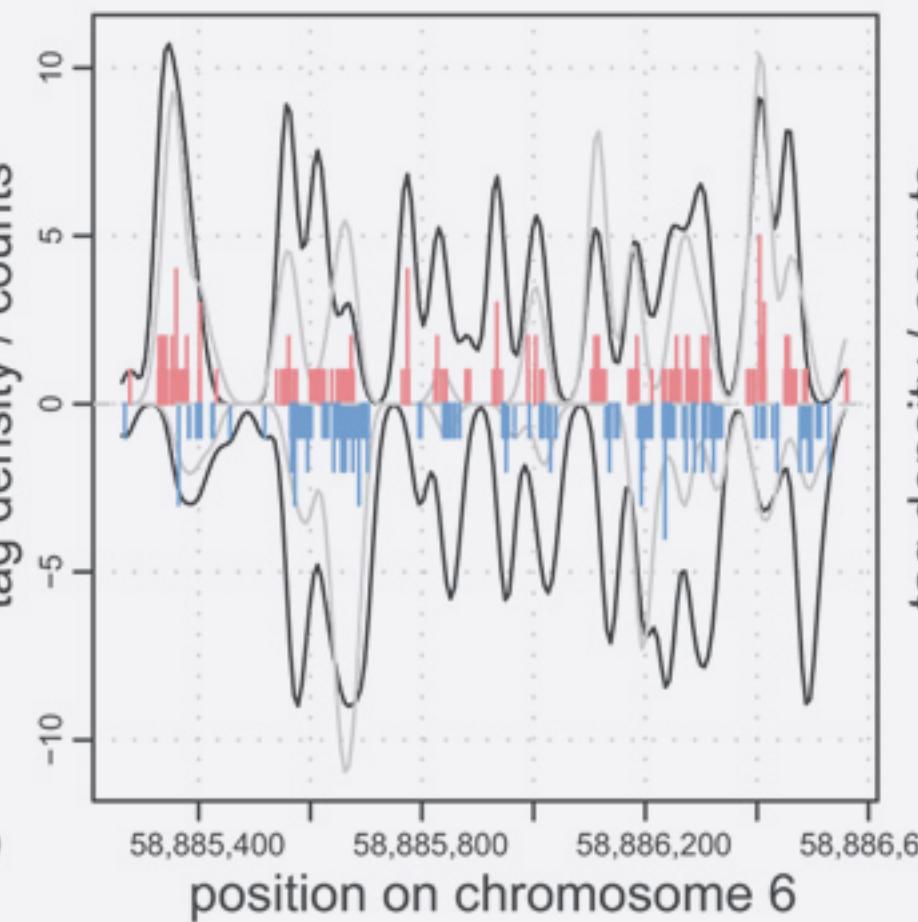
## Density of tags from ChIP and input samples showing three types of anomalies

a.



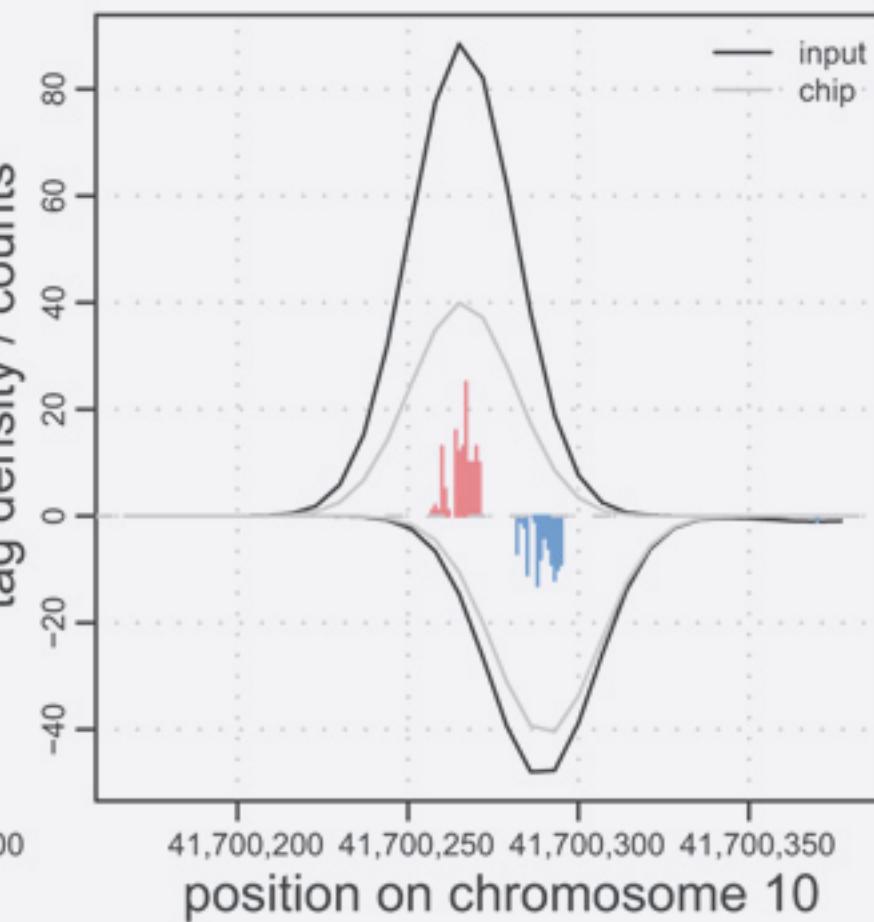
Singular positions with extremely high tag count.

b.



Larger (>1000bp), non-uniform regions of increased background tag density.

c.

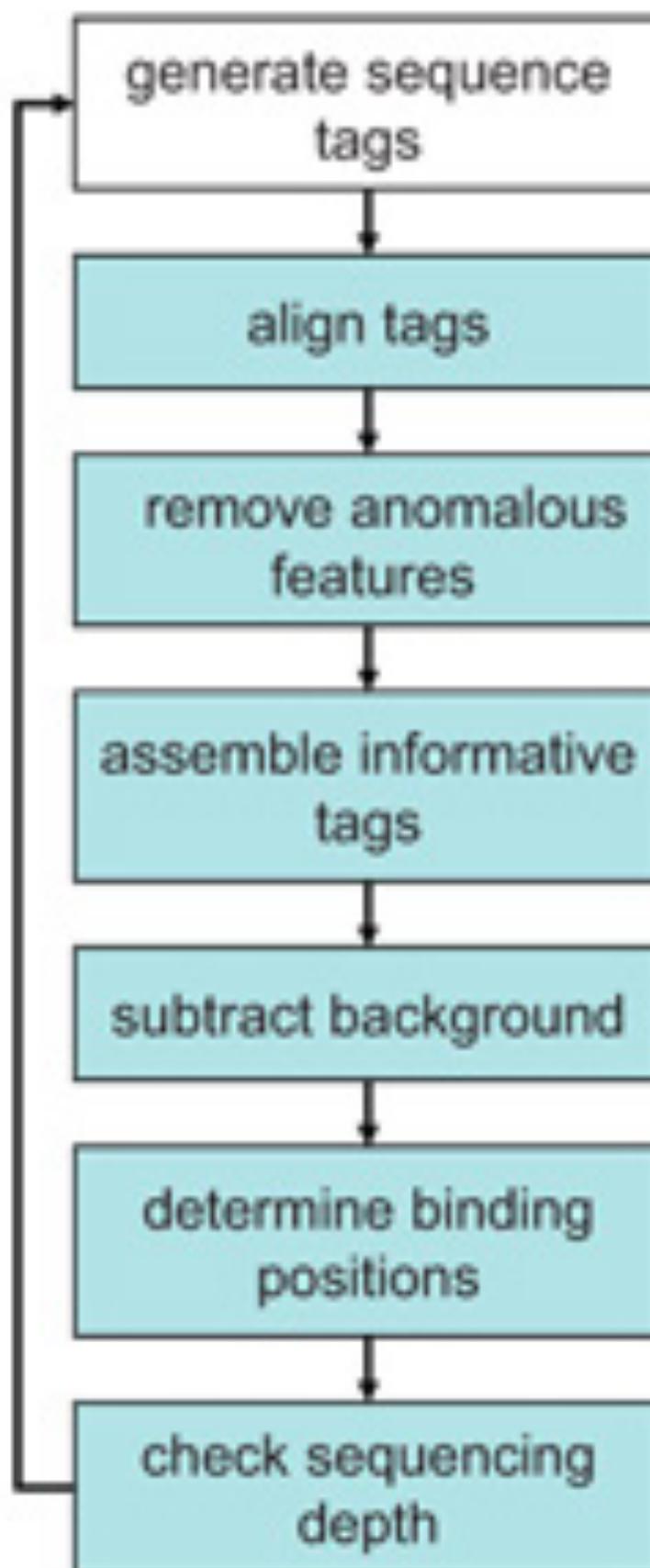


Background tag density patterns resembling true protein binding positions.

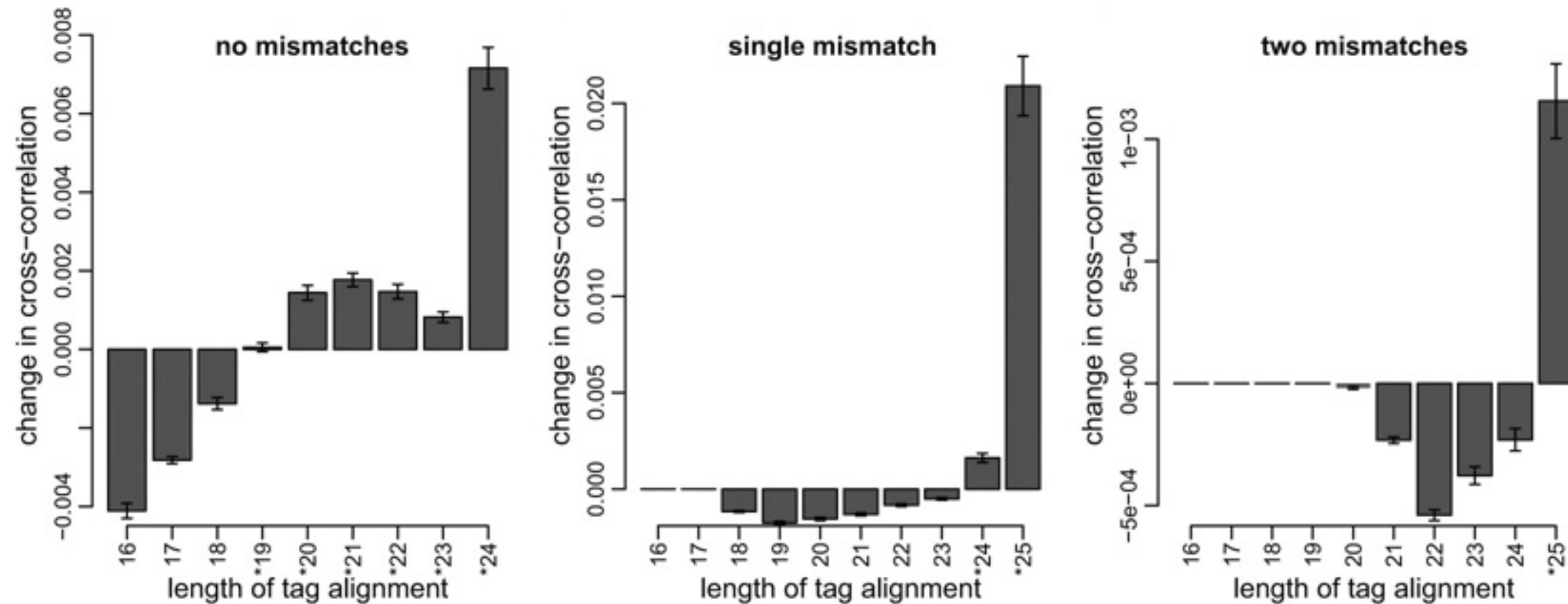
Positions with number of mapped sequence tags (5' ends) with Z-score >10. All tags mapping to such anomalous position (on either strand) are omitted.

*Kharchenko, Nature Biotech, 2008*

# SPP: Remove anomalous features



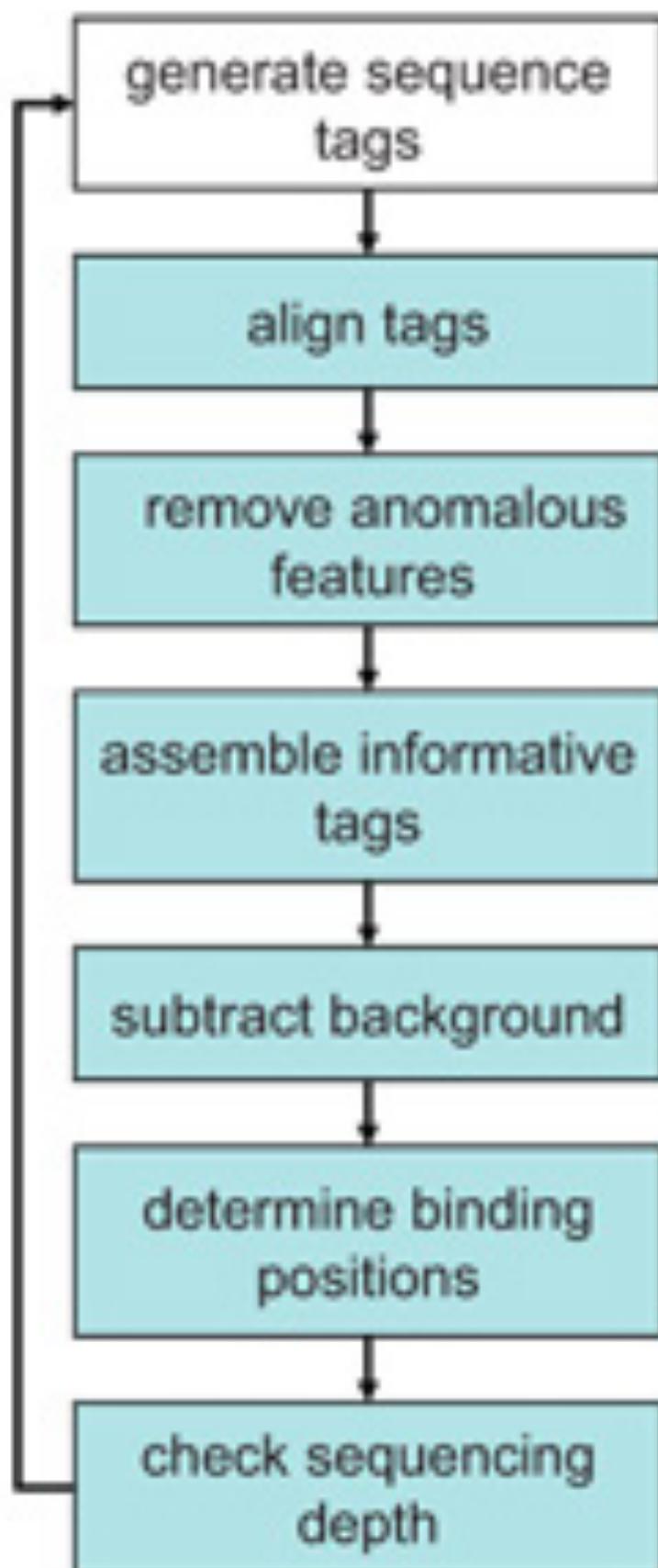
### Select informative tag classes based on change in strand cross-correlation magnitude



