# Epigenomics overview

# Epigenomics overview

**Common Targets**
- DNA-associated factors
- RNA-associated factors
- histone modifications
- chromatin accessibility
- nucleosome positioning
- DNA modifications

**Common Methods**
- immunoprecipitation
- transposase insertion
- nuclease digestion
- conversion resistance

**Mapping**
- standard or splice-aware (RNA)
- standard
- standard
- standard w/ converted genome

**Tag Location(s)**
- IP
- IP-exo

**Measurement**
- peaks or regions
- and/or more complex analyses
- peaks
- and/or more complex analyses
- peaks
- % of reads un-converted at given position

**Further Goals**
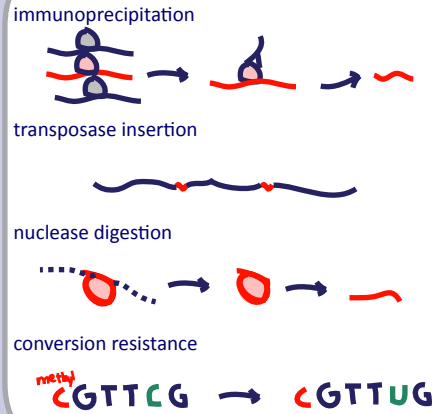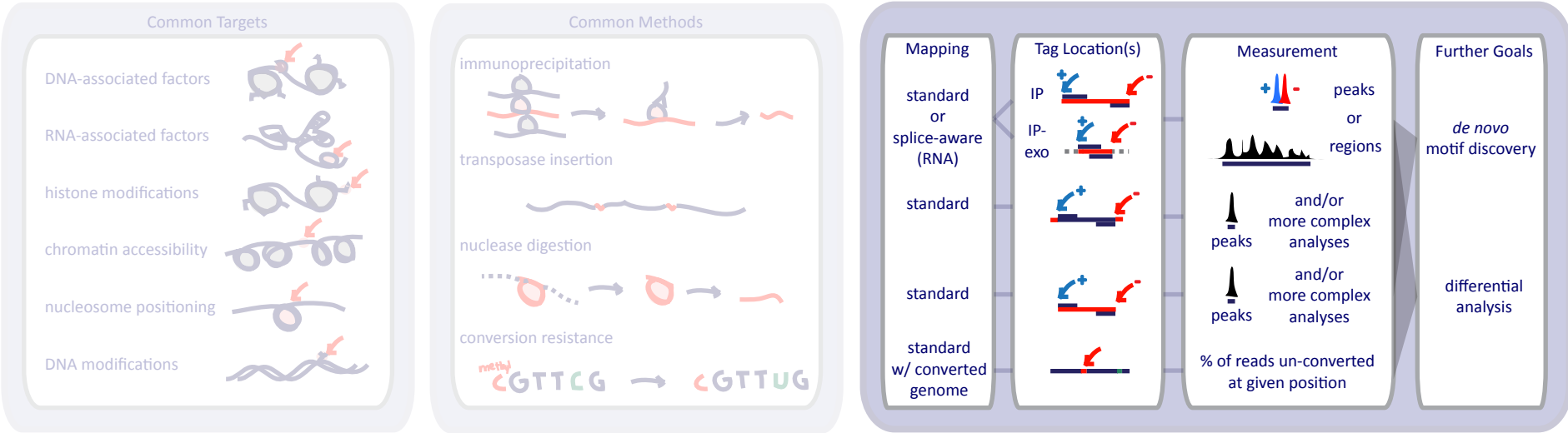- *de novo* motif discovery
- differential analysis

# Epigenomics overview

# Epigenomics overview

## Common Targets

DNA-associated factors

RNA-associated factors

histone modifications

chromatin accessibility

nucleosome positioning

DNA modifications

## Common Methods

immunoprecipitation

transposase insertion

nuclease digestion

conversion resistance

methyl
CGTTCG → CGTTUG

## Mapping

standard
or
splice-aware
(RNA)

standard

standard

standard
w/ converted
genome

## Tag Location(s)

IP

IP-
exo

## Measurement

peaks
or
regions

and/or
more complex
analyses
peaks

and/or
more complex
analyses
peaks

% of reads un-converted
at given position

## Further Goals

*de novo*
motif discovery

differential
analysis

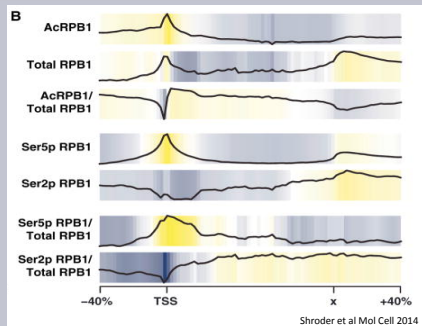# DNA-associated factors

## sequence-specific factors



+ tags
- tags

## chromatin associated



T2

## transcription machinery



B

AcRPB1

Total RPB1

AcRPB1/
Total RPB1

Ser5p RPB1

Ser2p RPB1

Ser5p RPB1/
Total RPB1

Ser2p RPB1/
Total RPB1

−40%    TSS         x      +40%

Shroder et al Mol Cell 2014

## methodology

### immunoprecipitation



### library construction

### sequencing

### mapping to genome reference (BWA or bowtie2)

### identify tag locations



### calculate median insert size (single end)

### determine fragment midpoint (paired end)



### estimate actual binding location per tag



### calculate genomic density of shifted tags
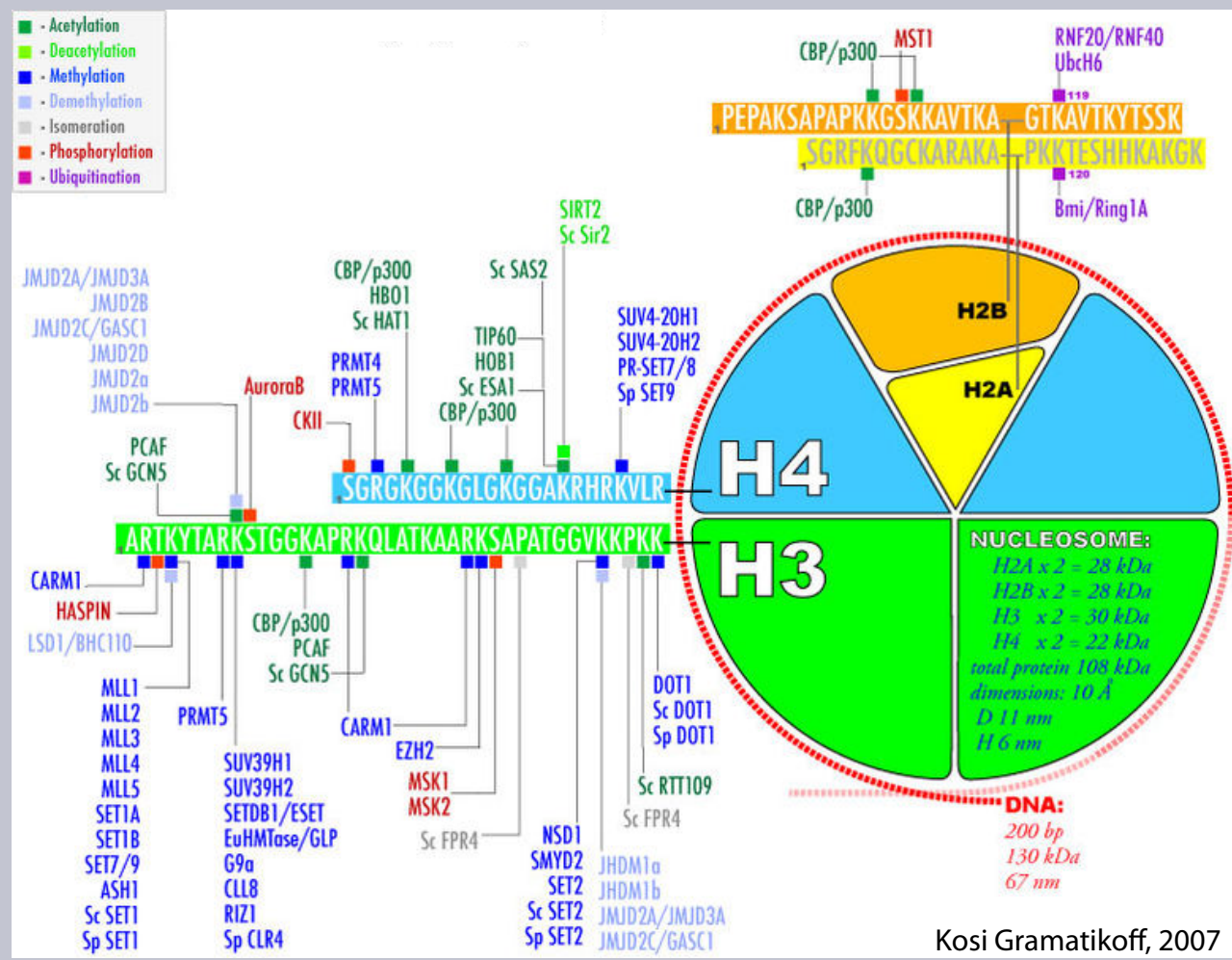


### measure background density



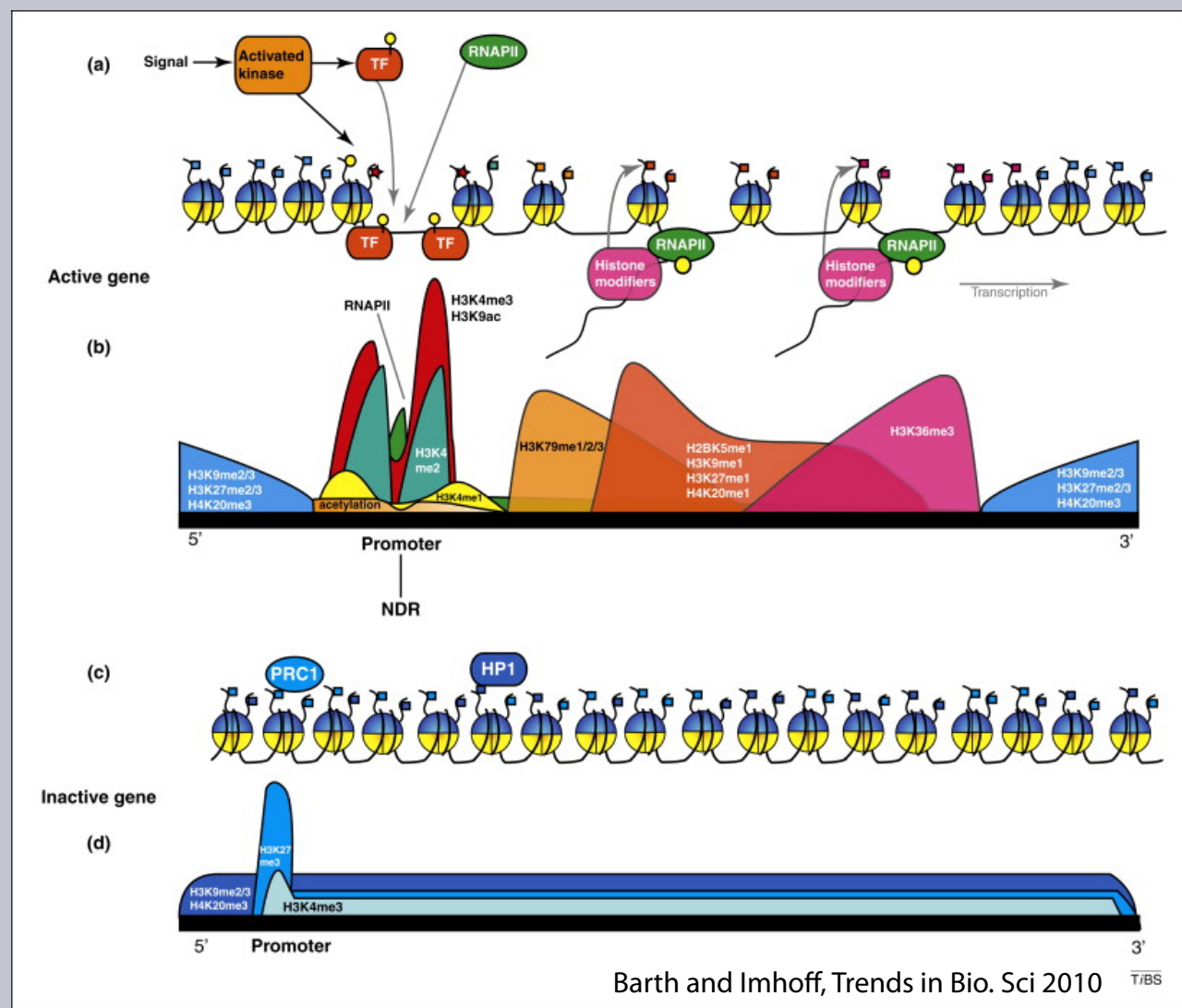### calculate background-normalized binding density



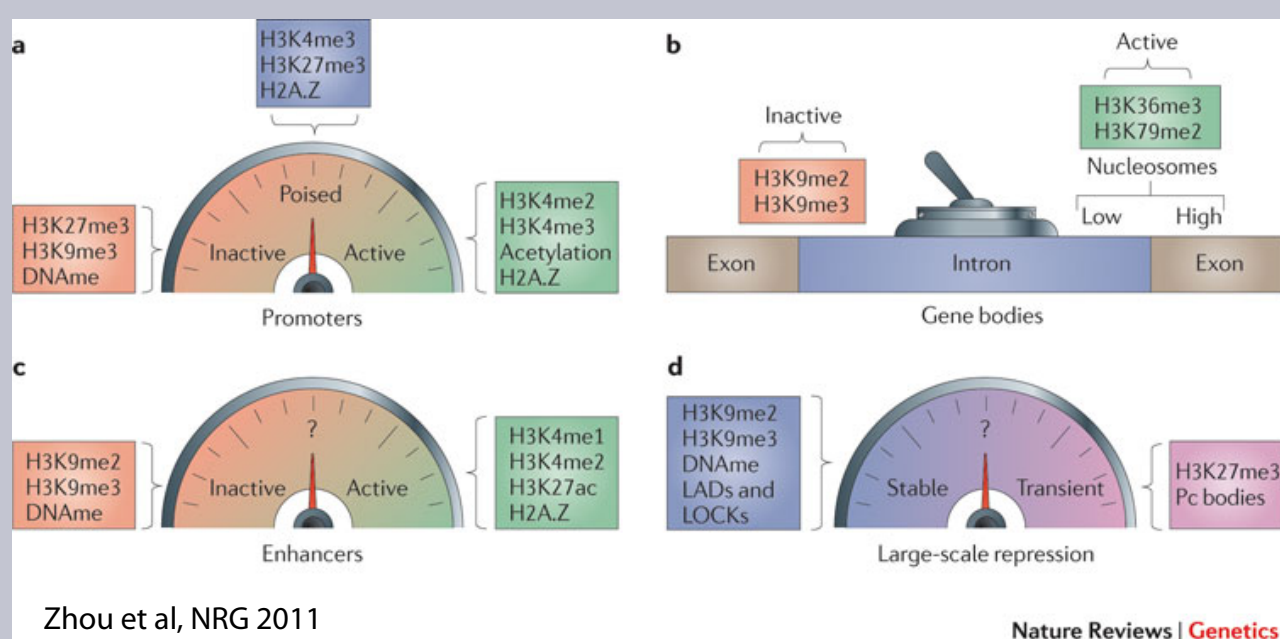### call peaks (punctate marks) or regions (broad marks)

