

Analyzing data from capture-based next generation sequencing assays

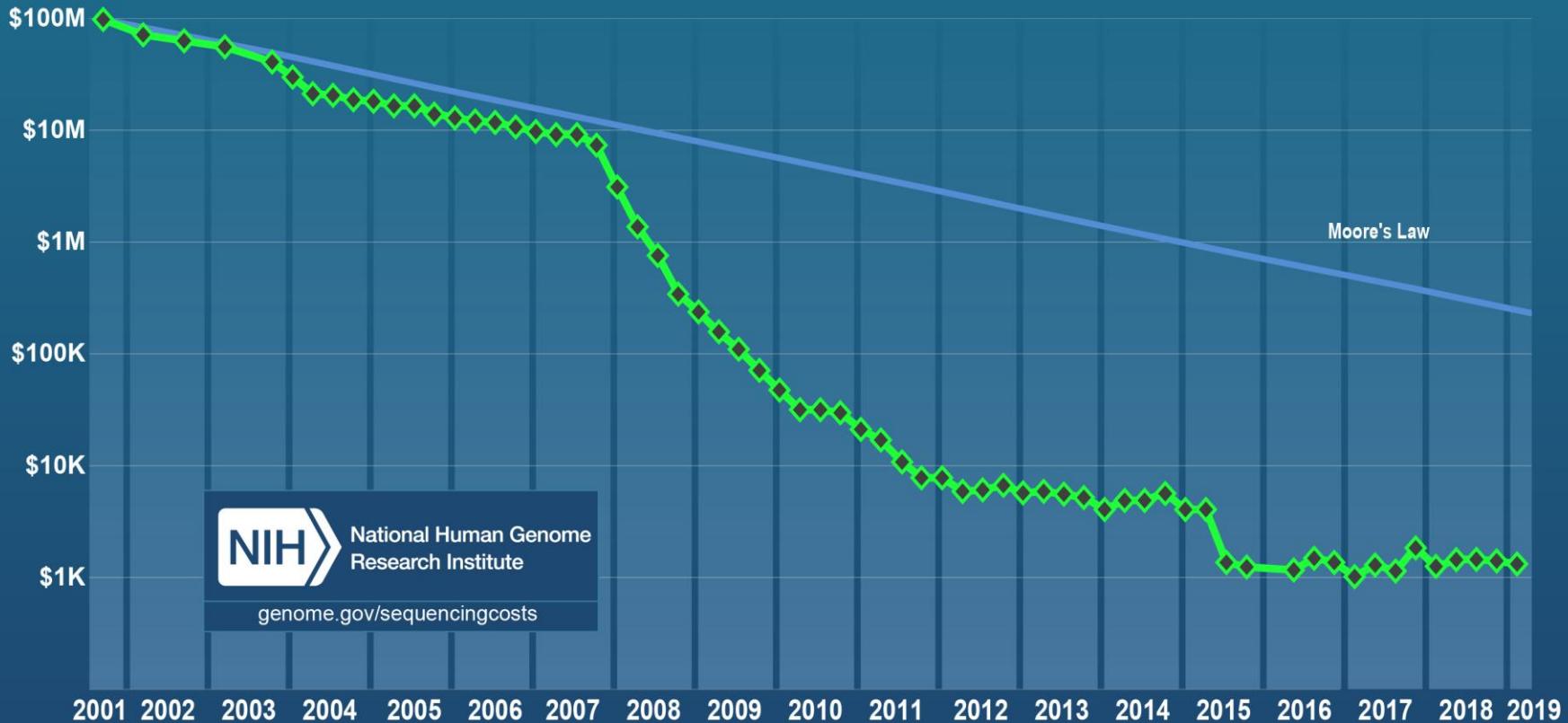
Steve Qin

Department of Biostatistics
and Bioinformatics

Rollins School of Public Health
Emory University

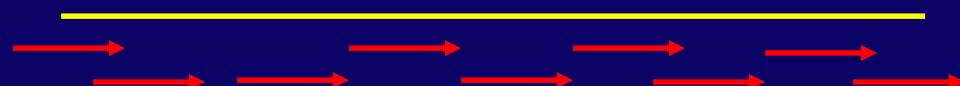


Cost per Genome



Different strategies of using sequencing technologies

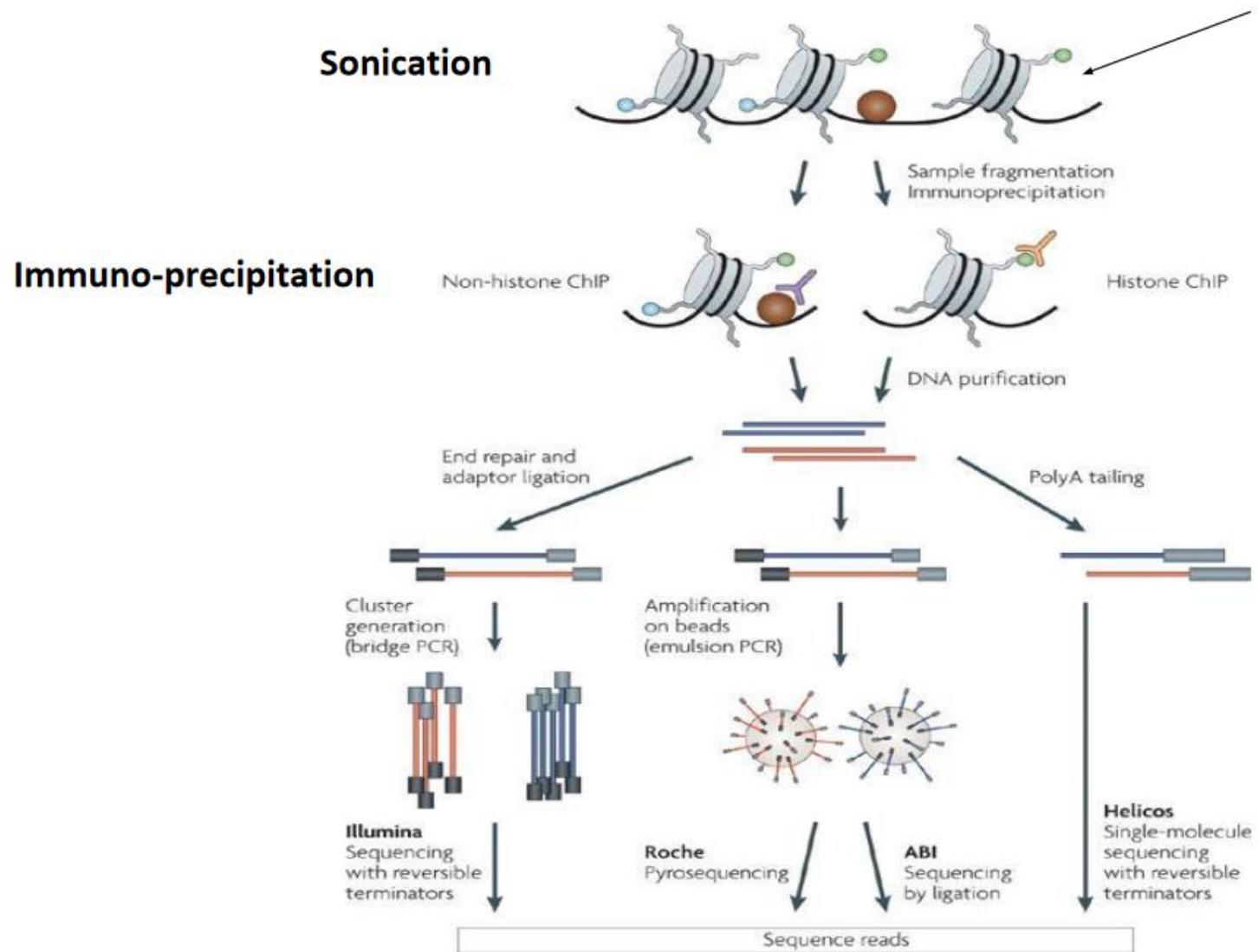
- DNA-seq:
 - Whole genome sequencing
 - Uniform coverage
 - Flat
- ChIP-seq, Dnase-seq, ATAC-seq...
 - Capture-based
 - Selected “genome-wide” sequencing
 - Subset
 - Peaky



What is ChIP-seq

- Combination of chromatin immunoprecipitation (ChIP) with ultra high-throughput massively parallel sequencing.
- Allows mapping of protein–DNA interactions *in vivo* on a genome scale.
- Enables mapping of transcription factors binding, DNA binding proteins (HP1, Lamins, HMGA etc), RNA Pol II occupancy or Histone modification marks at genome scale.
- The typical ChIP assay usually take 4–5 days, and require approx. 1 -10 million cells.

ChIP-seq methodology



ChIP sequencing

Resource **Cell**

High-Resolution Profiling of Histone Methylation in the Human Genome

Artem Barski,^{1,3} Suresh Cuddapah,^{1,3} Kairong Cui,^{1,3} Tae-Young Roh,^{1,3} Dustin E. Schones,^{1,3} Zhibin Wang,^{1,3} Gang Wei,^{1,3} Iouri Chepelev,² and Keji Zhao^{1,*}

¹Laboratory of Molecular Immunology, National Heart, Lung, and Blood Institute, NIH, Bethesda, MD 20892, USA
²Department of Human Genetics, Gonda Neuroscience and Genetics Research Center, University of California, Los Angeles, Los Angeles, CA 90095, USA
³These authors contributed equally to this work and are listed alphabetically.
*Correspondence: zhao@nhlbi.nih.gov
DOI 10.1016/j.cell.2007.05.009

Genome-Wide Mapping of in Vivo Protein-DNA Interactions

David S. Johnson,^{1,*} Ali Mortazavi,^{2,*} Richard M. Myers,^{1,†} Barbara Wold^{2,3,†}

www.sciencemag.org SCIENCE VOL 316 8 JUNE 2007

Genome-wide profiles of STAT1 DNA association using chromatin immunoprecipitation and massively parallel sequencing

Gordon Robertson¹, Martin Hirst¹, Matthew Bainbridge¹, Misha Bilenky¹, Yongjun Zhao¹, Thomas Zeng¹, Ghia Euskirchen², Bridget Bernier¹, Richard Varhol¹, Allen Delaney¹, Nina Thiessen¹, Obi L Griffith¹, Ann He¹, Marco Marra¹, Michael Snyder² & Steven Jones¹

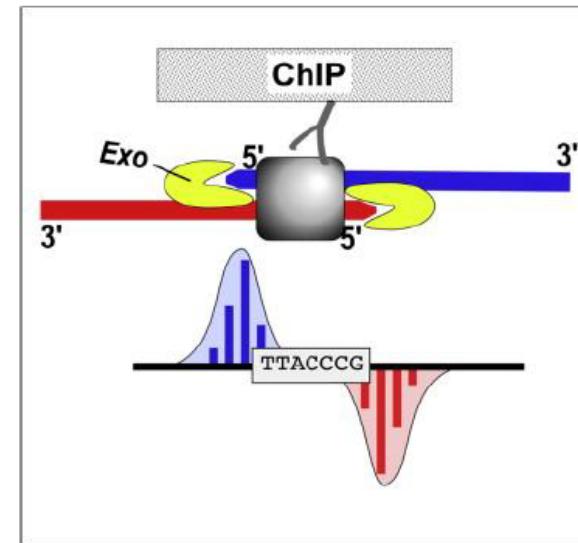
ature.com/naturemethods

¹British Columbia Cancer Agency Genome Sciences Centre, 675 West 10th Avenue, Vancouver, British Columbia V5Z 4S6, Canada. ²Department of Molecular, Cellular and Developmental Biology, Yale University, New Haven, Connecticut 06520, USA. Correspondence should be addressed to S.J. (sjones@bcgsc.ca).

RECEIVED 11 MAY; ACCEPTED 5 JUNE; PUBLISHED ONLINE 11 JUNE 2007; DOI:10.1038/NMETH1068

Advances in technologies for protein-DNA interaction

- ChIP-chip : combines ChIP with microarray technology.
- ChIP-PET : ChIP with paired end tag sequencing
- ChIP-exo : ChIP-seq with exonuclease digestion
- CLIP-seq / HITS-CLIP : cross-linking immunoprecipitation high throughput sequencing
- ATAC-seq : Assay for Transposon Accessible Chromatin
- Sono-seq : Sonication of cross linked chromatin sequencing.
- Hi-C: High throughput long distance chromatin interactions



Statistical aspects and best practices

These guidelines address :

- Antibody validation
- Experimental replication
- Sequencing depth
- Data and metadata reporting
- Data quality assessment
- Replicates

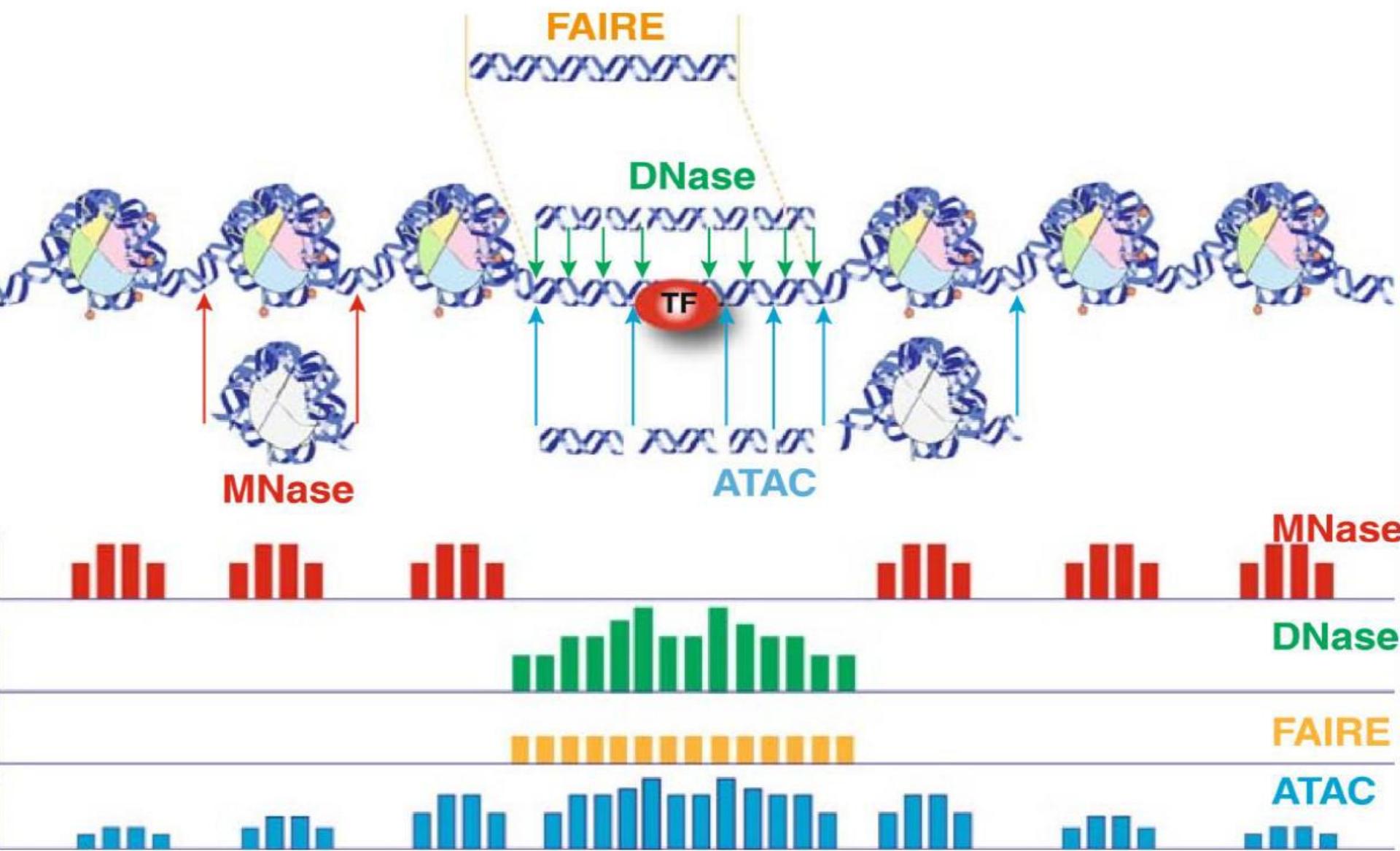
Experimental guidelines:

- Landt *et al.*, "ChIP-seq guidelines and practices of the ENCODE and modENCODE consortia." *Genome Res.* 2012.
- Marinov *et al.*, "Large-scale quality analysis of published ChIP-seq data." 2014 *G3*
- Rozowsky *et al.*, "PeakSeq enables systematic scoring of ChIP-seq experiments relative to controls" *Nat Biotechnol.* 2009

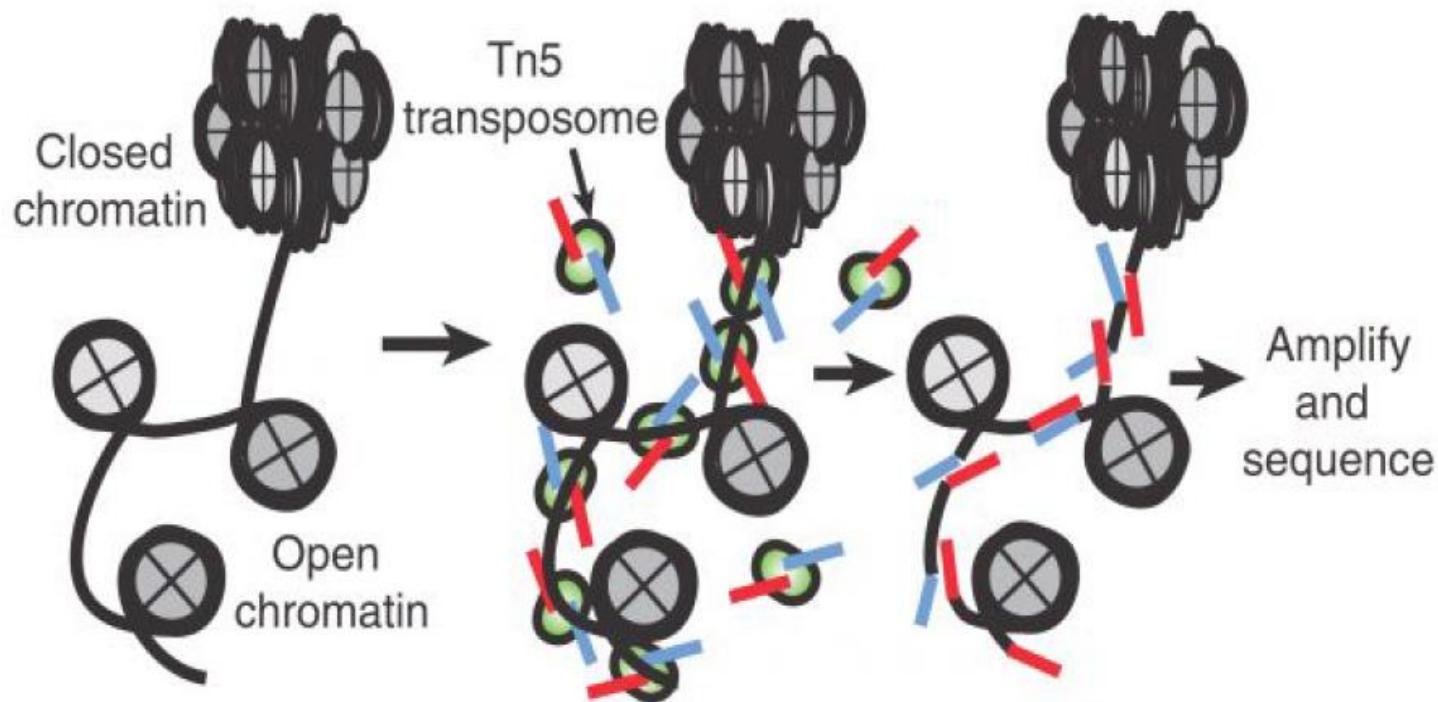
Statistical aspects:

- Cairns *et al.*, "Statistical Aspects of ChIP-Seq Analysis." *Adv. in Stat Bioinf.*, 2013.
- Carroll TS *et al.*, "Impact of artifact removal on ChIP quality metrics in ChIP-seq and ChIP-exo data." *Front Genet.* 2014
- Bailey *et al.*, "Practical guidelines for the comprehensive analysis of ChIP-seq data." *PLoS Comput Biol.* 2013.
- Sims *et al.*, "Sequencing depth and coverage: key considerations in genomic analyses." *Nat. Rev. Genet.* 2014.

Detecting Chromatin Accessibility



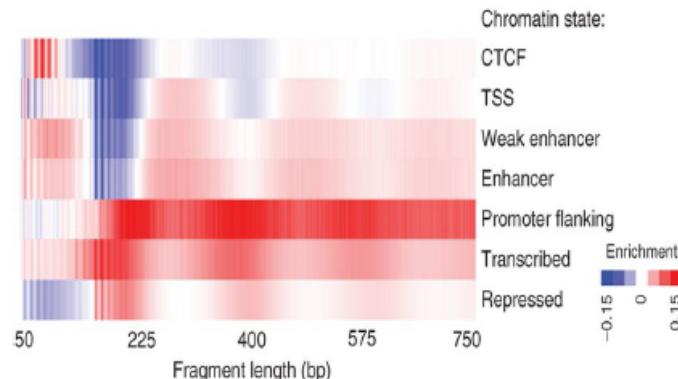
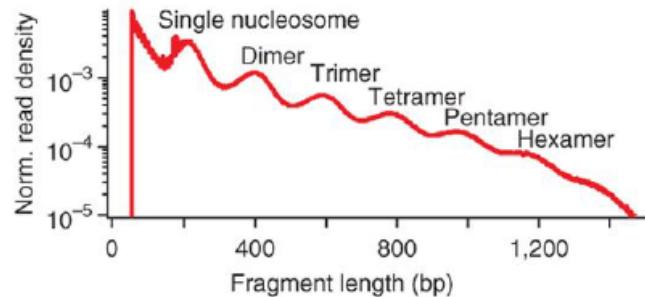
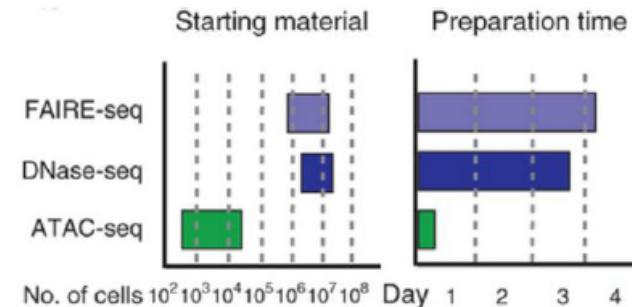
ATAC-seq



- Enables measurement of chromatin structure modifications (nucleosome free regions) on gene regulation.
- Does not require antibodies or tags that can introduce potential bias.
- Hyperactive Tn5 transposase is used to fragment DNA and integrate into active regulatory regions.
- During ATAC-seq, 500–50,000 unfixed nuclei are tagged *in vitro* with sequencing adapters by purified Tn5 transposase.
- Can also detect nucleosome packing, positioning and TF footprints.

