



# Quick intro to the MIT CompBio Group



# Who we are



# What we do: Research synopsis

- Why biology in a computer science group?
- Fundamental biological questions:
  1. Interpreting the human genome.
  2. Revealing the logic of gene regulation.
  3. Principles of evolutionary change.
- Algorithmic/machine learning methods:
  - Comparative genomics: evolutionary signatures
  - Regulatory genomics: motifs, networks, models
  - Epigenomics: chromatin states, dynamics, disease
  - Phylogenomics: evolution at the genome scale
- Defining characteristics of our group:
  - Learn genomic rules, exploit nature of problems
  - Interdisciplinary collaborations, high biology impact

# (1) Comparative genomics: evolutionary signatures

<i>D. mel.</i>	CATTTA <del>T</del> ATT-----T-----ATT-----AAATAATGGCGTT-----TCGCAGC-GGCTGG-C-----TGT <del>T</del> TTTATTAAACATTATT-----
<i>D. sim.</i>	CATTTA <del>T</del> ATT-----T-----ATT-----AAATAATGGCGTT-----TCGCAGC-GCTGG-C-----TGT <del>T</del> TTTATTAAACATTATT-----
<i>D. sec.</i>	CATTTA <del>T</del> ATT-----T-----ATT-----AAATAATGGCGTT-----TCGCAGC-GCTGG-C-----TTT <del>T</del> TTTATTAAACATTATT-----
<i>D. yak.</i>	CATTTA <del>T</del> ATT-----T-----ATT-----AAATAATGGCGTT-----TCGCAGC-GCTGG-CTG-----TGT <del>T</del> TTTATTAAACATTATT-----
<i>D. ere.</i>	CGTTTATTAT-----T-----ATC-----AAATAATGGCGTT-----TCGCAGC-GGTGG-C-----TGT <del>T</del> TTTATTAAACATTACTA-----
<i>D. ana.</i>	CAT <del>T</del> TTTATT-----T-----AAATAATGGTATT-----TCTTGACTGGCTGC-CTGCC-----TGCCGTAA-TTGTGT <del>T</del> TTTATTAAACATTATT-----
<i>D. pse.</i>	CATTTTATT-----T-----GAT-----AAATAATGGAACTTGGTCAGTT-----TTGCTGCGCTGCC-----TRGCTGCTGCC <del>T</del> TTGTGT <del>T</del> TTTATTAAACATTATT-----
<i>D. per.</i>	CATTTTTCT-----T-----GAT-----AAATAATGGAAATTTGGTCAC <del>T</del> TTTACTGCCTGCCG-----CACCTCTCGCTTCGTGT <del>T</del> TTTATTAAACATTATT-----
<i>D. will.</i>	CATTTTATT-----TATTTATATT-----AAATAATGAAGTT-----TCGTTTC-----G-T-----TTCGTATGGT-----TCGTT-----
<i>D. moj.</i>	TATTAATTATG <del>T</del> -----ATAATAATTAAATGAAGTT-----TT-----C-----GCTTAT-----CGTTTATGCAGC <del>T</del> TTTTTTAA-----
<i>D. vir.</i>	CATTAATTAT-----T-----ATA-----AAATAATGAAGTT-----GCGTT-----T-----CGTTTATGCAGC <del>T</del> TTTTTTAA-----
<i>D. gri.</i>	CATTAATTATGAGT-----ATT-----AAATAATGAAGTT-----T-----GCTCT-----T-----CGCTCACCGATAG <del>C</del> TTTTTTAAACAC-----

## Protein-coding signatures

- 1000s new coding exons
  - Translational readthrough
  - Overlapping constraints

# Non-coding RNA signatures

- Novel structural families
  - Targeting, editing, stability
  - Structures in coding exons

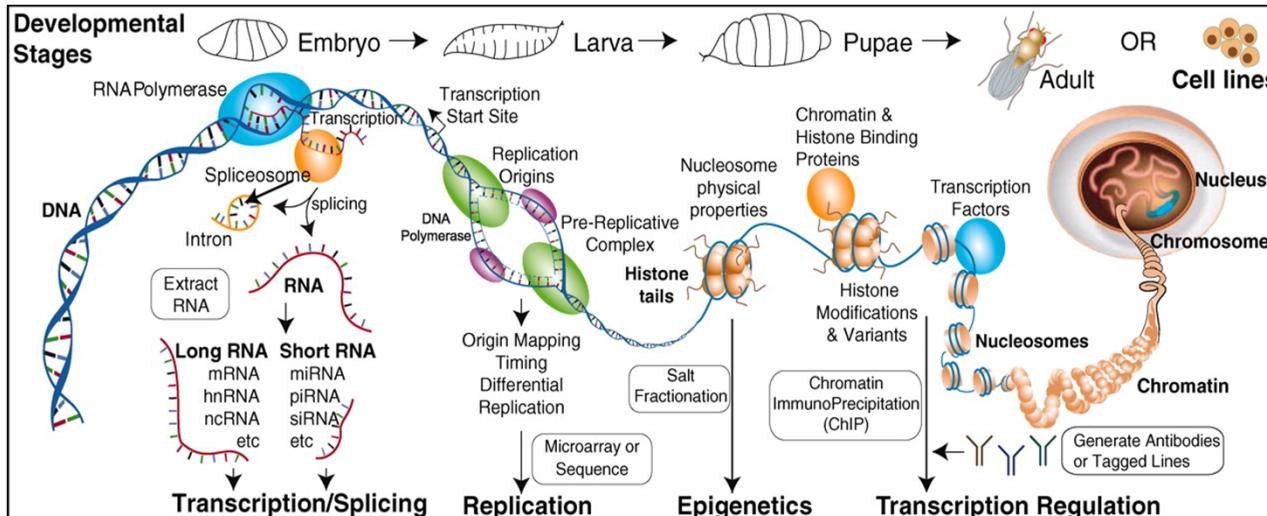
# **microRNA signatures:**

- Novel/expanded miR families
  - miR/miR\* arm cooperation
  - Sense/anti-sense switches

# Regulatory motif signatures

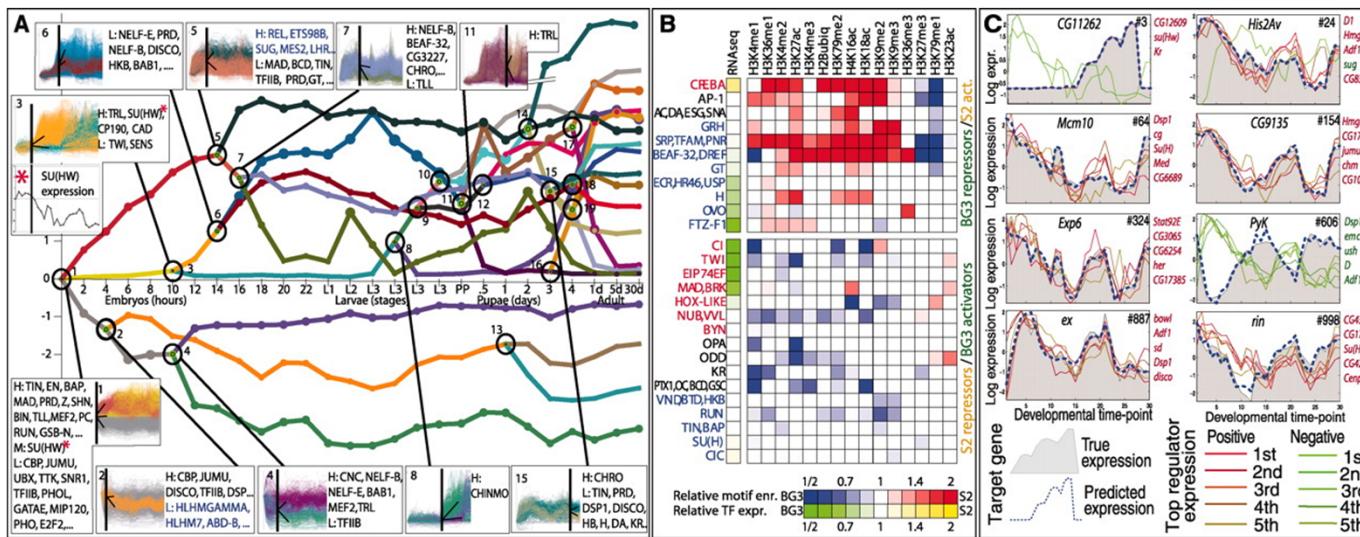
- Systematic motif discovery
  - Regulatory motif instances
  - TF/miRNA target networks
  - Single binding-site resolution