# histone modifications

## catalogue



Kosi Gramatikoff, 2007

## genomic distribution



Barth and Imhoff, Trends in Bio. Sci 2010

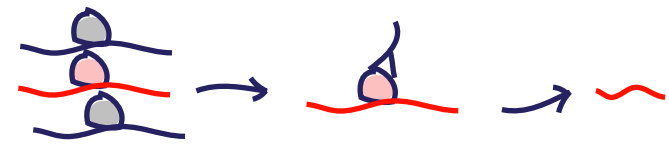## functional association



Zhou et al, NRG 2011

Nature Reviews | Genetics

## methodology
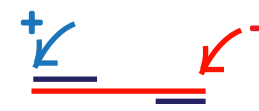
immunoprecipitation



library construction

sequencing

mapping to genome reference (BWA or bowtie2)

identify tag locations



calculate median insert size (single end)

determine fragment midpoint (paired end)



estimate actual binding location per tag



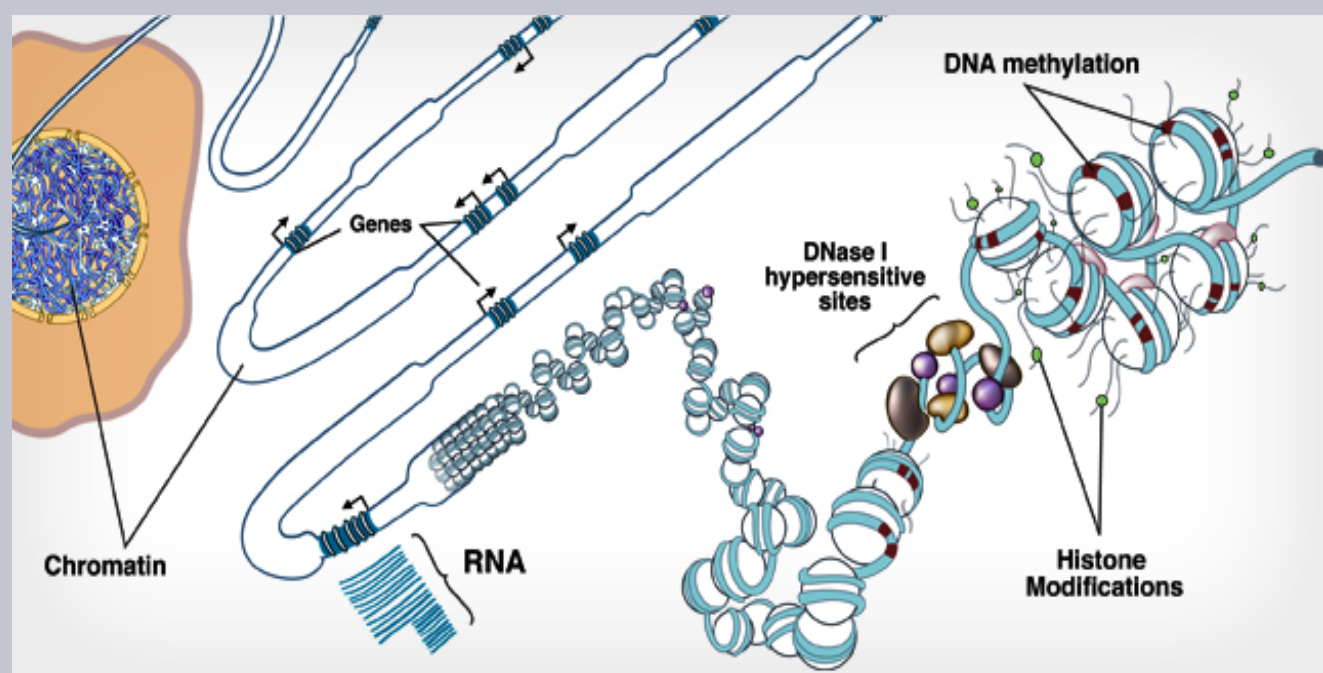calculate genomic density of shifted tags



measure background density



calculate background-normalized binding density



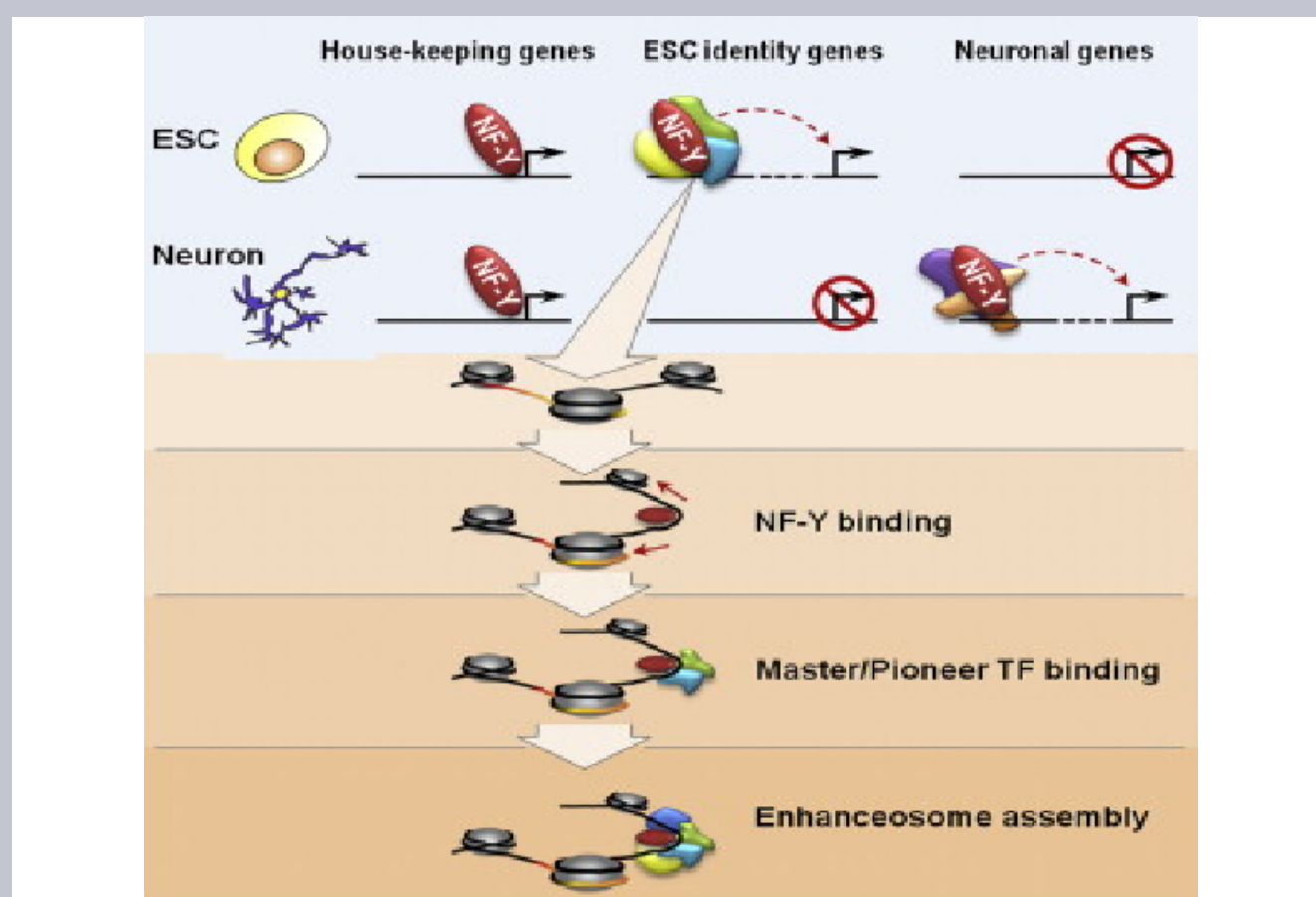call peaks (punctate marks) or regions (broad marks)

# chromatin accessibility

## functional regulatory elements are accessible



roadmapepigenomics.org

## formation of accessible chromatin



Oldfield et al, Mol Cell 2014

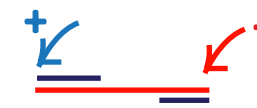## methodology (digestion or transposase)

DNase I or Tn5 release fragments of open chromatin



library construction and sequencing

mapping to genome reference (BWA or bowtie2)