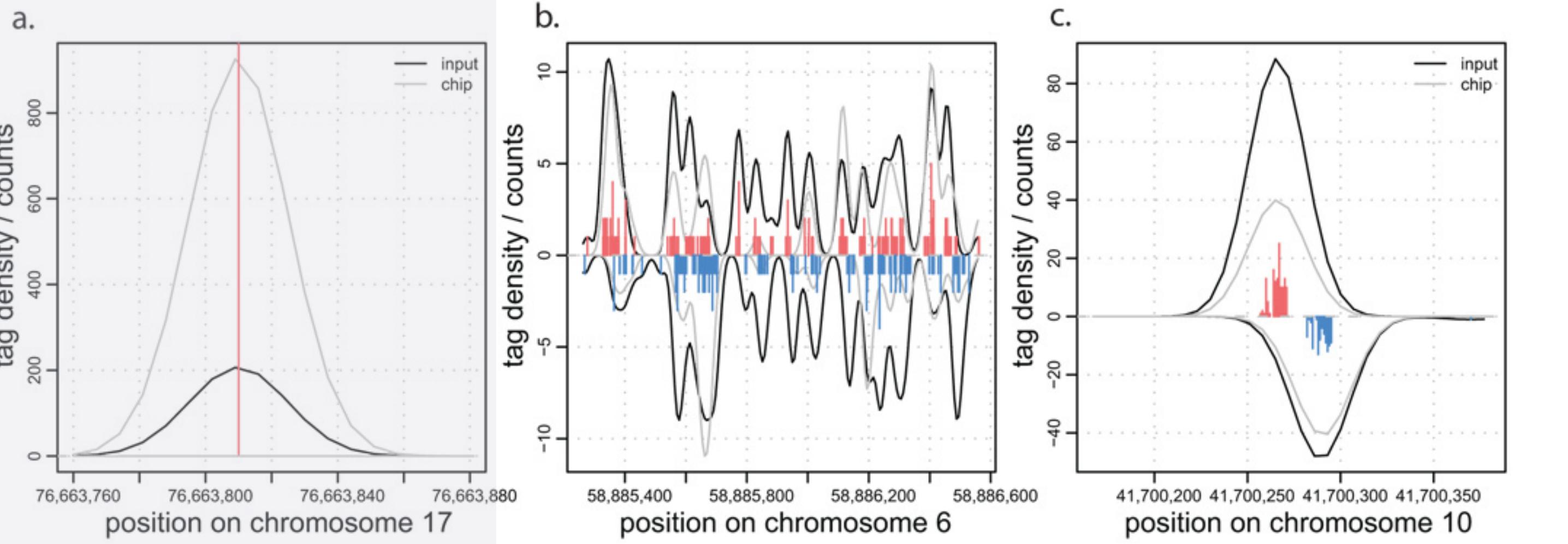
Informative tag classes improve cross-correlation (marked by \*), and are incorporated into the final tag set. The y-axis gives the mean change in cross-correlation profile within 40bp around the cross-correlation peak

*Kharchenko, Nature Biotech, 2008*

# SPP: Assemble informative tags



## Density of tags from ChIP and input samples showing three types of anomalies



Singular positions with extremely high tag count.

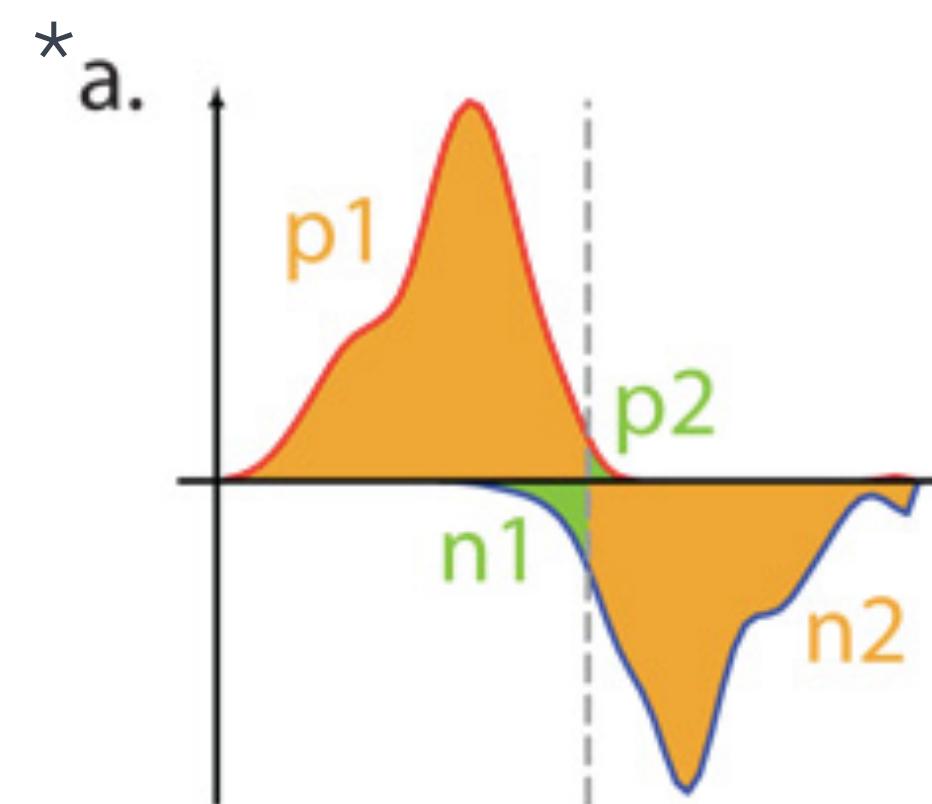
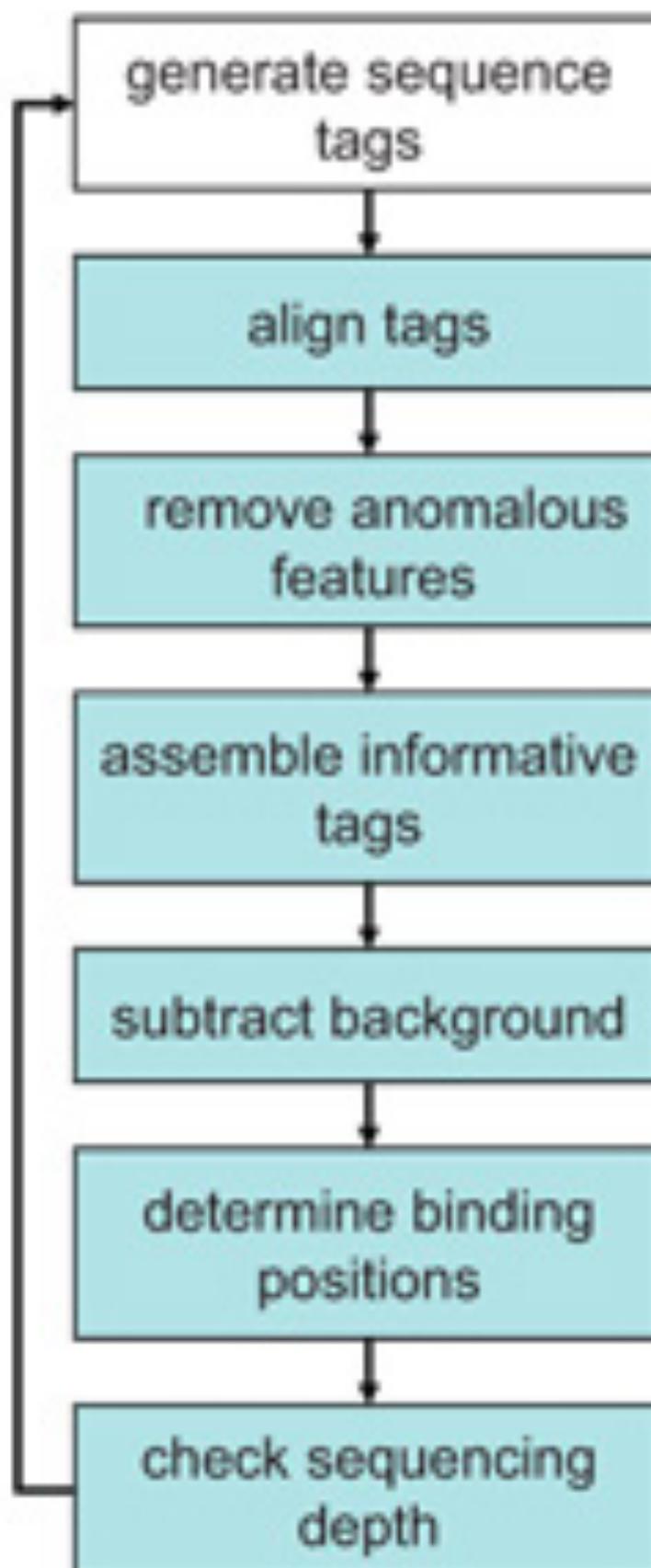
Larger (>1000bp), non-uniform regions of increased background tag density.

Background tag density patterns resembling true protein binding positions.

WTD and MTC methods were adjusted by subtracting the weighted number of background (input) tags occurring within that window.

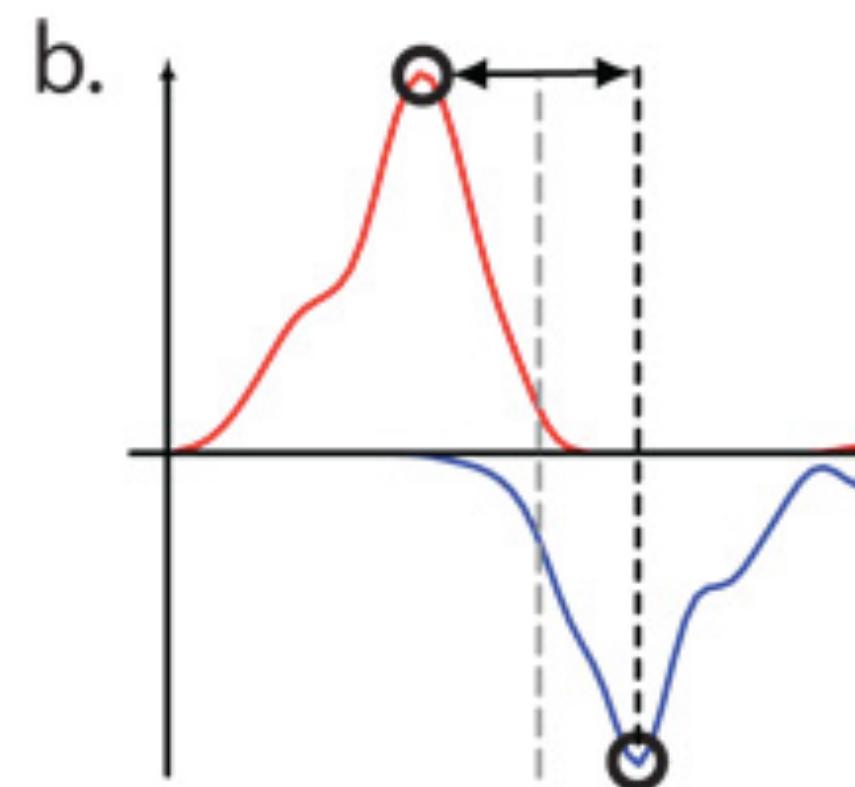
*Kharchenko, Nature Biotech, 2008*

# SPP: Subtract background



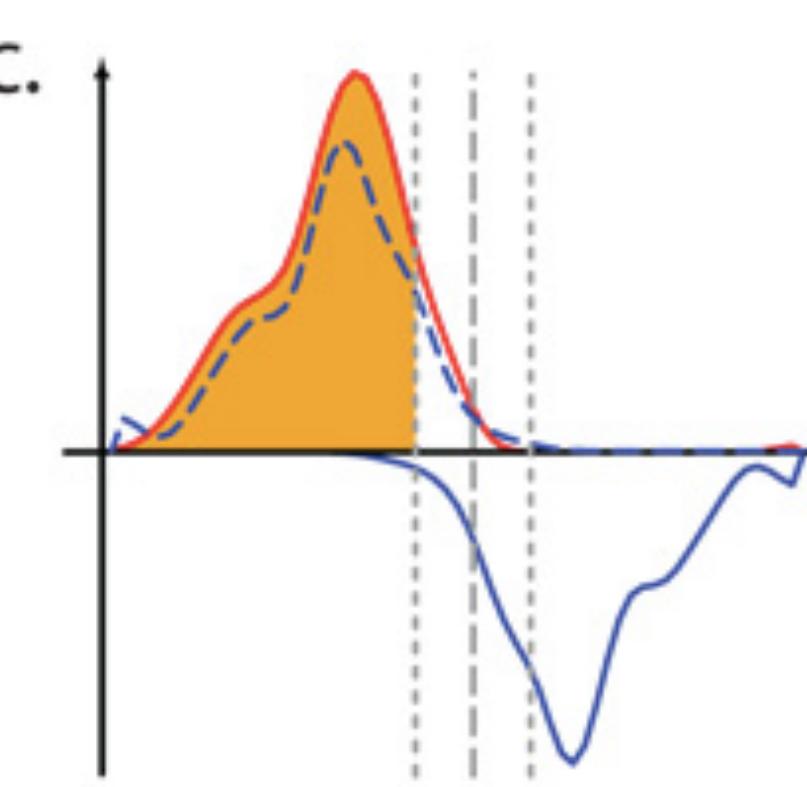
#### Window Tag Density (WTD)

Calculates the difference between geometric average of the tag counts within the regions marked by orange color (p1 and n2), and the average tag count within the regions marked by green color (n1 and p2). Window size based on average binding tag pattern (estimated from CC plot)



#### Matching Strand Peaks (MSP)

Identifies local maxima on positive and negative strands and then determines positions where such two peaks are present in the right order, with the expected separation (e.g. 20bp) and comparable magnitude (based on a likelihood ratio test)



#### Mirror Tag Correlation (MTC)

Similar to WTD. Based on the mirror correlation of the positive and negative strand tag densities. The mirror image of the negative strand density is shown by the blue dashed line. Uses the Pearson linear correlation coefficient.