## (2) Regulatory genomics: circuits, predictive models



### ENCODE/modENCODE

- 4-year effort, dozens of experimental labs
- Integrative analysis
- Systematic genome annotation
- Flagship NIH project

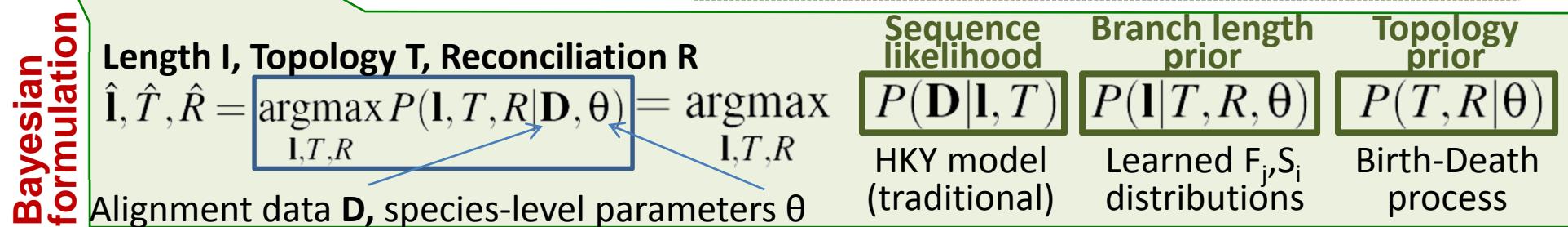
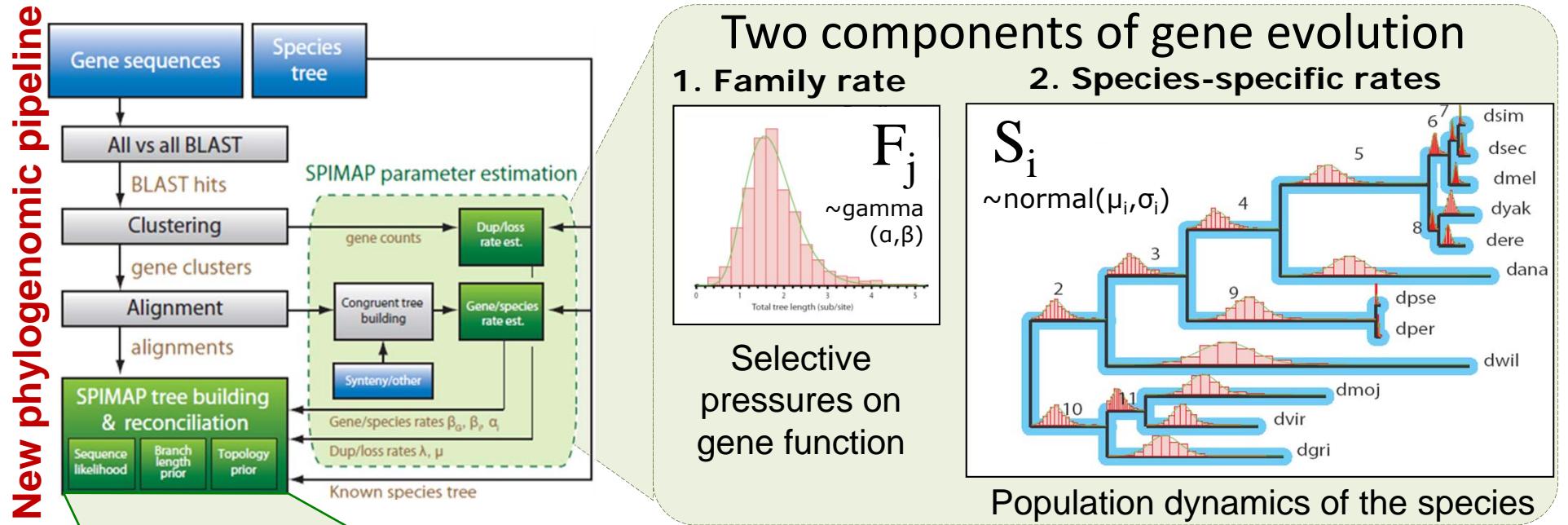
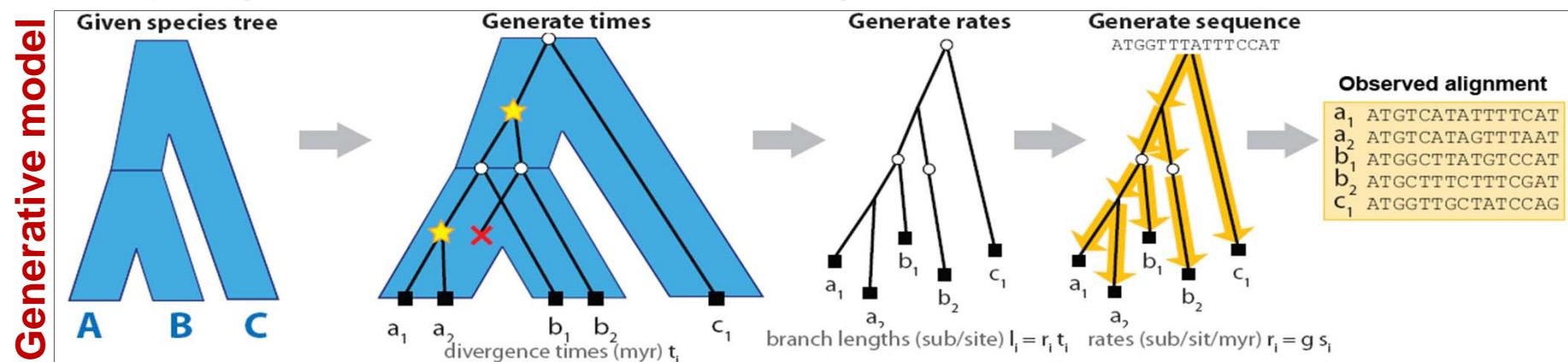


### Predictive models of gene regulation

- Infer networks
- Predict function
- Predict regulators
- Predict gene expression

- Initial annotation of the non-coding genome, from 20% to 70%
- Systems biology for an animal genome for the first time possible
- Students and postdocs are co-first authors, leadership roles

### (3) Phylogenomics: Bayesian gene-tree reconstruction





Jason Ernst

Discovery and characterization of chromatin states for systematic annotation of the human genome