identify cleavage/insertion location
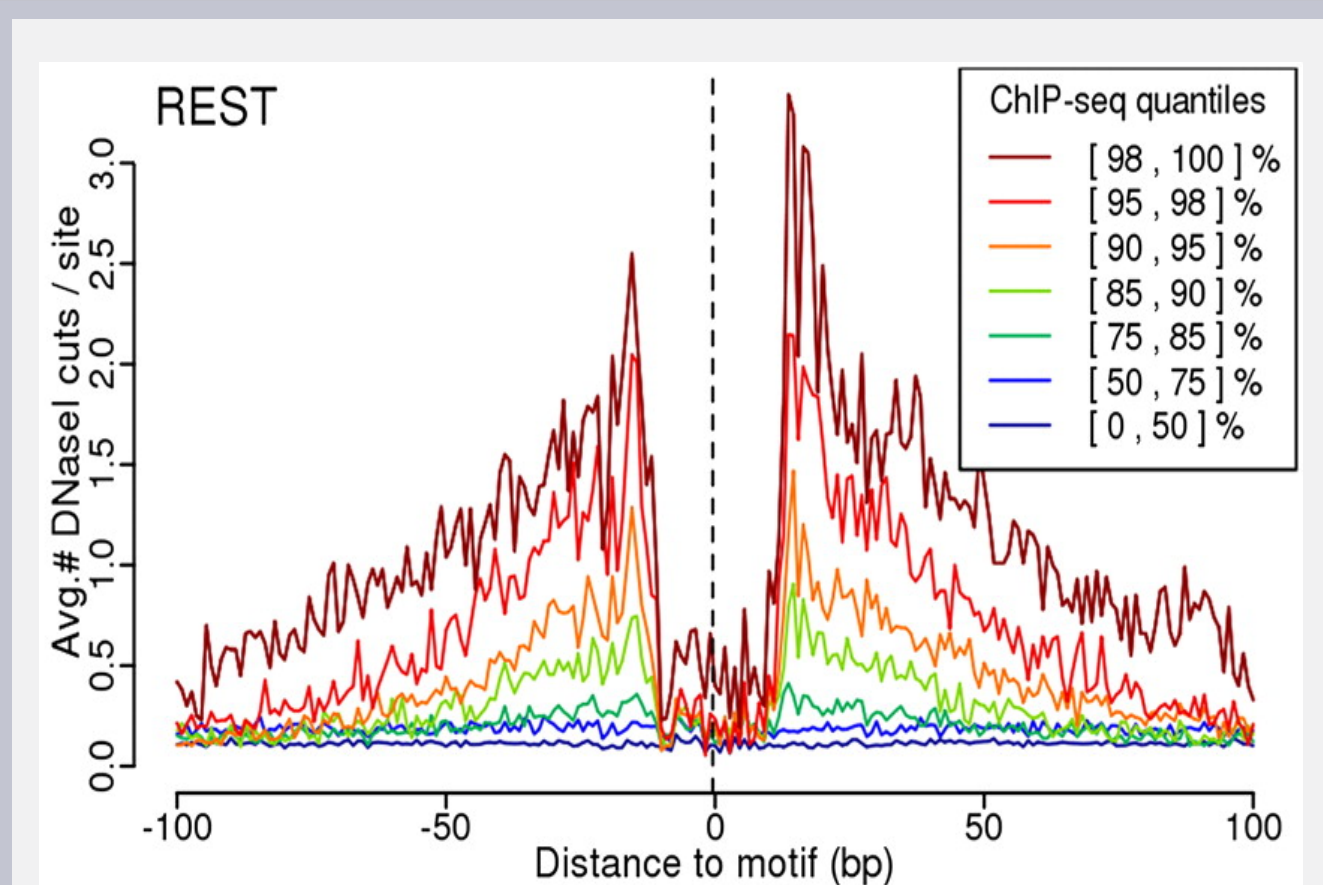


calculate density of cleavage/insertion events



call regions of accessibility
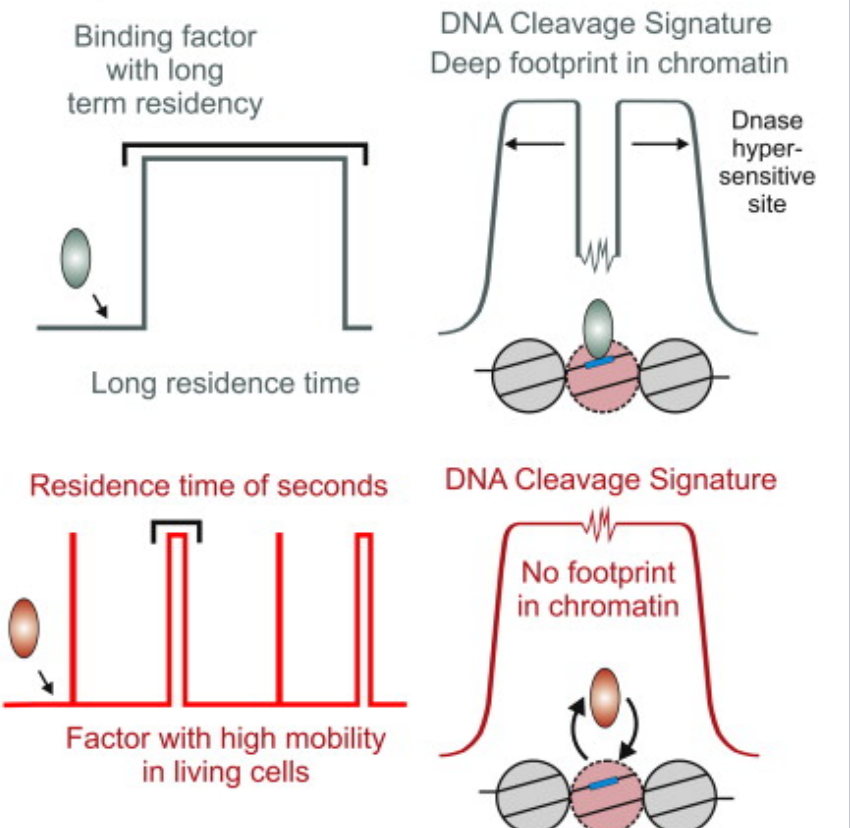


call footprints within accessible chromatin
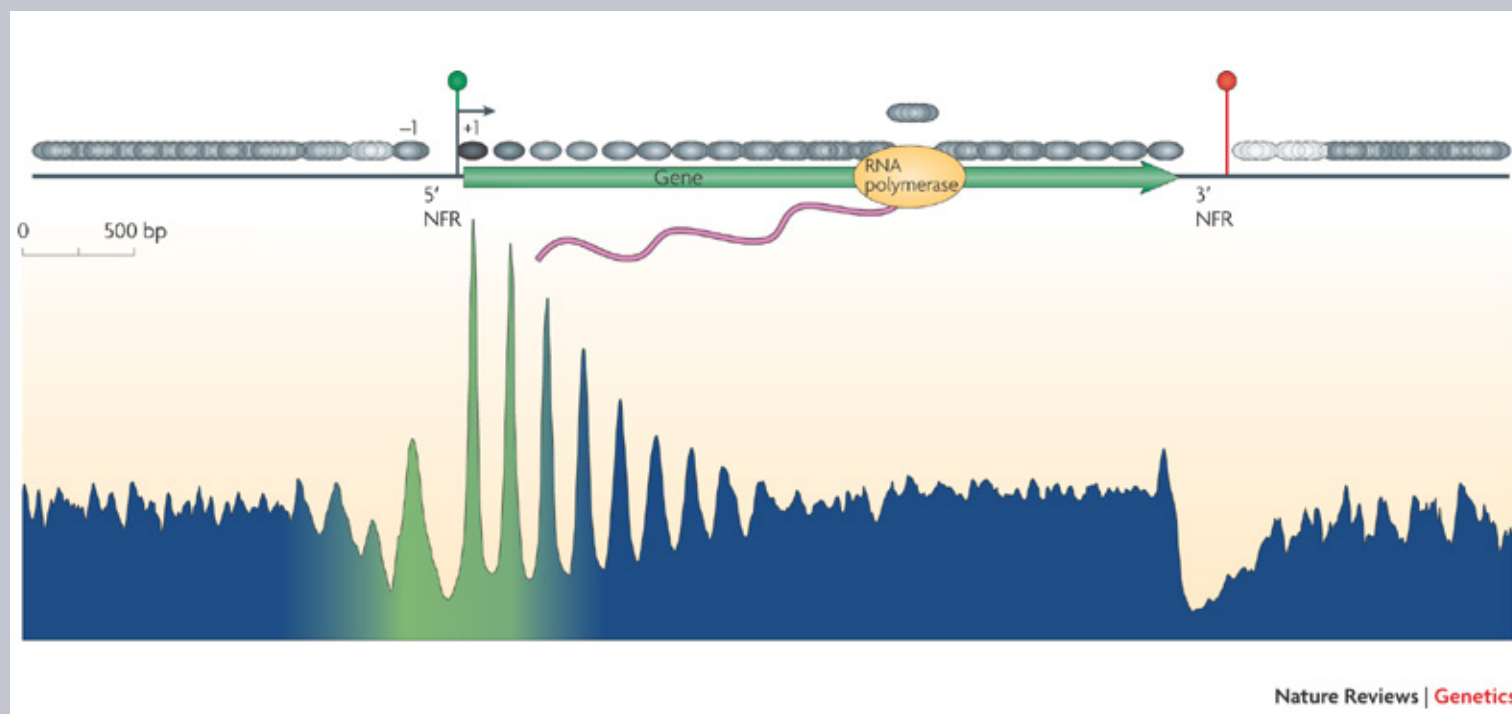


## footprinting



REST

ChIP-seq quantiles

[ 98 , 100 ] %
[ 95 , 98 ] %
[ 90 , 95 ] %
[ 85 , 90 ] %
[ 75 , 85 ] %
[ 50 , 75 ] %
[ 0 , 50 ] %

Avg.# DNaseI cuts / site

Distance to motif (bp)

Pique-Regi et al Genome Research 2011



Footprints reflect residence time in chromatin

Binding factor with long term residency

DNA Cleavage Signature Deep footprint in chromatin

Dnase hyper-sensitive site

Long residence time

Residence time of seconds

DNA Cleavage Signature

Factor with high mobility in living cells

No footprint in chromatin

Sung et al, Mol Cell 2014

## nucleosome positioning during transcription



Nature Reviews | Genetics
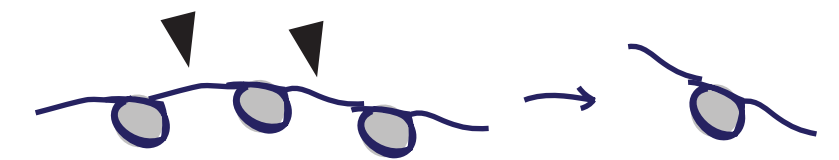
Jiang et al NRG 2009

## determinants of positioning stability



**A**

Highly positioned nucleosomes          "Fuzzy" nucleosomes

Poly(dA:dT) tract

−ATTGAGCTGCAATCTGGAATAACAGCCAGATAAGGAGCTACAGTACC−

Nucleosome favoring sequence:
A/T dinucleotide every 10 basepairs
G/C dinucleotide every 10 basepairs, in antiphase with A/T dinucleotides

**B**

ATP    ADP + P$_i$          ATP    ADP + P$_i$

ATP-consuming           ATP-consuming
chromatin remodeling factor    chromatin remodeling factor

HAT    HDAC

Acetyl

## methodology (MNase digestion)

MNase I digests linker DNA, releasing multisomes

purify mononucleosome-bound DNA

library construction and sequencing

mapping to genome reference (BWA or bowtie2)

identify cleavage location

calculate nucleosomal density

call positioned nucleosomes

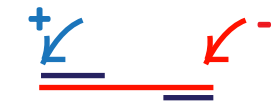## methodology (transposon insertion)

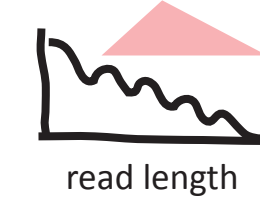Tn5 integrates into accessible chromatin, and releases multisomes

library construction and sequencing

mapping to genome reference (BWA or bowtie2)

identify cleavage location

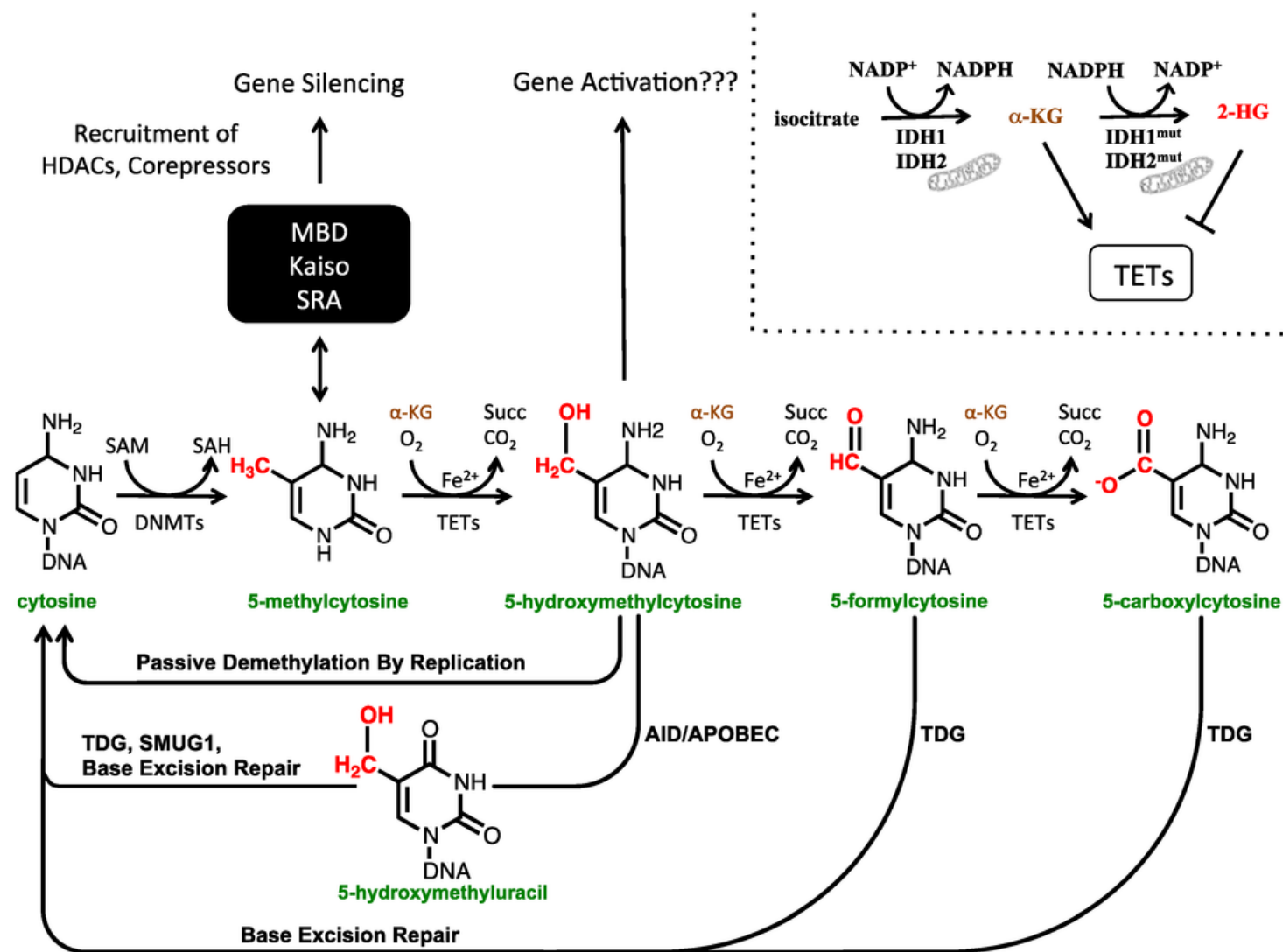identify reads that span at least one nucleosome

read length

calculate density of tags / number of insertion events
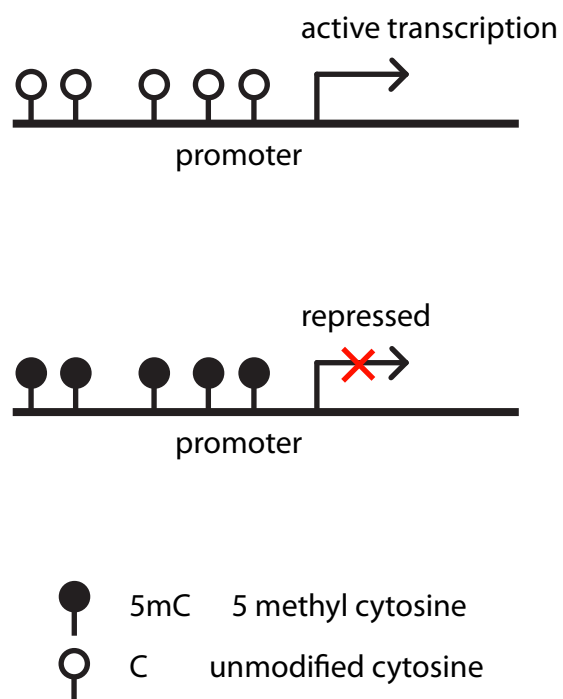
call positioned nucleosomes

# DNA modifications

## catalogue of common base modifications (in mammals)



Mariani et al, Cancers 2013

## 5mC repression of gene expression



- ● 5mC   5 methyl cytosine
- ○ C   unmodified cytosine

## bisulfite conversion / protection



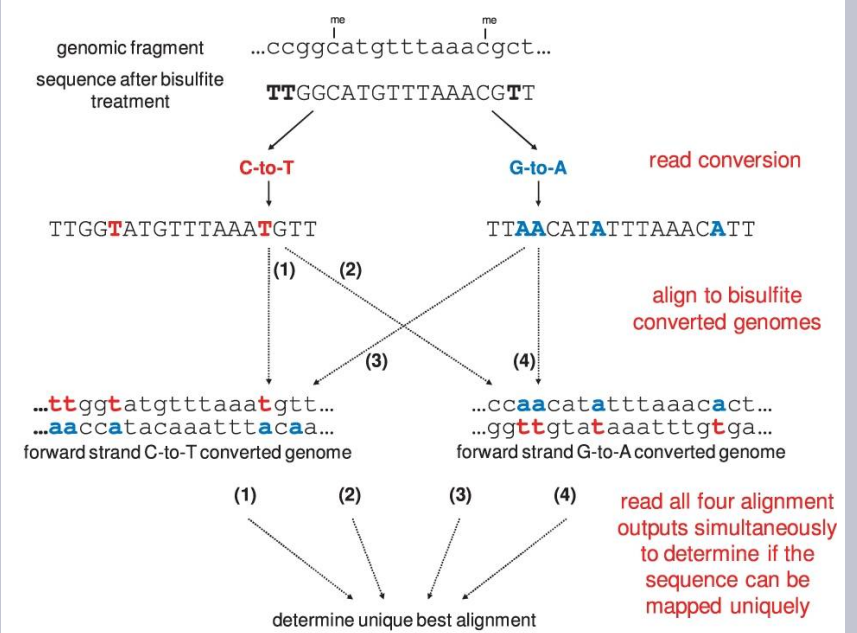bisulfite

PCR

daigenode.com

## methodology

### convert DNA with bisulfite



### fragmentation
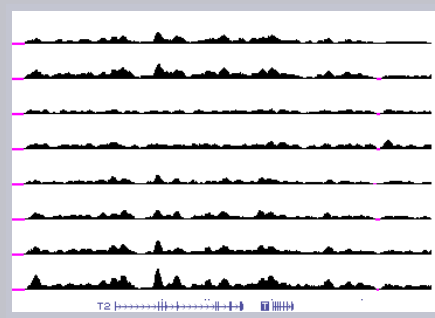
### library construction and sequencing

### mapping



Krueger et al Bioinformatics 2011

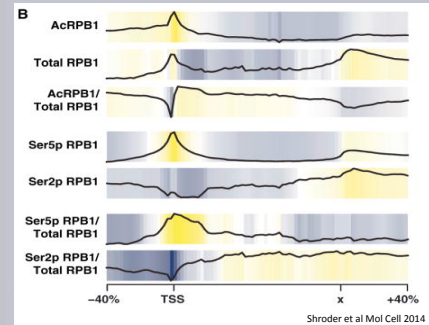### calculate % of tags methylated per genomic position