*Kharchenko, Nature Biotech, 2008*

# SPP: Determine binding positions

# Estimation of the False Discovery Rate (FDR)

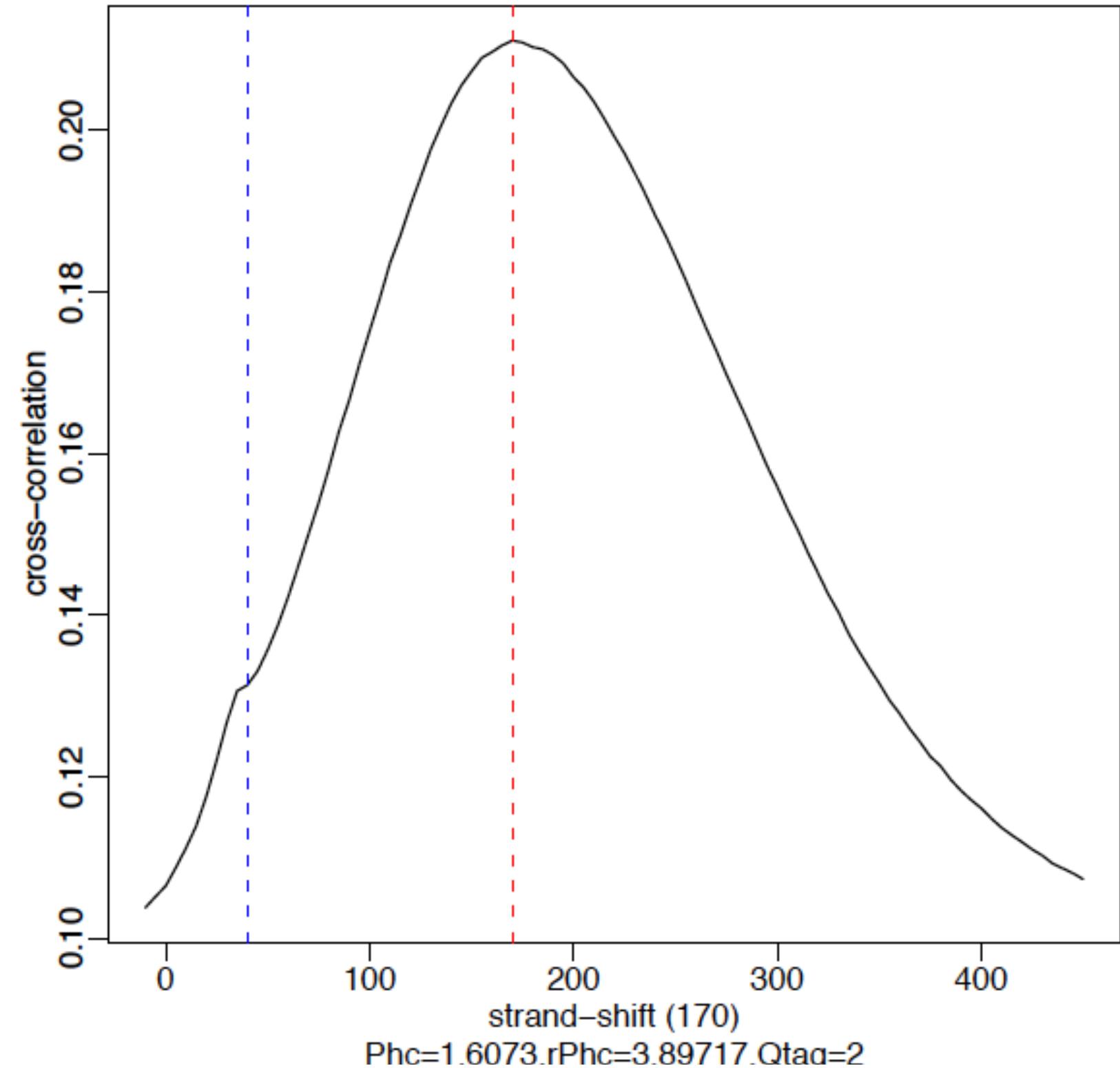
- ▶ Expected proportion of false discoveries among total rejections of the null hypothesis (i.e. how many false positive peaks of the total set of peaks)
- ▶ FDR determined empirically by exchanging ChIP and control
- ▶ Given by

$$\text{FDR}(s) = \frac{N_r(s) + 0.5}{N_c(s) + 0.5}$$

where  $N_r(s)$  is the number of binding positions with score  $s$  or higher found in the real dataset, and  $N_c(s)$  is the number found in a control dataset

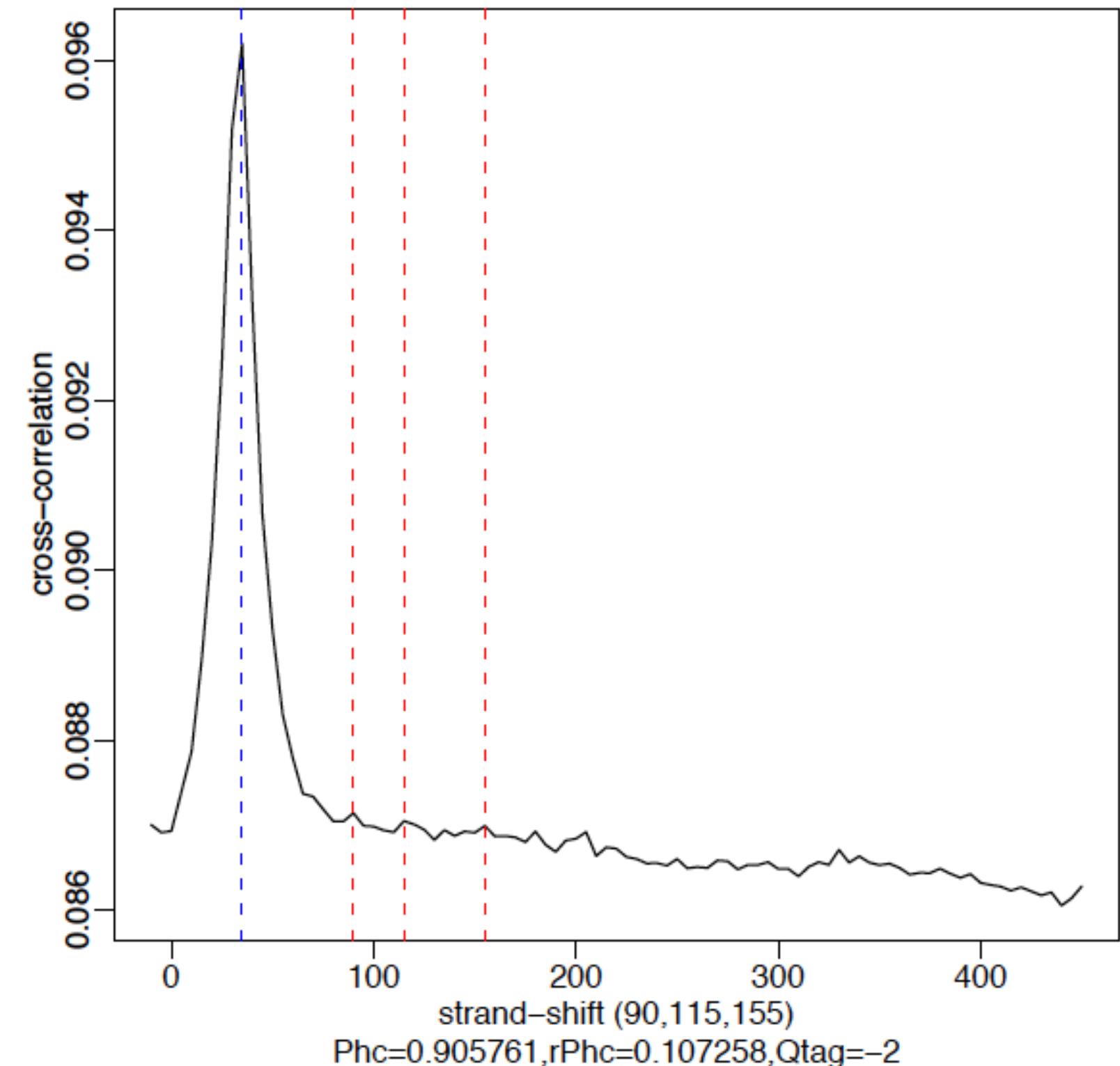
# Strong dataset: CTCF in human cells

- ▶ Great antibody and 45-60K peaks typically
- ▶ Red vertical line shows the dominant peak at the true peak shift
- ▶ Small bump at the blue vertical line is at read-length



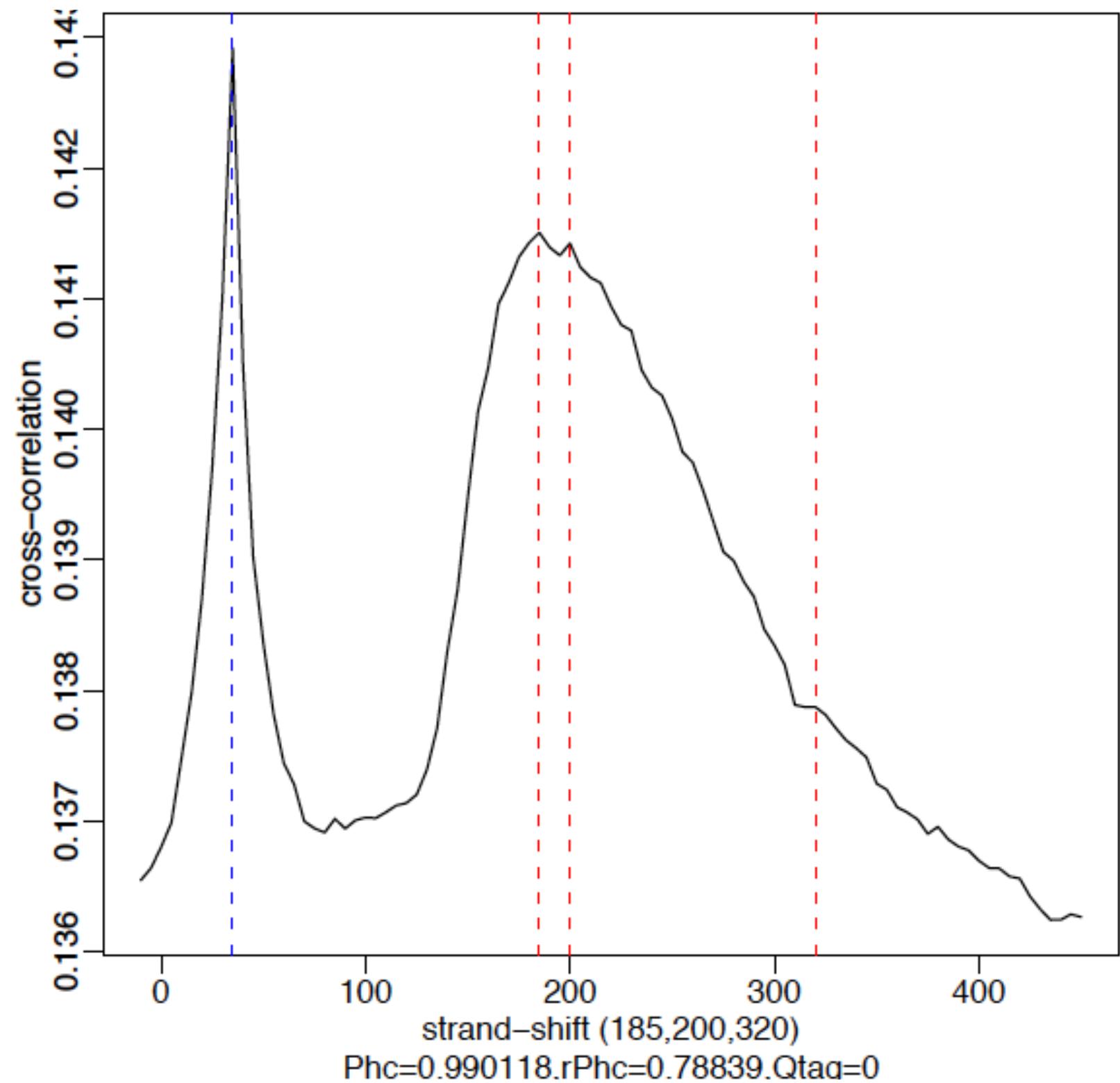
# Control dataset (input DNA)

- ▶ Note the strongest peak is the blue line (read length) and there is basically almost no other significant peak in the profile
- ▶ The absence of a peak is expected since there should be no significant clustering of fragments around specific target sites (except potentially weak biases in open chromatin regions depending on the protocol used)
- ▶ The read-length peak occurs due to unique mappability properties of the mapped reads



# Weaker dataset: POL2

- ▶ This particular antibody is not very efficient and these are broad scattered peaks
- ▶ Has few peaks (~3000 detectable)
- ▶ Two peaks in the cross-correlation profile: one at the true peak shift (~185-200 bp) and the other at read length. For such weaker datasets, the read-length peak starts to dominate.