ATAC-seq

- Two-step protocol
 - Insertion of Tn5 transposase with adaptors
 - PCR amplification
- Needs ~500-50,000 cells
- Paired-end reads produce information about nucleosome positioning.
- Insert size distribution of fragments has a periodicity of ~200bp, suggesting that fragments are protected by multiples of nucleosomes
- Different fragmentation patterns can be associated with different functional states (eg. TSSs are more accessible than promoter flanking or transcribed regions)



[1] JD Buenrostro et al, *Nature Methods*, 2013.

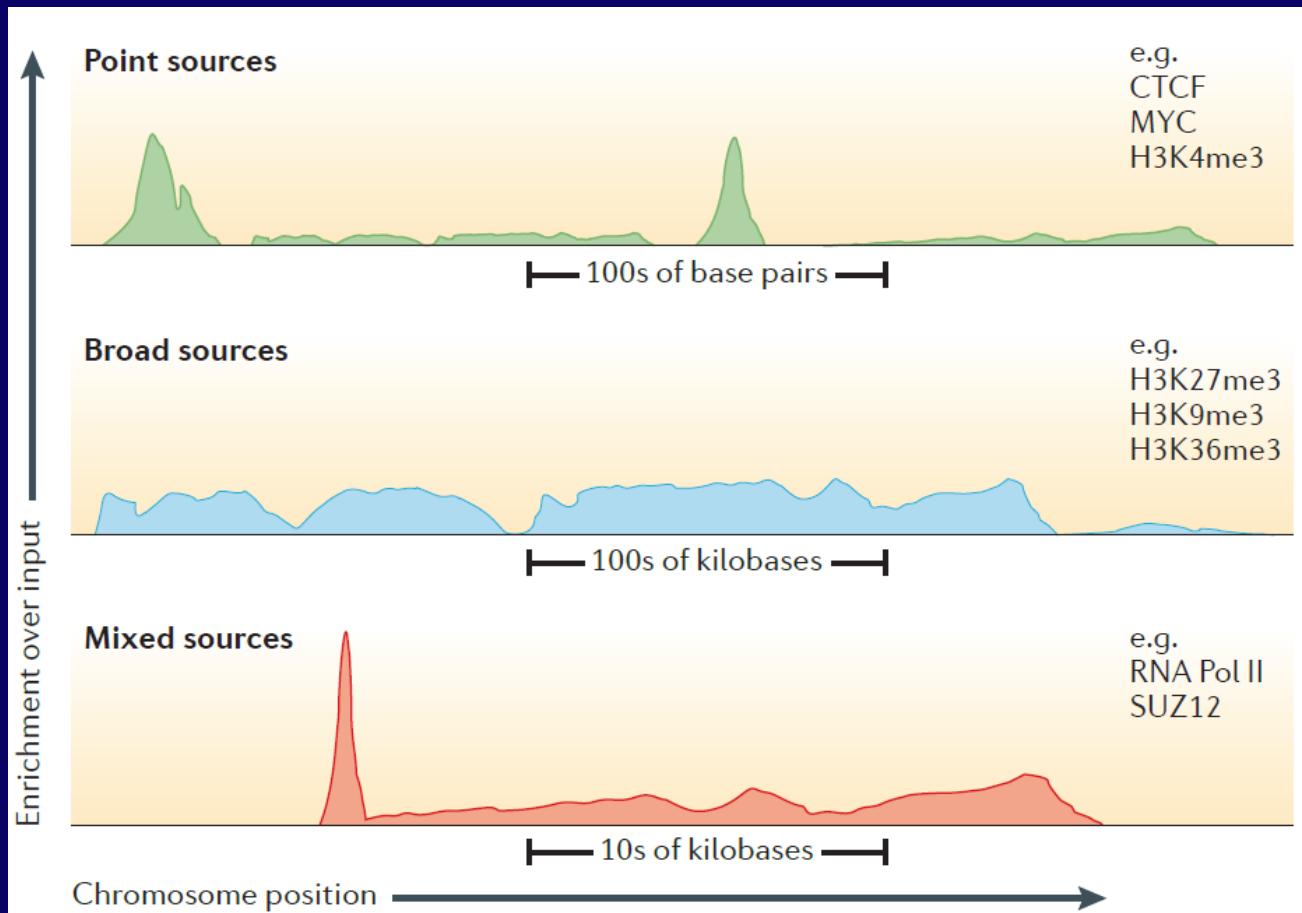
Using model-based methods to analyze ChIP-seq data

Outline

- Hidden Markov model for peak detection
- Hierarchical Hidden Markov model for combining ChIP-seq and ChIP-chip data, or analyze multiple ChIP-seq data
- Hybrid Monte Carlo strategy for Motif finding

Peak calling tools

- MACS
- HOMER
- CisGenome
- PeakSeq
- HPeak
- etc.



HPeak algorithm

Align reads to genome, get summary statistics, estimate model parameters.



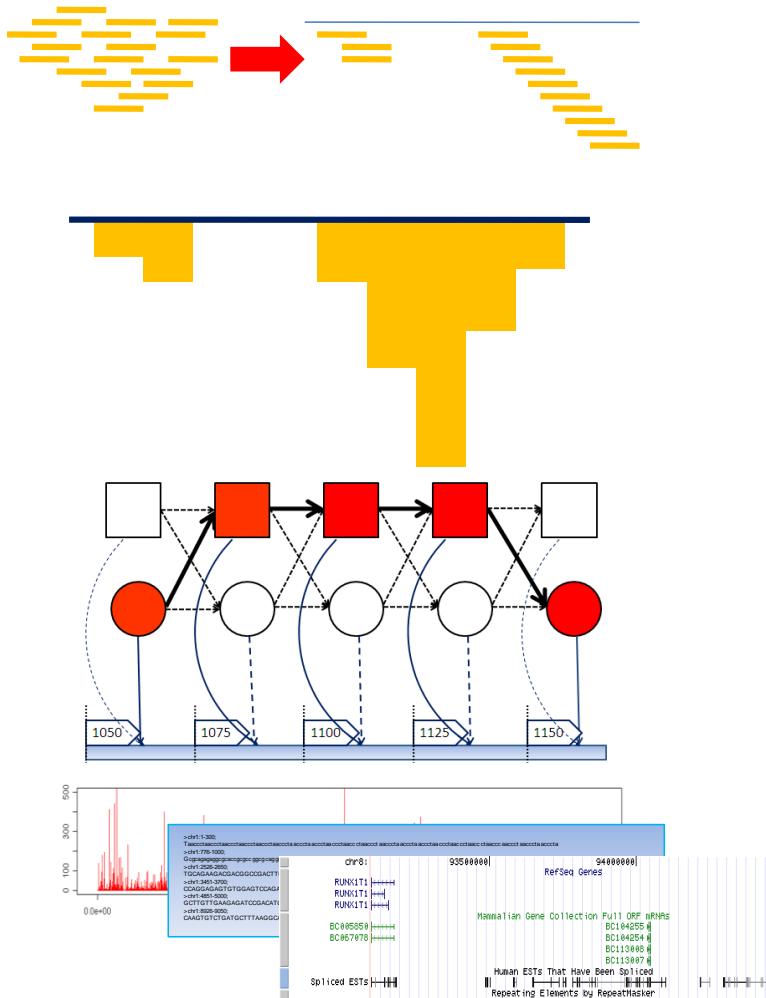
Get read coverage for each bin on all chromosomes.



Build HMM to infer whether a bin belongs to peak or background.



Post-processing on identified peaks.



GP and ZIP distribution

- Do not require mean equal to variance which is useful to model over-dispersion and under-dispersion.

$$P(Y = y | \lambda, \phi) = \left(\frac{\lambda}{1 + \phi\lambda} \right)^y \frac{(1 + \phi\lambda)^{y-1}}{y!} \exp \left\{ \frac{-\lambda(1 + \phi\lambda)}{1 + \phi\lambda} \right\}$$

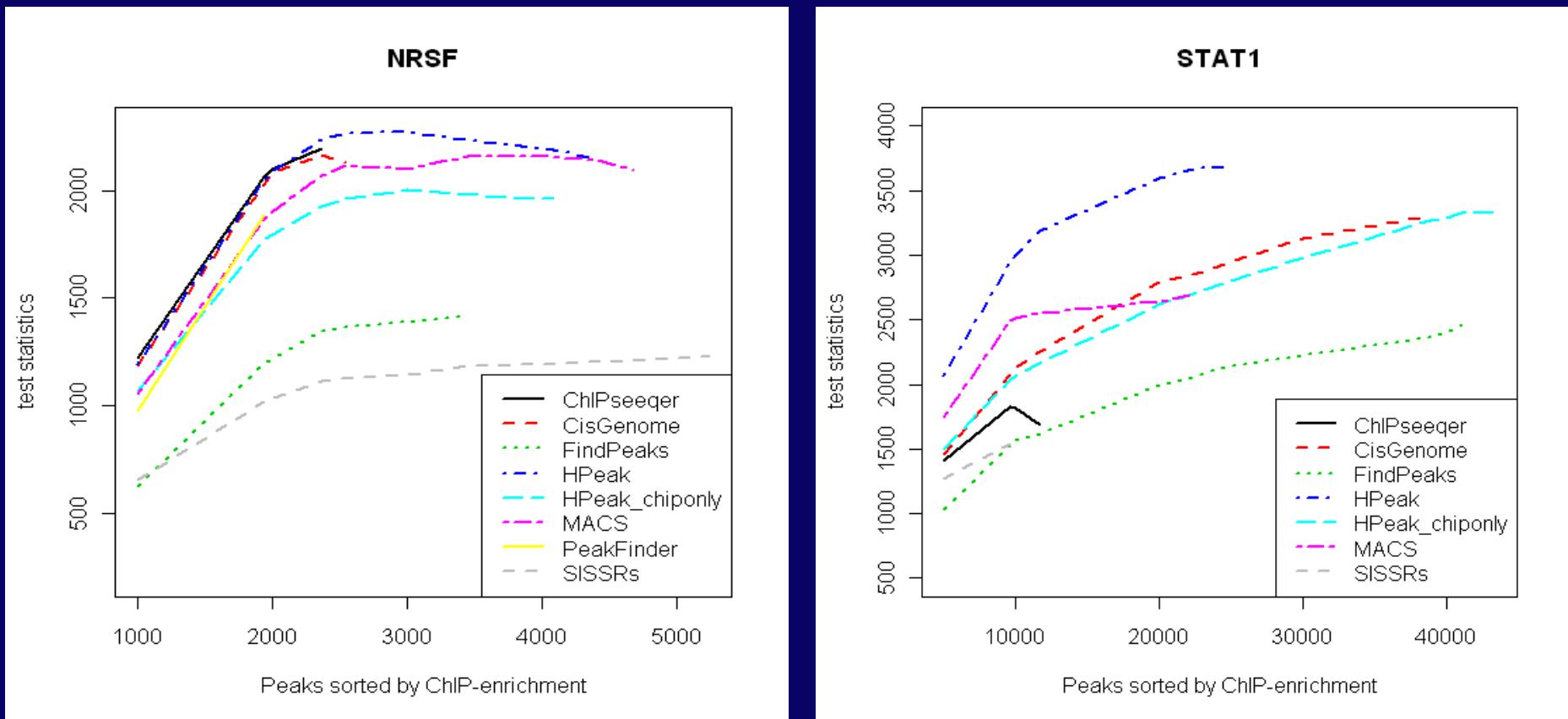
$$E(Y) = \lambda$$

$$Var(Y) = \lambda(1 + \phi\lambda)^2$$

- Zero-inflated Poisson distribution

$$f(Y | \pi, \mu) = \begin{cases} (1 - \pi) + \pi e^{-\mu} & \text{if } x = 0 \\ \frac{\pi e^{-\mu} \mu^x}{x!} & \text{if } x \neq 0 \end{cases}$$

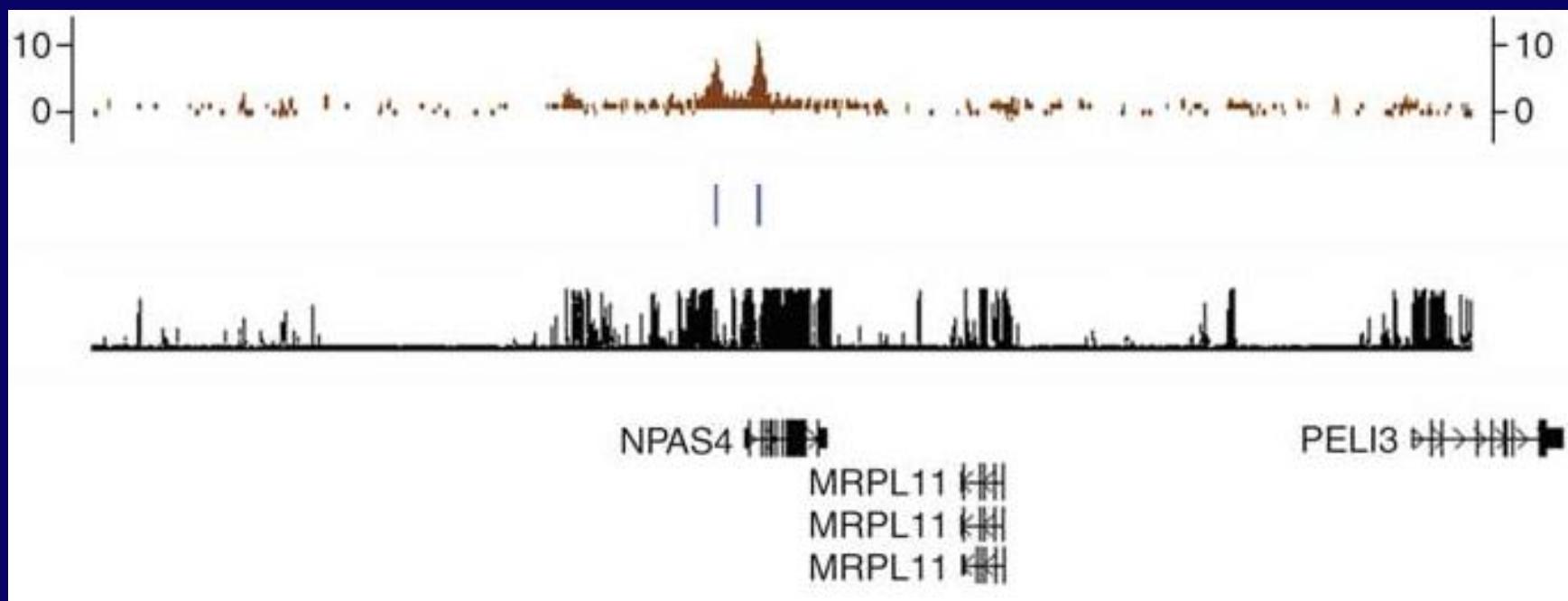
Motif enrichment results for NRSF and STAT1 data



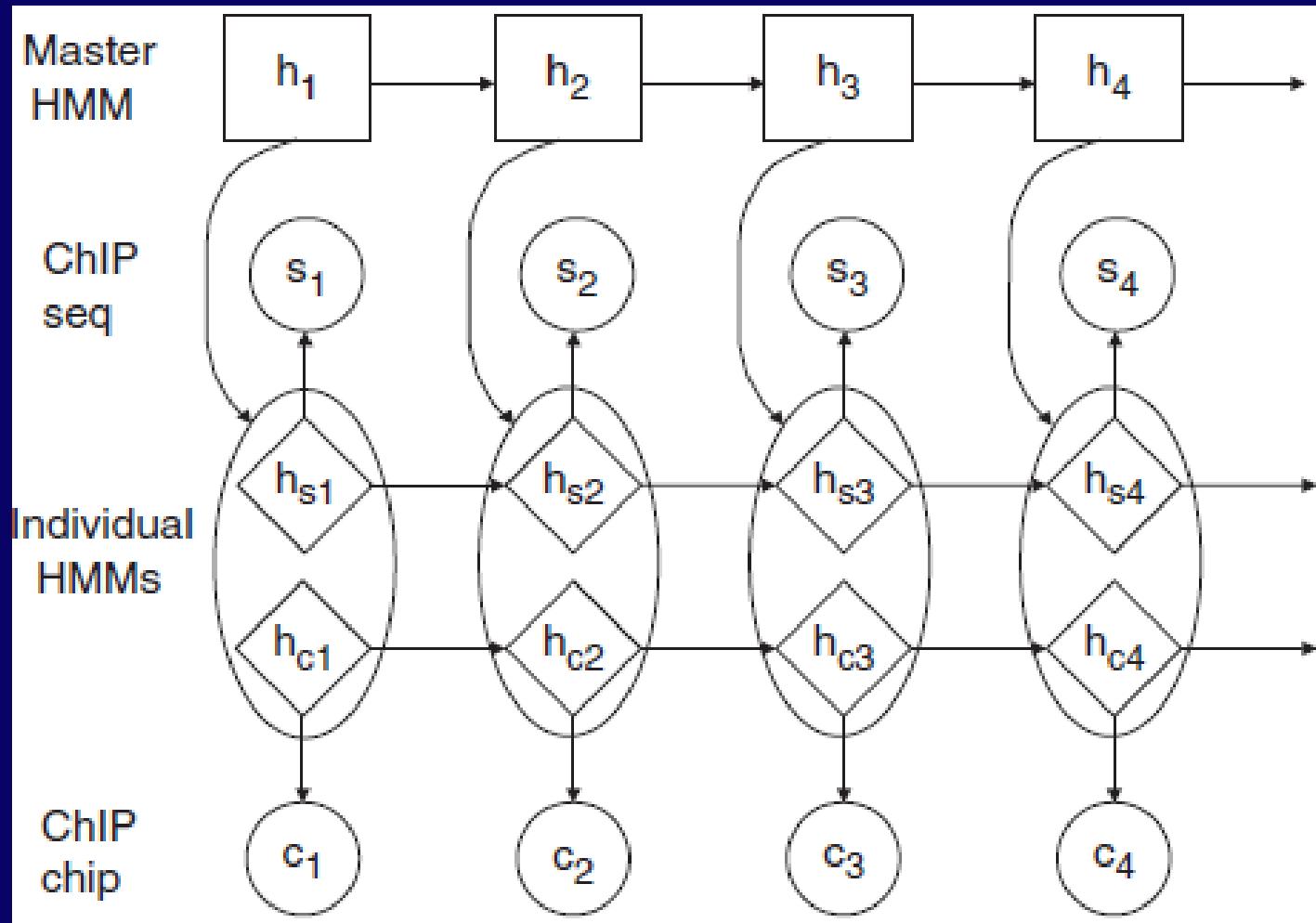
Outline

- Hidden Markov model for peak detection
- Hierarchical Hidden Markov model for combining ChIP-seq and ChIP-chip data
- Hybrid Monte Carlo strategy for Motif finding

Joint analysis of ChIP-chip and ChIP-seq



Hierarchical HMM

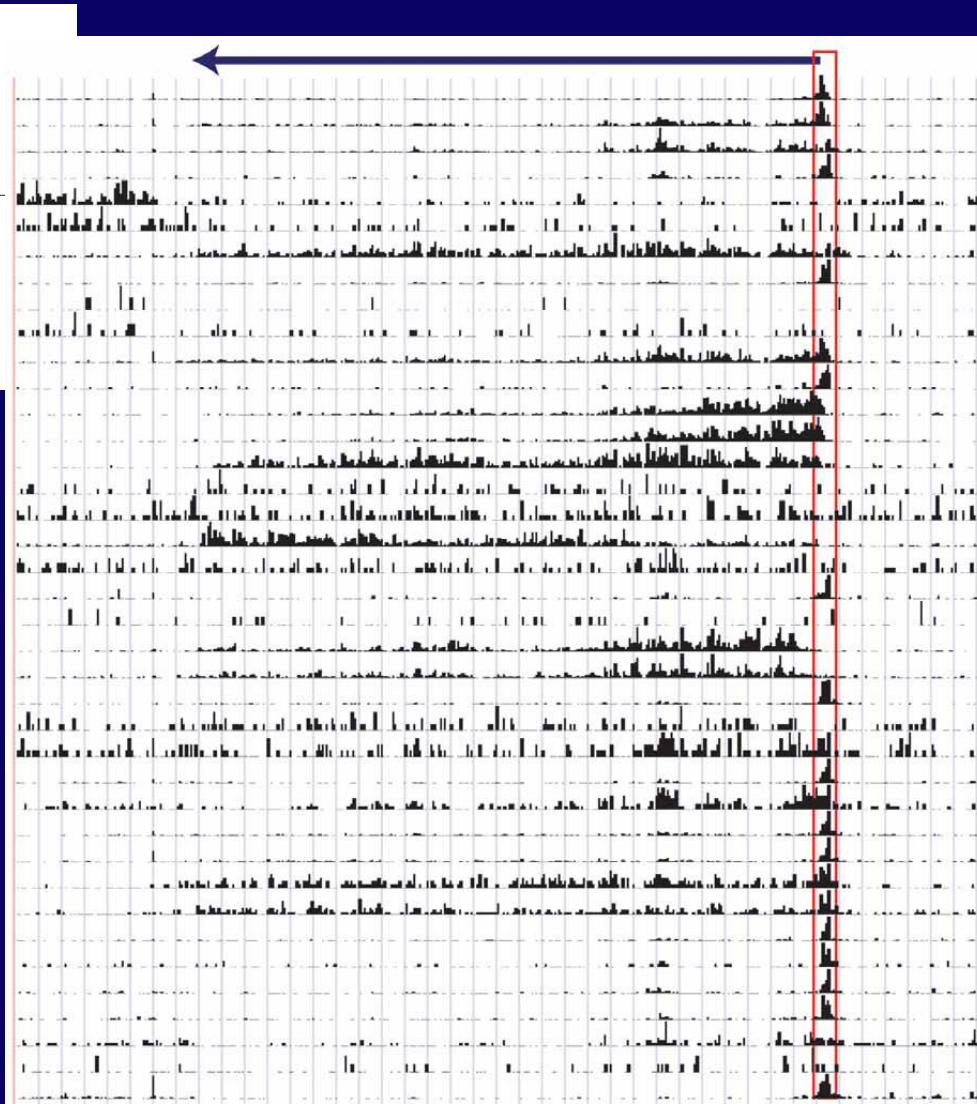


ChIP-Seq compendium

nature
genetics

Combinatorial patterns of histone acetylations and methylations in the human genome