# DNA-associated factors

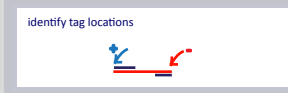## sequence-specific factors



+ tags
- tags

## chromatin associated



T2

## transcription machinery



B
AcRPB1
Total RPB1
AcRPB1/ Total RPB1
Ser5p RPB1
Ser2p RPB1
Ser5p RPB1/ Total RPB1
Ser2p RPB1/ Total RPB1

−40%    TSS    x    +40%

Shroder et al Mol Cell 2014
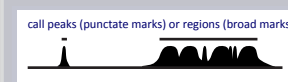
## methodology

### immunoprecipitation



### library construction

### sequencing

### mapping to genome reference (BWA or bowtie2)

### identify tag locations



### calculate median insert size (single end)

### determine fragment midpoint (paired end)



### estimate actual binding location per tag



### calculate genomic density of shifted tags



### measure background density



### calculate background-normalized binding density



### call peaks (punctate marks) or regions (broad marks)

# histone modifications

## catalogue



Kosi Gramatikoff, 2007

## genomic distribution



Barth and Imhoff, Trends in Bio. Sci 2010

## functional association



Zhou et al, NRG 2011

Nature Reviews | Genetics

## methodology

immunoprecipitation



library construction

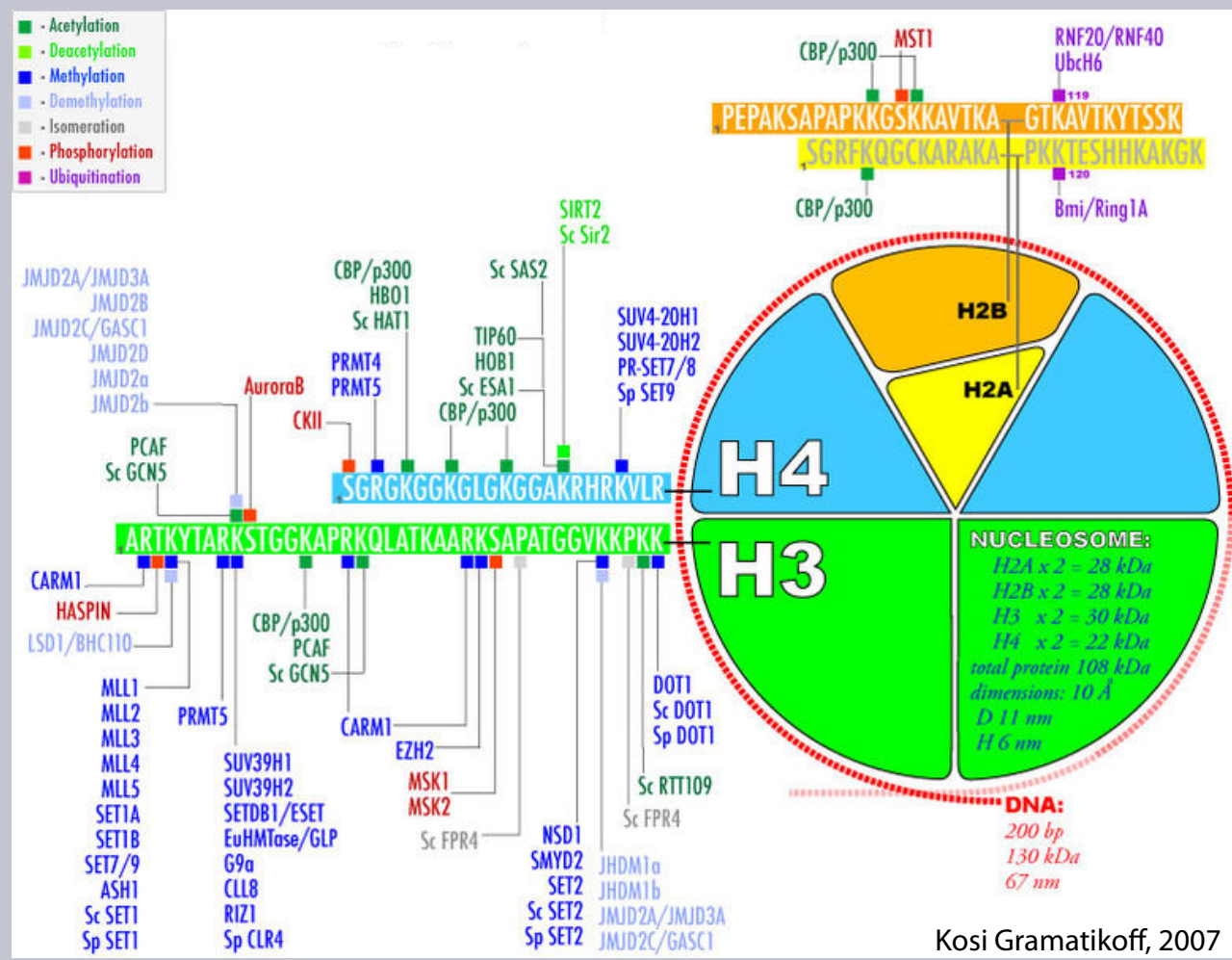sequencing

mapping to genome reference (BWA or bowtie2)

identify tag locations



calculate median insert size (single end)
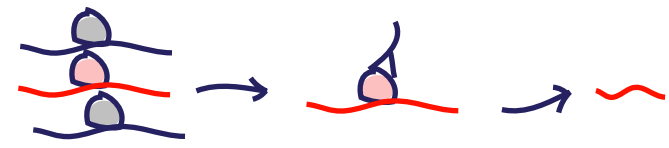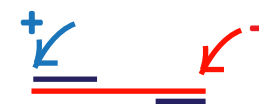
determine fragment midpoint (paired end)



estimate actual binding location per tag



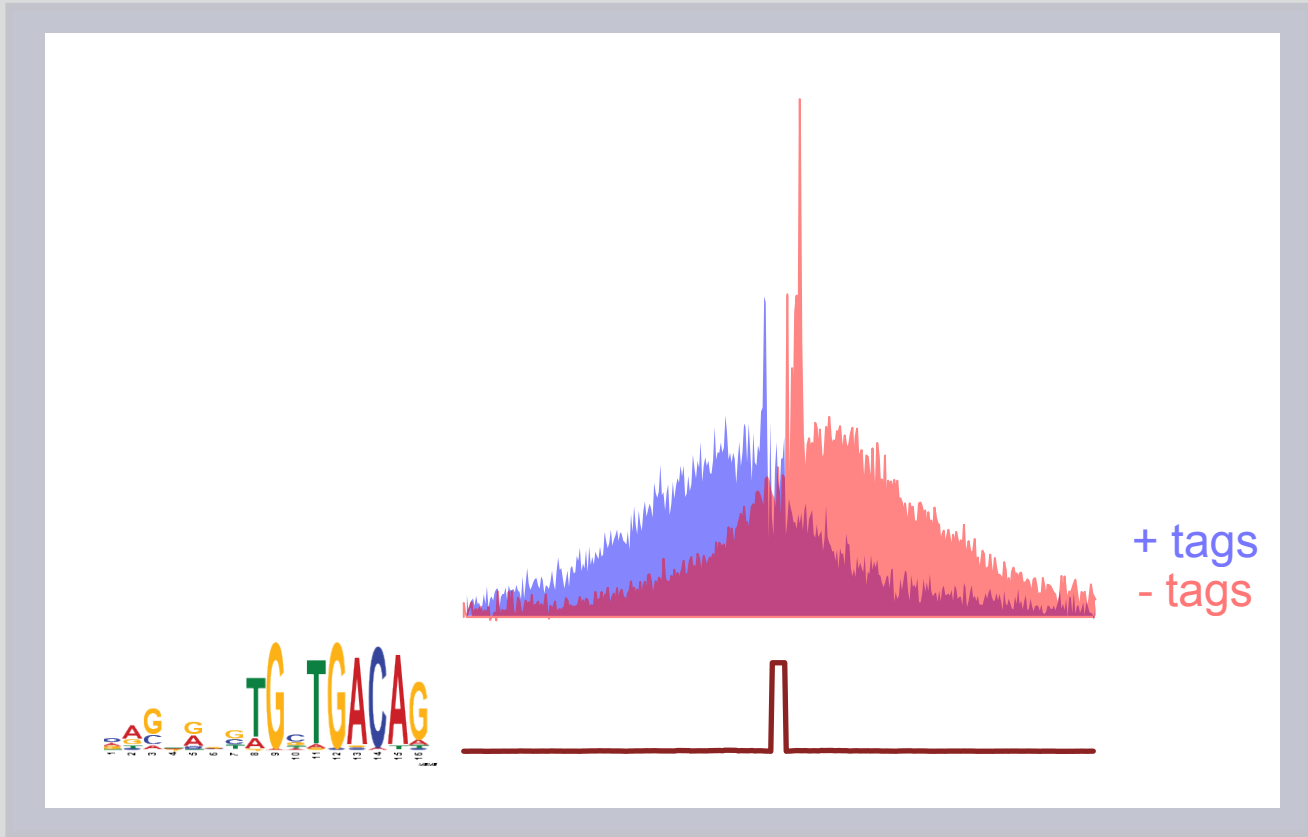calculate genomic density of shifted tags



measure background density
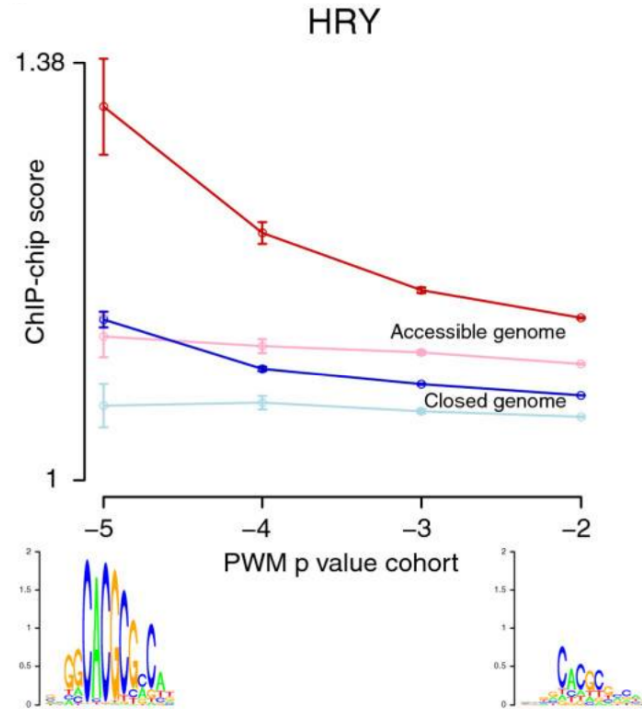


calculate background-normalized binding density



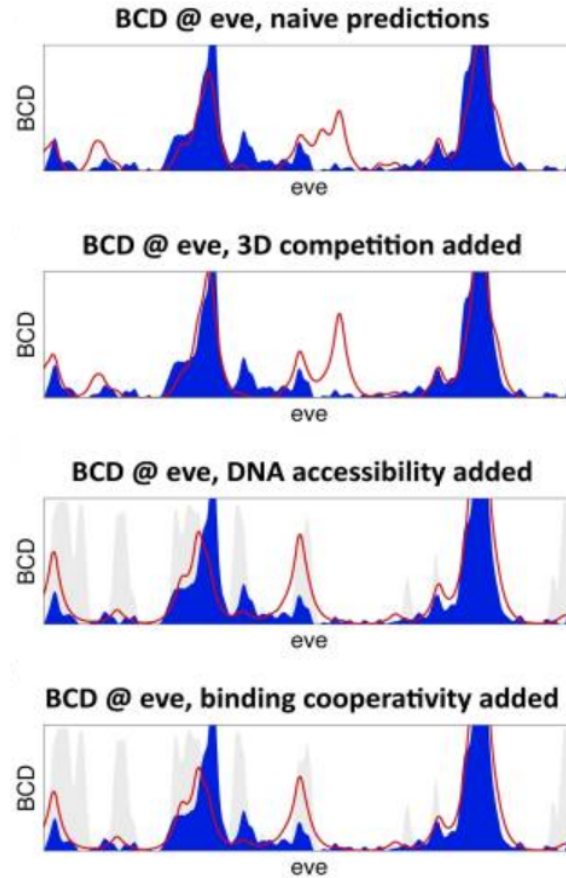call peaks (punctate marks) or regions (broad marks)

## sequence-specific factors



+ tags
- tags

## sequence-specific factors
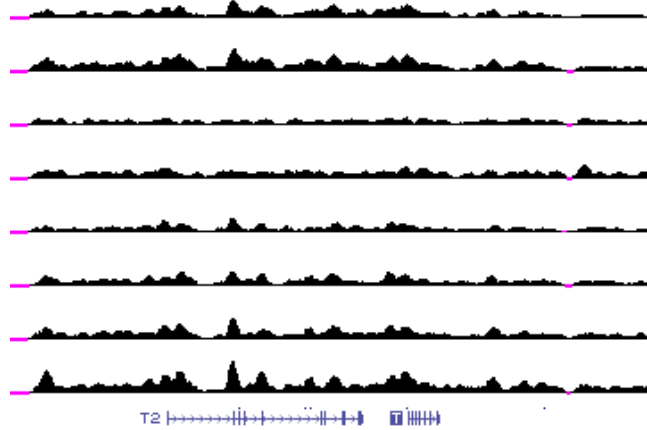


Li et al Genome Biology, 2011

## sequence-specific factors



BCD @ eve, naive predictions

BCD @ eve, 3D competition added

BCD @ eve, DNA accessibility added

BCD @ eve, binding cooperativity added

modeled ChIP tag density

actual ChIP tag density

Kaplan et al PLoS Genetics, 2011

# chromatin-associated factors

## chromatin remodellers



Kato et al IBMS Bonekey, 2010

## transcription machinery

e.g. RNA polymerase II

( and various post-translational modifications )