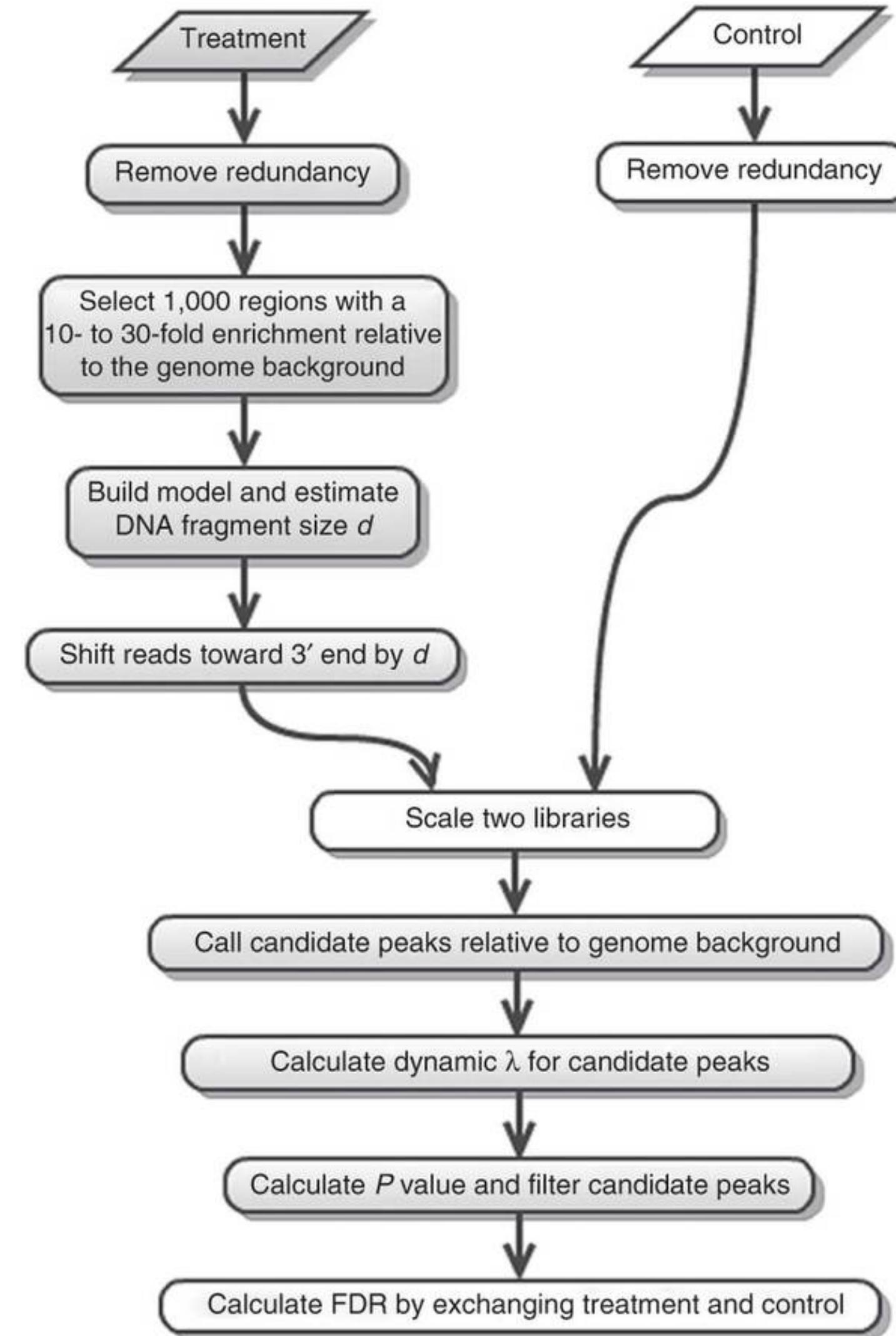


# MACS peak calling

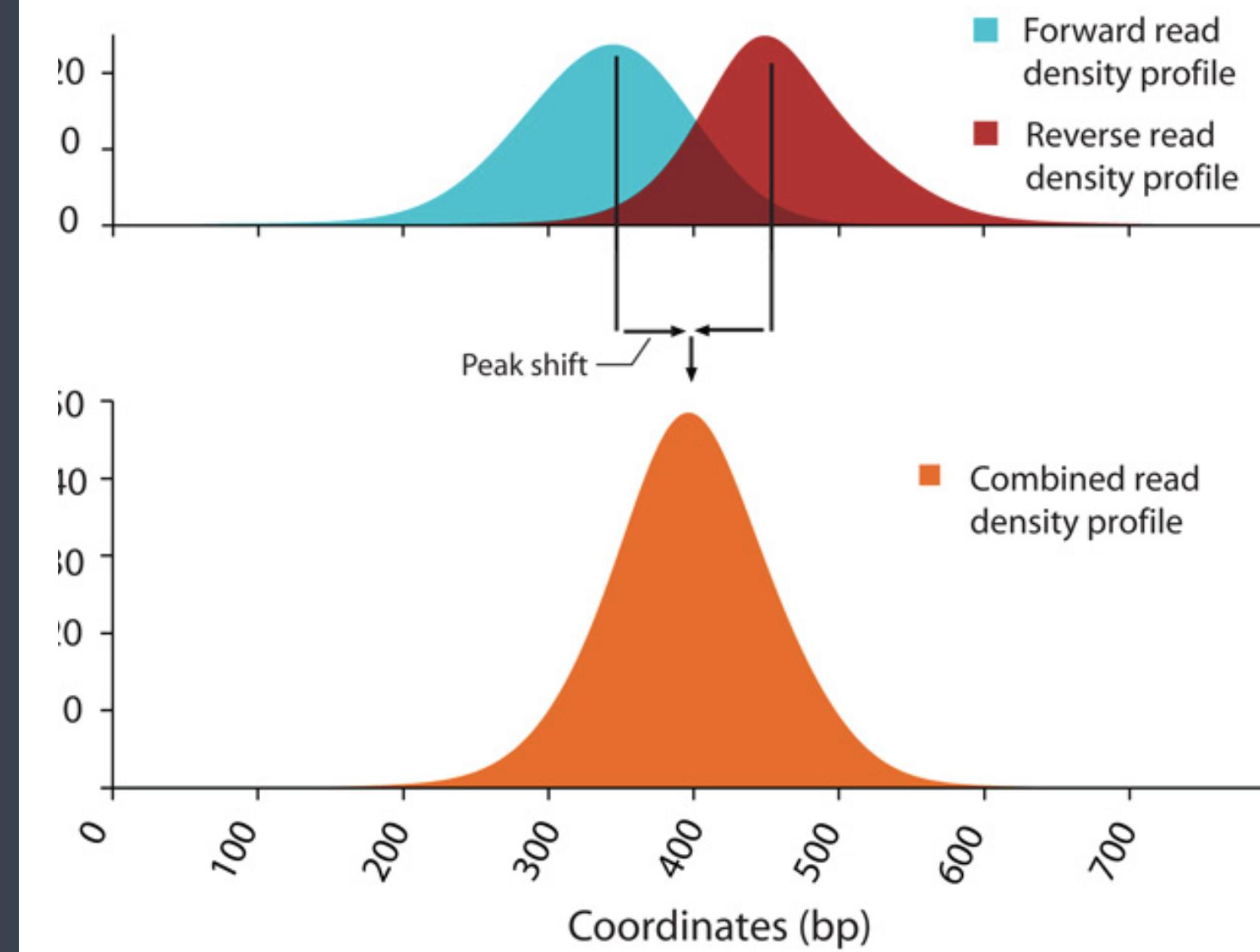
Zhang et al. Model-based analysis of ChIP-Seq,  
Genome Biol. (2008)

Developed for detection of transcription factor  
binding sites

Also suited for larger regions

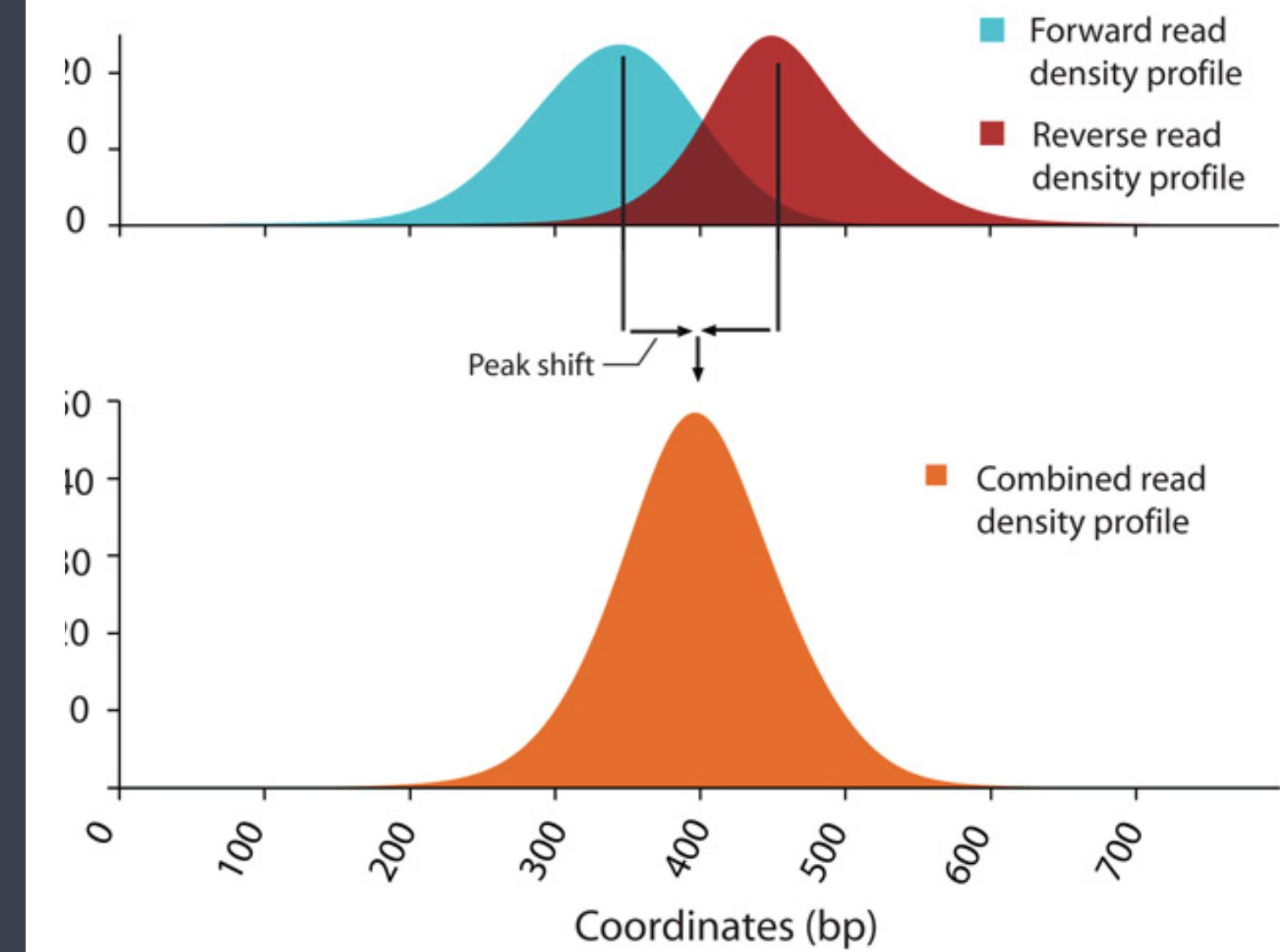


# Model read distribution and calculate the shift size, $d/2$



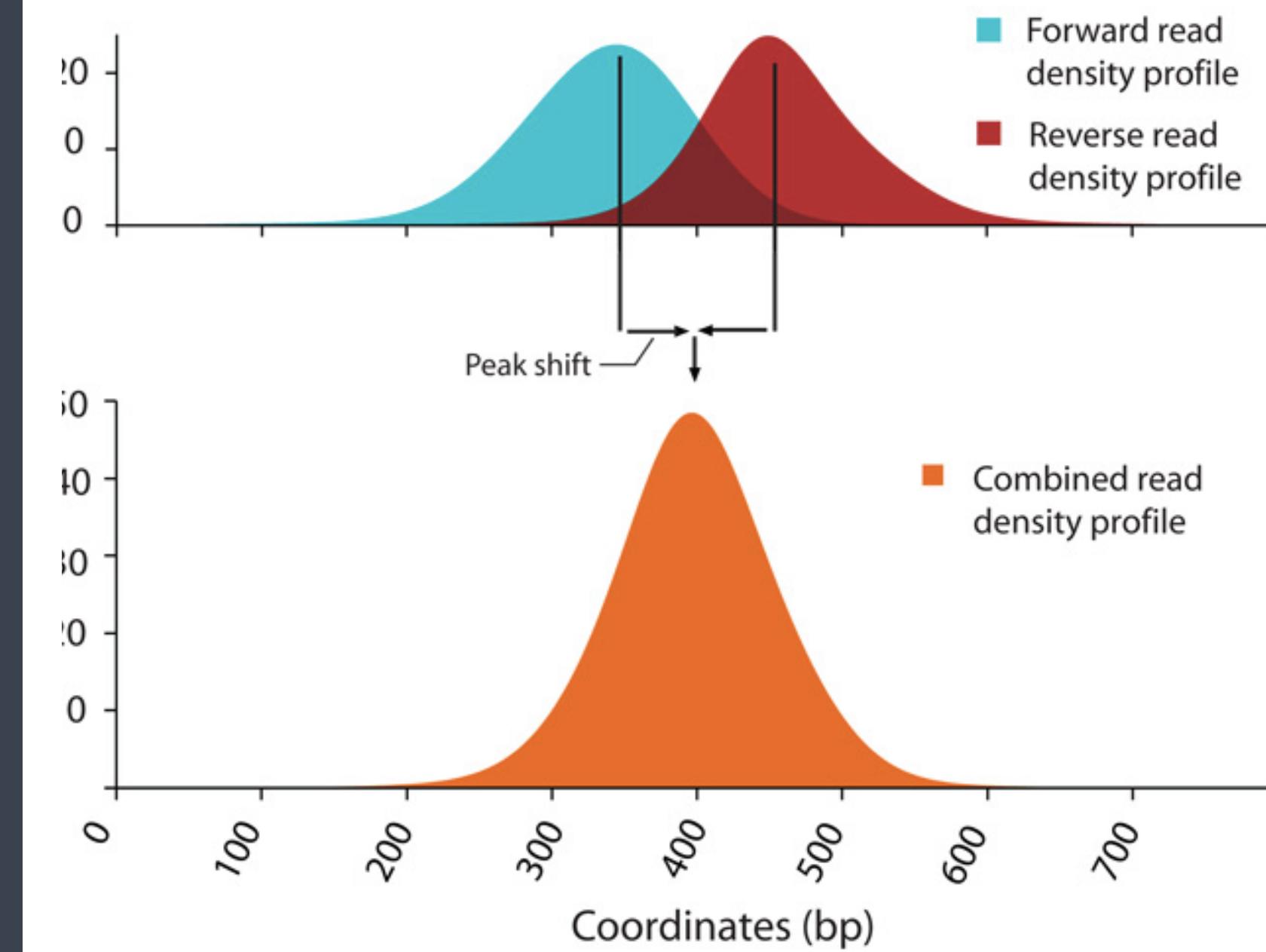
# Model read distribution and calculate the shift size, $d/2$

- ▶ Find paired peaks using a sliding window  $2^*$ size of sequence fragments
  - > Search for regions where reads are enriched more than  $MFOLD$