Jason Ernst & Manolis Kellis

nature  
biotechnology

## (4) Vignette on Epigenomics

Using chromatin information  
to understand human diseases

Pouya  
Kheradpour



**Mapping and analysis of chromatin state dynamics in nine human cell types**

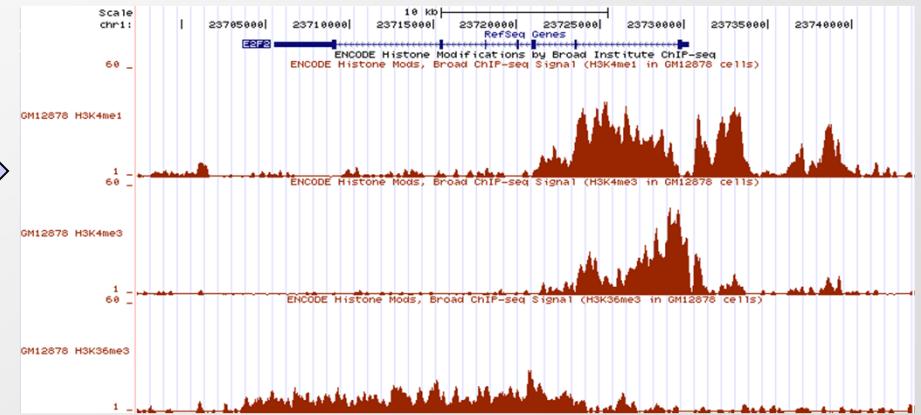
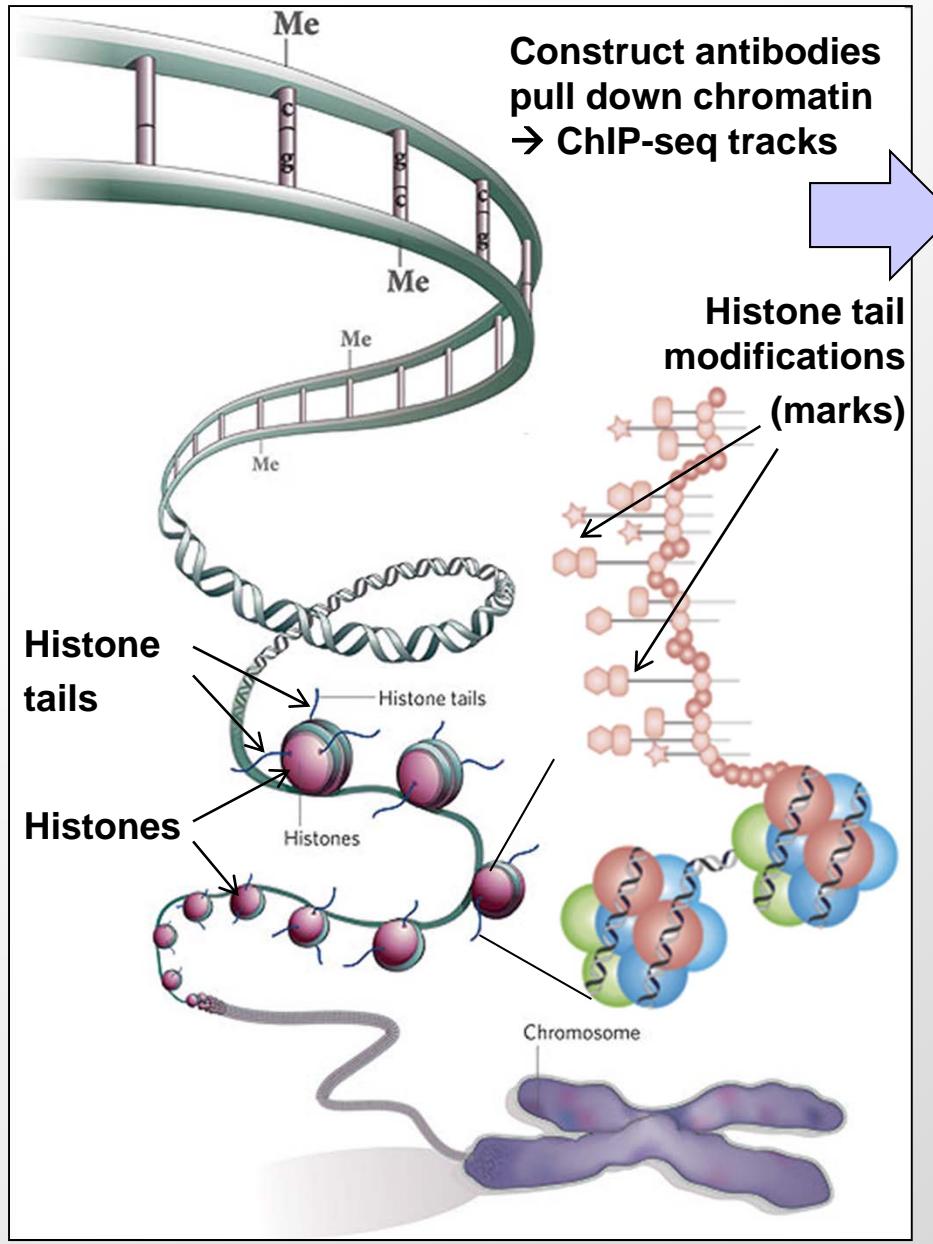
Jason Ernst, Pouya Kheradpour, Tarjei S. Mikkelsen, Noam Shoresh, Lucas D. Ward, Charles B. Epstein, Xiaolan Zhang, Li Wang, Robbyn Issner, Michael Coyne, Manching Ku, Timothy Durham, Manolis Kellis & Bradley E. Bernstein

doi:10.1038/nature09906

[Abstract](#) | [Full Text](#) | [PDF](#)