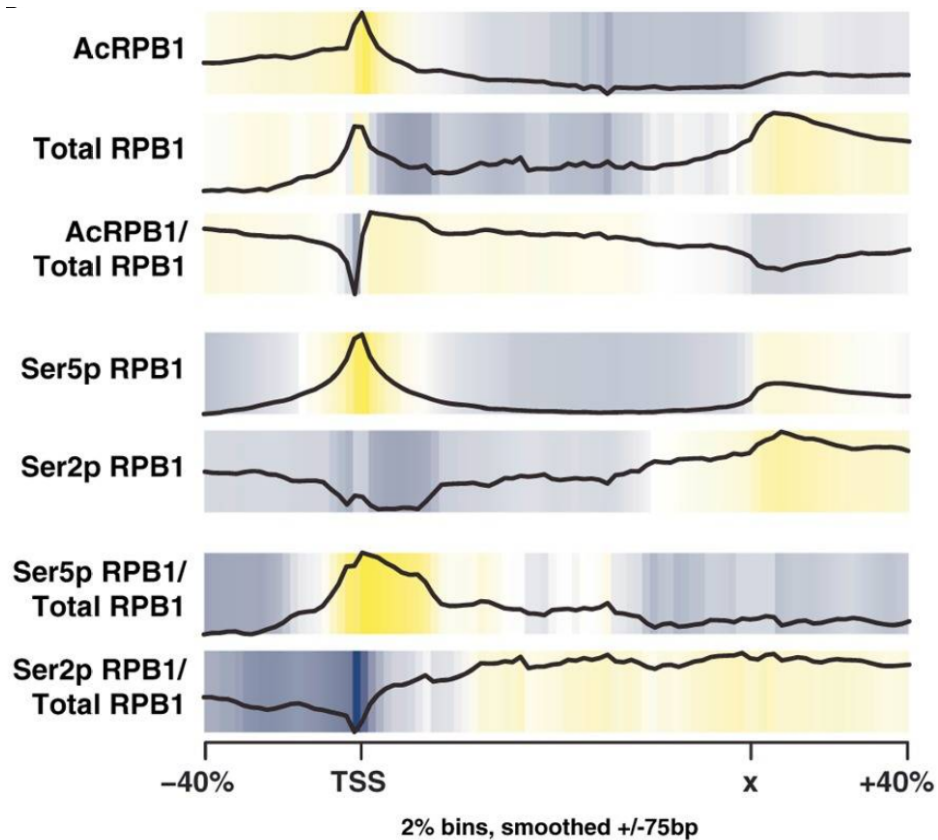Schroder et al Mol Cell 2012

# histone modifications



most commonly studied:

H3K4me3
H3K27ac
H3K27me3

Kosi Gramatikoff, 2007

## histone modifications

### wide range of binding patterns



Barth and Imhoff, Trends in Bio. Sci 2010

H3K4me3 - transcription start sites
H3K27ac - active enhancers and TSS
H3K27me3 - domain around TSS
H3K36me3 - active gene bodies

# histone modifications

## wide range of functional associations



Nature Reviews | Genetics

Zhou et al, NRG 2011

immunoprecipitation

library construction

sequencing

mapping to genome reference (BWA or bowtie2)

identify tag locations

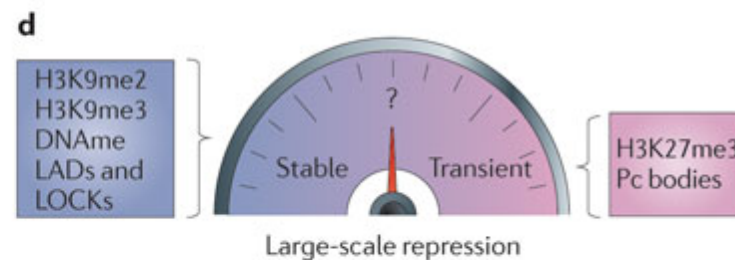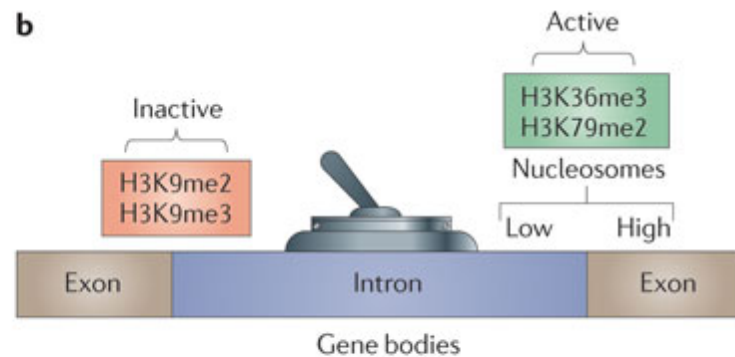calculate median insert size (single end)
determine fragment midpoint (paired end)

estimate actual binding location per tag

calculate genomic density of shifted tags

measure background density

calculate background-normalized binding density

call peaks (punctate marks) or regions (broad marks)

# methodology

## immunoprecipitation



## library construction

## sequencing

## mapping to genome reference (BWA or bowtie2)

## methodology

immunoprecipitation

library construction

sequencing

mapping to genome reference (BWA or bowtie2)

identify tag locations

calculate median insert size (single end)

determine fragment midpoint (paired end)

estimate actual binding location per tag

calculate genomic density of shifted tags
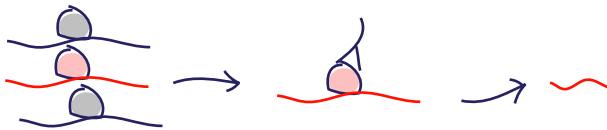
measure background density

calculate background-normalized binding density
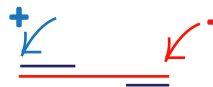
call peaks (punctate marks) or regions (broad marks)

calculate genomic density of shifted tags

measure background density

calculate background-normalized binding density

call peaks (punctate marks) or regions (broad marks)

# chromatin accessibility
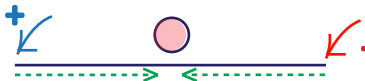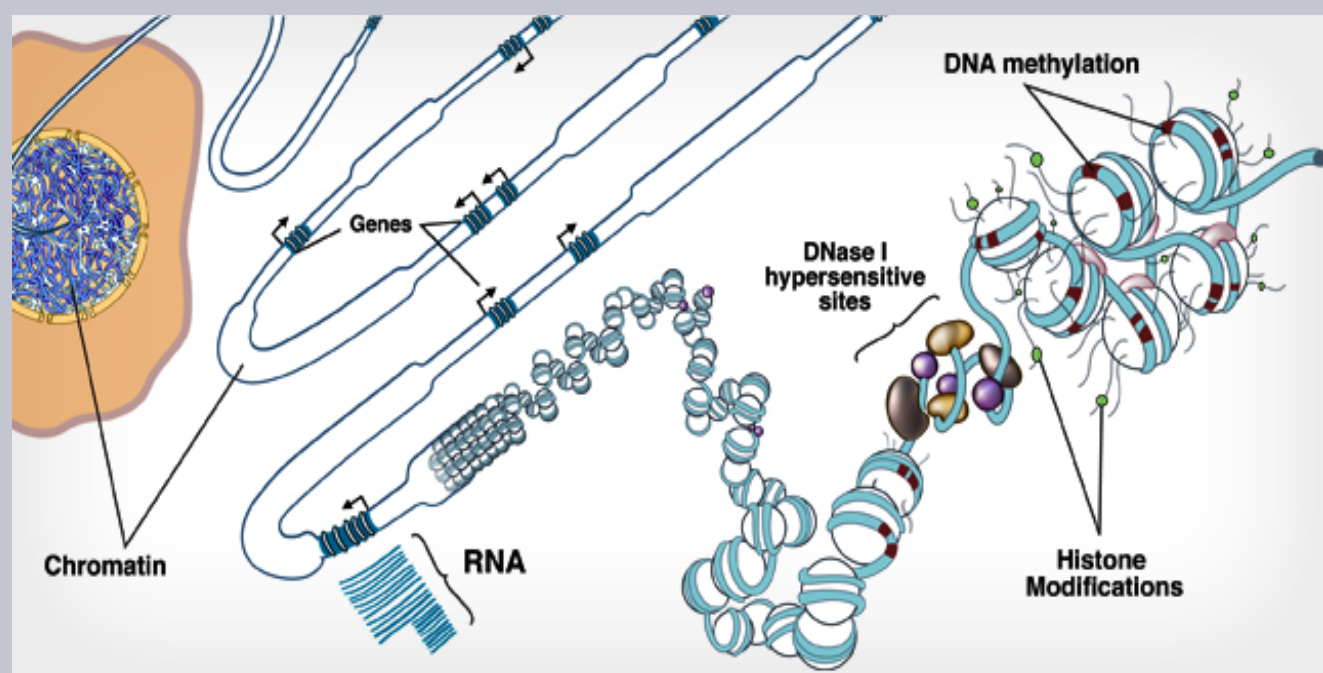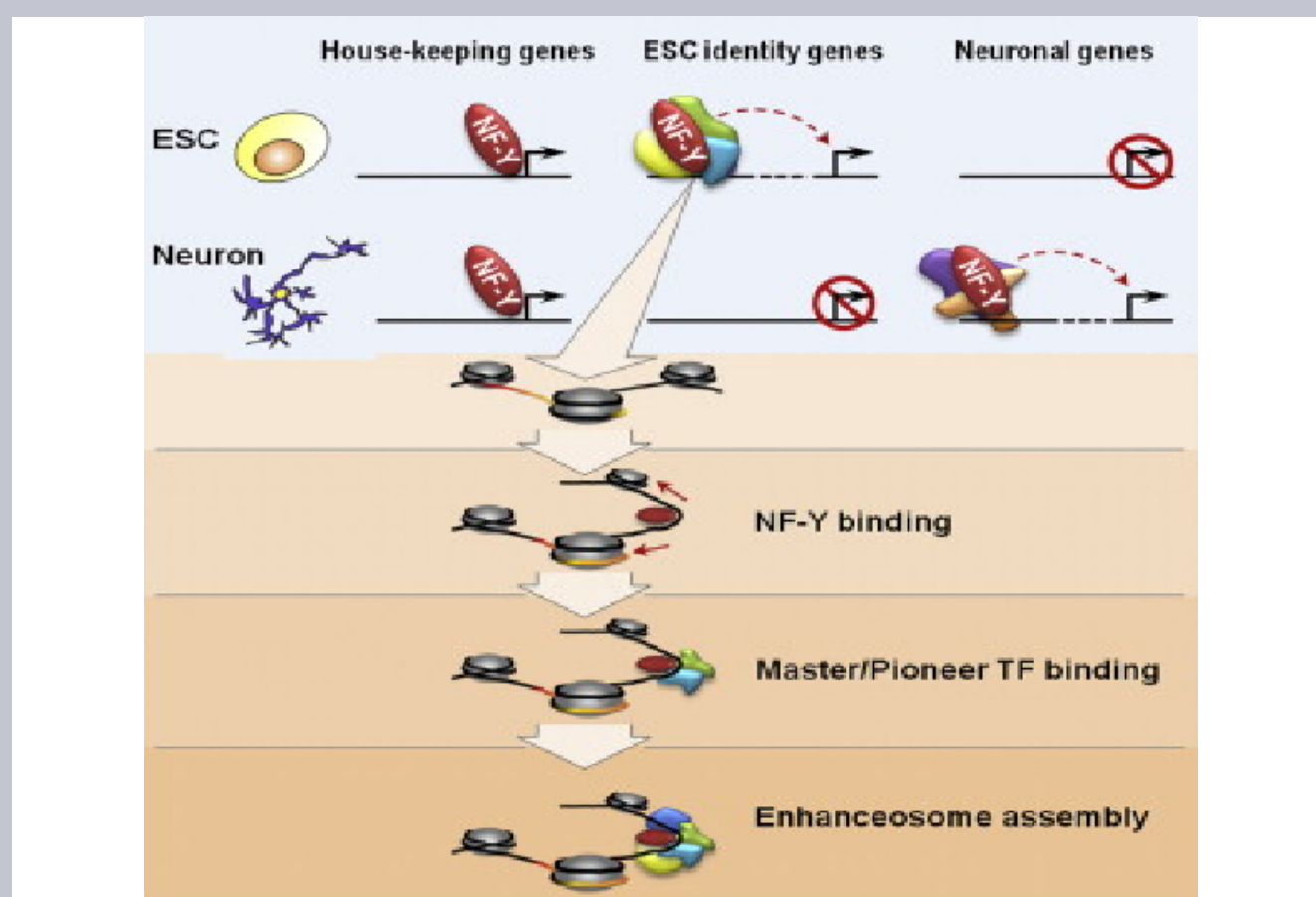
## functional regulatory elements are accessible



DNA methylation

Genes

DNase I hypersensitive sites

Chromatin

RNA

Histone Modifications

roadmapepigenomics.org

## formation of accessible chromatin



House-keeping genes    ESC identity genes    Neuronal genes

ESC

Neuron

NF-Y binding

Master/Pioneer TF binding

Enhanceosome assembly

Oldfield et al, Mol Cell 2014

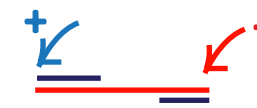## methodology (digestion or transposase)

DNase I or Tn5 release fragments of open chromatin



library construction and sequencing

mapping to genome reference (BWA or bowtie2)

identify cleavage/insertion location

calculate density of cleavage/insertion events

call regions of accessibility

call footprints within accessible chromatin

## footprinting



REST

ChIP-seq quantiles
- [ 98 , 100 ] %
- [ 95 , 98 ] %
- [ 90 , 95 ] %
- [ 85 , 90 ] %
- [ 75 , 85 ] %
- [ 50 , 75 ] %
- [ 0 , 50 ] %

Avg. # DNaseI cuts / site

Distance to motif (bp)