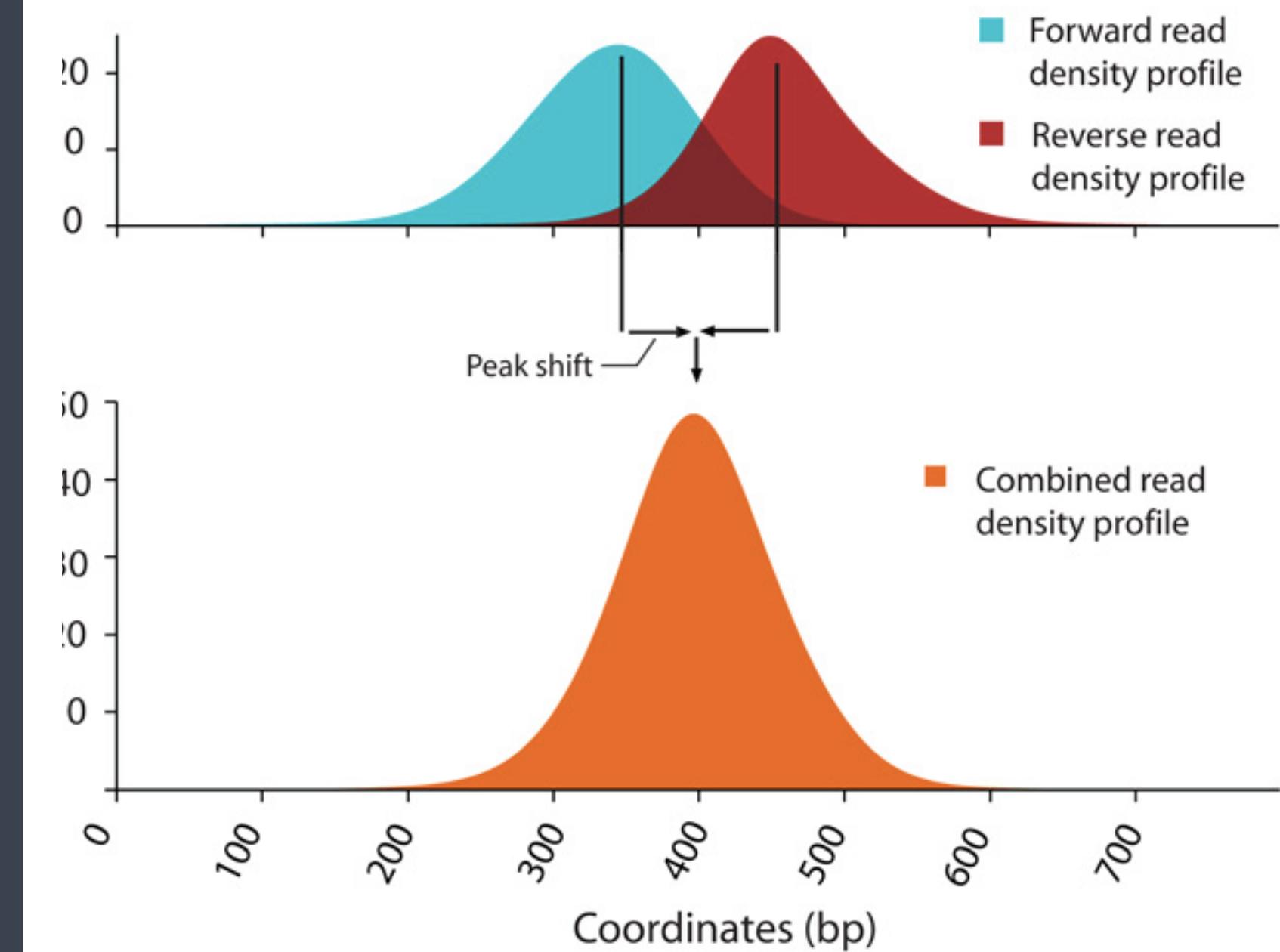
# Model read distribution and calculate the shift size, $d/2$

- ▶ Find paired peaks using a sliding window  $2^*$ size of sequence fragments
  - > Search for regions where reads are enriched more than  $MFOLD$
- ▶ Estimate the fragment length, ' $d$ '
  - > Distance between the modes of the positive and negative strand peaks



# Model read distribution and calculate the shift size, $d/2$

- ▶ Find paired peaks using a sliding window  $2^*$ size of sequence fragments
  - > Search for regions where reads are enriched more than  $MFOLD$
- ▶ Estimate the fragment length, ' $d$ '
  - > Distance between the modes of the positive and negative strand peaks
- ▶ Peak calling
  - > Shift the tags by  $1/2$  the fragment length, ' $d$ '
  - > Scan for enriched peaks comparing to background



# Peak detection

MACS models the number of reads from a genomic region as a Poisson distribution

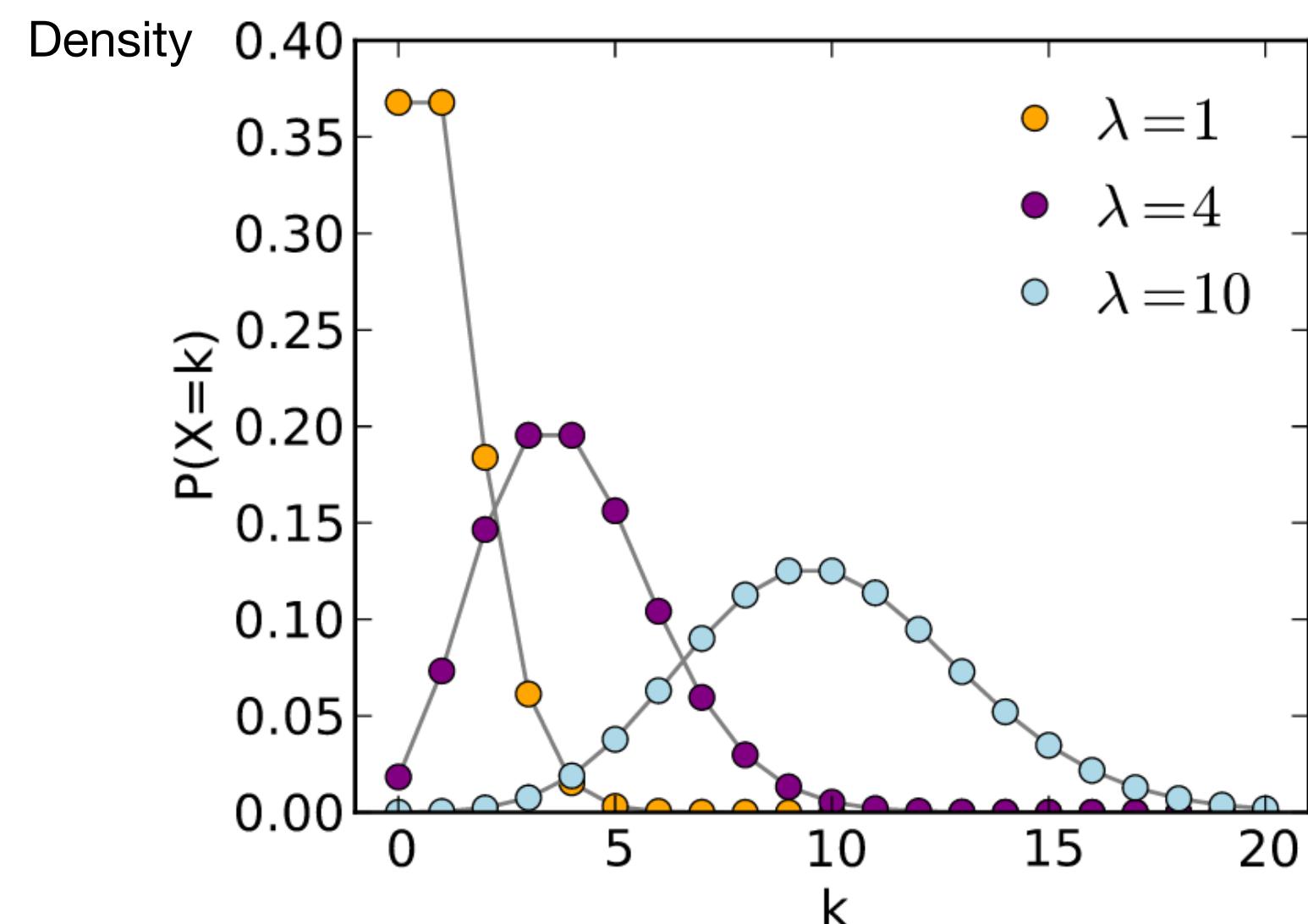
A **discrete probability distribution** that expresses the probability of a given number of events occurring in a fixed interval of time and/or space (with a known average rate and **independently** of the time since the last event)

$$P_{\lambda}(X=k) = \frac{\lambda^k}{k! * e^{-\lambda}}$$

$\lambda$  = mean = expected value = variance

$\lambda = \frac{\text{total number of events (k)}}{\text{number of units (n) in the data}}$

=  $\frac{\text{Read length (nt)} * \text{Total read number}}{\text{Effective genome length (nt)}}$



# Significance of enrichment

- ▶ MACS estimates  $\lambda$ , the expected number of reads, from the control to determine the significance level
- ▶ The probability distribution function is given by

$$F_\lambda(n) = P(X \leq n) = \sum P_\lambda(k) = e^{-\lambda} \sum \frac{\lambda^k}{k!}$$

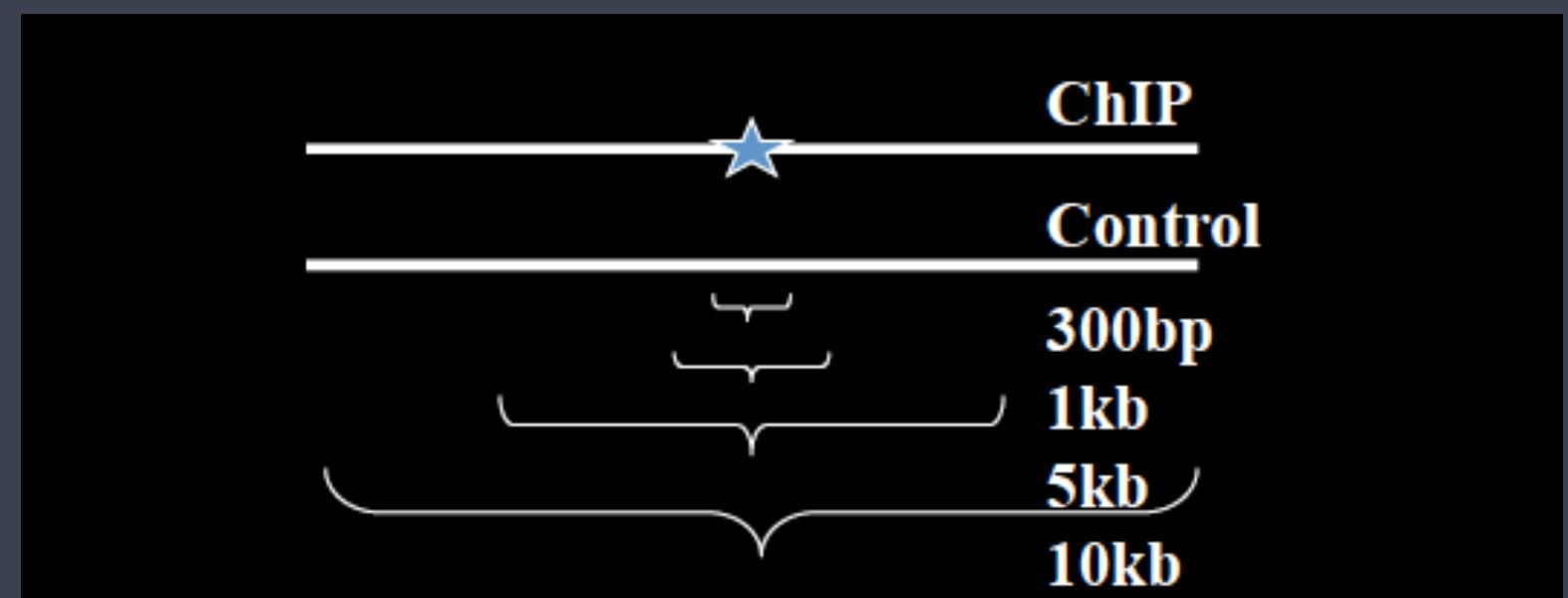
and the probability of observing more than  $n$  reads is  $1 - F_\lambda(n)$

# Background estimation

- ▶ Frequently observe more variance in the data than assumed by the Poisson distribution
  - > Local chromatin structure, PCR, sequencing bias, CNVs leads to false positive peaks
- ▶ Need to fit a distribution to smaller regions on the genome using sliding/discrete windows
- ▶ MACS uses a dynamic  $\lambda_{\text{local}}$ , determined as the maximum value of  $\lambda$  from the background and 1 kb, 5 kb and 10 kb windows centered at the peak location in the control sample

$$\lambda_{\text{BG}} = \text{total tags} / \text{genome size}$$

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



# Library scaling

- ▶ When ChIP-seq and control samples are sequenced at different depths, MACS either linearly scales down the larger sample (default behavior) or scales up the smaller sample.
- ▶ If the total number of reads in the control sample is greater than the number of reads obtained from ChIP-seq by a factor of  $r$  ( $r > 1$ ), then when calculating the P value  $\lambda_{local}$  will be divided by  $r$  by default.

# Effective genome length (Mappability)

- ▶ Not possible to unambiguously assign reads to all genomic regions
- ▶ ‘Mappability’ or uniqueness influences the average mapped depth
- ▶ Mappability improves with increased read length

**Table 1.** Proportions of unique start sites for nucleotide-space short tag alignments

Species	25 (1) (%)	30 (1) (%)	35 (1) (%)	50 (2) (%)	60 (3) (%)	75 (4) (%)	90 (5) (%)
<i>Homo sapiens</i> <sup>a</sup>	66.0	70.9	74.1	76.9	77.5	79.3	80.8
<i>Mus musculus</i> <sup>b</sup>	69.9	74.4	77.1	79.1	79.4	80.7	81.7
<i>Caenorhabditis elegans</i> <sup>c</sup>	85.3	87.7	89.0	89.8	89.9	90.6	91.1
<i>Drosophila melanogaster</i> <sup>d</sup>	67.5	68.4	69.0	69.2	69.2	69.5	69.8

Columns shown are length of tag matched; numbers in parentheses represent the number of mismatches allowed.

<sup>a</sup>Build hg19.

<sup>b</sup>Build mm9.

<sup>c</sup>Build ce6.

<sup>d</sup>Build dm3.