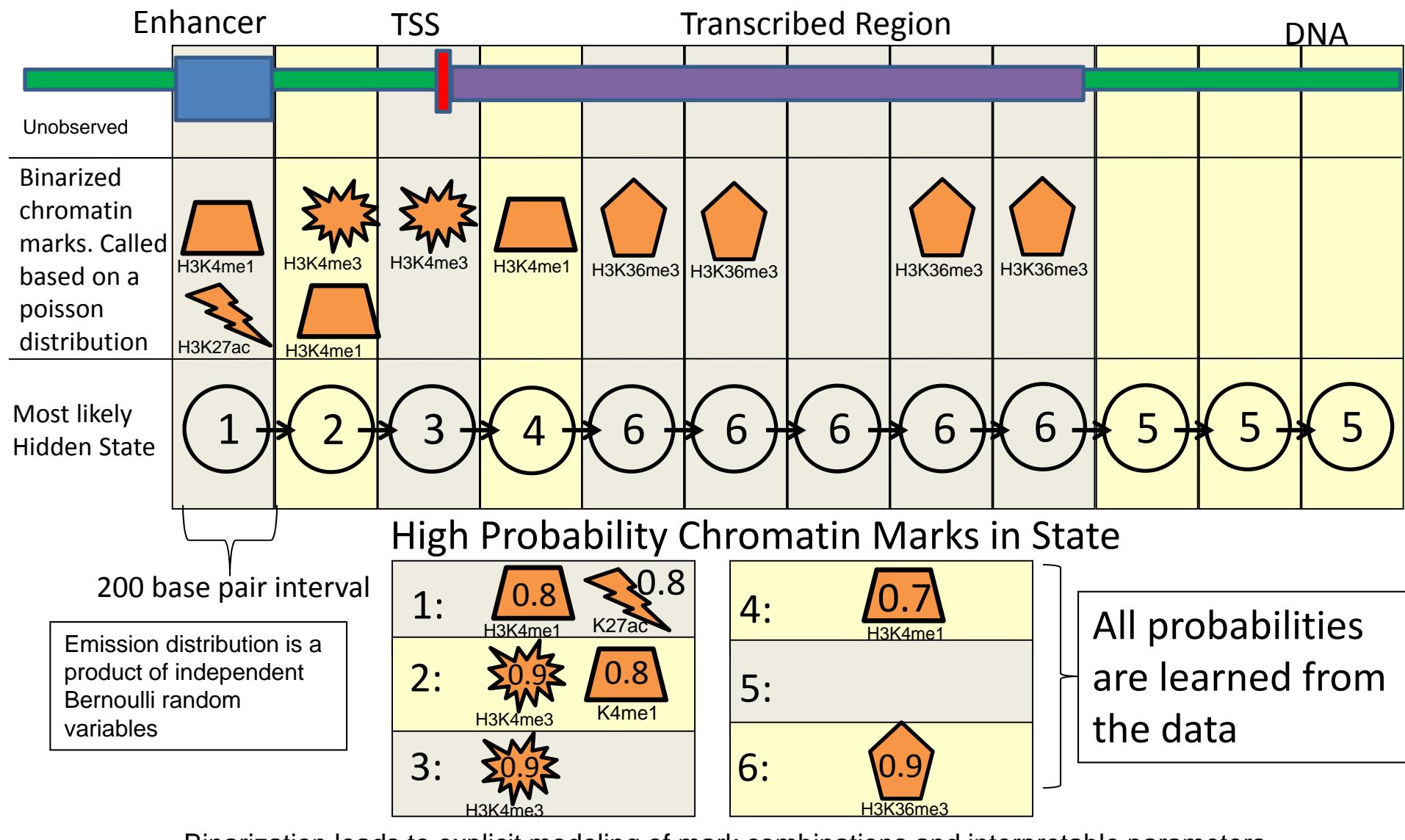
nature

# Challenge of data integration in many marks/cells



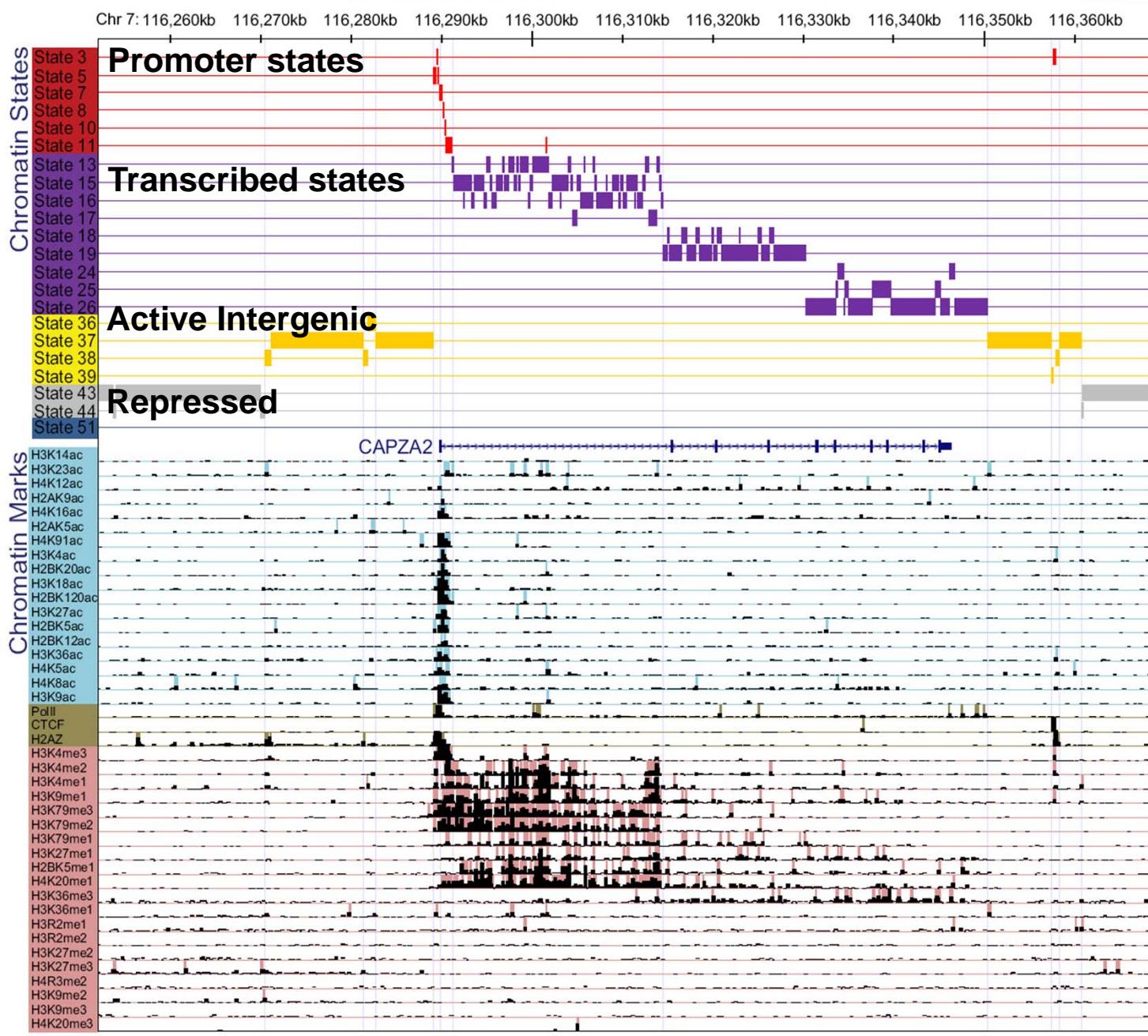
- **Dozens of chromatin tracks**
  - Understand their function
  - Reveal their combinations
  - Annotate systematically
- **Our approach: learn common chromatin states**
  - Explicitly model combinations
  - Unsupervised approach, probabilistic model

# Our approach: Multivariate Hidden Markov Model (HMM)



9

# From ‘chromatin marks’ to ‘chromatin states’



- Learn *de novo* significant combinations of chromatin marks
- Reveal functional elements, even without looking at sequence
- Use for genome annotation
- Use for studying regulation dynamics in different cell types

# ENCODE: Study nine marks in nine human cell lines

9 marks

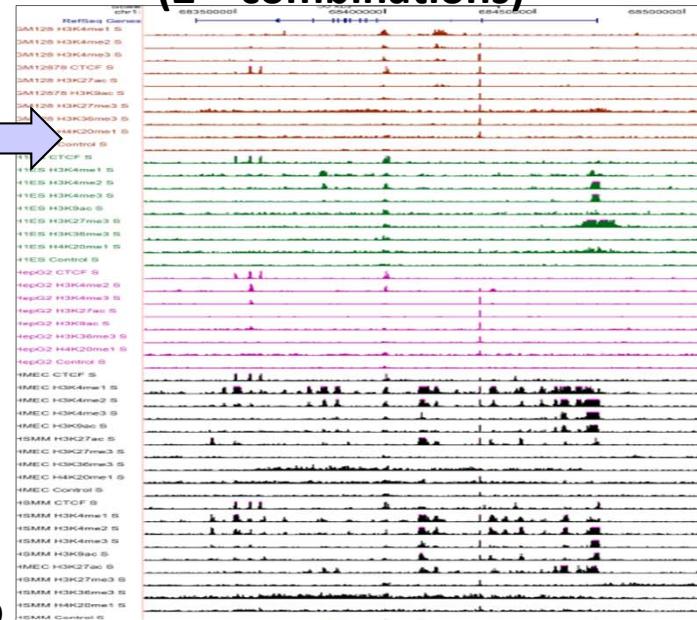
H3K4me1
H3K4me2
H3K4me3
H3K27ac
H3K9ac
H3K27me3
H4K20me1
H3K36me3
CTCF
+WCE
+RNA