Pique-Regi et al Genome Research 2011

Footprints reflect residence time in chromatin



Binding factor with long term residency

DNA Cleavage Signature Deep footprint in chromatin

Dnase hyper-sensitive site

Long residence time

Residence time of seconds

DNA Cleavage Signature

Factor with high mobility in living cells

No footprint in chromatin

Sung et al, Mol Cell 2014

# functional regulatory elements are accessible

# formation of accessible chromatin

pioneering factors make important genes accessible during development



Oldfield et al, Mol Cell 2014

## footprinting



Pique-Regi et al Genome Research 2011

# footprinting



distinct for
different factors

not consistently
present for all factors

(effect seen in naked DNA for many factors)

methodology (digestion or transposase)

library construction and sequencing

mapping to genome reference (BWA or bowtie2)

identify cleavage/insertion location
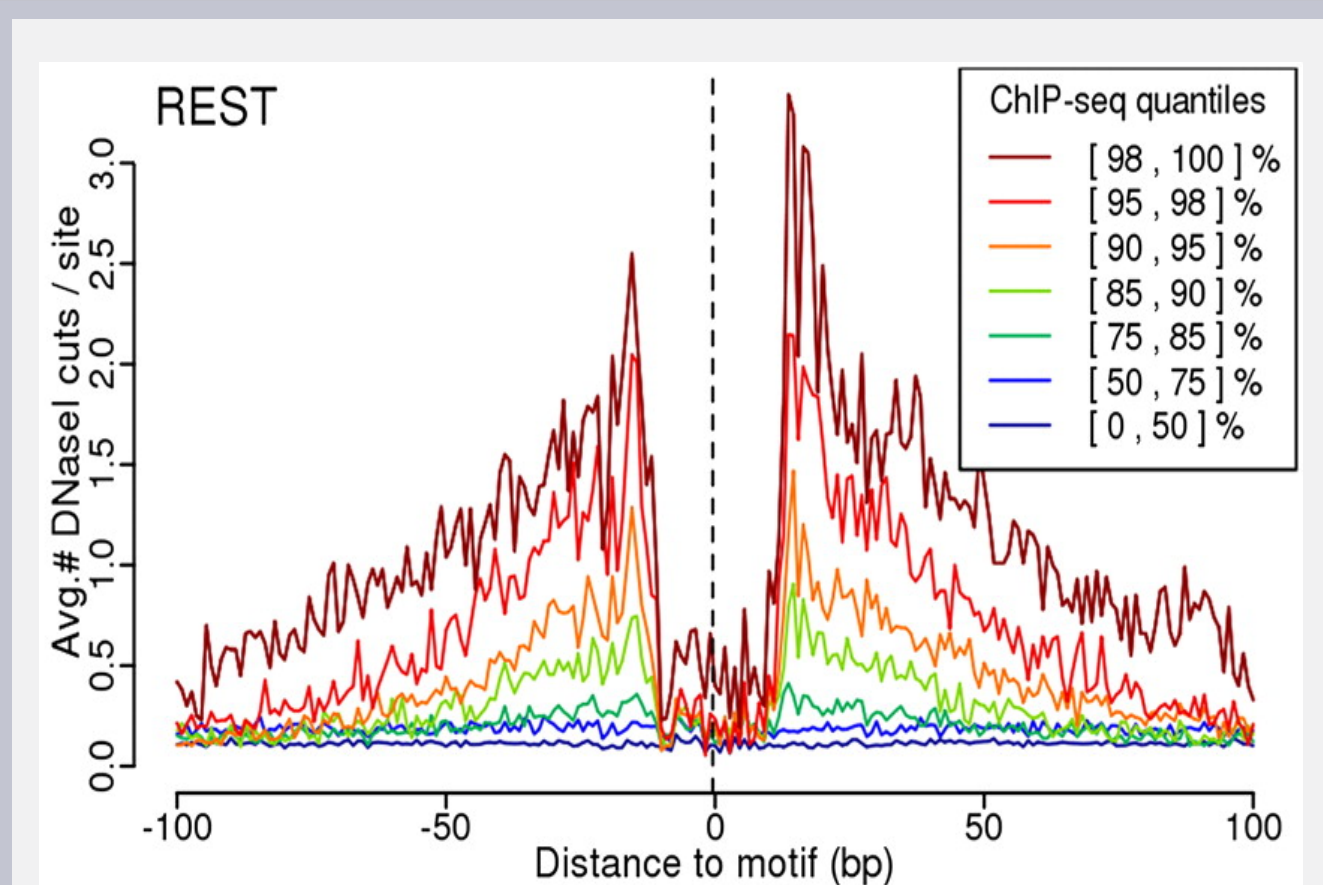
calculate density of cleavage/insertion events

call regions of accessibility

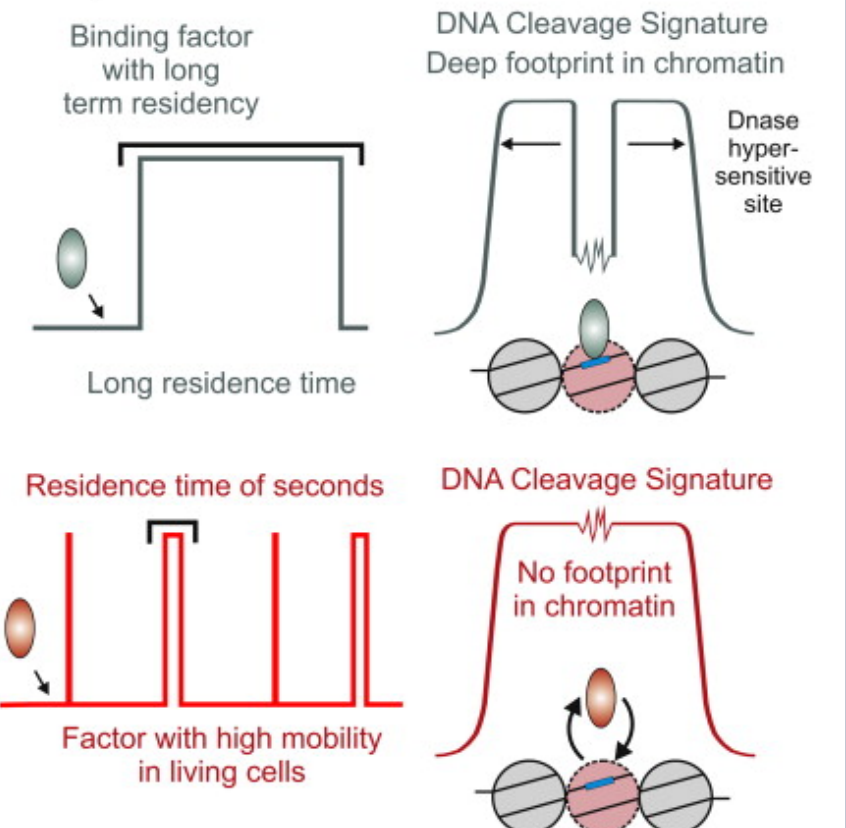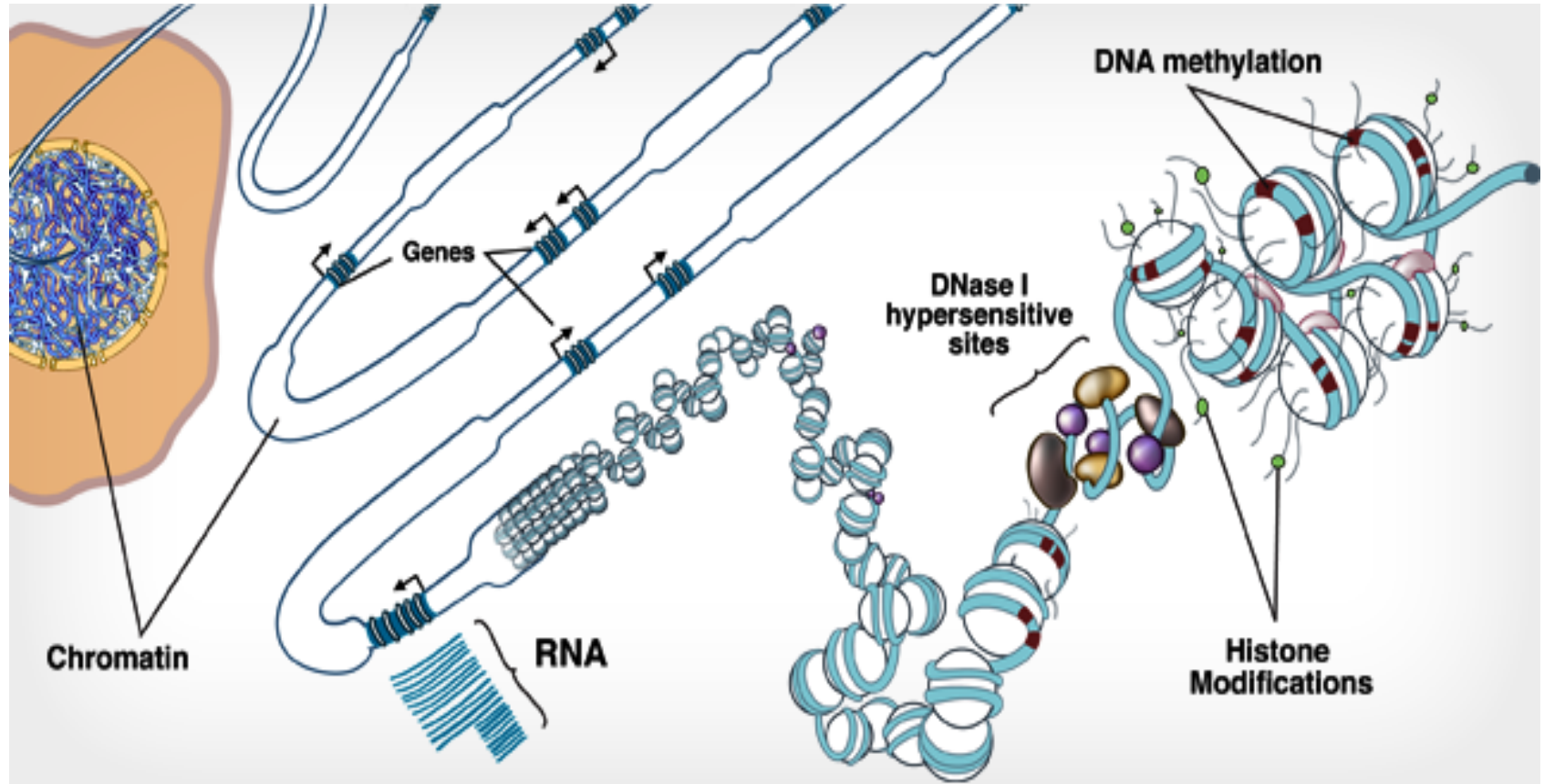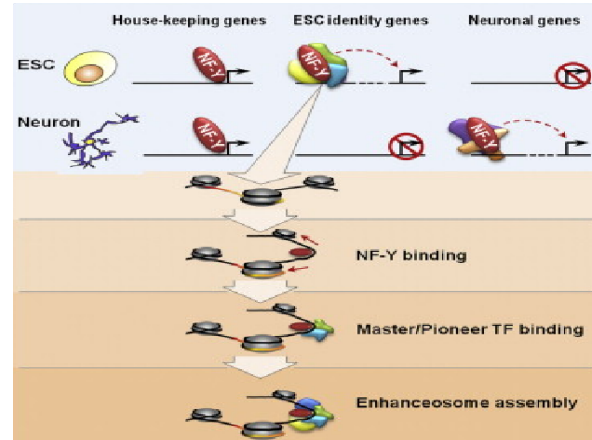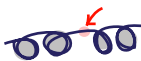call footprints within accessible chromatin

# methodology

DNase I digestion
(DNase-seq)

transposon integration
(ATAC-seq)

DNase I  or Tn5 release short fragments of open chromatin

methodology (digestion or transposase)

DNase I or Tn5 release short fragments of open chromatin

library construction and sequencing

mapping to genome reference (BWA or bowtie2)

identify cleavage/insertion location

calculate density of cleavage/insertion events

call regions of accessibility

call footprints within accessible chromatin

# methodology

library construction and sequencing

mapping to genome reference (BWA or bowtie2)

identify cleavage/insertion location

## methodology (digestion or transposase)

DNase I or Tn5 release short fragments of open chromatin

library construction and sequencing

mapping to genome reference (BWA or bowtie2)

identify cleavage/insertion location

calculate density of cleavage/insertion events
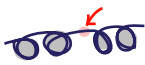
call regions of accessibility

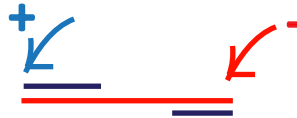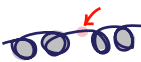call footprints within accessible chromatin

# methodology

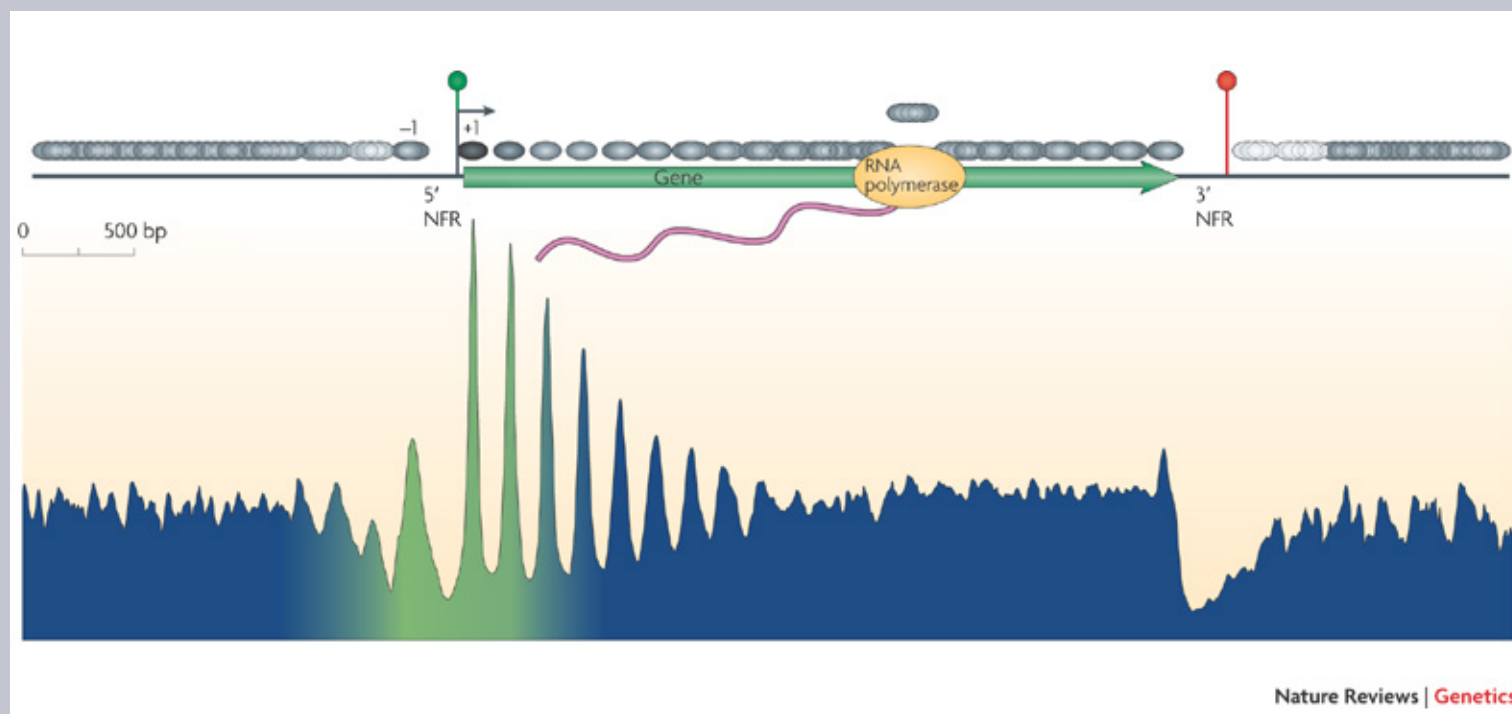### calculate density of cleavage/insertion events