Koehler et al, Bioinformatics (2011)

# Estimation of the false discovery rate (FDR)

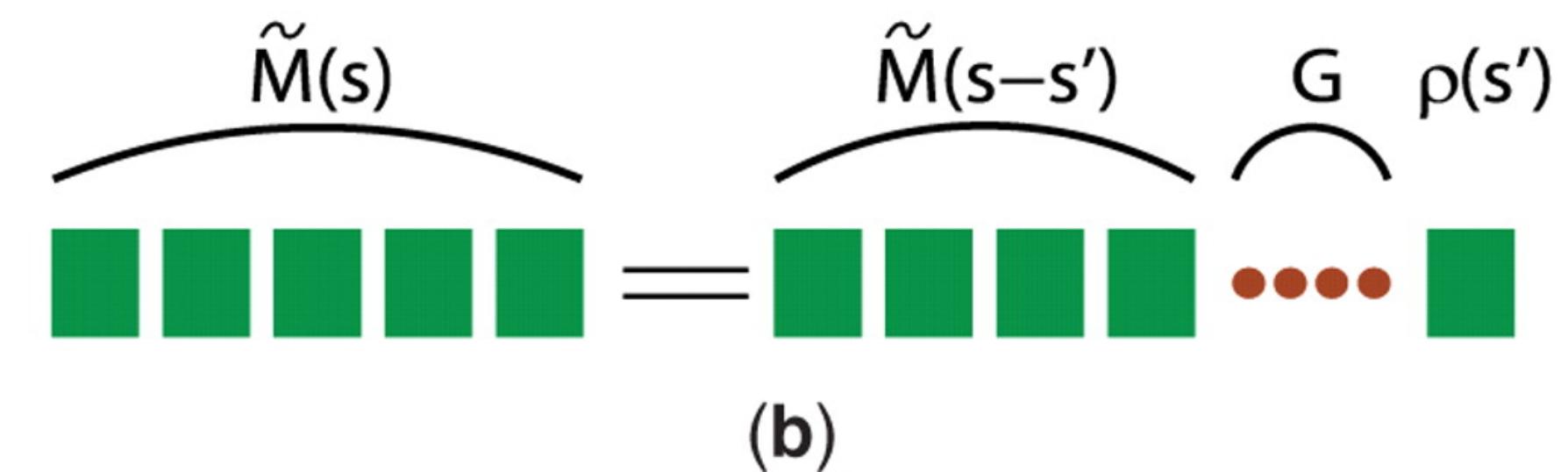
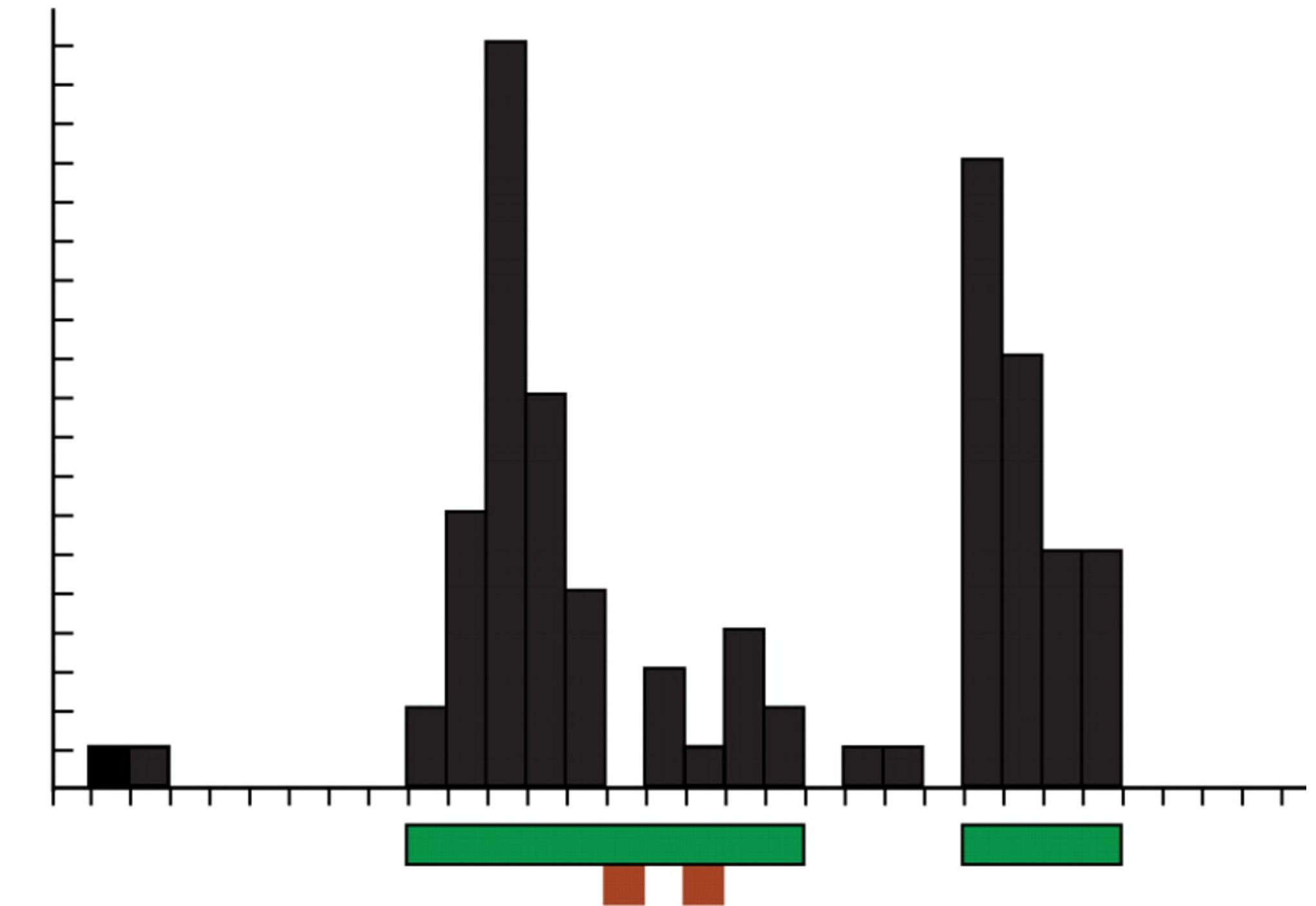
- ▶ Expected proportion of false discoveries among total rejections of the null hypothesis (i.e. how many false positive peaks of the total set of peaks)
- ▶ In MACSv1.4, FDR determined empirically by exchanging ChIP and control
- ▶ In MACSv2, p-values are calculated at every basepair in the genome and then corrected for multiple comparison using the Benjamini-Hochberg correction
- ▶ P-values and FDR are affected by sequencing depth with greater sequencing depth leading to lower p-values and FDRs

# SICER

Hstose modifications tend to cluster to form domains

Identifies spatial clusters of signals unlikely to appear by chance

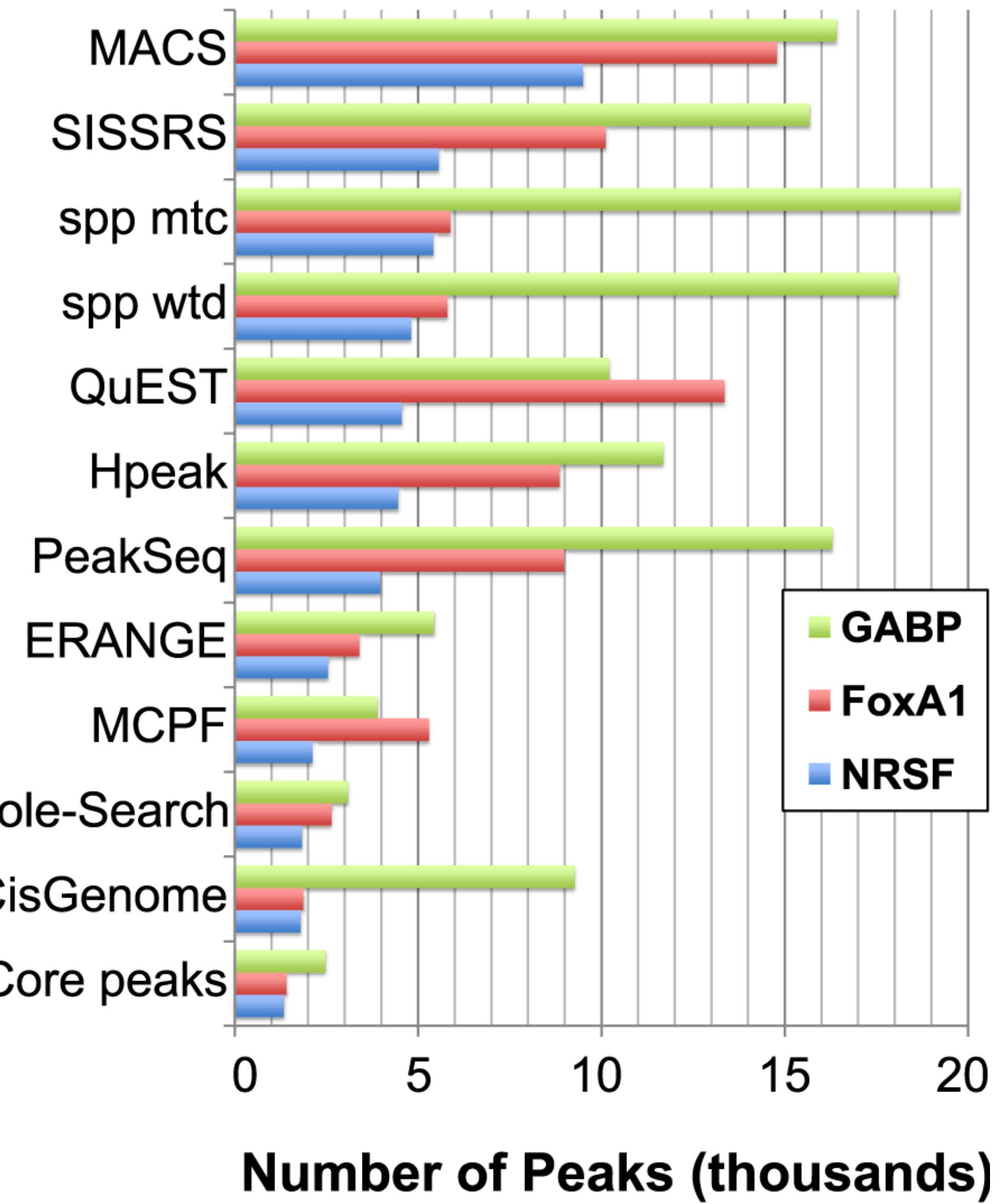
Pools together enrichment information from neighboring nucleosomes to increase sensitivity and specificity.



# Peak callers agree on the strongest signals

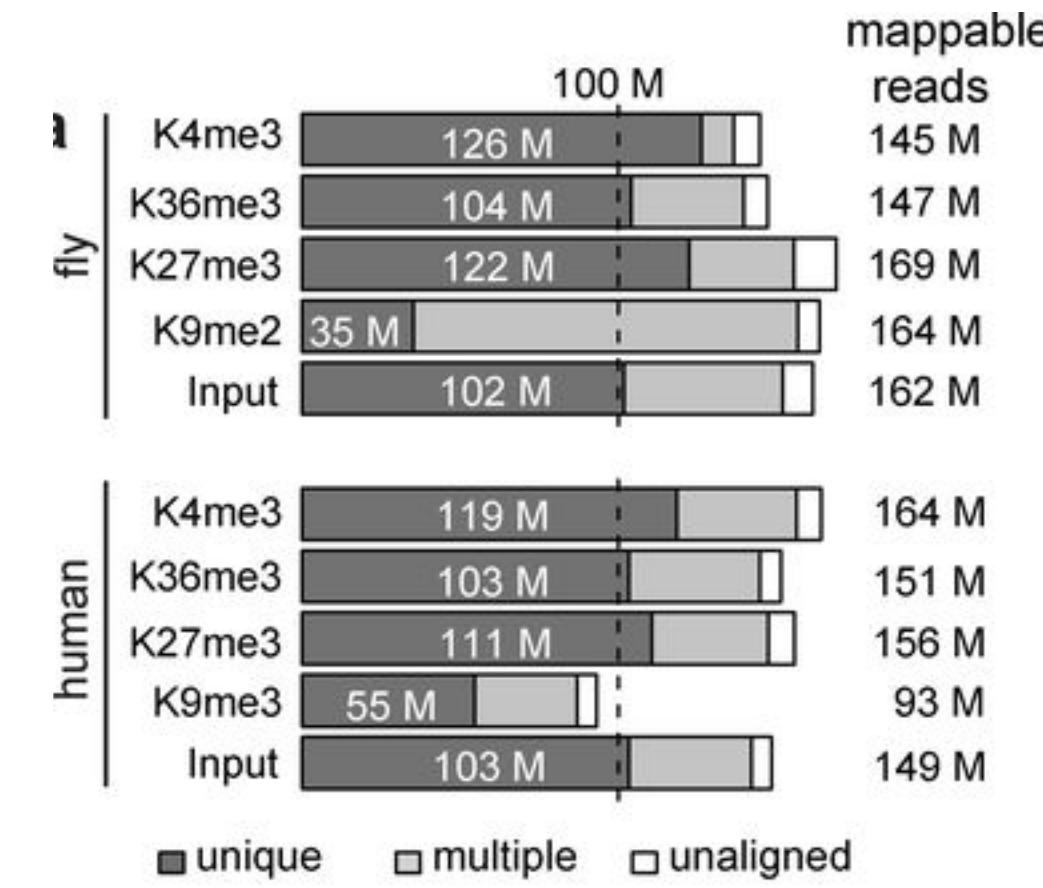
Wilbanks EG, Facciotti MT (2010) Evaluation of Algorithm Performance in ChIP-Seq Peak Detection. PLoS ONE 5(7): e11471.