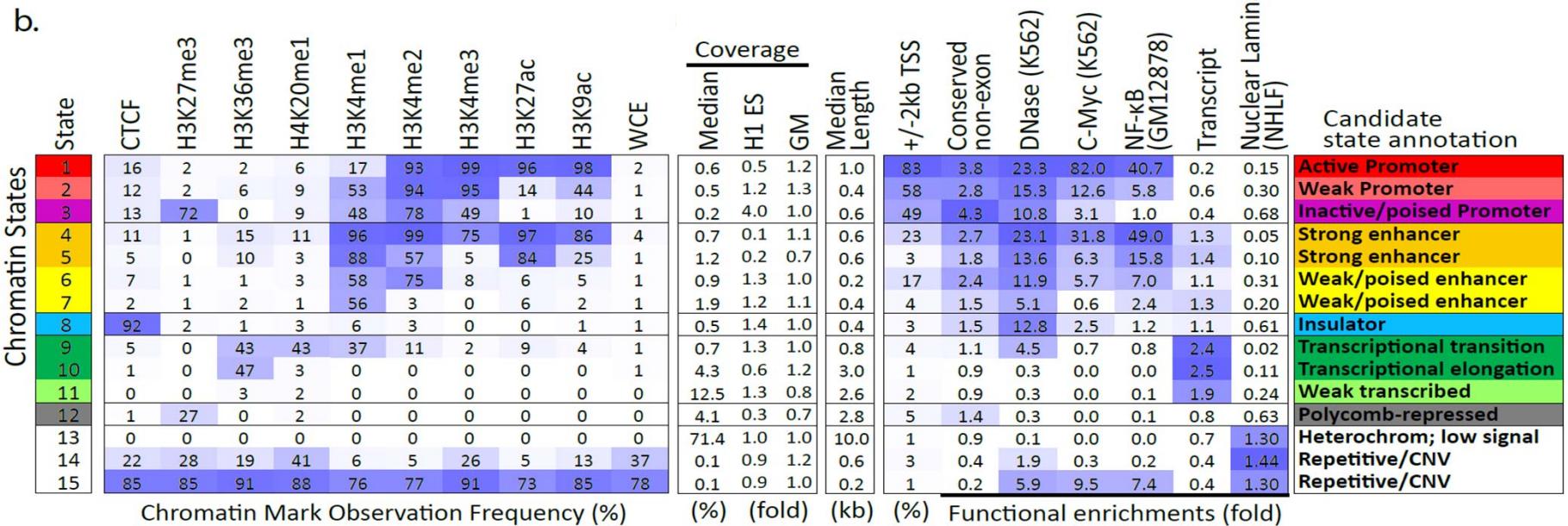
9 human cell types

HUVEC	Umbilical vein endothelial
NHEK	Keratinocytes
GM12878	Lymphoblastoid
K562	Myelogenous leukemia
HepG2	Liver carcinoma
NHLF	Normal human lung fibroblast
HMEC	Mammary epithelial cell
HSMM	Skeletal muscle myoblasts
H1	Embryonic

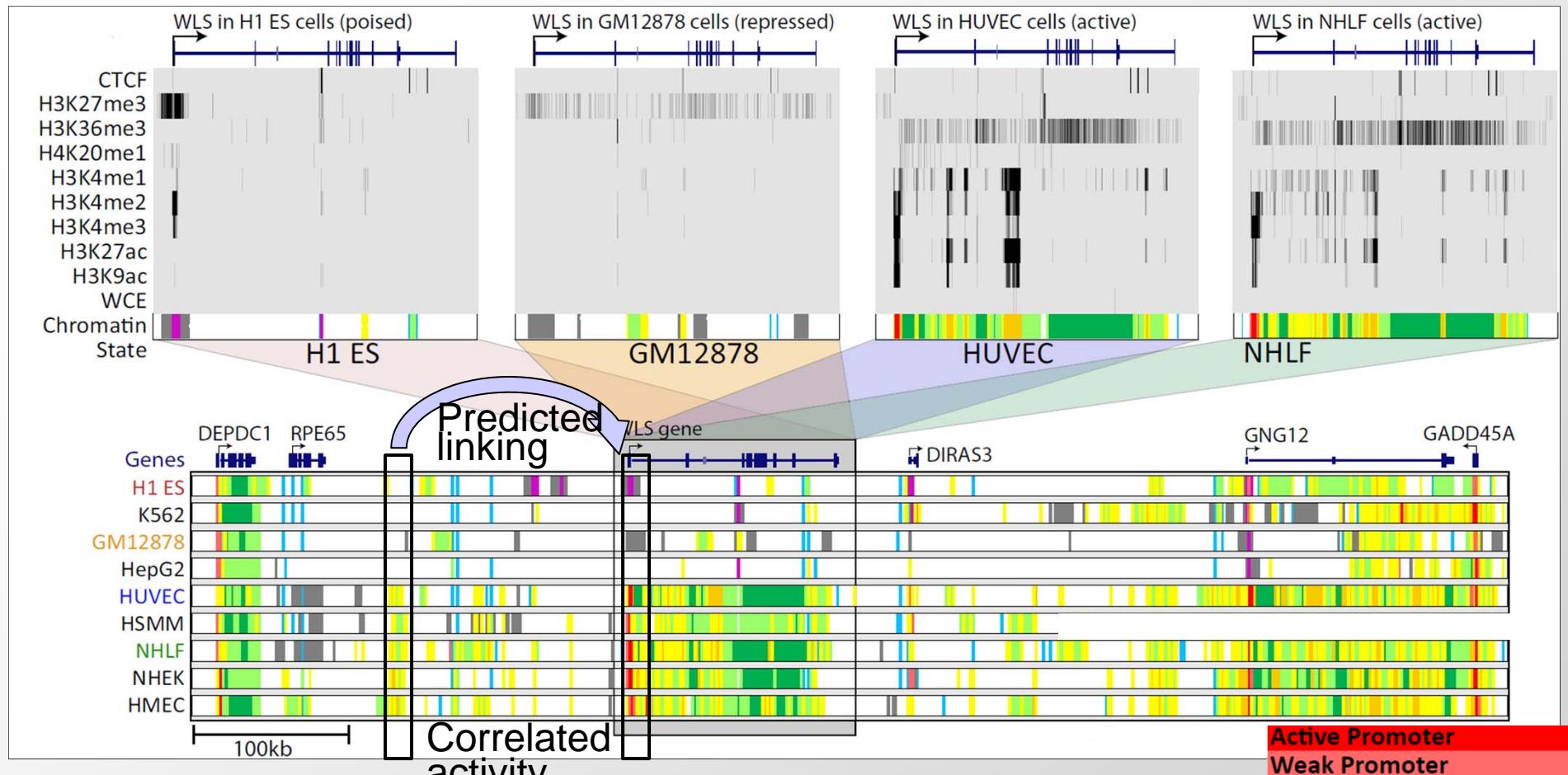
81 Chromatin Mark Tracks  
( $2^{81}$  combinations)