### call regions of accessibility

### call footprints within accessible chromatin

# nucleosome position

## nucleosome positioning during transcription



Nature Reviews | Genetics

Jiang et al NRG 2009

## determinants of positioning stability



A

Highly positioned nucleosomes — "Fuzzy" nucleosomes

Poly(dA:dT) tract

-ATTGAGCTGCAATCTGGAATAACAGCCAGATAAGGAGCTACAGTACC-

Nucleosome favoring sequence:
A/T dinucleotide every 10 basepairs
G/C dinucleotide every 10 basepairs, in antiphase with A/T dinucleotides

B

ATP        ADP + P_i          ATP        ADP + P_i

ATP-consuming
chromatin remodeling factor

ATP-consuming
chromatin remodeling factor

HAT        HDAC

Acetyl

## methodology (MNase digestion)

MNase I digests linker DNA, releasing multisomes

purify mononucleosome-bound DNA

library construction and sequencing

mapping to genome reference (BWA or bowtie2)

identify cleavage location

calculate nucleosomal density

call positioned nucleosomes

## methodology (transposon insertion)

Tn5 integrates into accessible chromatin, and releases multisomes

library construction and sequencing

mapping to genome reference (BWA or bowtie2)

identify cleavage location

identify reads that span at least one nucleosome

read length

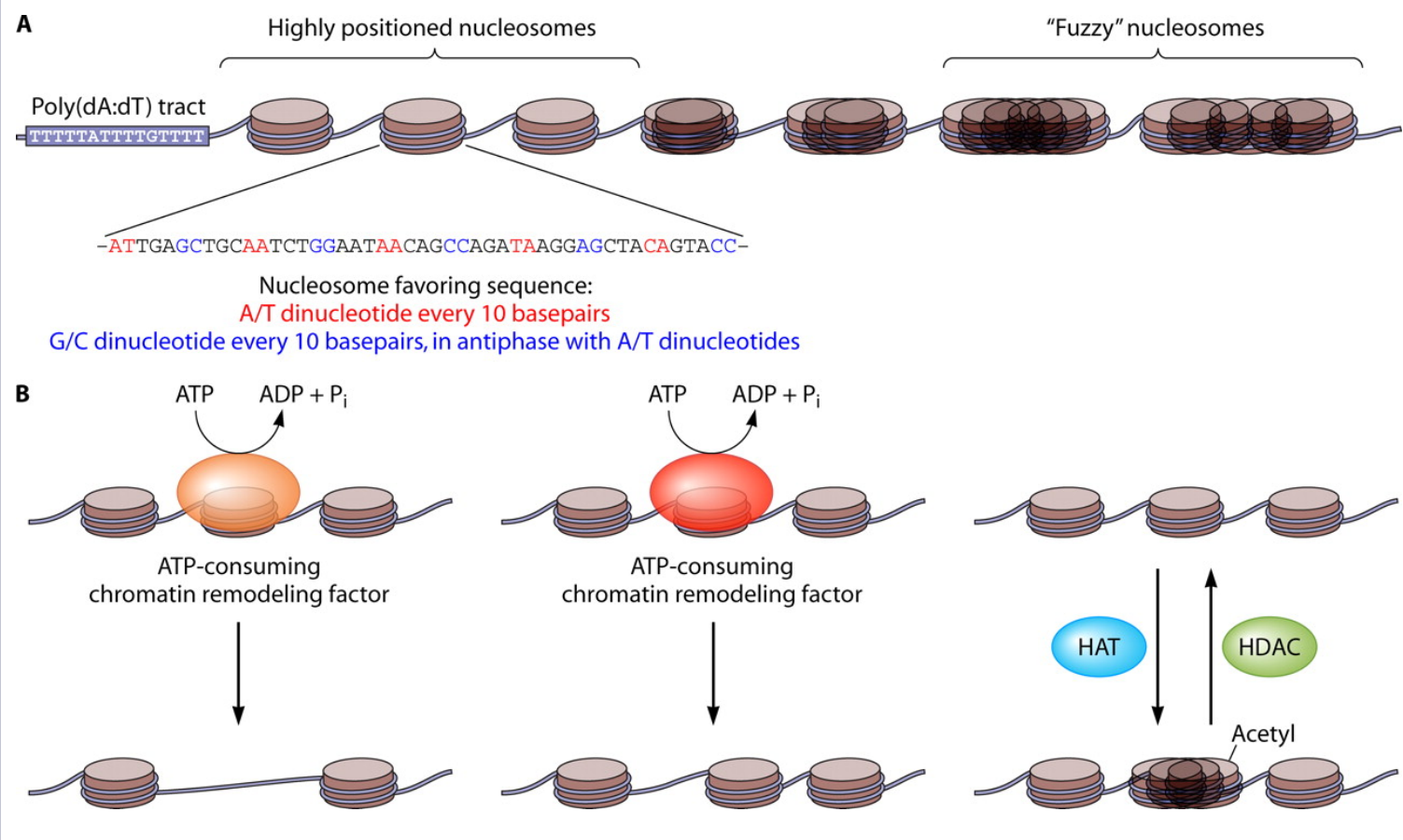calculate density of tags / number of insertion events

call positioned nucleosomes

# nucleosome positioning during transcription



Nature Reviews | Genetics

Jiang et al NRG 2009

# positioning stability
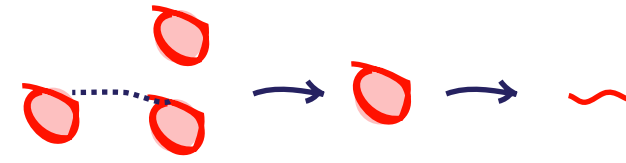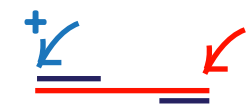


**A** Highly positioned nucleosomes      "Fuzzy" nucleosomes

Poly(dA:dT) tract
TTTTTATTTTGTTTT

−ATTGAGCTGCAATCTGGAATAACAGCCAGATAAGGAGCTACAGTACC−

Nucleosome favoring sequence:
A/T dinucleotide every 10 basepairs
G/C dinucleotide every 10 basepairs, in antiphase with A/T dinucleotides

**B** ATP → ADP + P$_i$     ATP → ADP + P$_i$

ATP-consuming chromatin remodeling factor

ATP-consuming chromatin remodeling factor

HAT    HDAC

Acetyl

# positioning stability



**A** Highly positioned nucleosomes — "Fuzzy" nucleosomes

Poly(dA:dT) tract
TTTTTATTTTGTTTT

−ATTGAGCTGCAATCTGGAATAACAGCCAGATAAGGAGCTACAGTACC−

Nucleosome favoring sequence:
A/T dinucleotide every 10 basepairs
G/C dinucleotide every 10 basepairs, in antiphase with A/T dinucleotides

**B**

ATP → ADP + P$_i$

ATP-consuming
chromatin remodeling factor

ATP → ADP + P$_i$

ATP-consuming
chromatin remodeling factor

HAT   HDAC

Acetyl

# two common methodologies (MNase, ATAC-seq)

## methodology (MNase digestion)

MNase I digests linker DNA, releasing multisomes

purify mononucleosome-bound DNA

library construction and sequencing

mapping to genome reference (BWA or bowtie2)

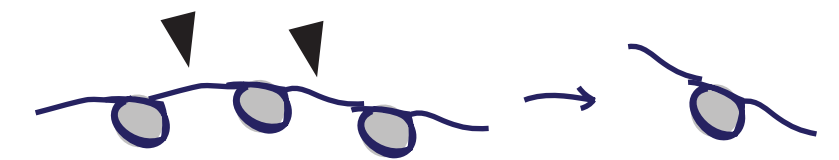identify cleavage location

calculate nucleosomal density

call positioned nucleosomes

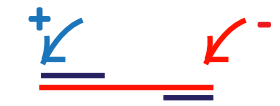## methodology (transposon insertion)

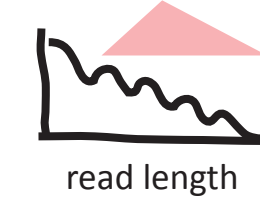Tn5 integrates into accessible chromatin, and releases multisomes

library construction and sequencing

mapping to genome reference (BWA or bowtie2)

identify cleavage location

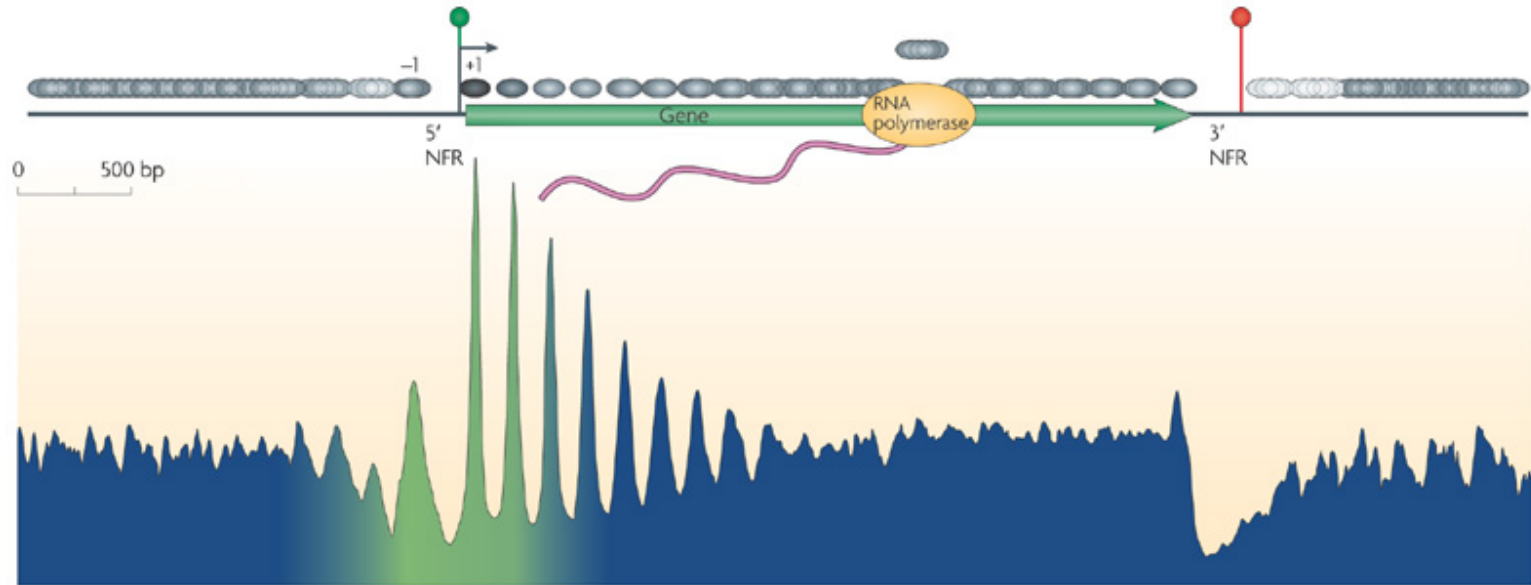identify reads that span at least one nucleosome

read length

calculate density of tags / number of insertion events
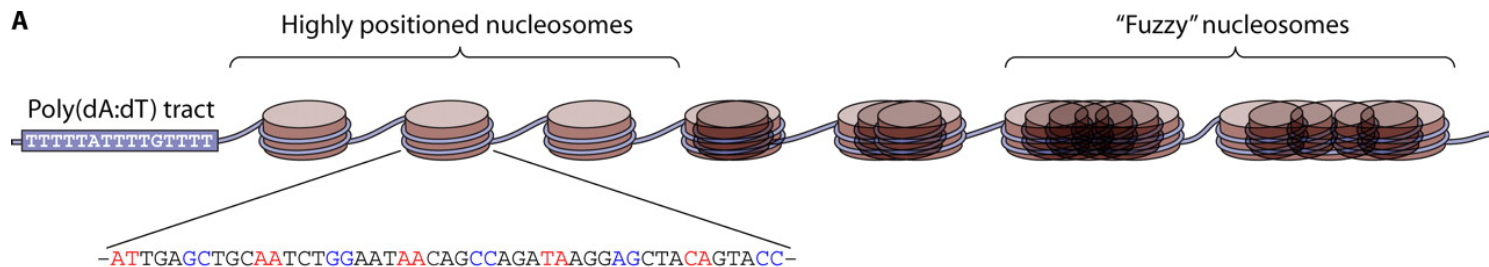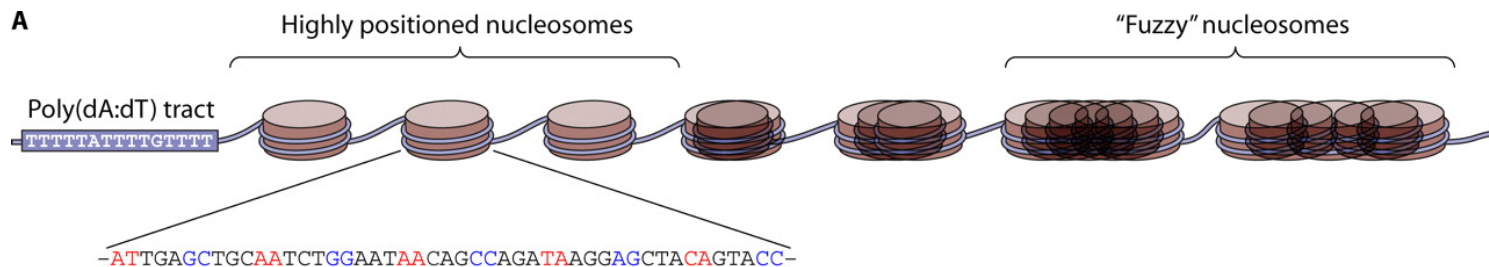
call positioned nucleosomes

# catalogue of common base modifications (in mammals)

## 5mC repression of gene expression



active transcription

promoter

repressed

promoter

● 5mC  5 methyl cytosine

○ C  unmodified cytosine

## bisulfite conversion / protection



T C T C G

bisulfite

T U T C G

PCR

T T T C G
A A A G C

daigenode.com

# methodology

### convert DNA with bisulfite



### fragmentation

### library construction and sequencing

### mapping



Krueger et al Bioinformatics 2011

### calculate % of tags methylated per genomic position

# methodology

## convert DNA with bisulfite



## fragmentation

## library construction and sequencing

Krueger et al Bioinformatics 2011

GEM (Guo et al, PLoS Comp. Bio 2012)



algorithm
1. predict protein-DNA binding events with sparse prior
2. discover set of enriched kmers at binding event
3. cluster set of enriched k-mers into k-mer classes
4. generate positional prior for discovery in most enriched class
5. predict improved binding event probabilities with pos. prior
6. repeat (2) and (3) to generate improved motif enrichments