Zhibin Wang^{1,5}, Chongzhi Zang^{2,5}, Jeffrey A Rosenfeld^{3–5}, Dustin E Schones¹, Artem Barski¹, Suresh Cuddapah¹, Kairong Cui¹, Tae-Young Roh¹, Weiqun Peng², Michael Q Zhang³ & Keji Zhao¹



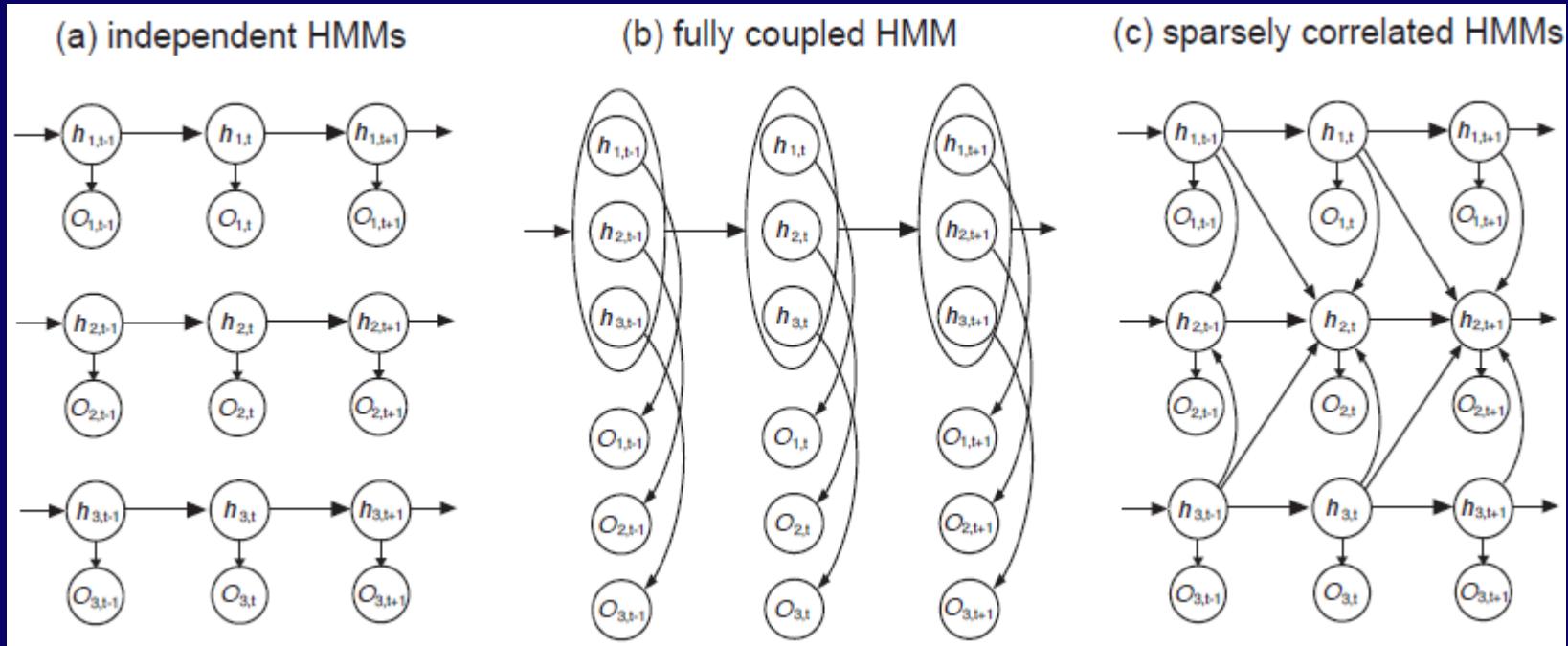
The problem

- N series of data, each can be modeled by an HMM,
- The goal is to infer the hidden states for all series,
- Suppose there are k states for each chain, then the total number of possible states for the whole datasets is k^N , the size of the transition matrix is k^{2N} ,
 - Independent: ignore correlation among the data series,
 - A single HMM to model all data together: intractable for large N .

Our goal

- Allow coupling among the chains,
- The goal is to borrow information across different experiments/datasets,
- Limit the amount of coupling allowed to reduce computation cost

Our scheme



Our learning plan

- Perform inference one series a time,
- Incorporate knowledge of hidden states in other series into the learning process,
- Assume sparsity in the correlation matrix.

Our model

- Use an inhomogeneous HMM to incorporate correlation,
- Define the transition kernel for series j and time t as:

$$K_j(t) = \begin{pmatrix} 1 - p_{jt} & p_{jt} \\ 1 - q_{jt} & q_{jt} \end{pmatrix}$$

$p_{jt} = \Pr(h_{j,t} = 1 | h_{j,t-1} = 0)$ and $q_{jt} = \Pr(h_{j,t} = 1 | h_{j,t-1} = 1)$.

$$\log\left(\frac{p_{jt}}{1 - p_{jt}}\right) = \beta_{j0}^p + \sum_{k \neq j} \left(\beta_{jk}^p h_{k,t-1} + \beta_{jk}^c h_{k,t} \right)$$

$$\log\left(\frac{q_{jt}}{1 - q_{jt}}\right) = \gamma_{j0}^p + \sum_{k \neq j} \left(\gamma_{jk}^p h_{k,t-1} + \gamma_{jk}^c h_{k,t} \right)$$

Our algorithm I

- Estimate regression parameters
 - Conditional on the current states, run penalized logistic regression to get model parameters,
 - LASSO penalty

$$\begin{aligned}y_t &= h_{j,t} \\x_t &= (h_{1,t-1}, \dots, h_{j-1,t-1}, h_{j+1,t-1}, \dots, h_{N,t-1}, h_{1,t}, \dots, h_{j-1,t}, h_{j+1,t}, \dots, h_{N,t})\end{aligned}$$

$$\min_{(\beta_{j0}, \vec{\beta}_j^p, \vec{\beta}_j^c)} \left\{ -\ell(\beta_{j0}, \vec{\beta}_j^p, \vec{\beta}_j^c) + \lambda P(\vec{\beta}_j^p, \vec{\beta}_j^c) \right\}$$

$$P(\vec{\beta}_j^p, \vec{\beta}_j^c) = \sum_{k \neq j} |\beta_{jk}^p| + \sum_{k \neq j} |\beta_{jk}^c|.$$

Our algorithm II

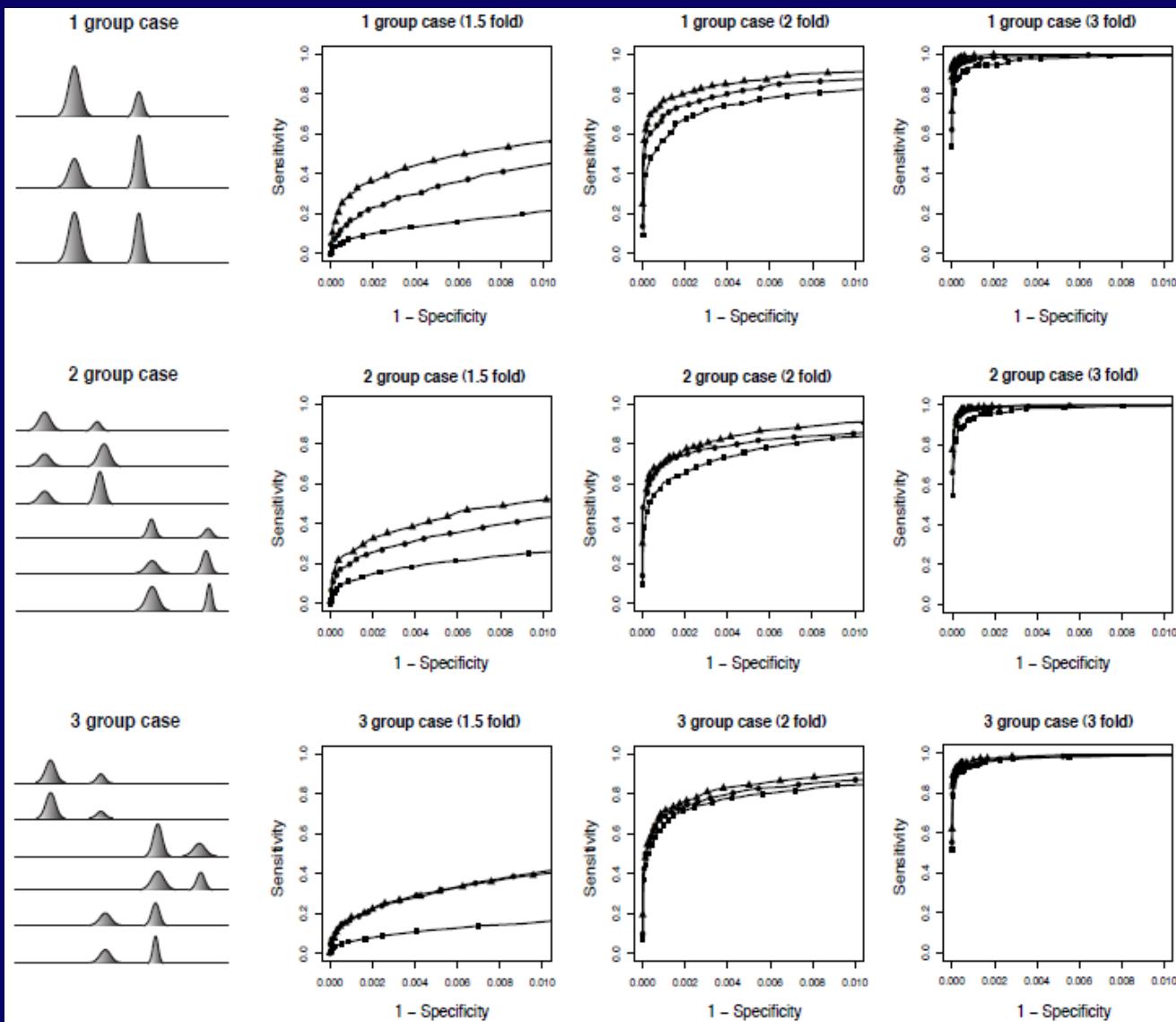
- Estimate transition kernel
 - Use the regression parameters estimated in step 1 and the current states of chains other than j , to get log odds for chain j at all time point t , then get estimated transition kernel.

$$\log \left(\frac{p_{jt}}{1 - p_{jt}} \right) = \beta_{j0}^p + \sum_{k \neq j} \left(\beta_{jk}^p h_{k,t-1} + \beta_{jk}^c h_{k,t} \right)$$
$$\log \left(\frac{q_{jt}}{1 - q_{jt}} \right) = \gamma_{j0}^p + \sum_{k \neq j} \left(\gamma_{jk}^p h_{k,t-1} + \gamma_{jk}^c h_{k,t} \right)$$

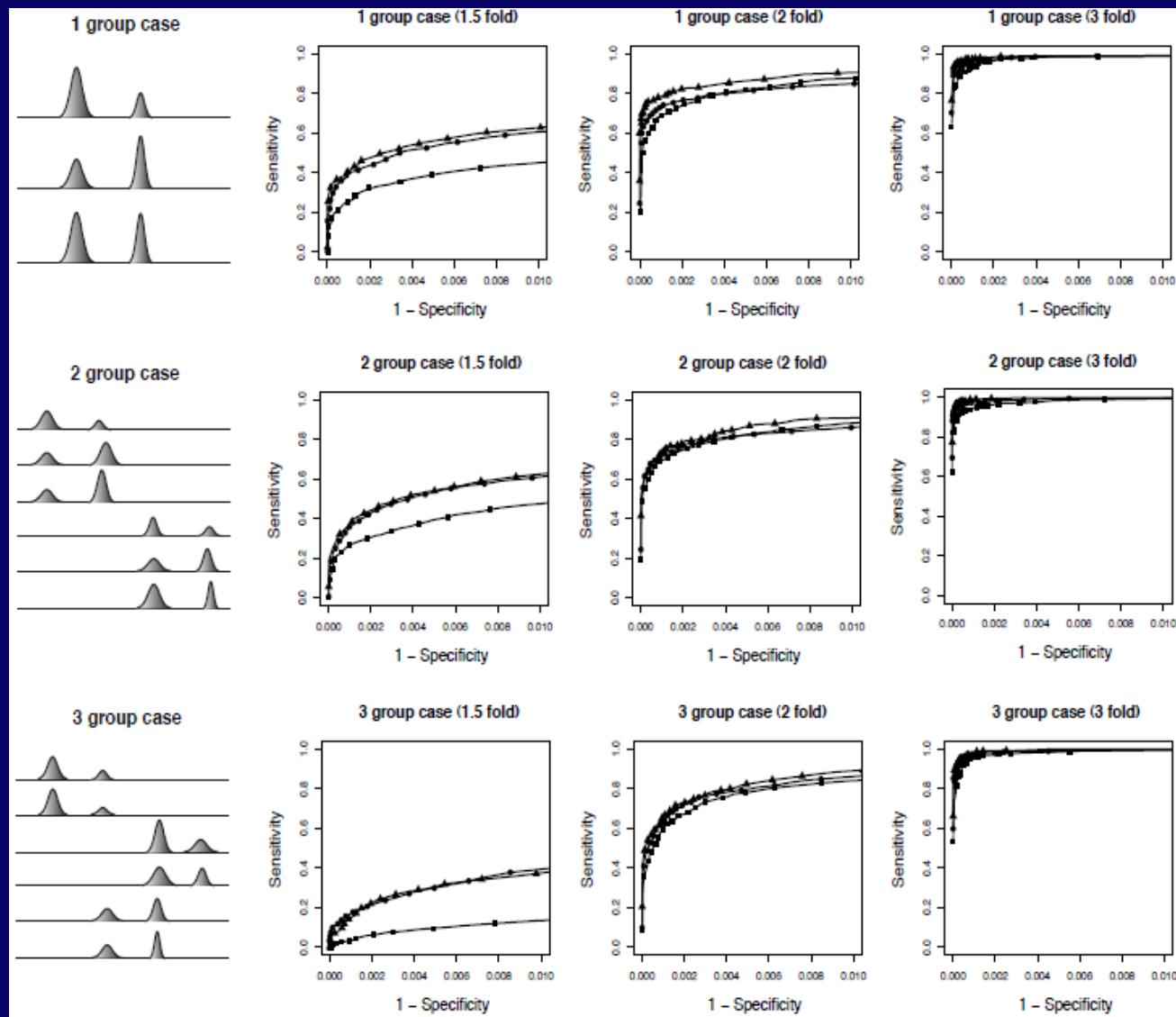
Our algorithm III

- Infer hidden states
 - Use the transition kernel estimated in step 2, current emission probabilities and observed data to run regular HMM (forward-backward algorithm) to get updated hidden states,
- Estimate the emission probabilities
 - Use the hidden states estimated in step 3 and observed data to update emission probabilities.

Simulation studies



Simulation studies



Joint inference of multiple ChIP-seq data

- JAMIE
 - Joint analysis of multiple ChIP-chip data
 - Wu, Ji Bioinformatics 2010
- HHMM
 - Joint analysis of ChIP-seq and ChIP-chip data
 - Choi et al. Bioinformatics 2009
- scHMM
 - Joint analysis of multiple ChIP-seq data
 - Choi et al. Bioinformatics 2013

Outline

- Hidden Markov model for peak detection
- Hierarchical Hidden Markov model for combining ChIP-seq and ChIP-chip data
- Hybrid Monte Carlo strategy for Motif finding

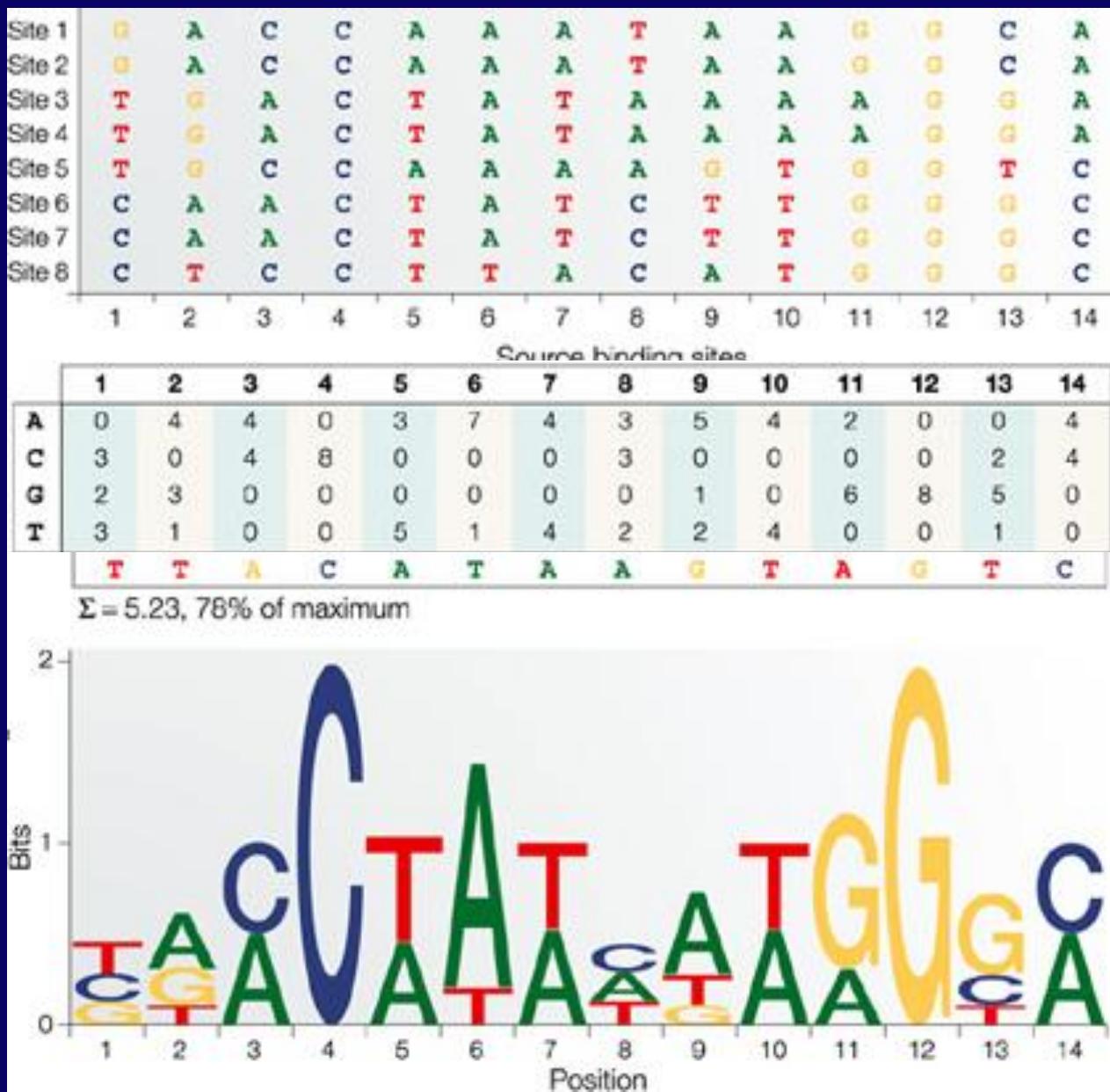
Example: cyclic receptor protein (CRP)

cole1	taatgtttgtctgtgtttgtggcatggcgagaatagcgcgtgtgaaagactgtttttgatcgccccacatgttgcac
ecoarabop	gacaaaaacgcgtAACAAAGTGTCTATAATCACGGCAGAAAAGTCACATTGATTATTCACGGCGTACACTTGCTATGCCATAGCATTATCCATAAG
ecobglr1	ACAAATCCCATAACTTAATTATGGGATTTGTATATATAACTTTATAAATTCTAAATTACACAAAGTTAATAACTGTGAGCATGGTCATATTTATAAT
ecocrp	CACAAAGCGAAAGCTATGCTAAACAGTCAGGATGCTACAGTAATAACATTGATGTACTGCATGTGCAAAGGACGTACATTACCGTGCAGTACAGTTGATAGC
ecocya	ACGGTGCCTACACTTGTATGCTAGCGCATCTTCTTACGGTCAATCAGCATGGTAAATTGATCACGTTAGACCATTTCGTCGTGAAACTAAAAAAAC
ecodecop	AGTGAATTATTGAACCAAGATCGCATTACAGTGATGCAAACCTGTAAGTAGATTCCCTTAATTGTGATGTGATCGAAGTGTGTTGCGGAGTAGATGTTAGAATA
ecogale	GCGCATAAAAAACGGCTAAATTCTGTGAAACGATTCCACTAATTATTCCATGTACACCTTCGCATCTTGTATGCTATGGTTATTCTACCCATAAGCC
ecoilvbpr	GCTCCGGCGGGTTTTGTTATCTGCATTCTAGTACAAAACGTGATCAACCCCTCAATTTCCTTGTGAAAATTTCATTGTCTCCCTGTTAAAGCTGT
ecolac	AACGCAATTATGTGAGTTAGCTCACTCATTAGGCACCCAGGCTTACACTTATGTTCCGGCTGTATGTTGTGTTGAATTGTGAGCGGATAACAAATTCAAC
ecomale	ACATTACCGCCAATTCTGTAAACAGAGATCACACAAAGCGACGGTGGGCGTAGGGGCAAGGAGGATGGAAAGAGGTTGGCTATAAAAGAAACTAGAGTCGTTA
ecomalk	GGAGGAGGCGGGAGGATGAGAACACGGCTTCTGTGAACTAAACCGAGGTCTGTGATGTTGGCTGTTGCTGCAAAATCGTGGCGATTTATGTGCGCA
ecomalt	GATCAGCGTCCTTAGGTGAGTGTAAATAAGATTGGATTGTGACACAGTCAGACACATAAAAACGTACATCGCTTCGATTAGAAAGGTTCT
ecoompa	GCTGACAAAAAGATTAAACACACCTTACACAGACTTTTCTATGCTGACGGAGTTACACTTGAAGTTCAACTACGTTAGACTTACATCGCC
ecotnaa	TTTTTAAACATTAAATTCTACGTAAATTATACTTTAAAGCATTAATATTGTCCTTGAACGATTGTGATTGATTCACTTAAACAATTTCAGA
ecouxul	CCCATGAGAGTGAATTGTTGTGATGTGGTTAACCAATTAGAATTCTGGGATTGACATGCTTACCAAAAGGTAGAACTTACGCCATCTCATCGATGCAAGC
pbr-p4	CTGGCTTAACTATGCGGCATCAGAGCAGATTGACTGAGAGTGCACCATATGCGGTGTGAAATACCGCACAGATGCGTAAGGAGAAAATACCGCATCAGCGCTC
trn9cat (tdc)	CTGTGACGGAAGATCCTCGCAGAATAAAATCTGGTGTCCCTGTTGATACCGGAAAGCCCTGGCCAACTTTGGCAAAATGAGACGTTGATCGGCACG GATTTTATACTTAACTTGTGATTTAAAGGTATTAACTTAAACGATACTCTGGAAAGTATTGAAAGTTAATTGTGAGTGGTCGACATATCCTGTT