Brad Bernstein ENCODE Chromatin Group

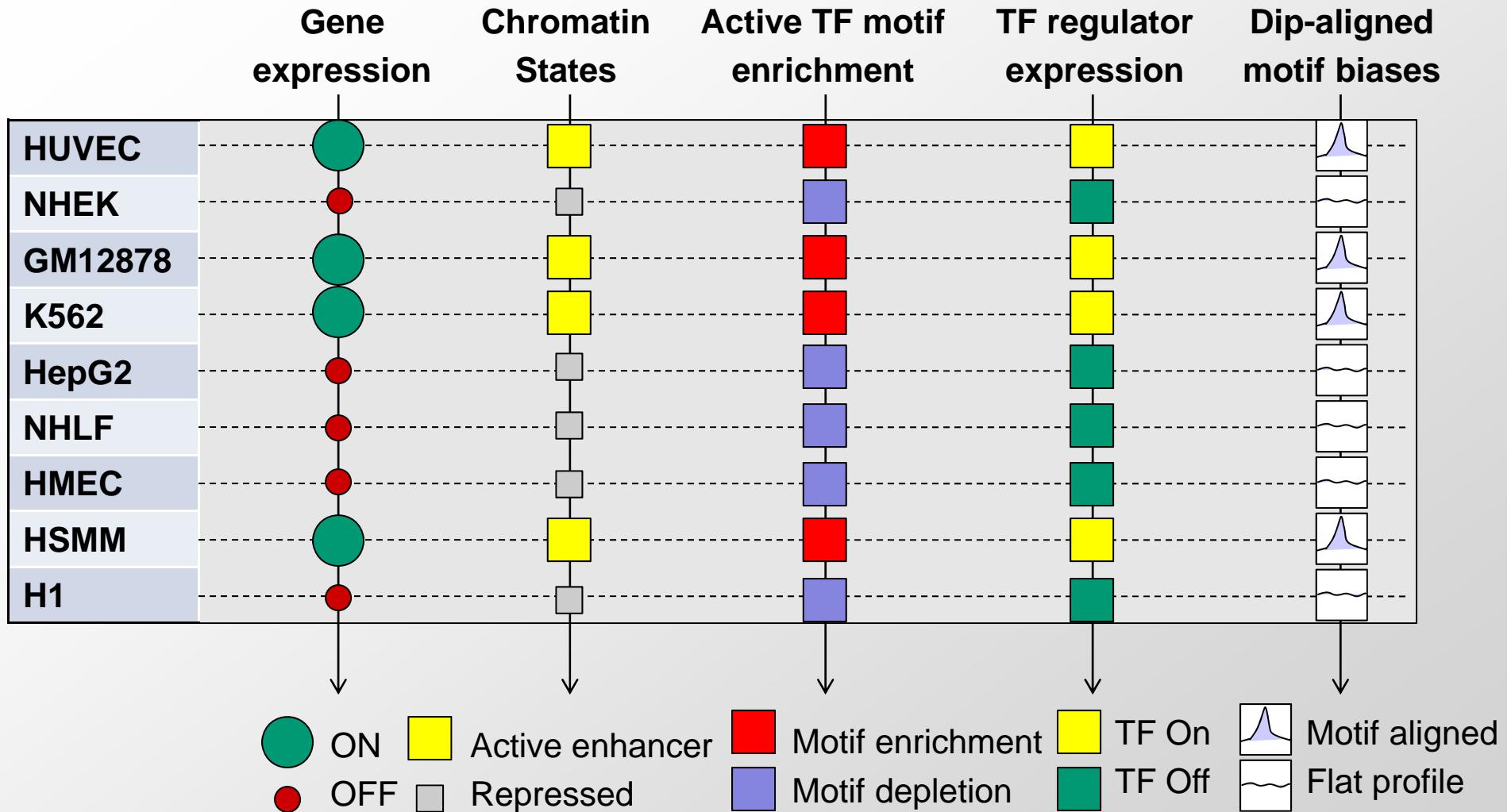


# Chromatin states dynamics across nine cell types



- Single annotation track for each cell type
- Summarize cell-type activity at a glance
- Can study 9-cell activity pattern across ↓

# Multi-cell activity profiles and their correlations



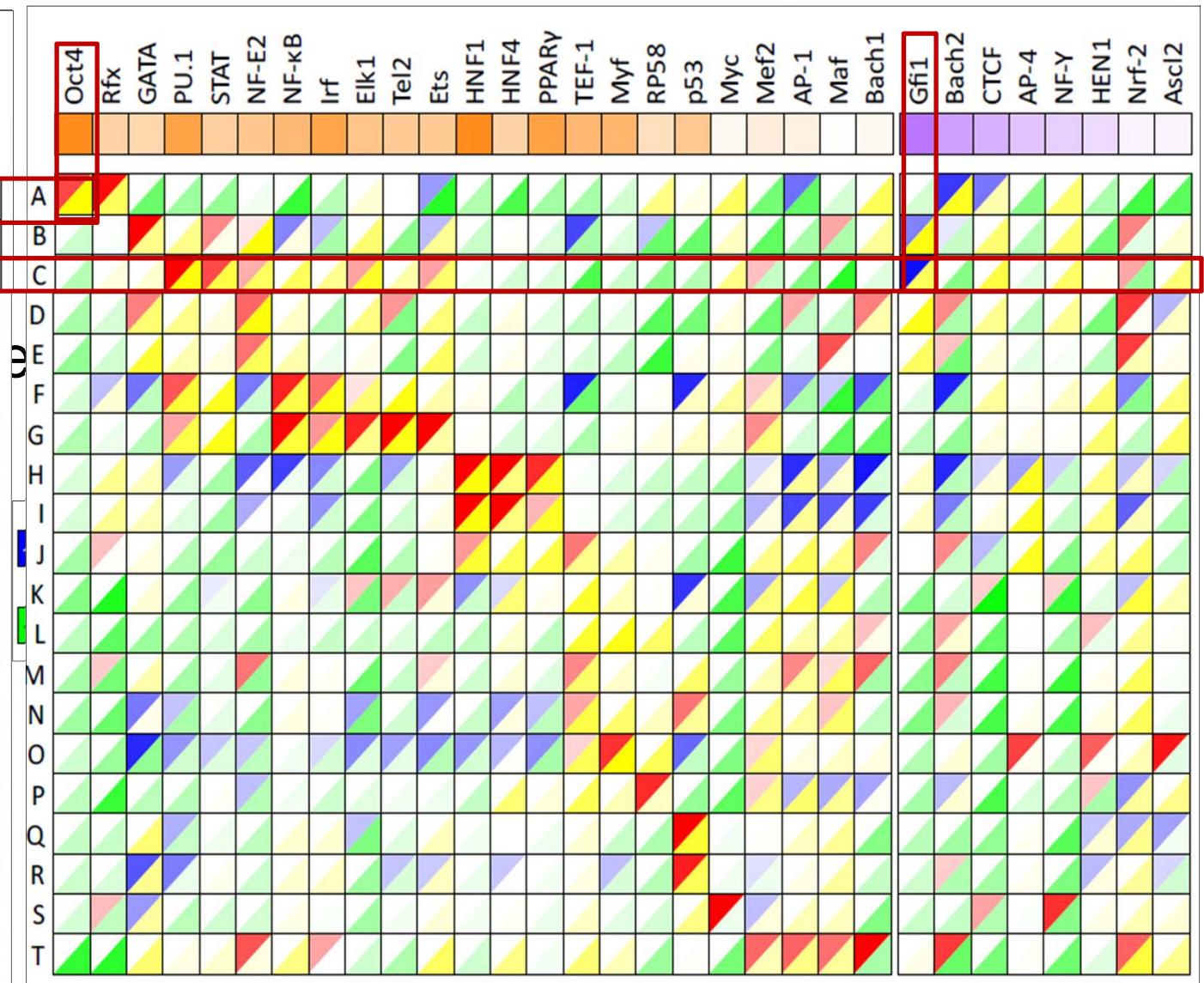
Chromatin state & gene expression → link enhancers and target genes  
 TF motif enrichment & TF expression → reveal activators / repressors

# Coordinated activity reveals activators/repressors

Enhancer activity

	H1 ES	K562	GM12878	HepG2	HUVEC	HSMM	NHLF	NHEK	HMEC
A	44	3	2	3	2	2	2	2	2
B	1	48	2	4	2	2	2	2	2
C	2	48	47	8	5	5	4	5	3
D	4	46	9	10	45	10	11	14	10
E	4	48	12	14	10	10	10	44	26
F	0	1	48	1	1	1	1	1	1
G	2	5	48	6	46	16	13	12	7
H	1	2	2	51	1	1	1	2	1
I	1	1	1	36	1	1	1	1	1
J	3	7	6	49	31	30	18	12	10
K	1	2	2	2	46	3	4	4	3
L	2	3	2	2	47	45	20	7	9
M	3	3	5	6	49	26	19	47	34
N	2	2	6	5	5	47	18	46	31
O	1	1	2	1	2	45	5	2	3
P	1	3	4	2	4	21	45	4	4
Q	2	1	5	5	3	3	4	49	45
R	1	1	3	2	3	3	3	45	8
S	2	4	4	4	8	6	5	8	43
T	11	40	27	25	45	41	39	45	41