GEM (Guo et al, PLoS Comp. Bio 2012)

algorithm

1. predict protein-DNA binding events with sparse prior
2. discover set of enriched kmers at binding event
3. cluster set of enriched k-mers into k-mer classes
4. generate positional prior for discovery in most enriched class
5. predict improved binding event probabilities with pos. prior
6. repeat (2) and (3) to generate improved motif enrichments

1. predict events with sparse prior



**(a)** Mixture of reads from joint events
**(b)** CTCF event expected read density
**(c)** GPS mixture model
**(d)** Deconvolved protein-DNA interaction events

Guo et al, Bioinformatics 2010

2. discover set of enriched kmers

set of predicted events from (1)
set of negative regions (+/-300bp)

| | positive region count | negative region count |
|---|---|---|
| kmer1 | 400 | 34 |
| kmer2 | 25 | 42 |
| kmer... | 17 | 25 |

hypergeometric test

Barash et al WABI 2001

3. cluster kmer classes

| K-mer | Offset | Pos Hit | Neg Hit |
|---|---|---|---|
| ----ATGCAAAT | -3 | 739 | 30 |
| ----TATGCAAA | -4 | 628 | 33 |
| ------TGCAAATG | -2 | 460 | 22 |
| ----ATGCTAAT | -3 | 382 | 12 |
| ---TTATGCAA | -5 | 358 | 13 |
| ------ATGCATAT | -3 | 320 | 21 |
| ------TGCAAATT | -2 | 222 | 18 |
| ... ... | ... | ... | ... |

Guo et al, PLoS Comp Bio 2012

4. generate positional prior for events in most enriched kmer class
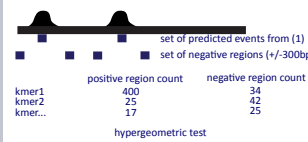
$$p(\pi) \propto \prod_{m=1}^{M} (\pi_m)^{-\alpha_S + \alpha_m}$$

5. predict binding event probabilities

$$\hat{\pi}_m^{(t)} = \frac{\max(0, N_m - \alpha_S + \alpha_m)}{\sum_{m'=1}^{M} \max(0, N_{m'} - \alpha_S + \alpha_m)}, \quad N_m = \sum_{n=1}^{N} \gamma(z_n = m)$$

6. Redo steps (2) and (3) to improve motif quality

# GEM (Guo et al, PLoS Comp. Bio 2012)

## 1. predict events with sparse prior



**(a)** Mixture of reads from joint events

**(b)** CTCF event expected read density

**(c)** GPS mixture model

Possible events $b_j$

Observed reads $r_n$

**(d)** Deconvolved protein-DNA interaction events

GEM (Guo et al, PLoS Comp. Bio 2012)

## 2. discover set of enriched kmers



set of predicted events from (1)

set of negative regions (+/-300bp)

| | positive region count | negative region count |
|---|---|---|
| kmer1 | 400 | 34 |
| kmer2 | 25 | 42 |
| kmer... | 17 | 25 |

hypergeometric test

---

GEM (Guo et al, PLoS Comp. Bio 2012)

algorithm
1. predict protein-DNA binding events with sparse prior
2. discover set of enriched kmers at binding event
3. cluster set of enriched k-mers into k-mer classes
4. generate positional prior for discovery in most enriched class
5. predict improved binding event probabilities with pos. prior
6. repeat (2) and (3) to generate improved motif enrichments

1. predict events with sparse prior

(a) Mixture of reads from joint events
(b) CTCF event expected read density
(c) GPS mixture model
(d) Deconvolved protein-DNA interaction events

Guo et al, Bioinformatics 2010

2. discover set of enriched kmers

set of predicted events from (1)
set of negative regions (+/-300bp)

| | positive region count | negative region count |
|---|---|---|
| kmer1 | 400 | 34 |
| kmer2 | 25 | 42 |
| kmer... | 17 | 25 |

hypergeometric test

Barash et al WABI 2001

3. cluster kmer classes

| K-mer | Offset | Pos Hit | Neg Hit |
|---|---|---|---|
| -----ATGCAAAT | -3 | 739 | 30 |
| ----TATGCAAA | -4 | 628 | 33 |
| ------TGCAAATG | -2 | 460 | 22 |
| -----ATGCTAAT | -3 | 382 | 12 |
| ---TTATGCAA | -5 | 358 | 13 |
| -----ATGCATAT | -3 | 320 | 21 |
| ------TGCAAATT | -2 | 222 | 18 |
| ... ... | ... | ... | ... |

Guo et al, PLoS Comp Bio 2012

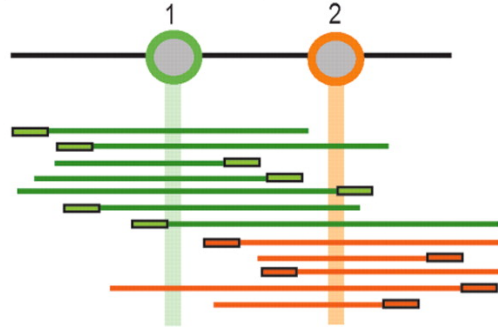4. generate positional prior for events in most enriched kmer class

$$p(\pi) \propto \prod_{m=1}^{M} (\pi_m)^{-\alpha_S + \alpha_m}$$
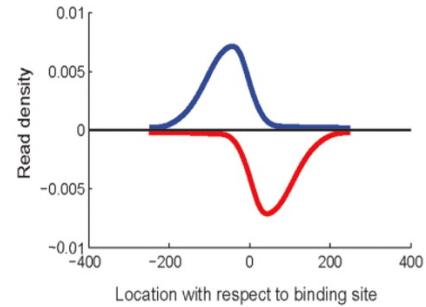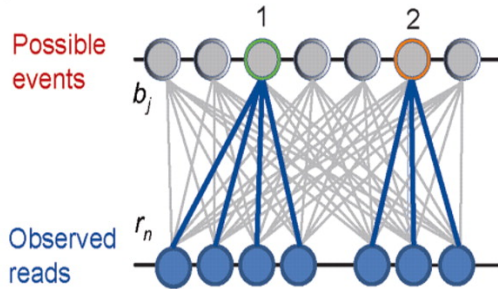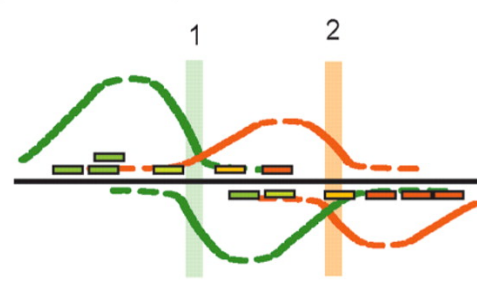
5. predict binding event probabilities

$$\hat{\pi}_m^{(t)} = \frac{\max(0, N_m - \alpha_S + \alpha_m)}{\sum_{m'=1}^{M} \max(0, N_{m'} - \alpha_S + \alpha_m)}, N_m = \sum_{n=1}^{N} \gamma(z_n = m)$$

6. Redo steps (2) and (3) to improve motif quality

GEM (Guo et al, PLoS Comp. Bio 2012)

## GEM (Guo et al, PLoS Comp. Bio 2012)

### algorithm

1. predict protein-DNA binding events with sparse prior
2. discover set of enriched kmers at binding event
3. cluster set of enriched k-mers into k-mer classes
4. generate positional prior for discovery in most enriched class
5. predict improved binding event probabilities with pos. prior
6. repeat (2) and (3) to generate improved motif enrichments

### 1. predict events with sparse prior

(a) Mixture of reads from joint events
(b) CTCF event expected read density
(c) GPS mixture model
(d) Deconvoluted protein-DNA interaction events

Guo et al, Bioinformatics 2010

### 2. discover set of enriched kmers

set of predicted events from (1)
set of negative regions (+/-300bp)

| | positive region count | negative region count |
|---|---|---|
| kmer1 | 400 | 34 |
| kmer2 | 25 | 42 |
| kmer... | 17 | 25 |

hypergeometric test

Barash et al WABI 2001

### 3. cluster kmer classes

| K-mer | Offset | Pos Hit | Neg Hit |
|---|---|---|---|
| -----ATGCAAAT | -3 | 739 | 30 |
| ----TATGCAAA | -4 | 628 | 33 |
| ------TGCAAATG | -2 | 460 | 22 |
| -----ATGCTAAT | -3 | 382 | 12 |
| ---TTATGCAA | -5 | 358 | 13 |
| -----ATGCATAT | -3 | 320 | 21 |
| ------TGCAAATT | -2 | 222 | 18 |
| ... ... | ... | ... | ... |

Guo et al, PLoS Comp Bio 2012

### 4. generate positional prior for events in most enriched kmer class
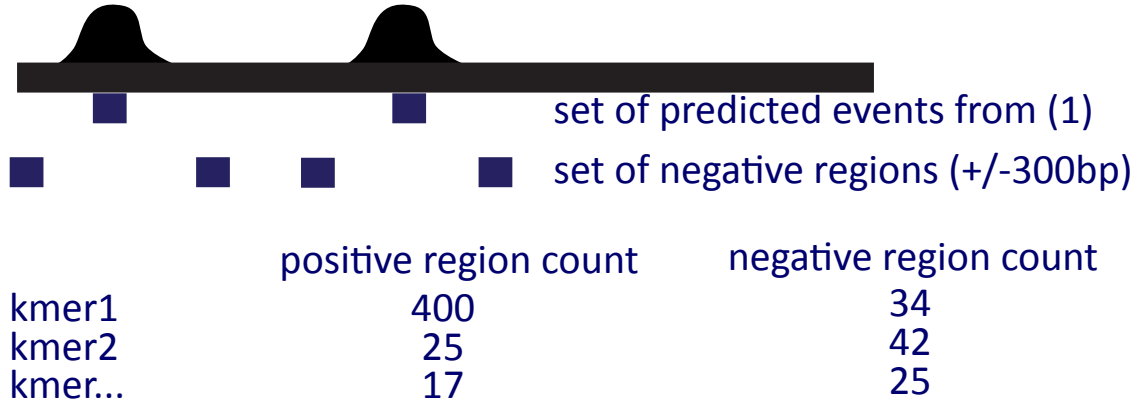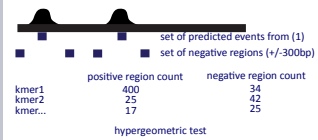
$$p(\pi) \propto \prod_{m=1}^{M} (\pi_m)^{-\alpha_S + \alpha_m}$$

### 5. predict binding event probabilities

$$\hat{\pi}_m^{(t)} = \frac{\max(0, N_m - \alpha_S + \alpha_m)}{\sum_{m'=1}^{M} \max(0, N_{m'} - \alpha_S + \alpha_m)}, \quad N_m = \sum_{n=1}^{N} \gamma(z_n = m)$$

### 6. Redo steps (2) and (3) to improve motif quality

# 3. cluster kmer classes

| K-mer | Offset | Pos Hit | Neg Hit |
|---|---|---|---|
| -----ATGCAAAT | -3 | 739 | 30 |
| ----TATGCAAA | -4 | 628 | 33 |
| ------TGCAAATG | -2 | 460 | 22 |
| -----ATGCTAAT | -3 | 382 | 12 |
| ---TTATGCAA | -5 | 358 | 13 |
| -----ATGCATAT | -3 | 320 | 21 |
| ------TGCAAATT | -2 | 222 | 18 |
| ... ... | ... | ... | ... |

GEM (Guo et al, PLoS Comp. Bio 2012)

**4. generate positional prior for events in most enriched kmer class**

$$p(\pi) \propto \prod_{m=1}^{M} (\pi_m)^{-\alpha_S + \alpha_m}$$

**5. predict binding event probabilities**

$$\hat{\pi}_m^{(i)} = \frac{\max(0, N_m - \alpha_S + \alpha_m)}{\sum_{m'=1}^{M} \max(0, N_{m'} - \alpha_S + \alpha_m)}, \; N_m = \sum_{n=1}^{N} \gamma(z_n = m)$$

GEM (Guo et al, PLoS Comp. Bio 2012)

GEM (Guo et al, PLoS Comp. Bio 2012)

algorithm
1. predict protein-DNA binding events with sparse prior
2. discover set of enriched kmers at binding event
3. cluster set of enriched k-mers into k-mer classes
4. generate positional prior for discovery in most enriched class
5. predict improved binding event probabilities with pos. prior
6. repeat (2) and (3) to generate improved motif enrichments

1. predict events with sparse prior

(a) Mixture of reads from joint events   (b) CTCF event expected read density

(c) GPS mixture model   (d) Deconvolved protein-DNA interaction events

Possible events

Observed reads

Guo et al, Bioinformatics 2010

2. discover set of enriched kmers

set of predicted events from (1)
set of negative regions (+/-300bp)

| | positive region count | negative region count |
|---|---|---|
| kmer1 | 400 | 34 |
| kmer2 | 25 | 42 |
| kmer… | 17 | 25 |

hypergeometric test

Barash et al WABI 2001

3. cluster kmer classes

| K-mer | Offset | Pos Hit | Neg Hit |
|---|---|---|---|
| -----ATNCAAAT | -3 | 739 | 30 |
| ----TATNCAAA | -4 | 628 | 33 |
| ------TNCAAATG | -2 | 460 | 22 |
| -----ATNCTAAT | -3 | 382 | 12 |
| ---TTATNCAA | -5 | 358 | 13 |
| -----ATNCATAT | -3 | 320 | 21 |
| ------TNCAAATT | -2 | 222 | 18 |
| … … | … | … | … |

Guo et al, PLoS Comp Bio 2012

4. generate positional prior for events in most enriched kmer class

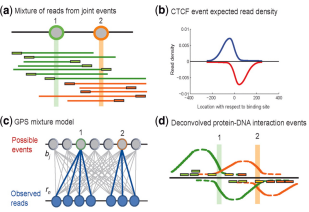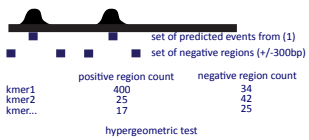$$p(\pi) \propto \prod_{m=1}^{M} (\pi_m)^{-\alpha_S + \alpha_m}$$

5. predict binding event probabilities

$$\hat{\pi}_m^{(i)} = \frac{\max(0, N_m - \alpha_S + \alpha_m)}{\sum_{m'=1}^{M} \max(0, N_{m'} - \alpha_S + \alpha_m)}, \quad N_m = \sum_{n=1}^{N} \gamma(z_n = m)$$

6. Redo steps (2) and (3) to improve motif quality

6. Redo steps (2) and (3) to improve motif quality