Peak calling program

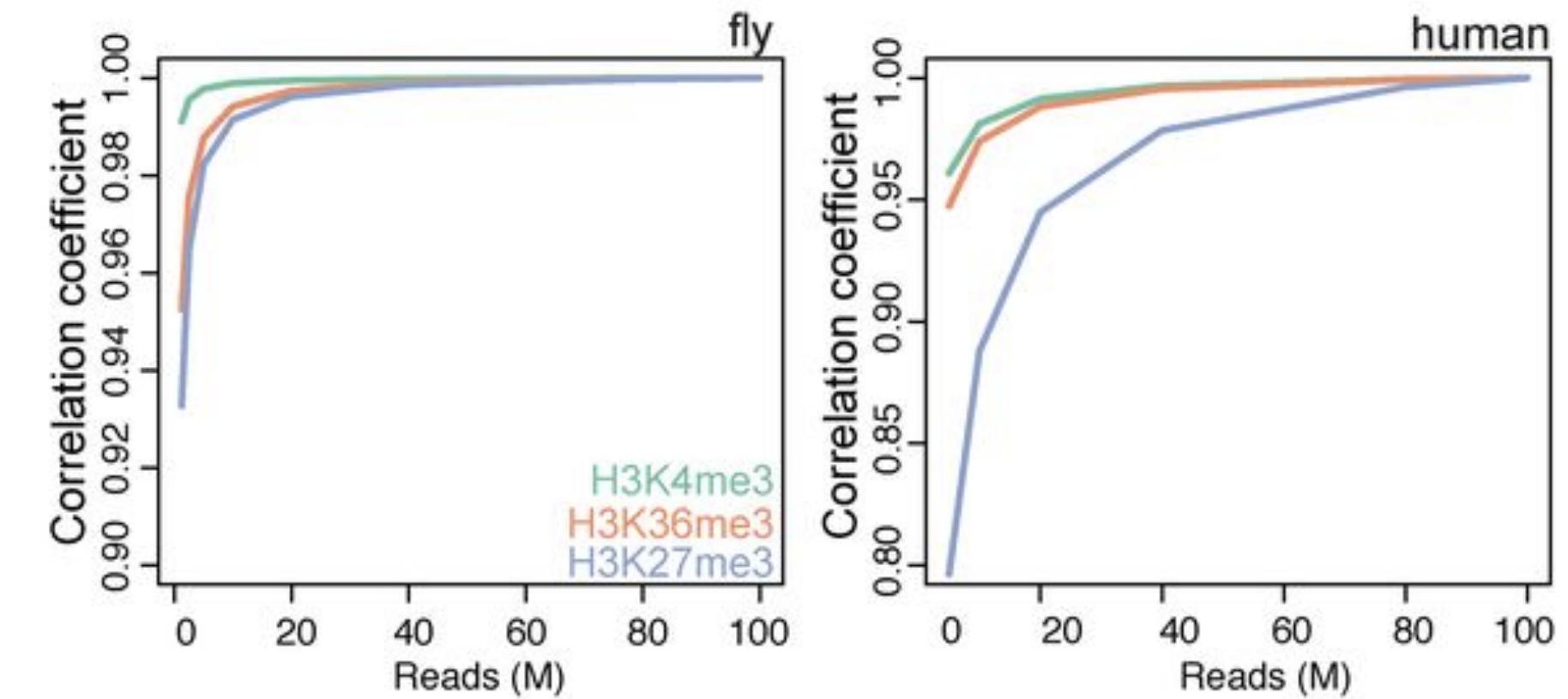


# Effects of sequencing depth and profile type

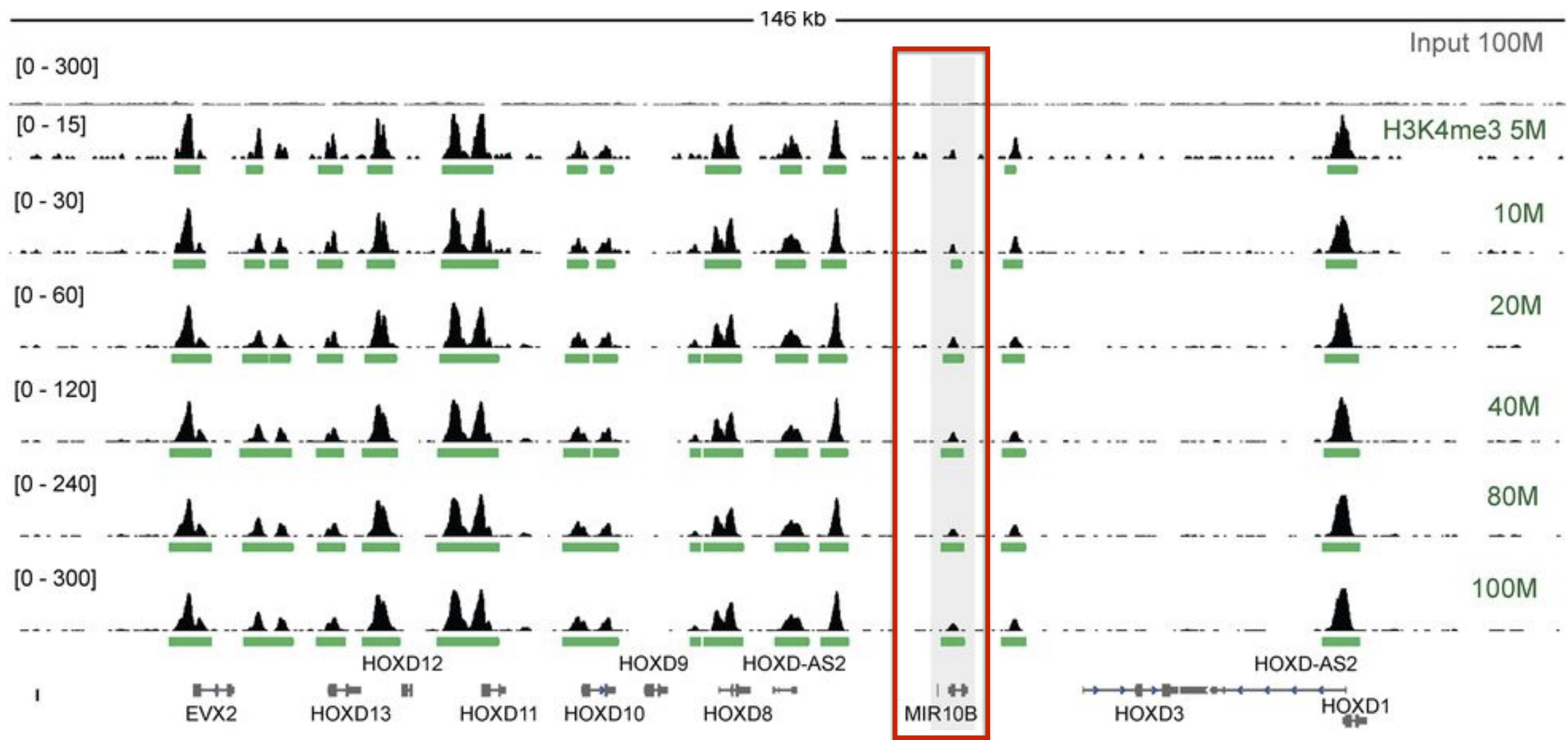
Jung et al., NAR (2014)



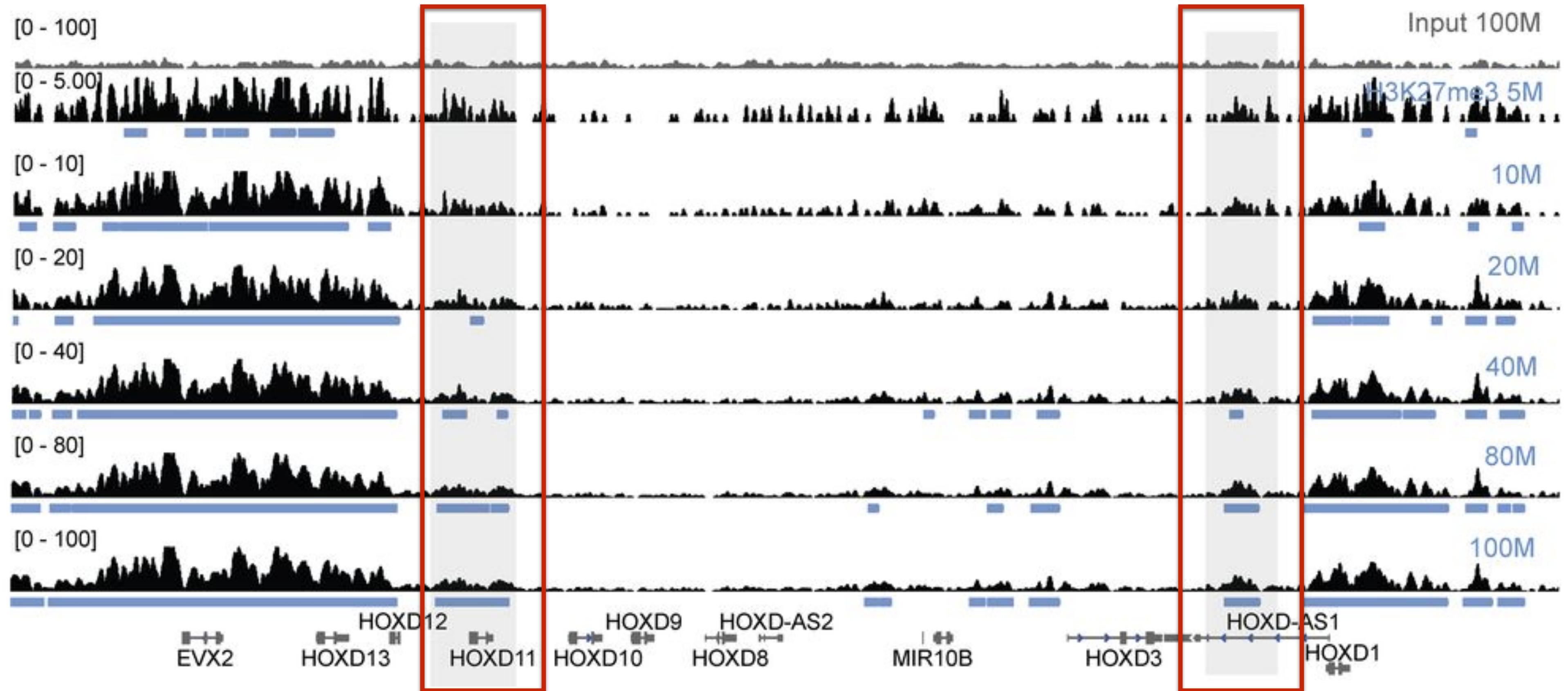
Number of uniquely aligned reads (dark gray), multiply aligned reads (light gray) and unaligned reads (white)



Genome-wide Pearson correlation coefficients between tag density profiles from the 100 million reads and those from subsampled data in fly (left) and human (right).



Impact of sequencing depth: H3K4me3

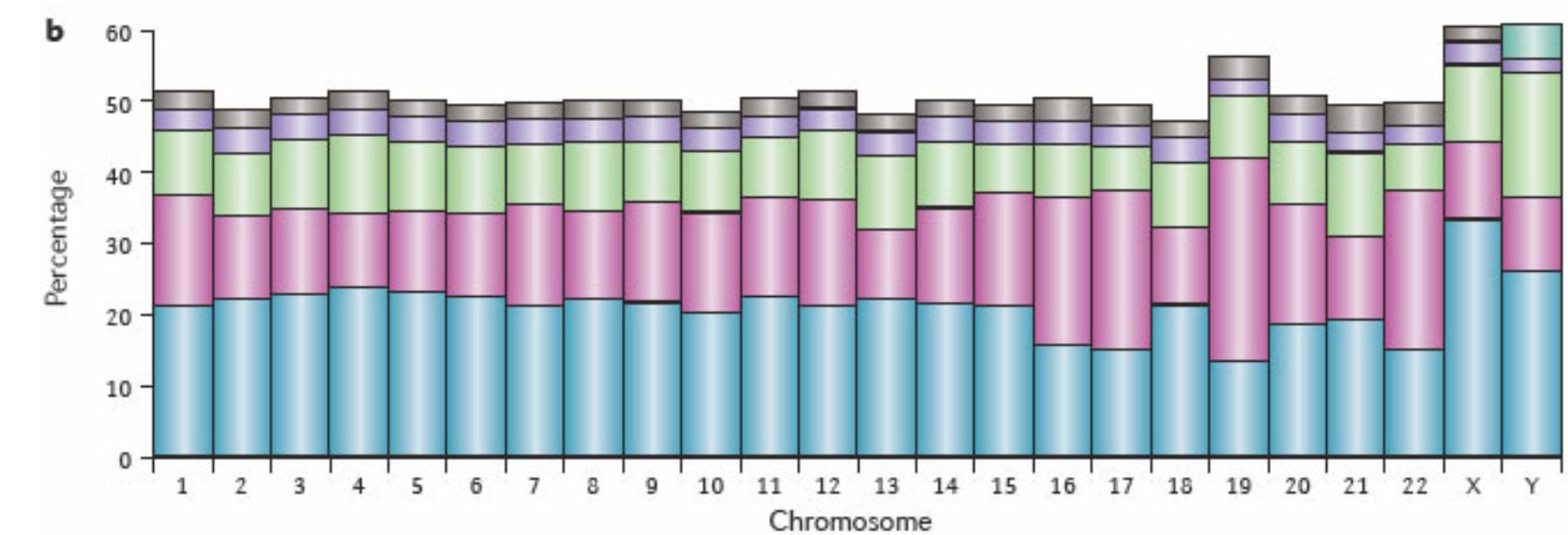


Impact of sequencing depth: H3K27me3

# ENCODE Blacklisted Regions

- ▶ Regions with anomalous, unstructured, high read counts in NGS experiments
- ▶ High ratio of multi-mapping to unique mapping reads
- ▶ Overlap repeat elements (satellites, centromeres, telomeres)
- ▶ Recommended to use this blacklist to filter regions

Repeat class	Repeat type	Number (hg19)	Cvg	Length (bp)
Minisatellite, microsatellite or satellite	Tandem	426,918	3%	2–100
SINE	Interspersed	1,797,575	15%	100–300
DNA transposon	Interspersed	463,776	3%	200–2,000
LTR retrotransposon	Interspersed	718,125	9%	200–5,000
LINE	Interspersed	1,506,845	21%	500–8,000
rDNA (16S, 18S, 5.8S and 28S)	Tandem	698	0.01%	2,000–43,000
Segmental duplications and other classes	Tandem or interspersed	2,270	0.20%	1,000–100,000



Nature Reviews | Genetics

Treangen, T.J. & Salzberg, S.L. (2012) Repetitive DNA and next-generation sequencing: computational challenges and solutions. Nature Reviews Genetics 13, 36–46.

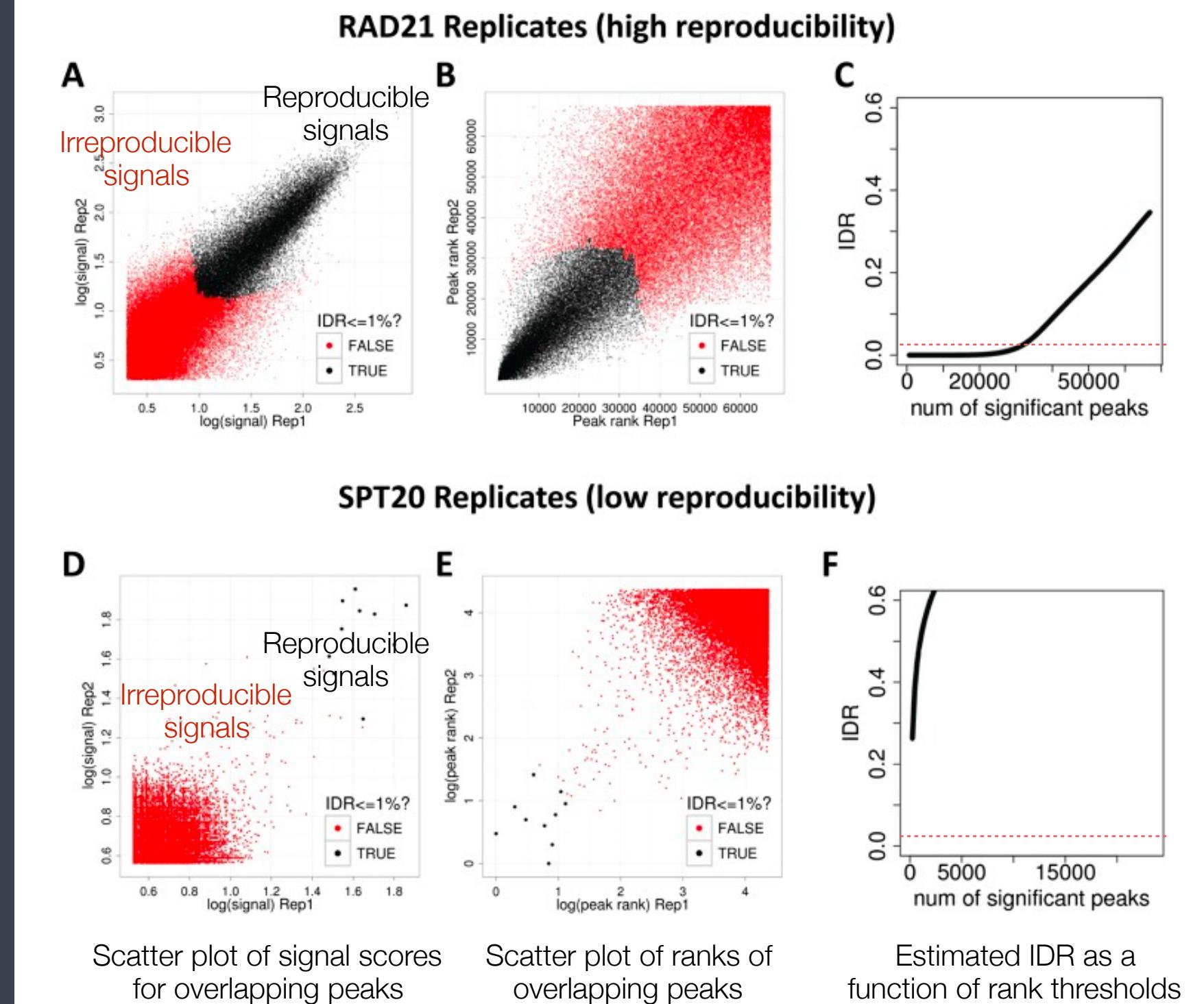
# Handling Replicates

- ▶ Technical replicates are usually merged before peak calling
- ▶ Biological replicates are analyzed separately and compared at the peak call level
- ▶ Historical ENCODE guidelines:
  - > either 80% of the top 40% of the targets identified from one replicate using an acceptable scoring method should overlap the list of targets from the other replicate, OR
  - > target lists scored using all available reads from each replicate should share more than 75% of targets in common
- ▶ Better approach: IDR

# Consistency between replicates

- Irreproducible discovery rate (IDR)

- ▶ Compares a pair of ranked lists of peaks for consistency over the replicates to separate signal from noise.
- ▶ The most significant peaks (i.e. genuine signals) are expected to have high consistency between replicates.
- ▶ Peaks with low significance (i.e. noise) are expected to have low consistency.