Activity signatures for each TF



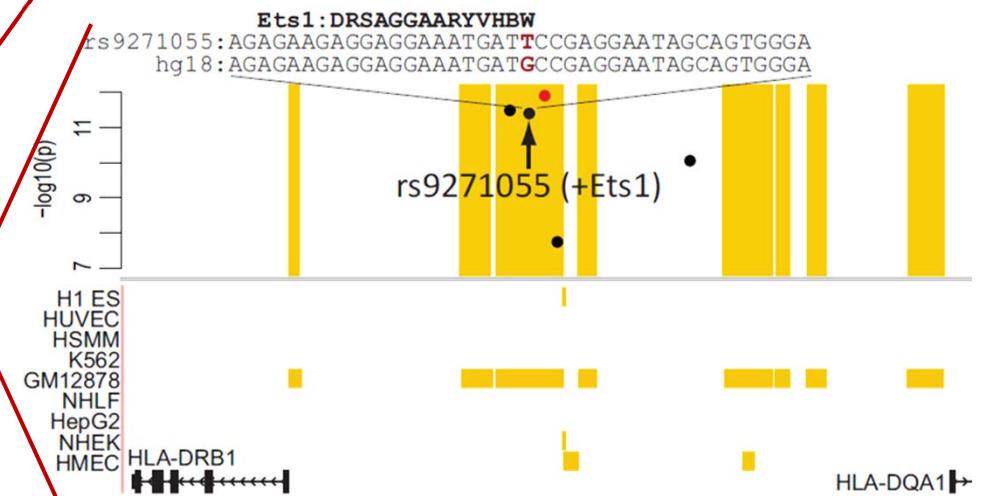
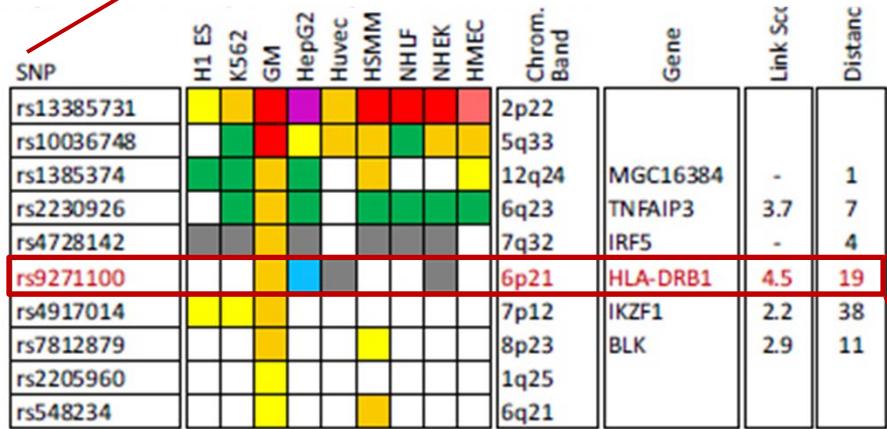
- Enhancer networks: Regulator → enhancer → target gene

# Revisiting disease-associated variants

## Phenotype

Erythrocyte phenotypes (Ref. 38)
Blood lipids (Ref. 39)
Rheumatoid arthritis (Ref. 40)
Primary biliary cirrhosis (Ref. 41)
Systemic lupus erythematosus (Ref. 42)
Lipoprotein cholesterol/triglycerides (Ref. 43)
Hematological traits (Ref. 44)
Hematological parameters (Ref. 45)
Colorectal cancer (Ref. 46)
Blood pressure (Ref. 47)

Top Cell Type	Total #SNPs from Study	#SNPs in enh. States 4 and 5	p-value	FDR	H1 ES	K562	GM12878	HepG2	HUVEC	HSMM	NHLF	NHEK	HMEC
K562	35	9	<10 <sup>-7</sup>	0.02	9	17	4	0	0	1	2	1	1
HepG2	101	13	<10 <sup>-7</sup>	0.02	3	5	0	11	2	3	3	4	3
GM12878	29	7	2.0 x 10 <sup>-7</sup>	0.03	0	0	15	0	2	0	0	2	3
GM12878	6	4	6.0 x 10 <sup>-7</sup>	0.03	0	11	41	0	0	0	0	8	8
GM12878	18	6	9.0 x 10 <sup>-7</sup>	0.03	0	4	21	0	5	8	0	3	5
HepG2	18	5	1.2 x 10 <sup>-6</sup>	0.03	17	8	0	24	3	6	4	3	3
K562	39	7	1.7 x 10 <sup>-6</sup>	0.03	0	12	10	2	1	0	0	1	0
K562	28	6	2.2 x 10 <sup>-6</sup>	0.03	0	15	7	0	5	7	7	3	2
HepG2	4	3	3.8 x 10 <sup>-6</sup>	0.03	0	0	0	66	0	12	0	12	12
K562	9	4	5.0 x 10 <sup>-6</sup>	0.04	0	30	14	0	10	6	7	5	11



- Disease-associated SNPs enriched for enhancers in relevant cell types
- E.g. **lupus SNP in GM enhancer disrupts Ets1 predicted activator**

# Regulatory roles revealed for many studies

Title	Author/Journal	Total #SNPs	Fold	Cell Type	# SNPs in enhancers	FDR
Multiple loci influence <b>erythrocyte</b> phenotypes in the CHARGE Consortium.	Ganesh et al Nat Genet 2009	35	17	K562	9	0.02
Biological, clinical and population relevance of 95 loci for <b>blood lipids</b>	Teslovich et al Nature 2010	101	11	HepG2	13	0.02
Genome-wide association study meta-analysis identifies seven new <b>rheumatoid arthritis</b> risk loci	Stahl et al Nat Genet 2010	29	15	GM12878	7	0.03
Genome-wide meta-analyses identify three loci associated with <b>primary biliary cirrhosis</b>	Liu et al Nat Genet 2010	6	41	GM12878	4	0.03
Chinese Han population identifies nine new susceptibility loci for <b>systemic lupus erythematosus</b> .	Han et al Nat Genet 2009	18	21	GM12878	6	0.03
Six new loci associated with blood low-density <b>lipoprotein cholesterol</b> , high-density <b>lipoprotein cholesterol</b> or <b>triglycerides</b> in humans.	Kathiresan et al Nat Genet 2008	18	24	HepG2	5	0.03
Genome-wide association study of <b>hematological</b> and biochemical traits in a Japanese population	Kamatani et al Nat Genet 2009	39	12	K562	7	0.03
A genome-wide meta-analysis identifies 22 loci associated with eight <b>hematological</b> parameters in the HaemGen consortium.	Soranzo et al Nat Genet 2009	28	15	K562	6	0.03
Meta-analysis of genome-wide association data identifies four new susceptibility loci for colorectal cancer.	Houlston et al Nat Genet 2008	4	66	HepG2	3	0.03
Genome-wide association study identifies eight loci associated with <b>blood</b> pressure.	Newton-Chen Nat Genet 2009	9	30	K562	4	0.04

# We work hard...

We aim to further our understanding of the **human genome** by computational integration of large-scale functional and comparative genomics datasets.

- We use comparative genomics of multiple related species to recognize **evolutionary signatures** of protein-coding genes, RNA structures, microRNAs, regulatory motifs, and individual regulatory elements.
- We use combinations of epigenetic modifications to define **chromatin states** associated with distinct functions, including promoter, enhancer, transcribed, and repressed regions, each with distinct functional properties.
- We develop **phylogenomic** methods to study differences between species and to uncover evolutionary mechanisms for the emergence of new gene functions

Our methods have led to numerous new insights on diverse regulatory mechanisms, uncovered evolutionary principles, and provide mechanistic insights for previously uncharacterized disease-associated SNPs

Science  
Nature  
Nature  
Nature

Nature  
Nature Biotech  
Nature  
Nature  
PLoS Genetics

MBE  
Genome Research  
Nature

Genome Research  
Nature  
Genome Research  
PLoS Comp. Bio.