Example: cyclic receptor protein (CRP)

cole1	taatgtttgtctgtgtttgtggcatggcgagaatagcgcgtgtgaaagactgttttttgatcgaaaaatggaaagtccacagtcttgcac
ecoarabop	gacaaaaacgcgtAACAAAGTGTCTATAATCACGGCAGAAAAGTCACATTGATTGGCATCGGCACACTTGCTATGCCATAGCATTATCCATAAG
ecobglr1	ACAAATCCCATAACTTAATTATGGGATTTGTATATATAACTTATAAATTCTAAATTACACAAAGTTAAACACGTTAAACAGTGTGAGCATGGTCAATTATCAAT
ecocrp	CACAAAGCGAAAGCTATGCTAAACAGTCAGGATGCTACAGTAATAACATTGATGTACTGCATGTATGCAAAGGACGTACATTACCGTGCAGTACAGTTGATAGC
ecocya	ACGGTGTACACTTGTATGTAGCGCATCTTCTTACGGTCAATCAGCATGGTTAAATGATCACGTTAGACCATTTCGTCGTGAAACTAAAAAAAC
ecodecop	AGTGAATTATTTGAACCAGATCGCATTACAGTGTGCAAACTGTAAAGTAGATTCCTTAATTGTGATGTATCGAAGTGTGTTGCGGAGTAGATGTTAGAATA
ecogale	GCGCATAAAAAACGGCTAAATTCTTGTAAACGATTCCACTAAATTATTCCATGTACACATTTCGCATCTTGTATGCTATGGTTATTCTACCCATAAGCC
ecoilvbpr	GCTCCGGCGGGTTTTGTTATCTGCAATTCACTGACAAACCGTGTCAACCCCTCAATTTCCTTTGCTGAAAATTTCATTGTCCTCCCTGTTAAAGCTGT
ecolac	AACGCAATTAAATGTGAGTTAGCTCACTCATAGGCACCCCCAGGCTTACACTTATGTTCCGGCTGTATGTTGTGTTGAATTGTGAGCGGATAACAAATTCAAC
ecomale	ACATTACCGCAATTCTGTAAACAGAGATCACACAAAGCGACGGTGGGCGTAGGGGCAAGGAGGATGGAAAGAGGTTGGCTATAAAAGAAACTAGAGTCGTTA
ecomalk	GGAGGAGGCGGGAGGATGAGAACACGGCTTCTGTGAACCAAACCGAGGTGTAGGAATTCTGTGATGTTGCTTGCACAAATCGTGGCGATTTATGTGCGCA
ecomalt	GATCAGCGTCCTTCTGTGAGTGTAAATAAAGATTGGATTGTGACACAGTCATGCAATTAGACACATAAAAACGTCATCGCTGCTTGTGCGCA
ecoompa	GCTGACAAAAAAAGATTAAACACACCTTACACAGACTTTCTCATATGCTGACGGAGTTACACTTGTAAAGTTCAACTACGTTGAGACTTACATGCC
ecotnaa	TTTTTAAACATTAAATTCTTACGTAAATTATACTTTAAAGCATTAATATTGTCCTTCAACGACGTTGTTAAACATTGTGATTGCTTGTGAGACTTACATGCC
ecouxul	CCCATGAGAGTGAATTGTGATGTTAACCCAAATTAGAATTCTGGGATTGACATGCTTACAAAGAGTGTGAAACTTACGCCATCTCATCGATGCAAGC
pbr-p4	CTGGCTTAACTATGCGGCATCAGAGCAGATTGACTGAGAGTGCACCATATGCGGTGTGAAATACCGCACAGATGCGTAAGGAGAAAATACCGCATCAGCGCTC
trn9cat (tdc)	CTGTGACGGAAGATCACTCGCAGAATAAAATCTGGTGTCCCTGTTGATACCGGAAAGCCCTGGCCAACTTTGGCAAAATGAGACGTTGATCGGCACG GATTTTATACTTAACTTGTGATTTAAAGGTATTAAACGATACTCTGGAAAGTATTGAAAGTTGTTGAGTGGTCGACATATCCTGTT

Transcription factor binding site (TFBS)



Existing *de novo* motif finding algorithms

- Consensus Hertz *et al.* 1990
 - Gibbs Motif Sampler Lawrence *et al.* 1993
 - MEME Bailey and Elkan 1994
 - AlignACE Roth *et al.* 1998
 - BioProspector Liu *et al.* 2001
 - MDScan Liu *et al.* 2002
 - Mobydick Bussemaker *et al.* 2000
- ...
- Review Tompa *et al.* 2005

Motif identification model

a_1
aaagg t cgag t agctactcgatcgataactagcaatcg t taccctagctcgatcgaaa
 a_2
acgtgagatcagctatgaccg t agctactcgataaccg
 a_3
gaat t agctactcgatcgataactagcaatcg t taccctagctcgagatggaaagactataa
...
 a_J
acgtgagatcagctatcgatcgattg t aactactcgatcgat

Alignment variable $A = \{a_1, a_2, \dots, a_J\}$

Posterior distributions

- The posterior conditional distribution for alignment variable \mathbf{A}

$$p(a_j = l | \boldsymbol{\theta}_0, \boldsymbol{\Theta}, \mathbf{R}_j, A_{-j}) \propto \prod_{k=1}^4 \theta_{0k}^{h_k(\mathbf{R}_j)} \prod_{i=1}^w \prod_{k=1}^4 \left(\frac{\theta_{ik}}{\theta_{0k}} \right)^{h_k(r_{j,l+i-1})} \propto \prod_{i=1}^w \prod_{k=1}^4 \left(\frac{\theta_{ik}}{\theta_{0k}} \right)^{h_k(r_{j,l+i-1})}$$

DNA sequence data $\mathbf{R} = (\mathbf{R}_1, \dots, \mathbf{R}_J)$

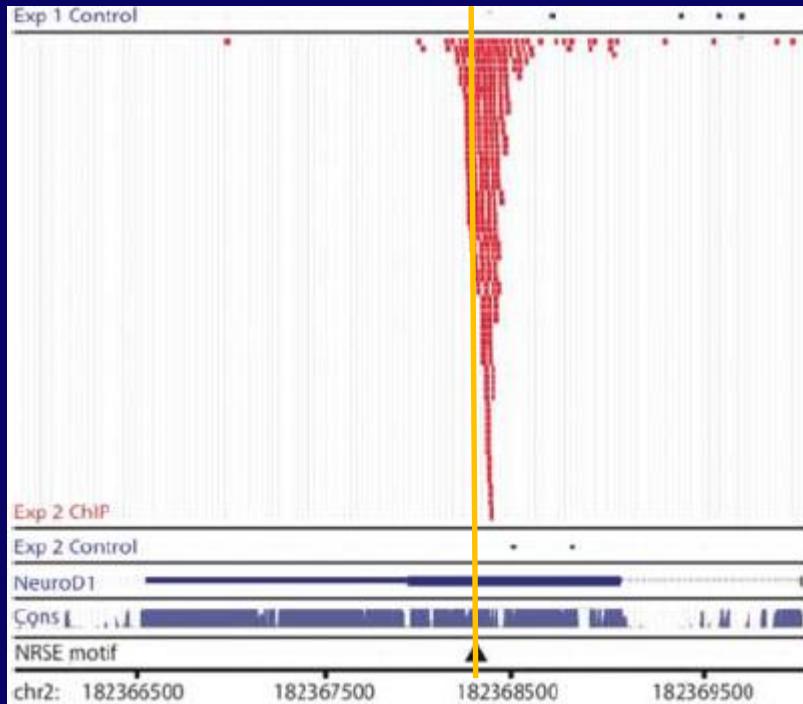
Lawrence *et al.* *Science* 1993, Liu *et al.* *JASA* 1995

Why de novo motif search

- The only option when the TF binding motif pattern is unknown.
- Reassuring to be able to rediscover the known TFBS motif.
- Many “known” motif patterns are biased and inaccurate.
- Multiple co-factors are often required in transcription regulation in eukaryotes.
- Binding specificity for some TFs may change under different conditions.

Challenges faced

- How to handle large number of input sequences?
- How to utilize sequencing depth information?



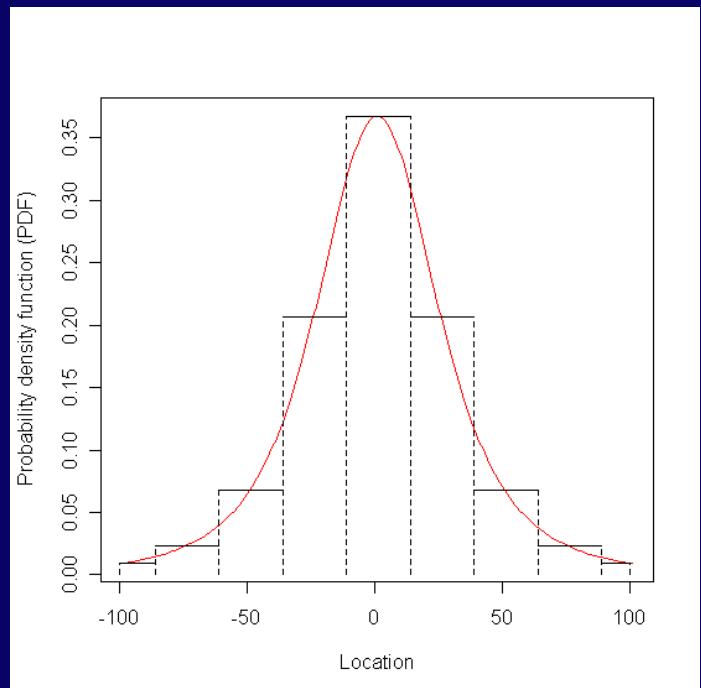
Features of our new algorithm

- Incorporate sequencing depth information in the statistical model.
- Generalize the product multinomial model to allow inter-dependent positions within the motif.
- Adopt a hybrid Monte Carlo strategy to speed up the traditional Gibbs sampler-based algorithm.

The informative prior

- The prior is symmetric and centered at the peak summit.
- The prior probabilities stem from Student's t -distribution with $\text{df}=3$.

$$p(a_j = l) \propto t_3 \left(\text{int} \left[\frac{|l + w/2 - s_j| + u/2}{u} \right] \right)$$

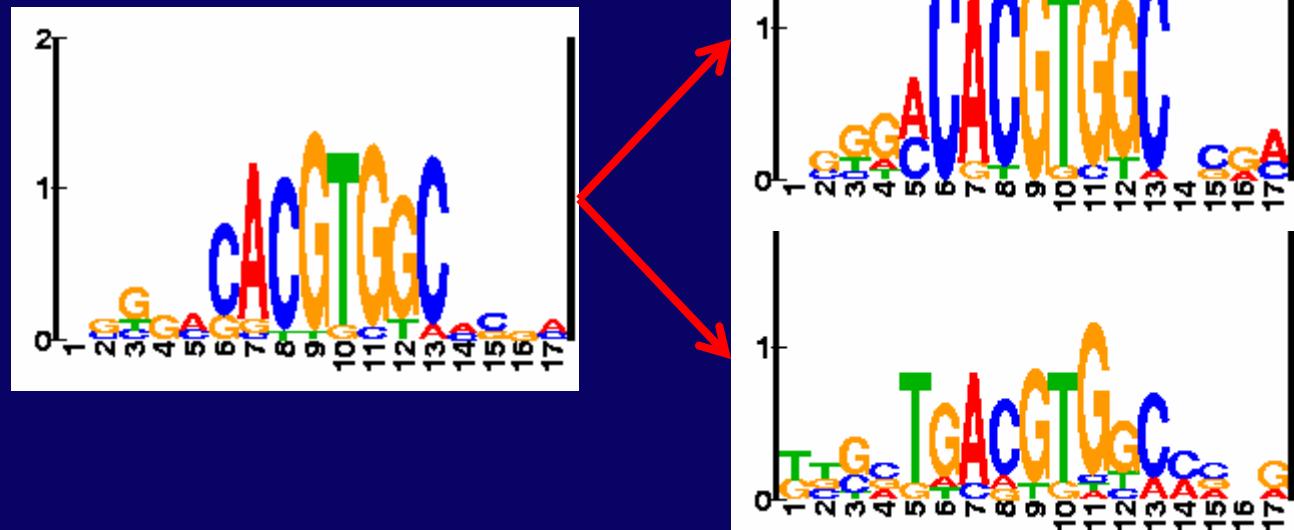


Modeling inter-dependent positions

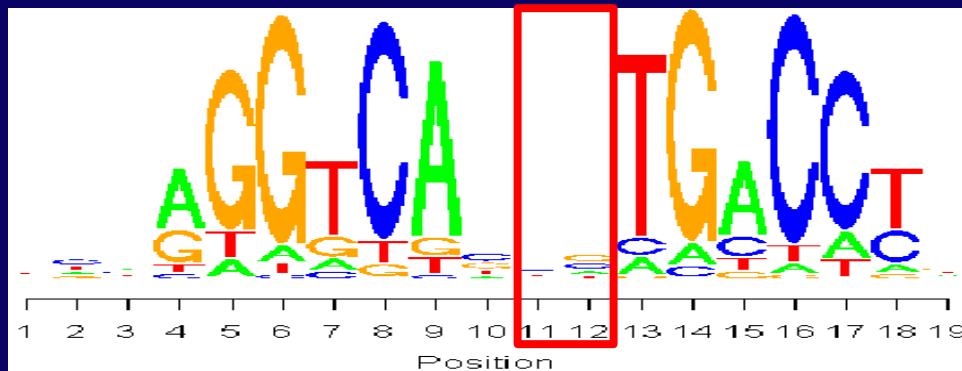
- Zhou and Liu
Bioinformatics 2005



- Barash *et al.*
RECOMB 2003

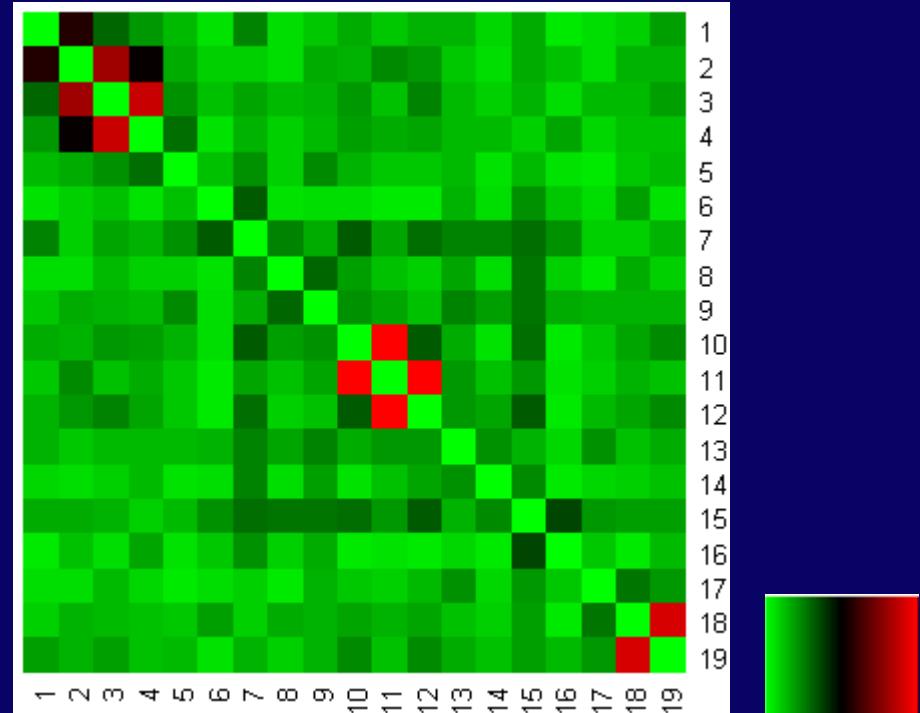


Detect intra-dependent position pairs



$$d_{ij} = \sum_{x=1}^4 \sum_{y=1}^4 |\hat{\eta}_{xy}(r_i, r_j) - \hat{\eta}_x(r_i)\hat{\eta}_y(r_j)|$$

	A	C	T	G	
A	0.03 (0.04)	0.15 (0.25)	0.28 (0.16)	0.03 (0.03)	0.49
C	0.00 (0.00)	0.01 (0.01)	0.00 (0.00)	0.00 (0.00)	0.01
T	0.05 (0.04)	0.34 (0.24)	0.06 (0.17)	0.03 (0.03)	0.48
G	0.00 (0.00)	0.02 (0.01)	0.00 (0.01)	0.00 (0.00)	0.02
	0.08	0.52	0.34	0.06	1



Prioritized hybrid Monte Carlo

- Subject each sequence to either stochastic sampling or greedy search.
- Input sequences are not created equal.
- ChIP-enrichment is indicative of binding affinity.

Implementation

- Hybrid Motif Sampler (HMS).
- Gibbs sampler type iterative procedure.
- Run multiple chains to avoid trapping in local mode.

Performance comparison

- Two established and popular motif discovery tools:
 - MEME (Bailey and Elkan 1994),
 - EM-based motif finding algorithm,
 - widely used.
 - MDscan (Liu *et al.* 2002),
 - designed to analyze ChIP-chip data,
 - combines word enumeration and probability matrix updating,
 - take into account ChIP-chip ranking,
 - very fast.

Real data analysis

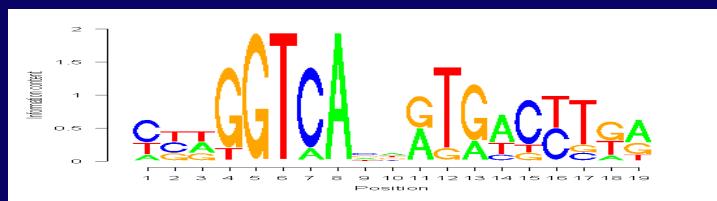
TF	Cell type	Antibody	# of peaks	Coverage	Reference
		Monoclonal			
NRSF	Jurkat T cell	12C11	4,982	1.4 MB	Johnson et al. (2007)
STAT1	HeLa S3 cell	Polyclonal	27,470	8.1 MB	Robertson et al. (2007)
CTCF	CD4+ T cell	Upstate 07-729	22,159	7.4 MB	Barski et al. (2007)
ER	MCF7 cell	ER α (HC-20)	10,072	2.5 MB	

Performance evaluation

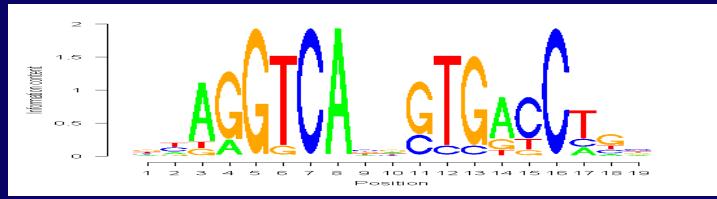
- Cross validation
 - Randomly separate all peaks into two halves: training and testing.
 - Run motif finding algorithms on the training data to predict the motif pattern.
 - Scan testing data using the identified motif pattern and compare to a set of control sequences.
- Testing
 - Using Chi-square test statistics to quantify motif enrichment .
 - Estimate FDR and plot FDR versus Chi-square test statistics.