**Main**[Home](#)[C.V.](#)[Publications](#)[News](#)[Positions](#)[Contact](#)[Sitemap](#)**Research**[Lab Members](#)[Projects](#)[Tutorials](#)[Datasets](#)[Code](#)[Lab Photos](#)[Interesting Papers](#)[Conferences](#)[Annals](#)[Projects >](#)

# (2012) ENCODE: TF ChIP-seq peak calling using the Irreproducibility Discovery Rate (IDR) framework

**Contents**[1 Mailing List](#)[2 Summary](#)[3 Peak callers tested with IDR](#)[4 Intuitive Explanation of IDR and IDR plots](#)[5 Code for IDR Analysis](#)[5.1 IDR CODE README](#)[6 IDR PIPELINE](#)[6.1 CALL PEAKS ON INDIVIDUAL REPLICATES](#)[6.2 CALL PEAKS ON POOLED REPLICATES](#)[6.3 FOR SELF-CONSISTENCY ANALYSIS CALL PEAKS ON PSEUDOREPLICATES OF INDIVIDUAL REPLICATES](#)[6.4 CREATE PSEUDOREPLICATES OF POOLED DATA AND CALL PEAKS](#)[6.5 INPUT TO IDR ANALYSIS](#)[6.6 IDR ANALYSIS ON ORIGINAL REPLICATES](#)[6.7 IDR ANALYSIS ON SELF-PSEUDOREPLICATES](#)[6.8 IDR ANALYSIS ON POOLED-PSEUDOREPLICATES](#)[6.9 GETTING THRESHOLDS TO TRUNCATE PEAK LISTS](#)

<https://sites.google.com/site/anshulkundaje/projects/idr>

Latest: <https://github.com/nboley/idr>

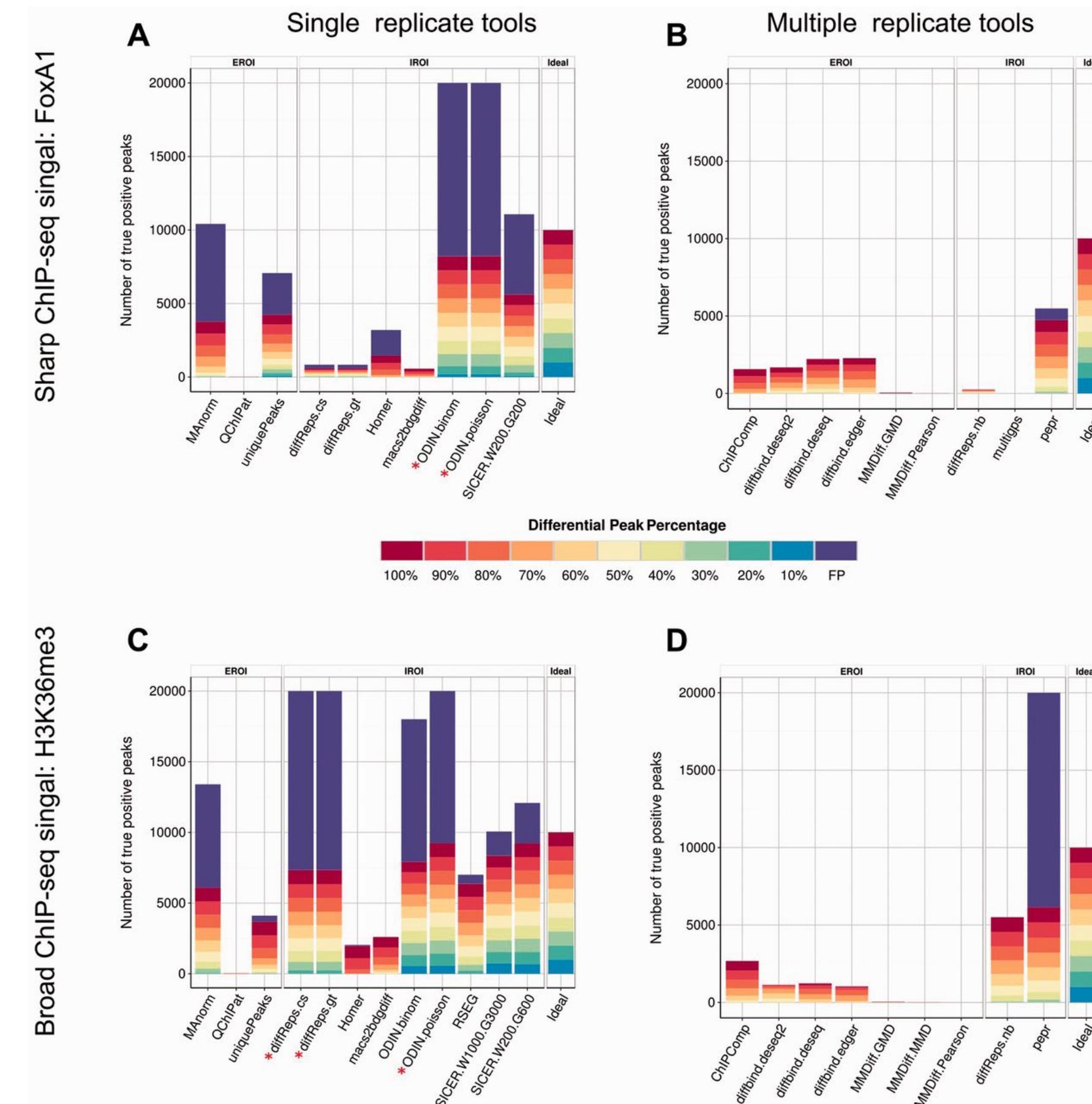
# Differential ChIP-seq analysis

- ▶ Detect differences in ChIP signal between two conditions
- ▶ Challenging compared to RNA-seq and whole genome bisulfite sequencing
- ▶ Search space is not limited to a particular region of the genome
- ▶ Range of the signal is not constrained to finite interval; requires transformation to apply standard statistical tools
- ▶ Amount of noise is considerable, making variations in the signal challenging to detect, especially when these differences are subtle
- ▶ The properties of the enriched regions (in particular their length) differ substantially depending on the protein or epigenetic modification targeted by the immunoprecipitation

# Software for Differential Binding

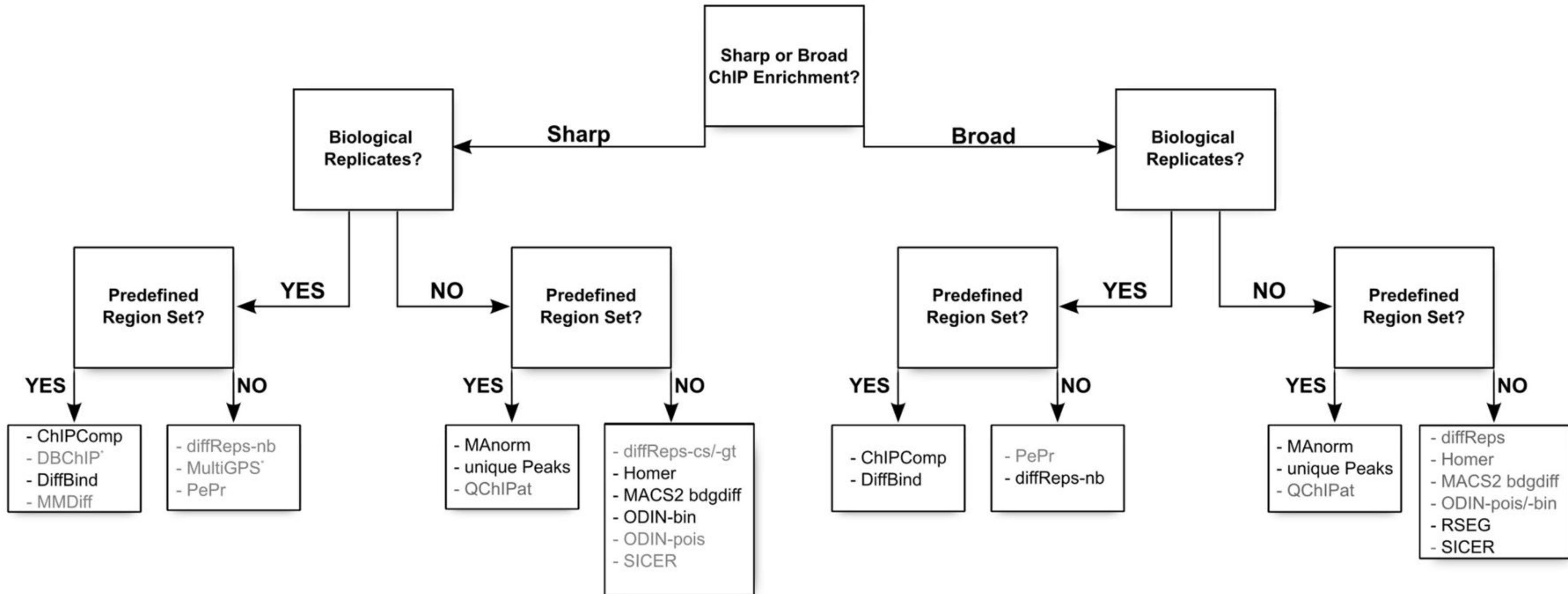
Tool	Peak Calling	Normalization	Statistical Test	Sharp Signal	Broad Signal	Biological Replicates	Significance Measure	PMID
<a href="#"><u>SICER</u></a>	Window based approach, merging of eligible clusters in proximity closer than defined gap size	Library size	Poisson distribution, ChIP norm. counts in possible island against Control norm. counts	n	y	n	FDR	19505939
<a href="#"><u>MACS2</u></a>	Not required	Library size	Computation of log10 likelihood ratios and setting a pre-defined cutoff for following comparisons: Cond1 > Cond2 and Cond1 > Control1 Cond1 < Cond2 and Cond2 > Control2	y	y	n	log10 likelihood ratio	18798982
<a href="#"><u>ODIN</u></a>	Not required	- SES normalization (ChIP / input) combined with input subtraction - Library size normalization (Cond1 / Cond2)	Hidden Markov Model (HMM) with a three state topology Emissions are calculated with a Binomial or a mixture of Poisson distribution	y	y	n	p-value	25371479
<a href="#"><u>RSEG</u></a>	Not required	-	Hidden Markov Model (HMM) with a three state topology NBDiff distribution is used to model read count differences between both conditions	n	y	n	-	21325299
<a href="#"><u>MAnorm</u></a>	Requires peak calling e.g. with MACS	Genome-wide MA plot combined with LOWESS Regression	Bayesian model approach	y	y	n	p-value	22424423
<a href="#"><u>HOMER</u></a>	Window based approach; Peak calling done by HOMER	Library size normalization	Fold-change thresholding combined with a Poisson distribution based enrichment analysis	y	y	n	FDR or p-value	20513432
<a href="#"><u>QChIPat</u></a>	Peak calling possible with BELT, MACS, SISSRs or FindPeaks	1) Nonparametric empirical Bayes correction normalization 2) Quantile normalization 3) Linear normalization	1) Wilcoxon rank sum test 2) Wilcoxon signed rank test	y	y	n	p-value	24564479
<a href="#"><u>diffReps</u></a>	Sliding window approach	Linear normalization	- Without replicates: G-test or Chi-square test - Replicates: exact negative binomial test Generalized linear model with negative Binomial distribution	n	y	y/n	p-value	23762400
<a href="#"><u>DBChip</u></a>	Requires peak calling e.g. with MACS	median ratio strategy (DESeq)	Generalized linear model with negative Binomial distribution	y	n	y/n	FDR	22057161
<a href="#"><u>ChIPComp</u></a>	Requires peak calling e.g. with MACS	Normalization with a Poisson distribution based model	Wald's test followed by probability calculation Using a Bayesian approach	y	n	y	Posterior probability	25682068
<a href="#"><u>MultiGPS</u></a>	Expectation maximization learning scheme		edgeR	y	n	y	p-value	24675637
<a href="#"><u>MMDiff</u></a>	Requires peak calling e.g. with MACS	DESeq	Kernel-based non-parametric test	y	n	y	p-value	24267901
<a href="#"><u>DiffBind</u></a>	Requires peak calling e.g. with MACS		Differential peak analysis can be performed with: 1) DESeq 2) DESeq2 3) edgeR	y	y	y	p-value or FDR	22217937
<a href="#"><u>PePr</u></a>	Window based approach	Trimmed Mean of M values (TMM) approach (edgeR)	Binomial distribution	y	y	y	p-value	24894502

**Proportion of true and false positives for each tool on the simulated FoxA1 data set (A, B) and H3K36me3 data set (C, D).**



Sebastian Steinhauser et al. Brief Bioinform  
2016;bib.bbv110

## Decision tree indicating the proper choice of tool depending on the data set



Sebastian Steinhauser et al. Brief Bioinform  
2016;bib.bbv110

© The Author 2016. Published by Oxford University Press.

# Ask us questions

[shosui@hsph.harvard.edu](mailto:shosui@hsph.harvard.edu)

[bioinformatics.hms.harvard.edu](mailto:bioinformatics.hms.harvard.edu)