Compare ER motif patterns

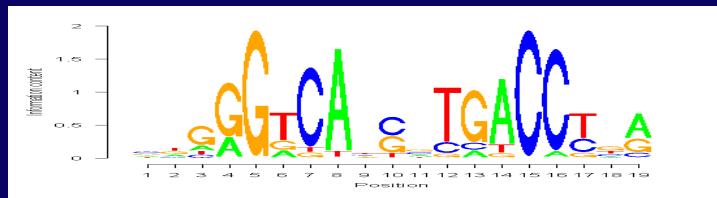
- V\$ER01*



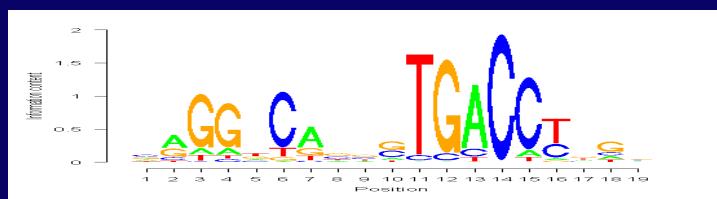
- V\$ER02*



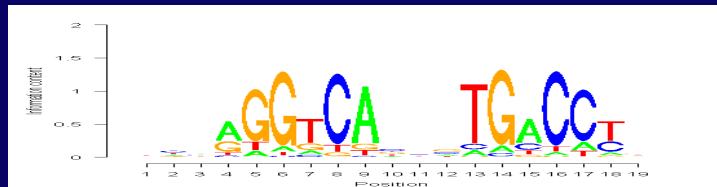
- V\$ER03*



- MEME

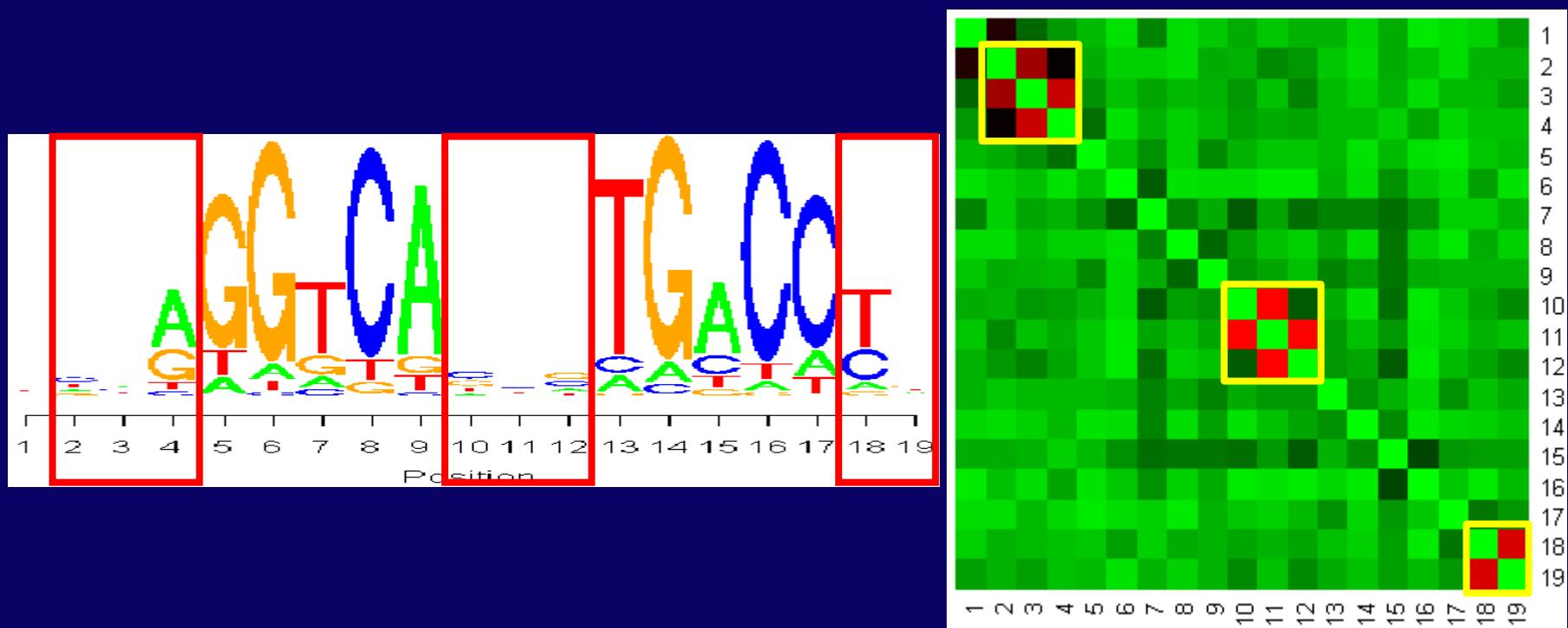


- HMS

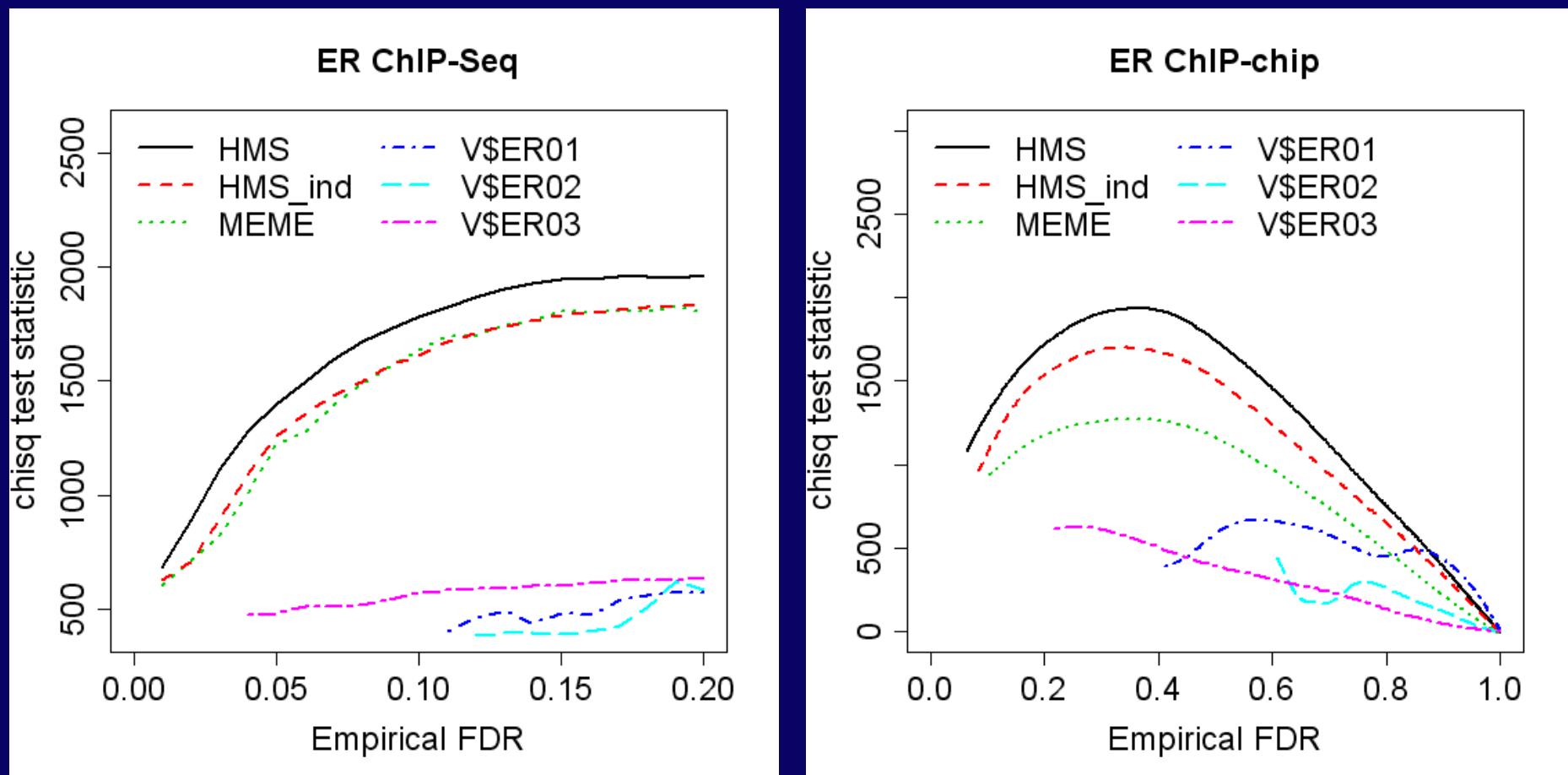


*

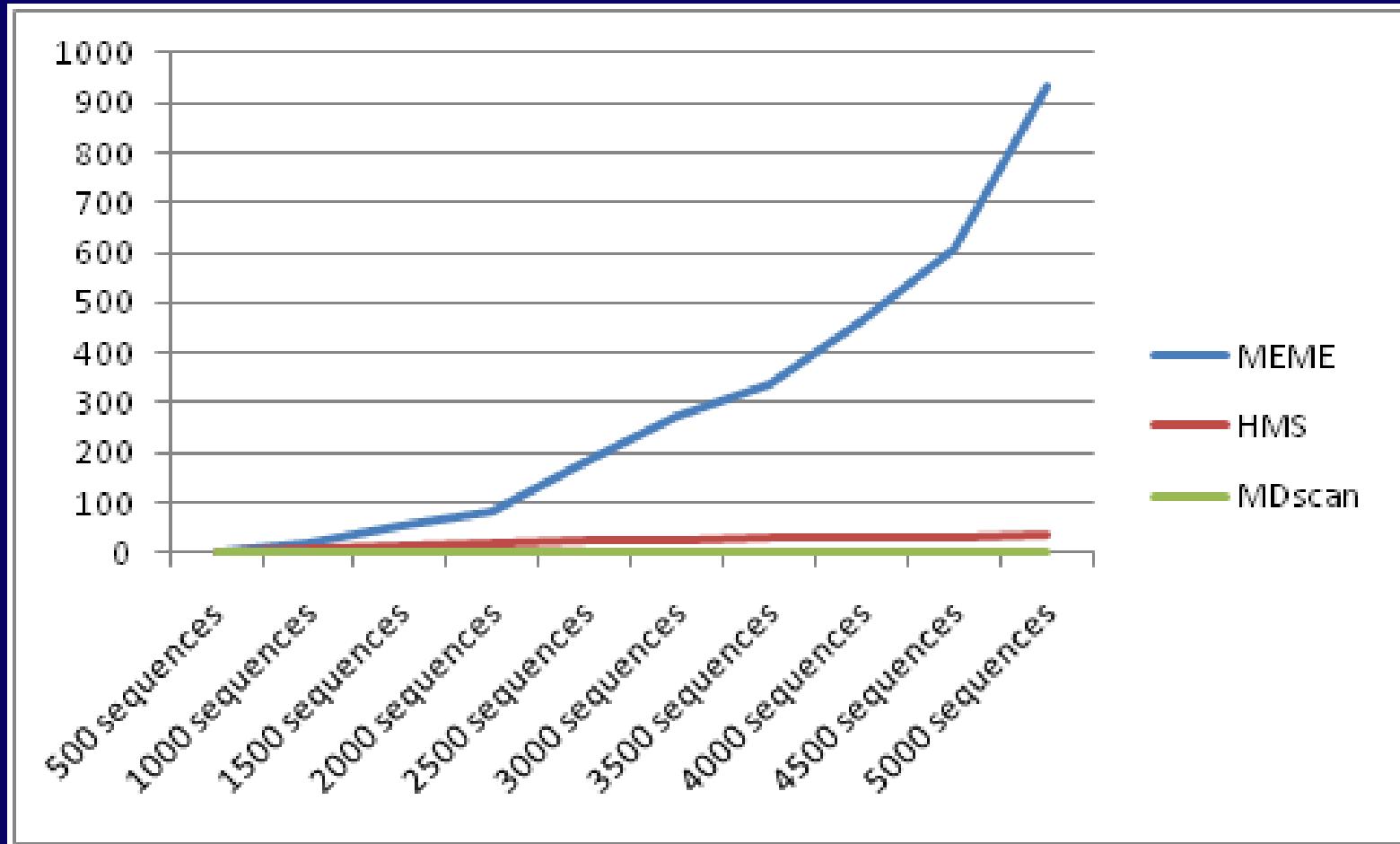
Positions show inter-dependency inside the ER motif



Compare ER motif enrichment



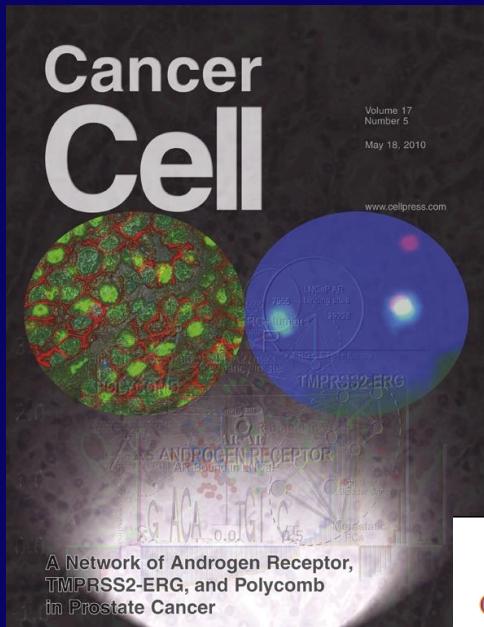
Computation time



Summary

- ChIP-Seq data offers abundant information and provides much improved opportunity for studying protein-DNA interaction.
- There are many biological and technical factors that affect the ChIP-Seq data we observe, careful modeling is critical in order to process ChIP-Seq data efficiently and thoroughly.
- New sequencing data are different from microarray, ChIP-chip data. Methods developed there do not work well for analyzing sequencing data, new models and algorithms need to be developed.

Apply to cancer genomics



Cancer Cell
Article

Cell
PRESS

An Integrated Network of Androgen Receptor, Polycomb, and TMPRSS2-ERG Gene Fusions in Prostate Cancer Progression

Jindan Yu,^{1,3,6,7} Jianjun Yu,^{1,3} Ram-Shankar Mani,^{1,3} Qi Cao,^{1,3} Chad J. Brenner,^{1,3} Xuhong Cao,^{1,2,3} Xiaoju Wang,^{1,3} Longtao Wu,⁷ James Li,^{1,3} Ming Hu,^{1,5} Yusong Gong,^{1,3} Hong Cheng,^{1,3} Bharathi Laxman,^{1,3} Adaikkalam Vellaichamy,^{1,3} Sunita Shankar,^{1,3} Yong Li,^{1,3} Saravana M. Dhanasekaran,^{1,3} Roger Morey,^{1,3} Terrence Barrette,^{1,3} Robert J. Lonigro,^{1,6} Scott A. Tomlins,^{1,3} Sooryanarayana Varambally,^{1,3,6} Zhaohui S. Qin,⁵ and Arul M. Chinnaiyan^{1,2,3,4,6,*}

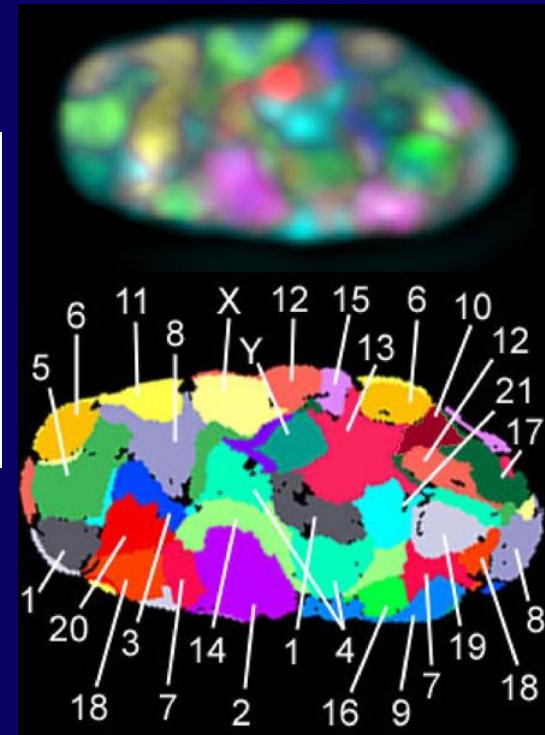
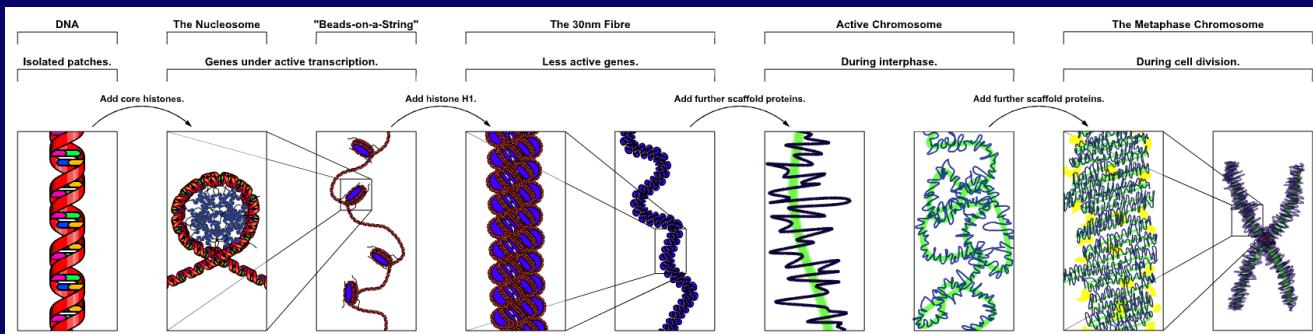
Reference

- Qin ZS, Yu J, Shen J, Maher CA, Hu M, Kalyana-Sundaram S, Yu J, Chinnaiyan AM. (2009) HPeak: An HMM-based Algorithm for Defining Read-enriched Regions in ChIP-Seq Data. *BMC Bioinformatics*. **11** 369.
<http://www.sph.umich.edu/csg/qin/HPeak/>
- Choi H, Nesvizhskii A, Ghosh D, Qin ZS. (2009) Hierarchical Hidden Markov Model with Application to Joint Analysis of ChIP-chip and ChIP-seq Data. *Bioinformatics* **25** 1715-1721.
<http://sourceforge.net/projects/chipmeta/>
- Hu M, Yu J, Taylor, JMG, Chinnaiyan AM, **Qin ZS**. (2010) On the Detection and Refinement of Transcription Factor Binding Sites Using ChIP-Seq Data. *Nucleic Acids Res.* **38** 2154-2167.
<http://www.sph.umich.edu/csg/qin/HMS/>
- Hu M, Zhu Y, Taylor JMG, Liu JS, Qin ZS (2011). Using Poisson mixed-effects model to quantify exon-level gene expression in RNA-seq. *Bioinformatics*. **28** 63-68.
<http://www.stat.purdue.edu/~yuzhu/pome.html>

Statistical model to infer chromosomal structures from Hi-C data

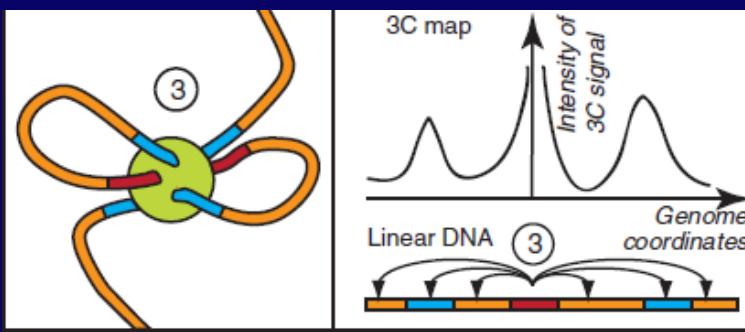
Chromosome folding

How can a two meter long polymer fit into a nucleus of ten micrometer (10^{-5} m) diameter?



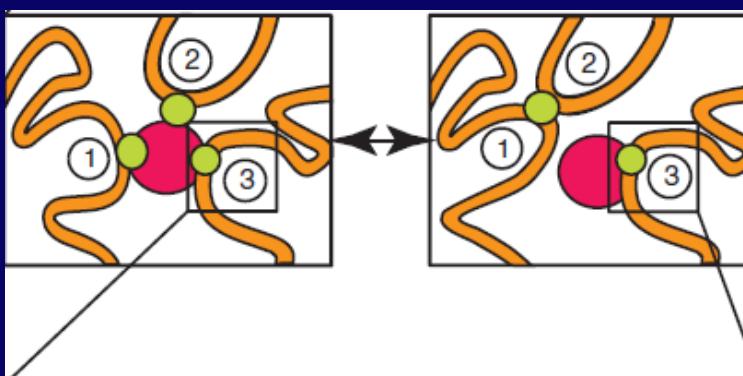
Chromosome Conformation Capture (3C)

Dekker et al. *Science* 2002

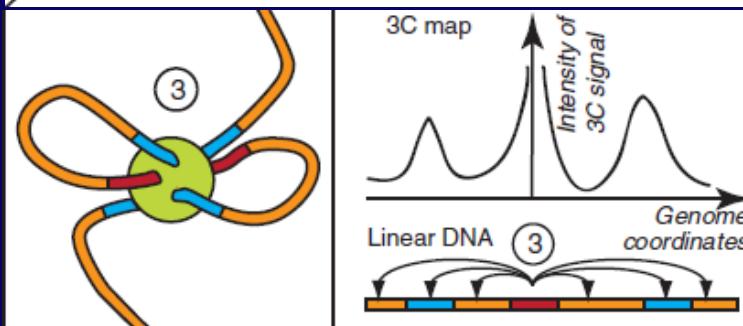


Fine scale: (0-kb)

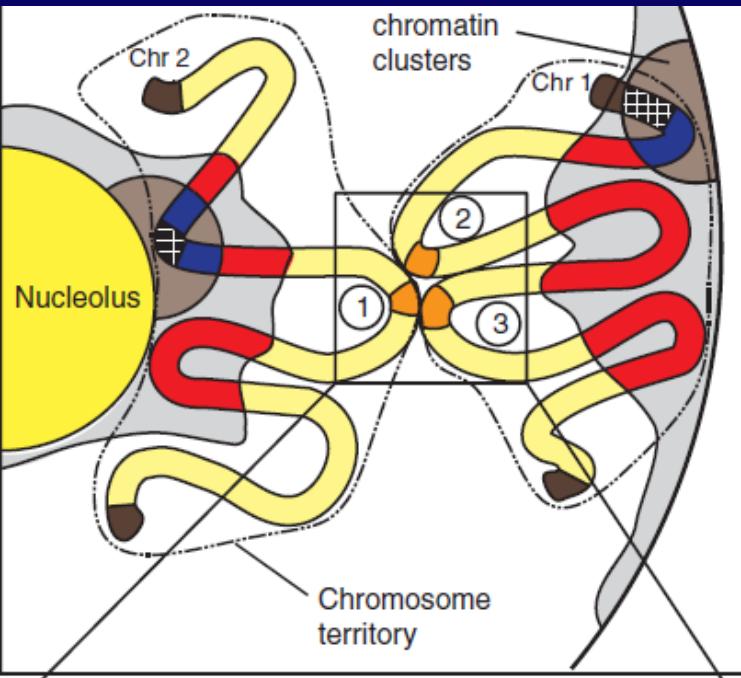
3C-on-chip/Circular 3C (4C) 5C



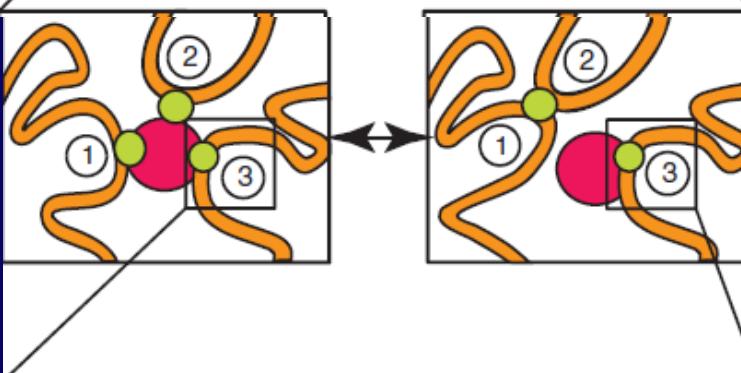
Intermediate: (0-Mb)



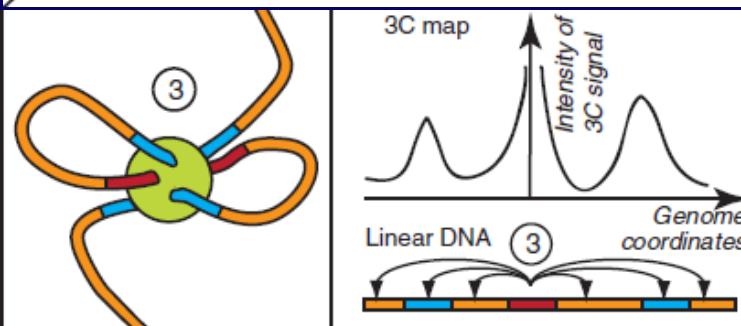
Fine scale: (0-kb)



Whole genome



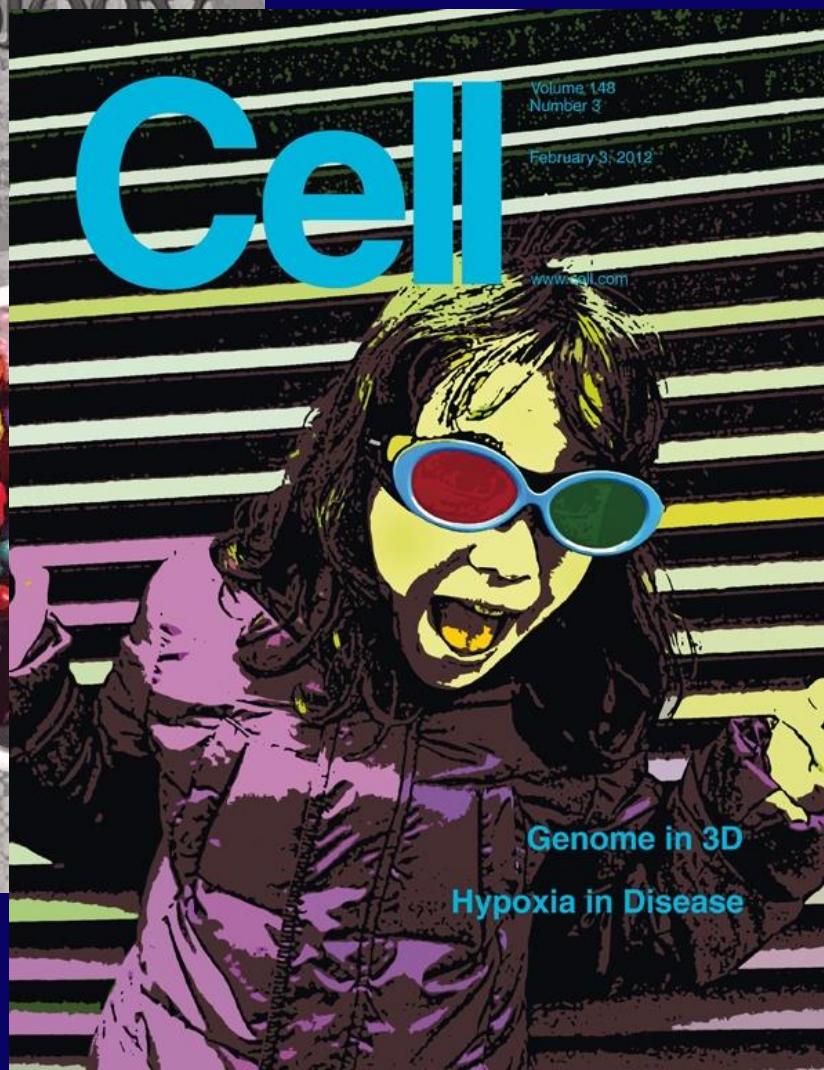
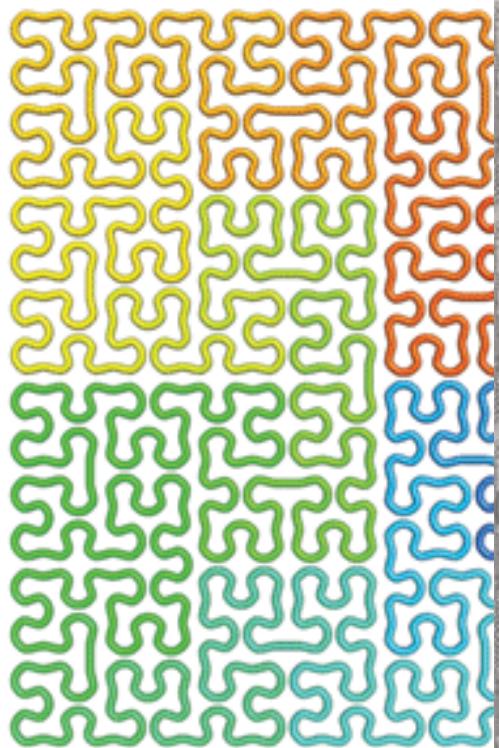
Intermediate: (0-Mb)



Fine scale: (0-kb)

9 October 2009 | \$10

Science



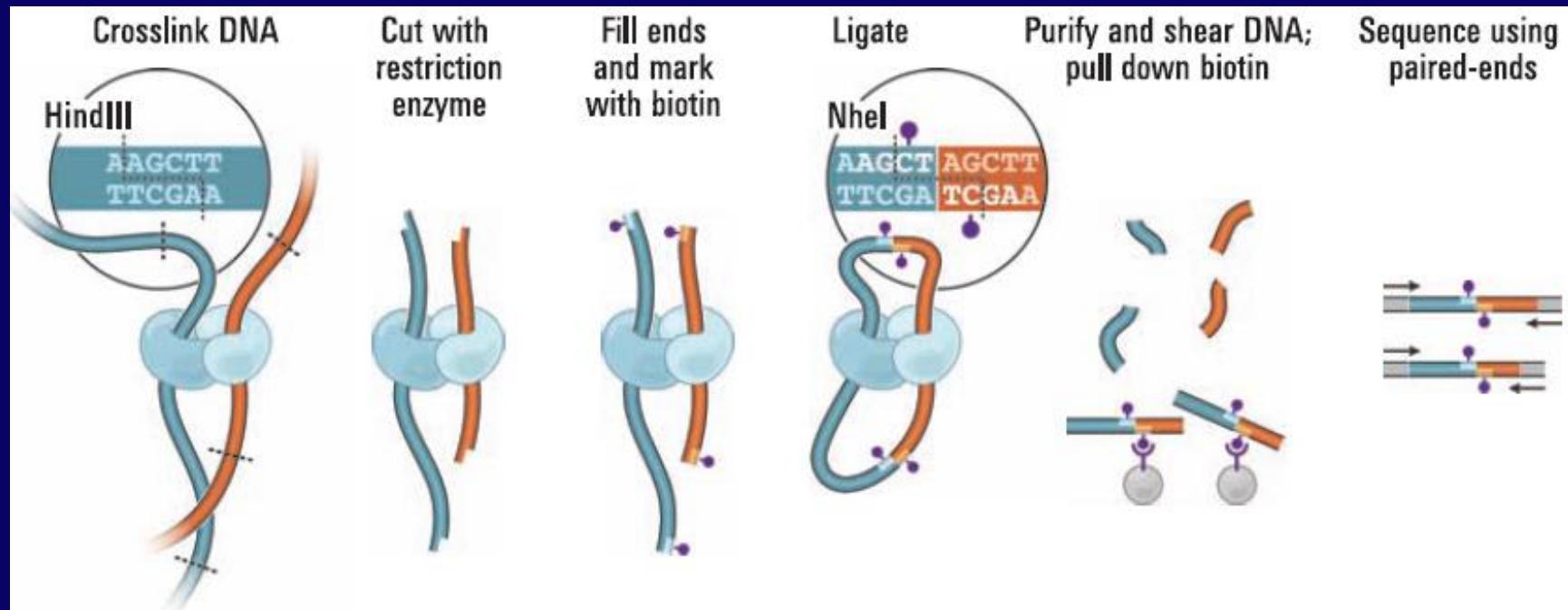
Comprehensive Mapping of Long-Range Interactions Reveals Folding Principles of the Human Genome

Erez Lieberman-Aiden,^{1,2,3,4*} Nynke L. van Berkum,^{5*} Louise Williams,¹ Maxim Imakaev,² Tobias Ragoczy,^{6,7} Agnes Telling,^{6,7} Ido Amit,¹ Bryan R. Lajoie,⁵ Peter J. Sabo,⁸ Michael O. Dorschner,⁸ Richard Sandstrom,⁸ Bradley Bernstein,^{1,9} M. A. Bender,¹⁰ Mark Groudine,^{6,7} Andreas Gnirke,¹ John Stamatoyannopoulos,⁸ Leonid A. Mirny,^{2,11} Eric S. Lander,^{1,12,13†} Job Dekker^{5†}

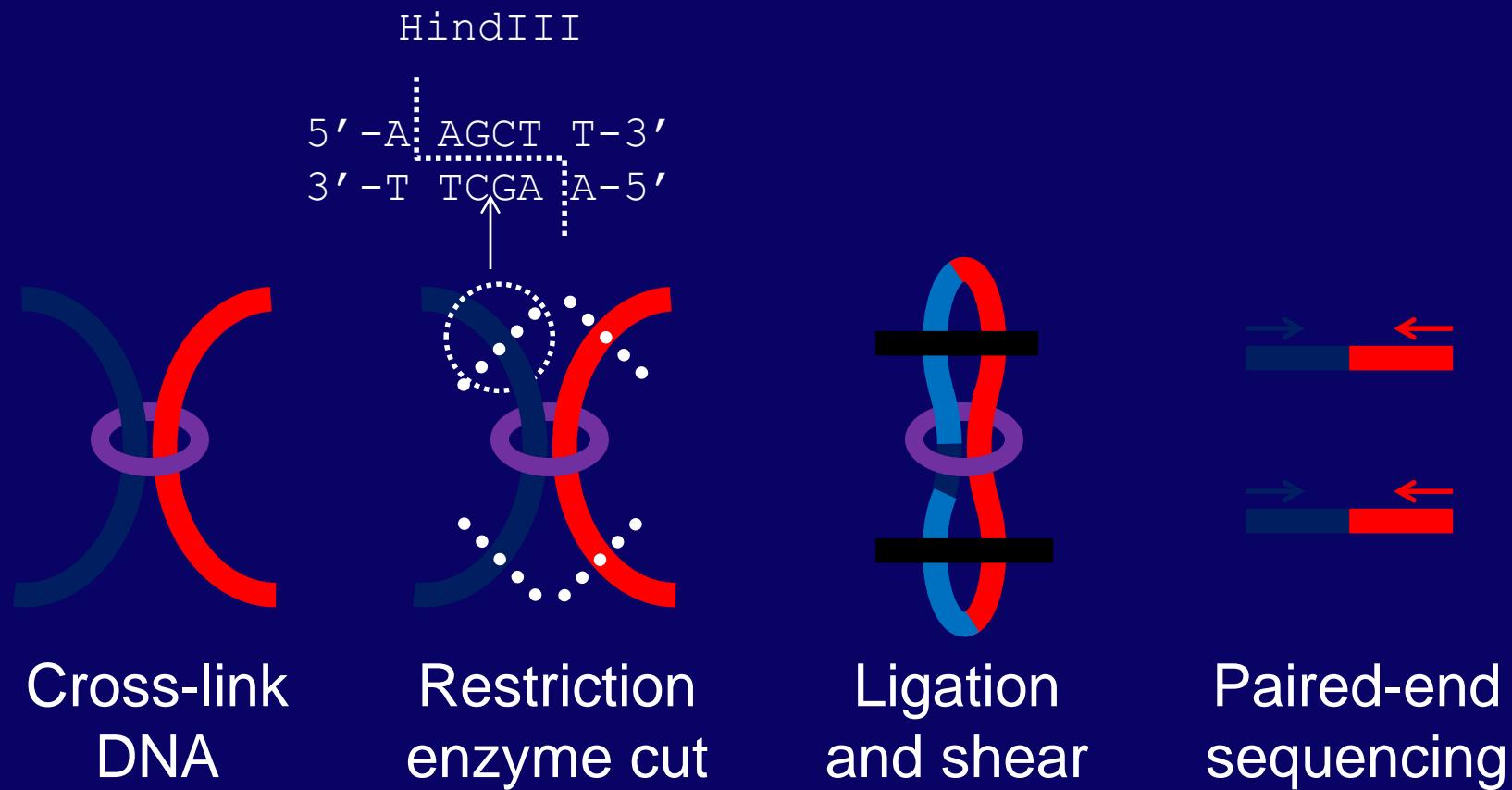
We describe Hi-C, a method that probes the three-dimensional architecture of whole genomes by

We created a Hi-C library from a karyotypically normal human lymphoblastoid cell line (GM06990) and sequenced it on two lanes of an Illumina Genome Analyzer (Illumina, San Diego, CA), generating 8.4 million read pairs that could be uniquely aligned to the human genome reference sequence; of these, 6.7 million corresponded to long-range contacts between segments >20 kb apart.

We constructed a genome-wide contact matrix M by dividing the genome into 1-Mb regions ("loci") and defining the matrix entry m_{ij} to be the number of ligation products between locus i and locus j (10). This matrix reflects an ensemble



Hi-C: one cell



Cross-link
DNA

Restriction
enzyme cut

Ligation
and shear

Paired-end
sequencing

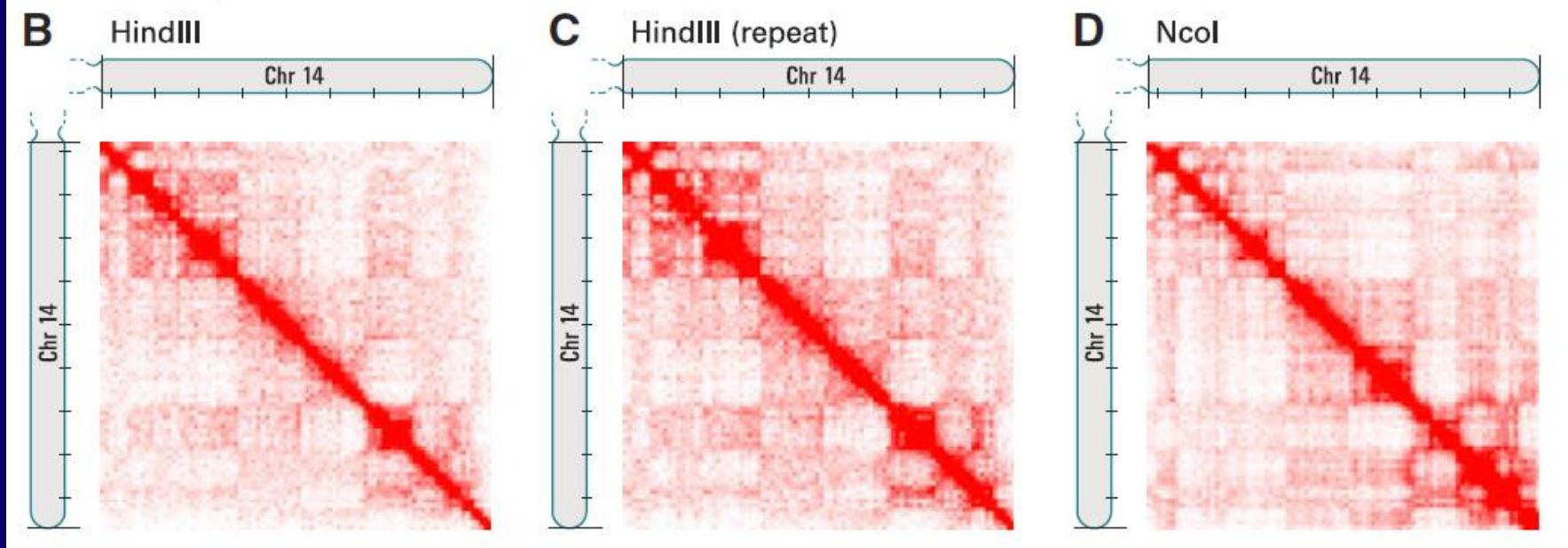
Comprehensive Mapping of Long-Range Interactions Reveals Folding Principles of the Human Genome

Erez Lieberman-Aiden,^{1,2,3,4*} Nynke L. van Berkum,^{5*} Louise Williams,¹ Maxim Imakaev,² Tobias Ragoczy,^{6,7} Agnes Telling,^{6,7} Ido Amit,¹ Bryan R. Lajoie,⁵ Peter J. Sabo,⁸ Michael O. Dorschner,⁸ Richard Sandstrom,⁸ Bradley Bernstein,^{1,9} M. A. Bender,¹⁰ Mark Groudine,^{6,7} Andreas Gnirke,¹ John Stamatoyannopoulos,⁸ Leonid A. Mirny,^{2,11} Eric S. Lander,^{1,12,13†} Job Dekker^{5†}

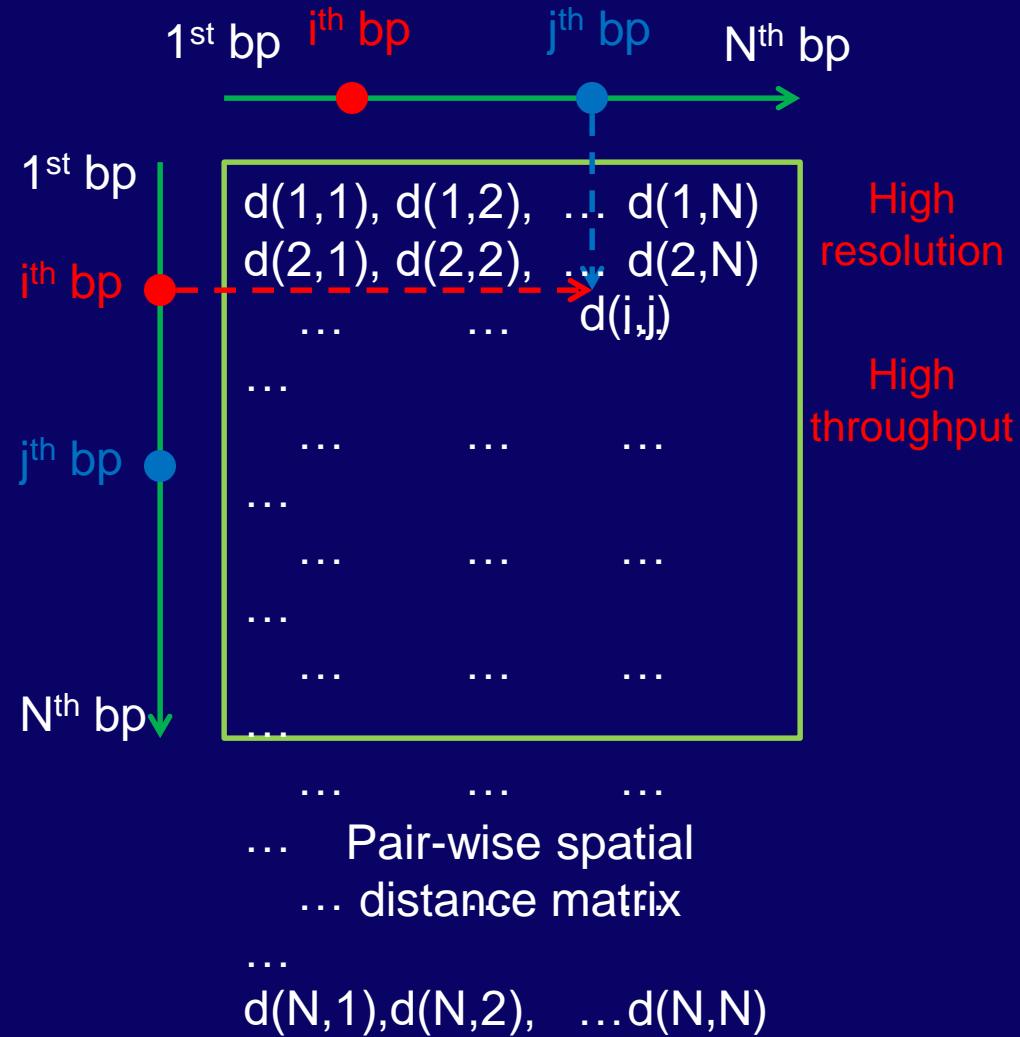
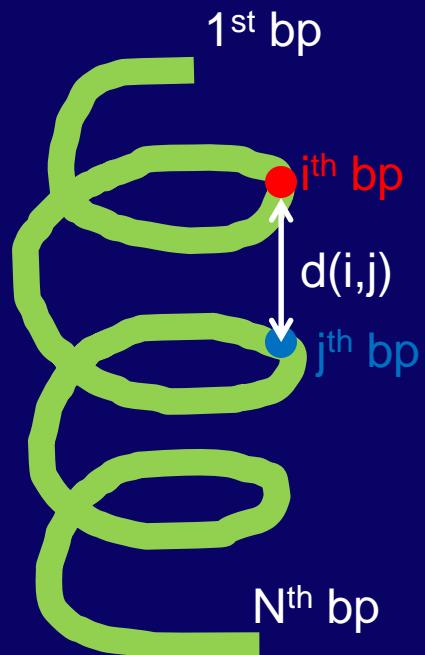
We describe Hi-C, a method that probes the three-dimensional architecture of whole genomes by

We created a Hi-C library from a karyotypically normal human lymphoblastoid cell line (GM06990) and sequenced it on two lanes of an Illumina Genome Analyzer (Illumina, San Diego, CA), generating 8.4 million read pairs that could be uniquely aligned to the human genome reference sequence; of these, 6.7 million corresponded to long-range contacts between segments >20 kb apart.

We constructed a genome-wide contact matrix M by dividing the genome into 1-Mb regions ("loci") and defining the matrix entry m_{ij} to be the number of ligation products between locus i and locus j (10). This matrix reflects an ensemble



Hi-C Data Representation

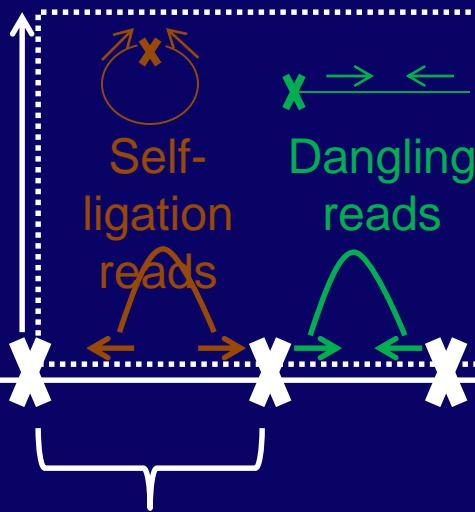


Challenges

- Quality control and pre-processing of the reads,
- Any bias in the data? and if so, how to normalize?
- Whether it is possible, and if so, how, to infer the 3-dimesnional chromosomal structure based on the Hi-C data?

Hi-C Data Preprocess

Restriction
enzyme
cutting site



PCR
amplificatio
n
 n
reads

Random
break

Random
break

Random
breaking
reads

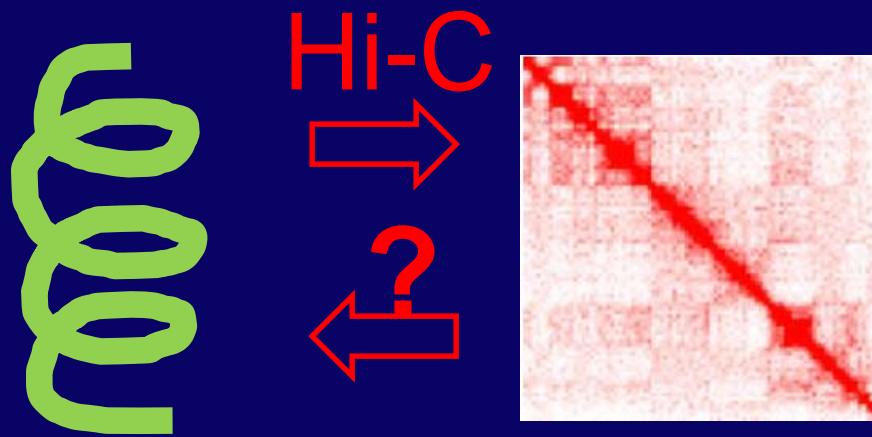
Valid reads

Downstream
analysis

Methods for Hi-C Bias Reduction

- Normalization (equal ‘visibility’, no assumption on biases)
- Iterative correction and eigenvector decomposition (ICE)
(Imakaev, et al, 2012)
- Sequential component normalization (SCN)
(Cournac, et al, 2012)
- Correction (posit a statistical model on biases)
- Yaffe & Tanay’s method (Yaffe & Tanay, 2011)
Fragment level (4KB, 10^{12}), 420 parameters
- **HiCNorm (Hu et al, 2012)**
Any resolution level
1MB, 10^6 , 3 parameters

3D structure prediction



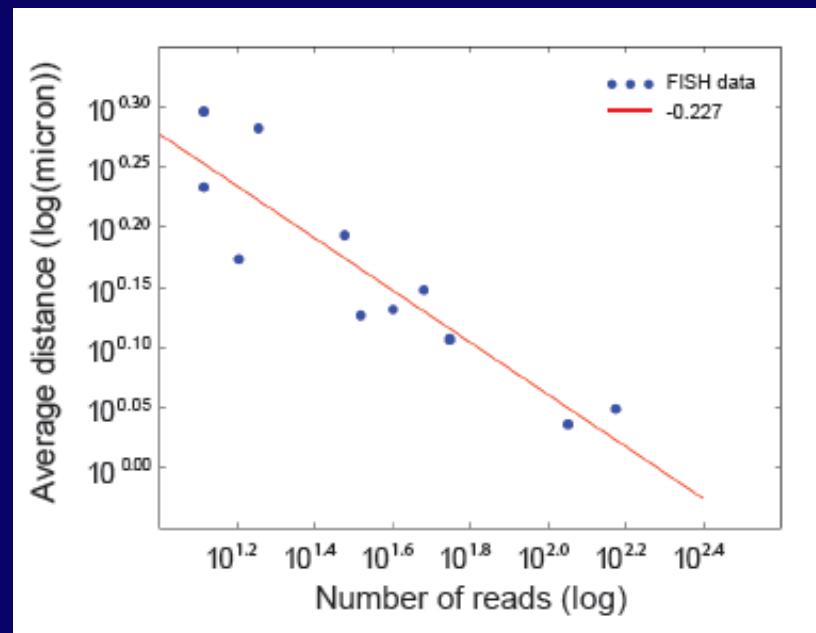
- Challenges:
 - Sequencing uncertainties
 - Biases: enzyme, GC content, mappability

What does the number mean?

- The Hi-C experiment is conducted on millions of cells,
- A captured pair-end read is from one cell,
- A number in the matrix (loci i and j) indicates the frequency of capture (link i and j) in the cell population,
- Do those numbers say anything about 3D distance?

Motivation and the key assumption

- Number of paired-end reads spanning the two loci is inversely proportional to the **3D spatial distance** between them (obtained from fluorescence in situ hybridization(FISH)).



Existing methods

- Optimizations-based method (Baù, et al, 2010, Duan, et al, 2010)
 - Biophysical properties of chromatin fiber.
 - No consideration of systematic biases.
 - No statistical inference.
- Statistical method: MCMC5C (Rousseau et al, 2011)
 - Normal model for count data.
 - No consideration of systematic biases.

Model

ACGTAGCTAGATACTGTAGTACATCGATAGCGTAGTTGGAACCTGAGGGTAAACC
TGGAGGGGATCATG

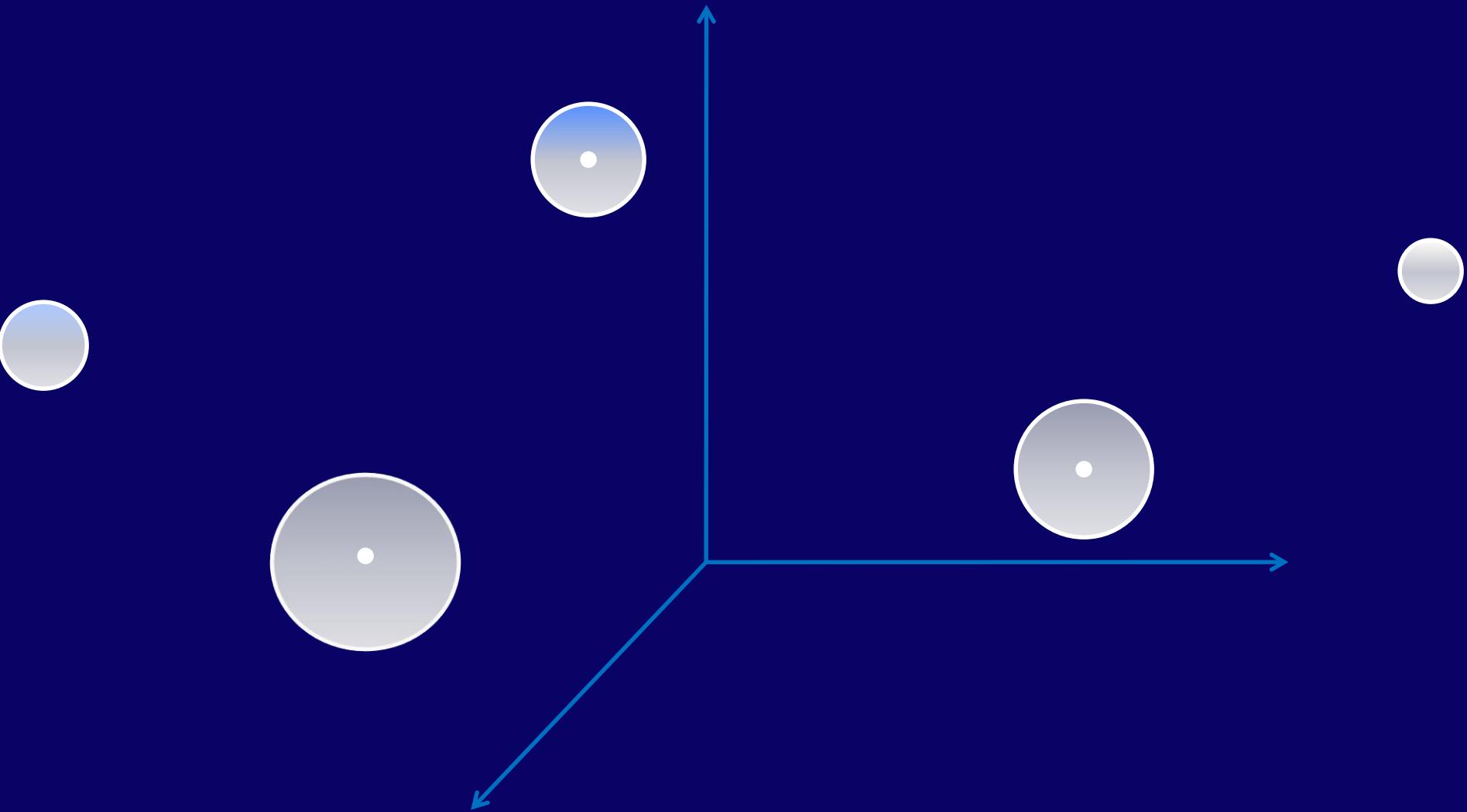
Model

**ACGTAGCTAGATACT GTAGTACATCGATAG CGTAGTTGGAACCT
GAGGGTAAACCTGG AGGGGAT**

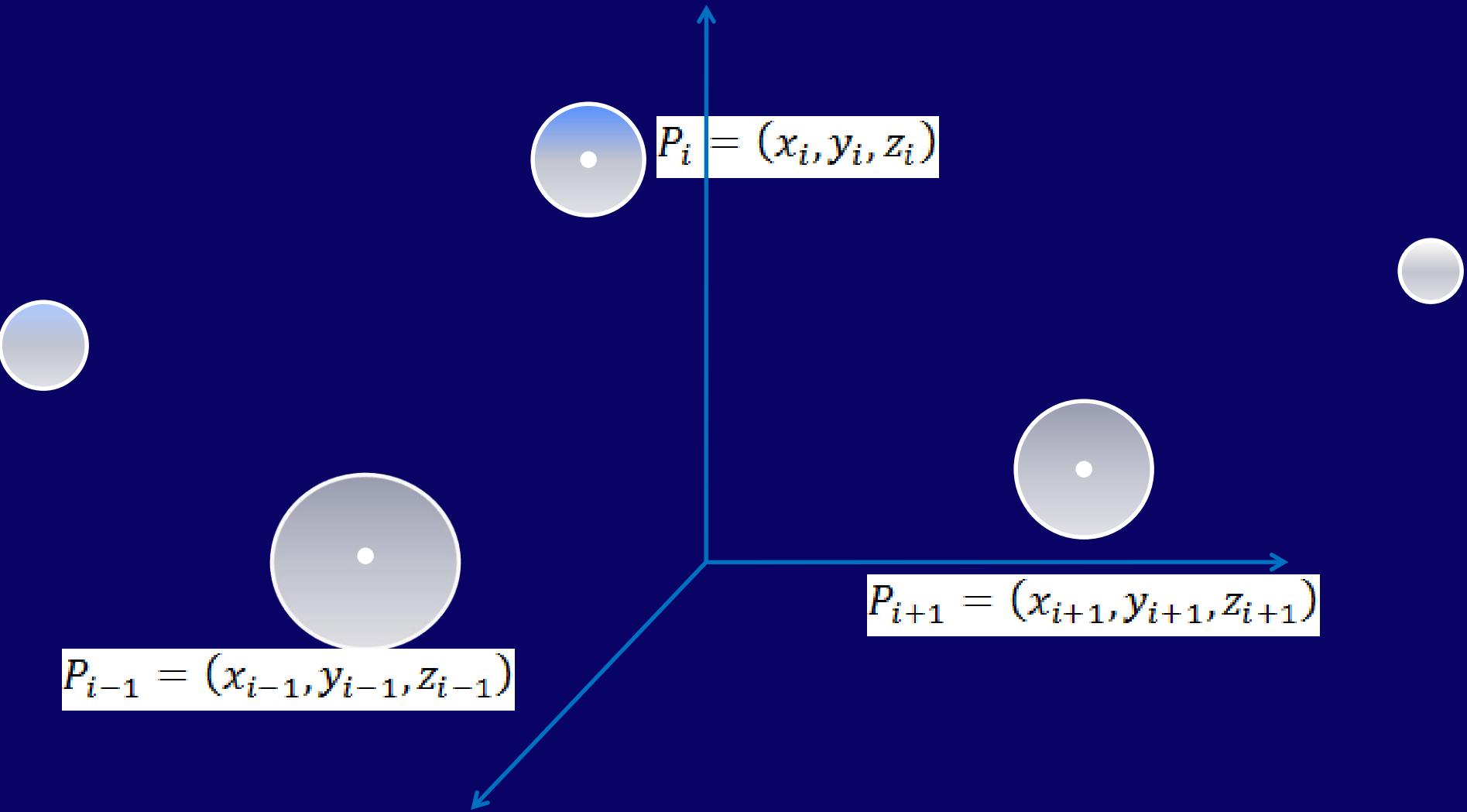
Model



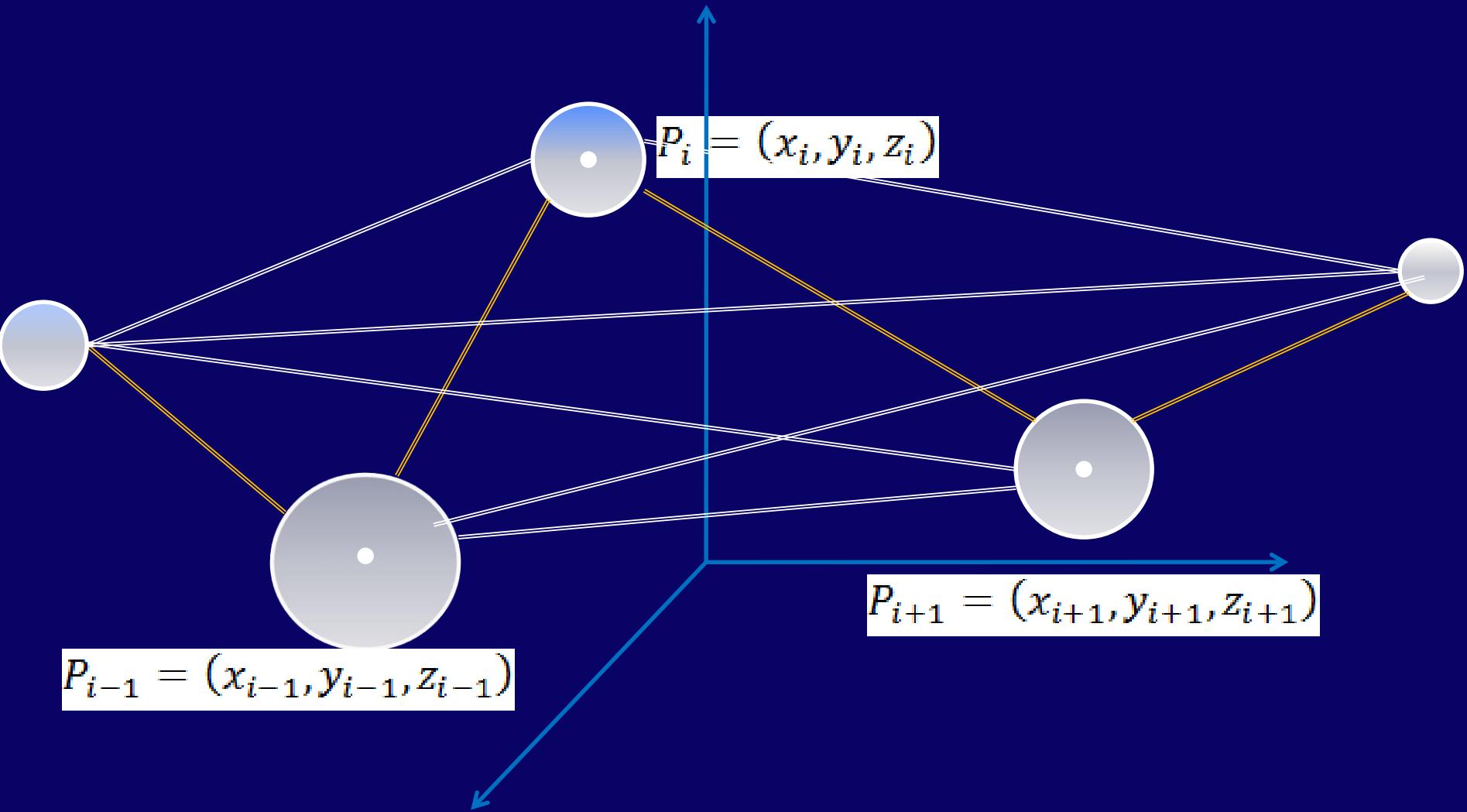
Beads-on-string



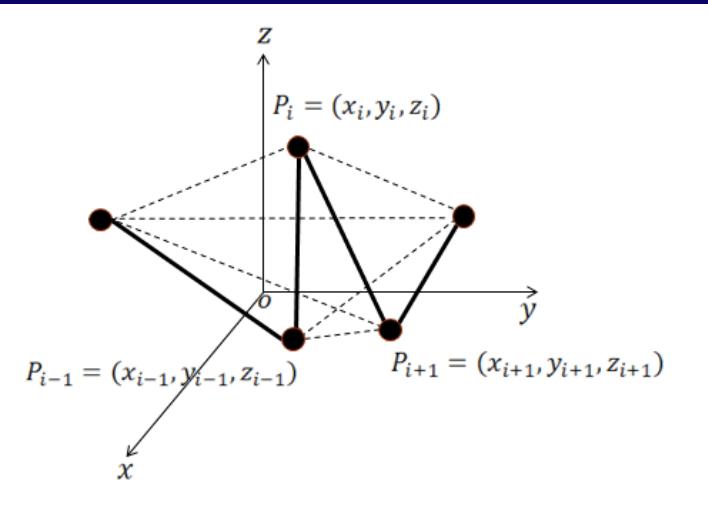
Beads-on-string



Beads-on-string



Bayesian statistical model



- u_{ij} : number of reads between loci i and j .
- d_{ij} : 3D Euclidian distance between loci i and j .
- enz_i : number of enzyme cut site in locus i .
- gcc_i : mean GC content in locus i .
- map_i : mean mappability score in locus i .

$$u_{ij} \sim Poisson(\theta_{ij})$$

$$\log(\theta_{ij}) = \beta_0 + \beta_1 \log(d_{ij}) + \beta_{enz} \log(enz_i enz_j)$$

$$+ \beta_{gcc} \log(gcc_i gcc_j) + \beta_{map} \log(map_i map_j)$$

Bayesian Statistical Model

- Likelihood: $\binom{N}{2}$ data points, $3N + 5$ parameters

$$L(u_{ij}, 1 \leq i < j \leq N | x_i, y_i, z_i, 1 \leq i \leq N, \beta_0, \beta_1, \beta_e, \beta_g, \beta_m) = \prod_{1 \leq i < j \leq N} \frac{e^{-\theta_{ij}} \theta_{ij}^{u_{ij}}}{u_{ij}!}$$

$$\begin{aligned} \log(\theta_{ij}) = & \beta_0 + \beta_1 \log \left(\sqrt{\left(x_i - x_j \right)^2 + \left(y_i - y_j \right)^2 + \left(z_i - z_j \right)^2} \right) \\ & + \beta_e \log(e_i e_j) + \beta_g \log(g_i g_j) + \beta_m \log(m_i m_j) \end{aligned}$$

Statistical Inference

- Algorithm: Bayesian 3D constructor for Hi-C data (**BACH**)
 - Initialization 1: use Poisson regression to obtain the initial values of model parameters.
 - Initialization 2: use sequential important sampling to get the initial 3D chromosomal structure .
 - Refinement: use Gibbs sampler with hybrid Monte Carlo to refine the initial values for parameters.

SIS in BACH: Outline

- Goal: use sequential importance sampling to **sequentially** put N loci into 3D space, i.e. sample from:

$$\pi(x_i, y_i, z_i, 1 \leq i \leq N | u_{ij}, 1 \leq i < j \leq N)$$

- Bridging distributions:

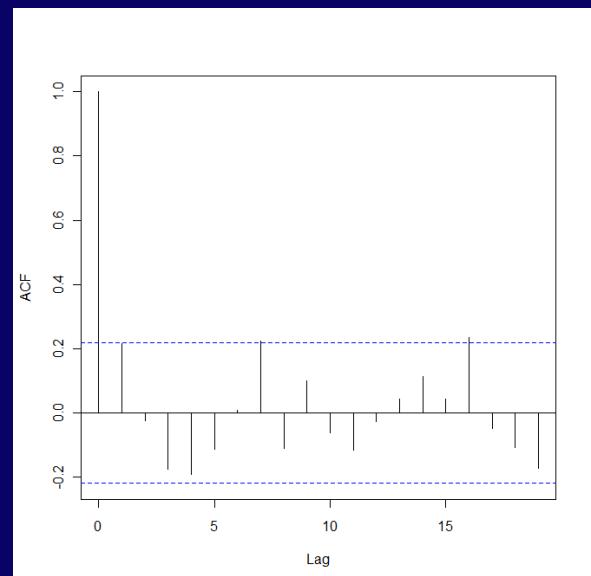
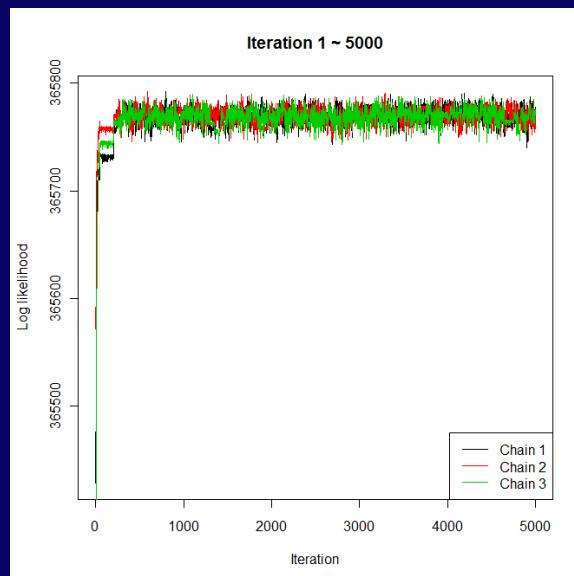
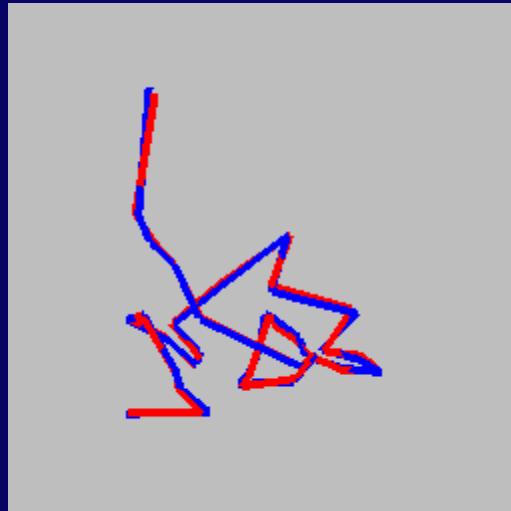
$$\pi_t(x_i, y_i, z_i, 1 \leq i \leq t | u_{ij}, 1 \leq i < j \leq t)$$

- Proposal distributions (given the previous $t-1$ loci, put the t th locus in to 3D space):

$$g_t(x_t, y_t, z_t | x_i, y_i, z_i, 1 \leq i \leq t-1, u_{ij}, 1 \leq i < j \leq t)$$

Simulation study

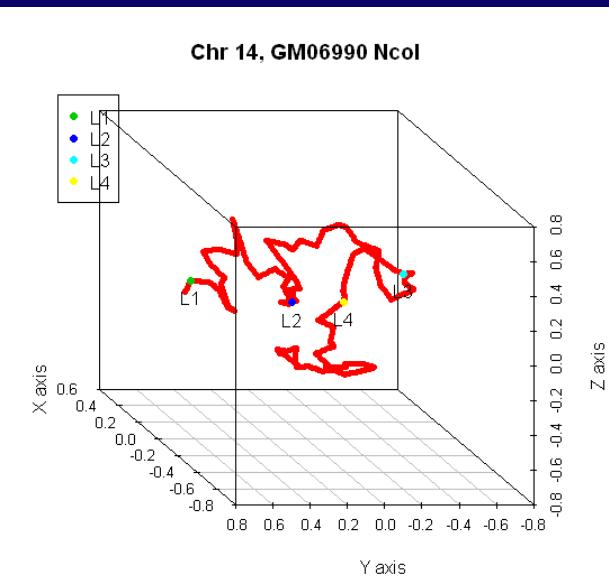
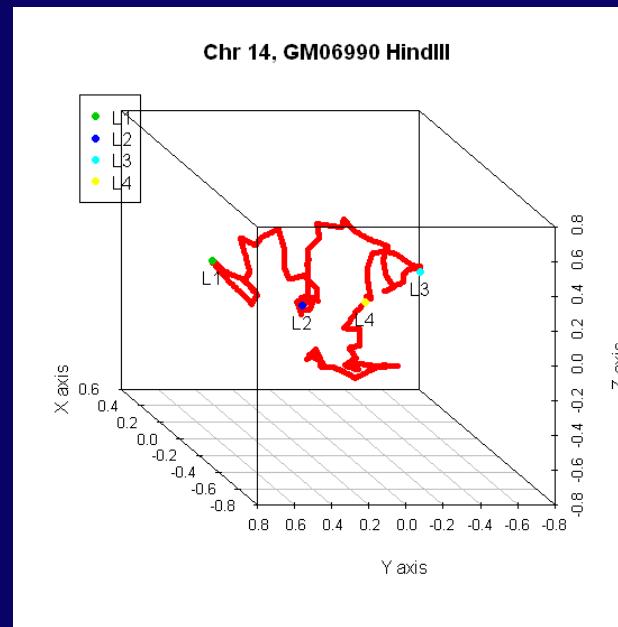
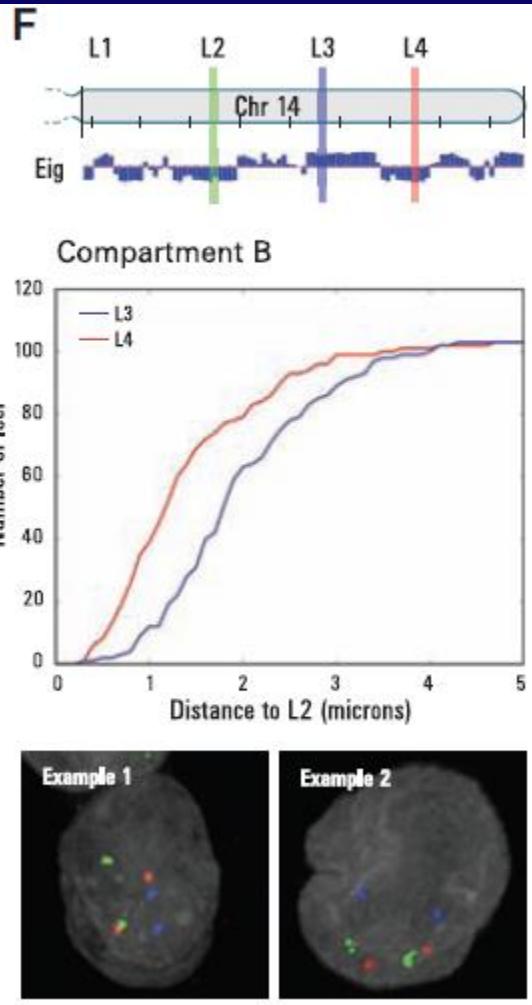
- Use random walk to simulate a 3D structure with 33 loci (**red lines**). Simulate Hi-C contact map from the posited model.
- Predicted 3D structure (**blue lines**) aligns well with true 3D structure (RMSD = 0.0091).



Human Hi-C data

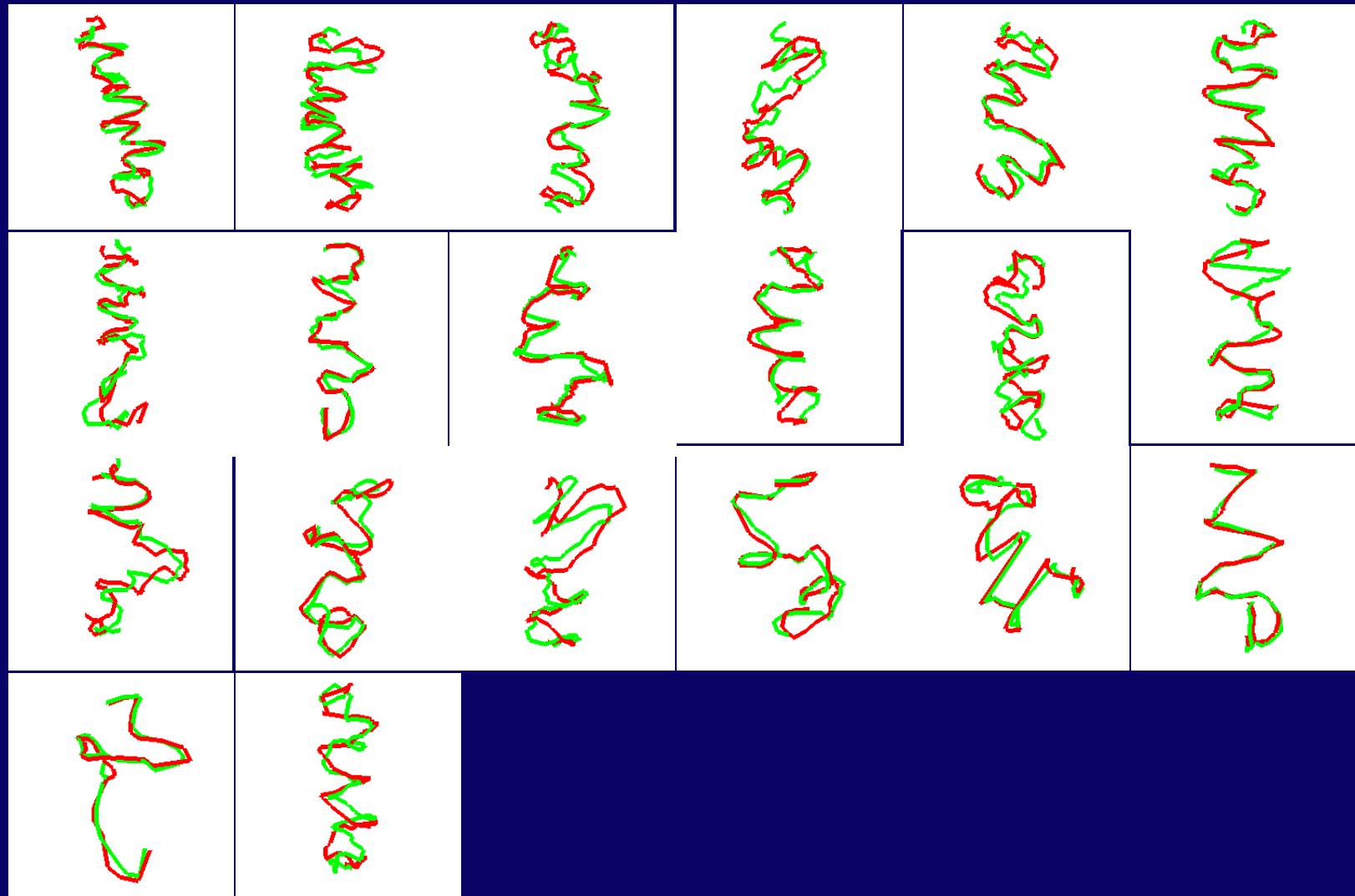
Cell line	Restriction enzyme	# of reads (million)
GM06990	HindIII	4.1
GM06990	HindIII	4.4
GM06990	HindIII	4.9
GM06990	HindIII	5.4
GM06990	NcoI	8.8
GM06990	NcoI	10.1
K562	HindIII	12.1
K562	HindIII	9.7

Real Hi-C data from Lieberman-Aiden et al. 2009



$d(L2, L4) = 1.4042, d(L2, L3) = 1.9755,$
significant

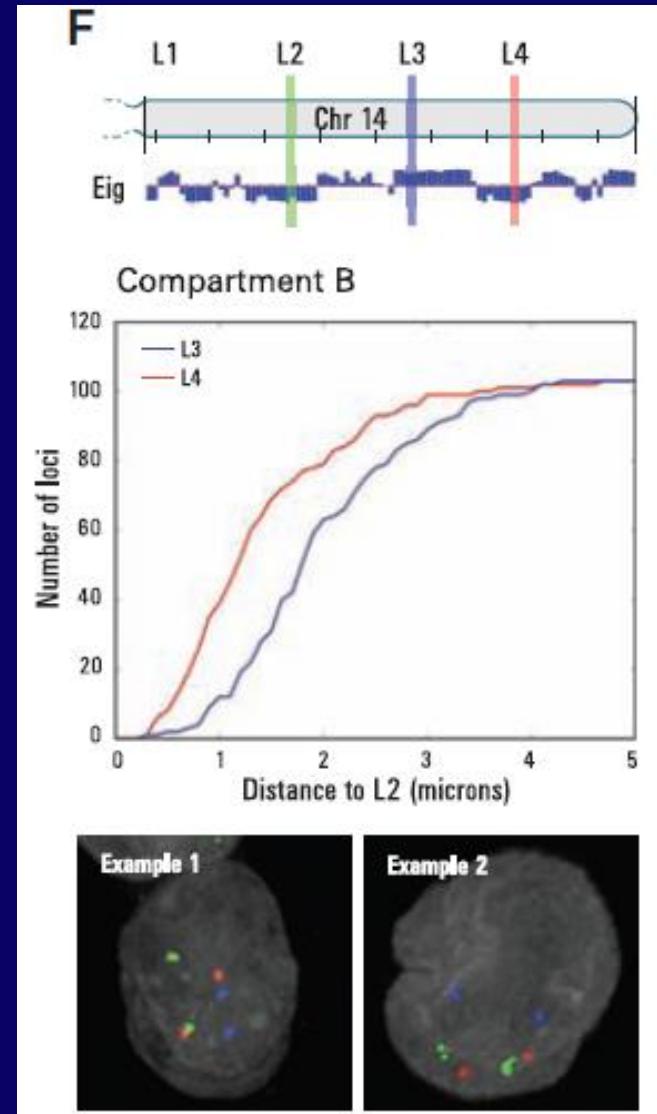
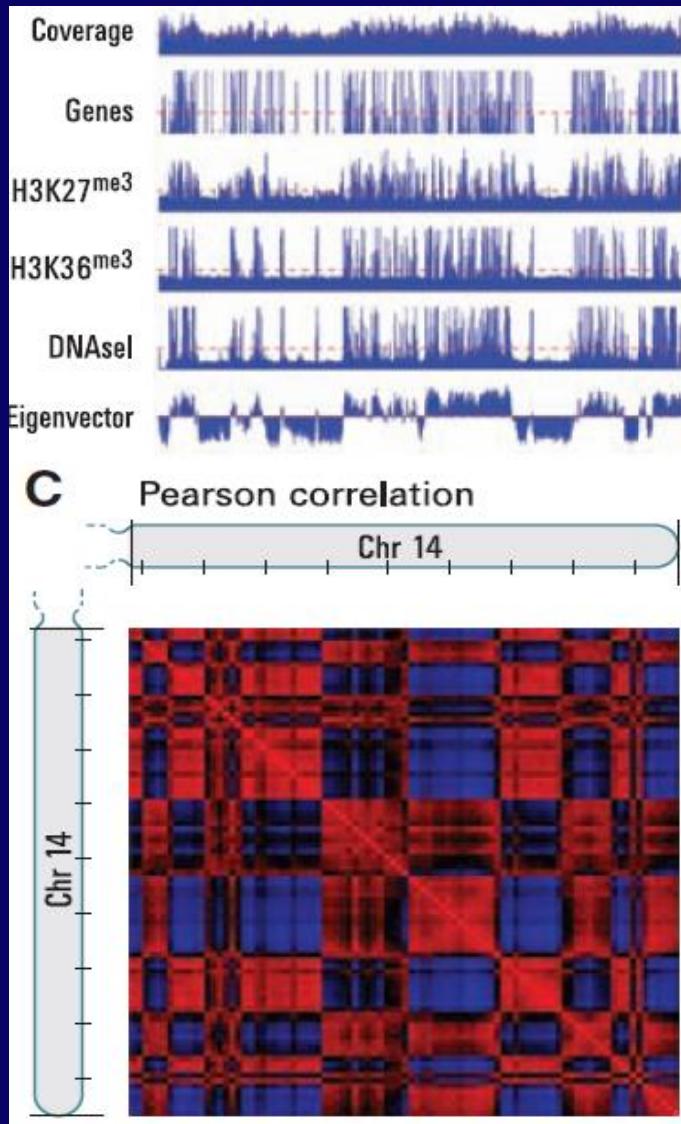
mESC: Hind3 vs. Nco1



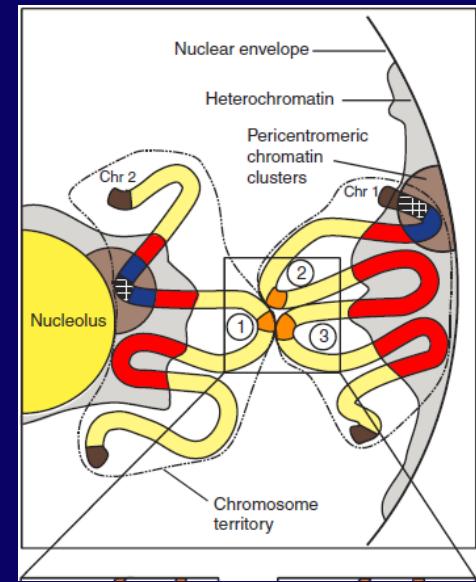
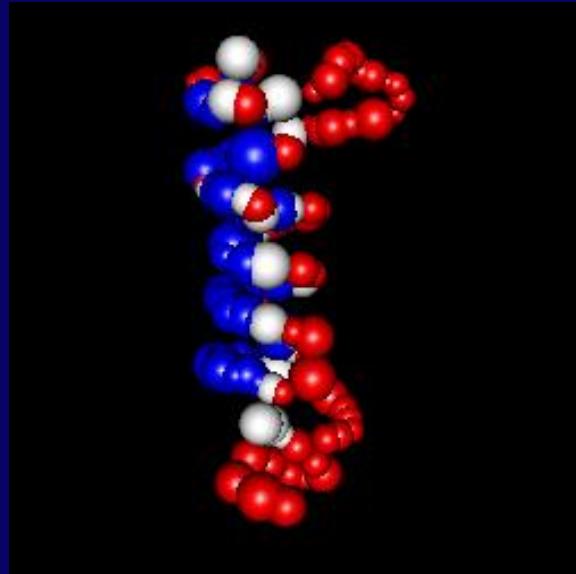
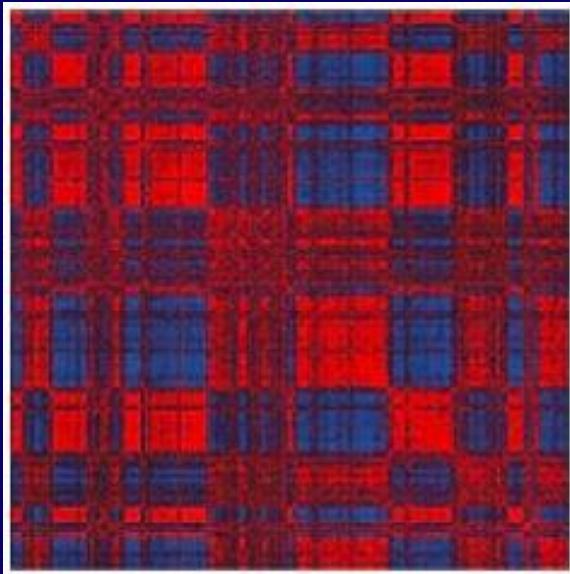
Whole Chromosome 3D Model

- Two compartments
 - Compartment A: gene rich, active transcription
 - Compartment B: gene poor, inactive transcription
- Same compartment: strong chromatin interactions, spatially close
- Different compartments: weak chromatin interaction, spatially isolated

Two compartment model

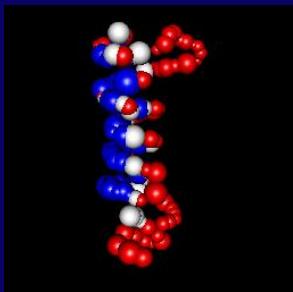


Whole Chromosome Model

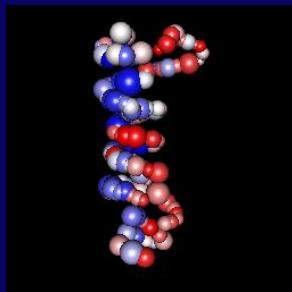


Other Features (Chromosome 2)

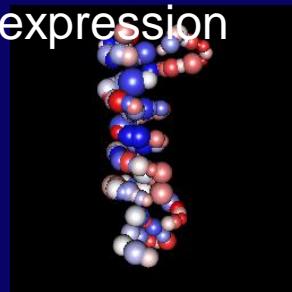
Compartment



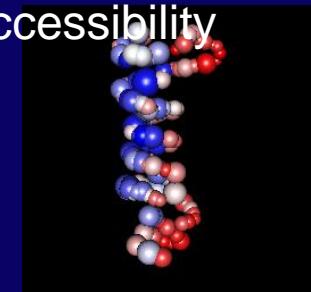
Gene density



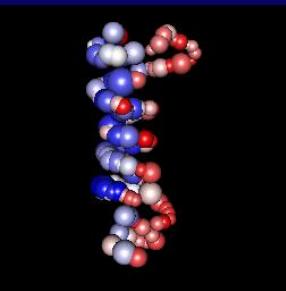
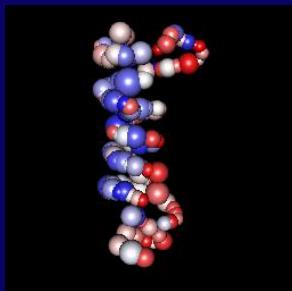
Gene expression



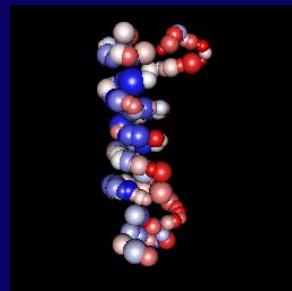
Chromatin accessibility



RNA polymerase II DNA replication time



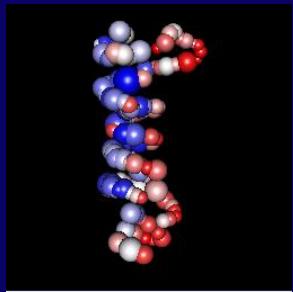
H3K36me3



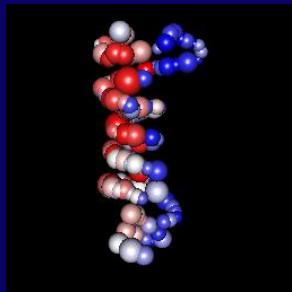
H3K27me3



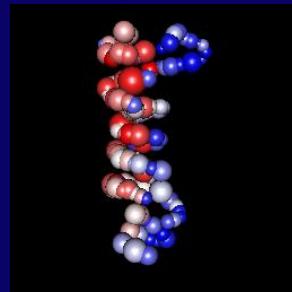
H3K4me3



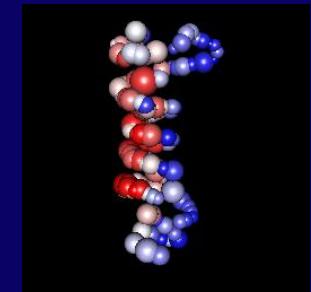
H3K9me3



H3K20me3



Lamina interaction



Conclusions

- BACH--Reconstruct chromosome 3D structures
- Remove systematic biases
- Consistent with FISH data
- Elongation of chromatin is highly associated with genetic/epigenetic features.
- Separation of compartments of A and B can be visualized.

References

- Hu M, Deng K, Selvaraj S, Qin ZS, Ren B, Liu JS. (2012) HiCNorm: removing biases in Hi-C data via Poisson regression. *Bioinformatics*. In press.
<http://www.people.fas.harvard.edu/~junliu/HiCNorm/>
- Hu M, Deng K, Qin ZS, Dixon J, Selvaraj S, Fang J, Ren B, Liu JS. (2012) Bayesian inference of three-dimensional chromosomal organization. *PLoS Computational Biology*. In press.
<http://www.people.fas.harvard.edu/~junliu/BACH/>
- Hou C, Li L, Qin ZS, Corces, VG. (2012) Gene Density, Transcription and Insulators Contribute to the Partition of the Drosophila Genome into Physical Domains. *Mol Cell*. **48** 471-484 (with preview article of Xu and Felsenfeld (2012) Order from Chaos in the Nucleus. *Mol Cell* 48. 327-328).
.
- Dixon JR, Selvaraj S, Yue F, Kim A, Li Y, Shen Y, Hu M, Liu JS and Ren B. (2012) Topological domains in mammalian genomes identified by analysis of chromatin interactions. *Nature* , 485, 376-380.

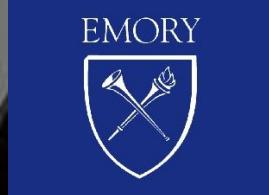


Hu M, Deng K, Qin ZS, Liu JS (2013) Understanding spatial organizations of chromosomes via statistical analysis of Hi-C data. *Quantitative Biology* 1. 156-174.

Acknowledgements



Ming Hu
Ke Deng
Jun S. Liu



Li Li
Chunhui Hou
Victor Corces



Jesse Dixon
Siddarth Selvara
Bing Ren





Thank You

Questions: zhaohui.qin@emory.edu