Genes & Development  
Genome Research

Nature  
PNAS  
BMC Evo. Bio.  
ACM TKDD

Genome Research  
RECOMB  
J. Comp. Bio.  
PNAS

## Computational Biology - Selected Publications

### Comparative and Integrative Genomics

modENCODE - Functional elements and regulatory circuits in fly  
12flies - Evolutionary signatures for systematic genome interpretation in fly  
4yeasts - Gene identification and motif discovery in yeast  
8candida - Revealing pathogenic gene families in fungal genomes  
Nature Nature Nature Nature Nature In review  
Other - Dog - Neurospora - Drosophila - Platypus - modENCODE - 29 mammals

### Epigenomics and Chromatin regulation

Chromatin States - Combinatorial chromatin patterns for genome annotation  
Fly landscape - Chromatin landscape of Drosophila genes and regulatory elements  
Human Enhancers - Tissue-specific enhancers in multiple human cell types  
Fly Insulators - Genome-wide map of fly boundary elements  
Nature Gen Genes&Dev Nature Nature Biotech  
Other - Polycomb - Stalling - Embryo patterning - Yeast code - Human Roadmap

### Gene and Genome Evolution

Phylogenomics - Bayesian gene-tree phylogenomic reconstruction  
Phylogeny - Learning common gene-tree properties to overcome sparse information  
Duplication - Proof and analysis of Whole-Genome Duplication (WGD) in yeast  
Nature Nature Nature Nature WBpress  
Other - FishDup - FlyEvo - Platypus - Candida - Birds

### Coding and non-coding genes

Uncommon genes - Unusual coding genes in fly: read-through, frameshifts, poly-cistronics  
lncRNAs - Chromatin reveals a new class of long intergenic non-coding RNAs  
Overlapping - Excess synonymous constraint in coding exons  
Signatures - Comparing evolutionary and single-species metrics  
PNAS Nature G.R. BioChem  
Other - Revisiting Human and Yeast - CCDS - MassSpec

### Small Regulatory RNAs

AS-miRNAs - Discovery of a new class of anti-sense functional microRNAs  
microRNAs - Comparative identification of microRNAs in fly  
Nature GenomRes Nature G.R. Science  
Other - endo-siRNAs - miR-targets - piRNAs - platypus - Tasmanian

### Regulatory Motifs

Human Motifs - *De novo* TF/miRNA motif discovery in human  
Long motifs - New class of long motifs in human CNEs, including insulators  
PSSMs - Position-specific constraints match motif information content  
Heart Beat - Motif discovery in cardiogram datasets to predict heart failure  
RECOMB RECOMB  
Other - Yeast Motifs - Motif Pairs

### Regulatory Networks

Motif Targets - Using motifs to infer regulatory networks in fly  
Network Motifs - Discovering large motifs using symmetry-breaking conditions  
SubMAP - Aligning large pathways with subnetwork mapping  
Network Evol - Post whole-genome duplication (WGD) network evolution

# But we also have fun



Kayaking



Whitewater Rafting



BBQs



Canadian  
Thanksgiving

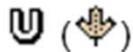


Sailing

# Want to get started?

## 6.047/6.878 – Computational Biology: Genomes, Networks, Evolution

### **6.047 Computational Biology: Genomes, Networks, Evolution**



(Subject meets with [6.878J](#), [HST.507J](#))

Prereq: [6.006](#), [6.041](#), [Biology \(GIR\)](#); or permission of instructor

Units: 3-0-9

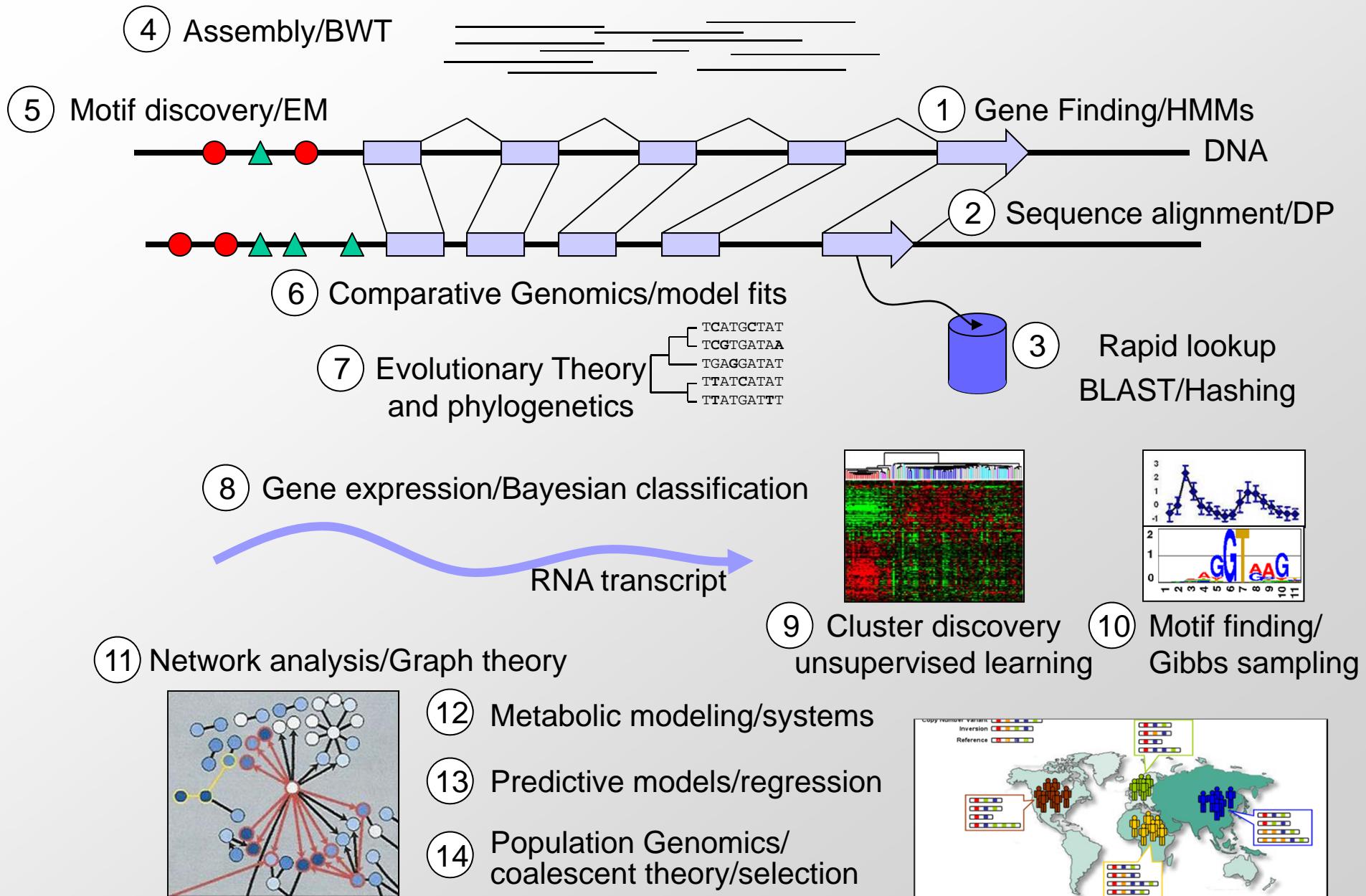


**Lecture:** TR9.30-11 (4-159) **Recitation:** F1 (36-156)

Covers the algorithmic and machine learning foundations of computational biology, combining theory with practice. Principles of algorithm design, influential problems and techniques, and analysis of large-scale biological datasets. Topics include (a) genomes: sequence analysis, gene finding, RNA folding, genome alignment and assembly, database search; (b) networks: gene expression analysis, regulatory motifs, biological network analysis; (c) evolution: comparative genomics, phylogenetics, genome duplication, genome rearrangements, evolutionary theory. These are coupled with fundamental algorithmic techniques including: dynamic programming, hashing, Gibbs sampling, expectation maximization, hidden Markov models, stochastic context-free grammars, graph clustering, dimensionality reduction, Bayesian networks.

*M. Kellis*

# Challenges in computational biology



# 4 modules: Genomes, Genes, Networks, Evolution

Psets	Week	Date	Topic	Category	Lec	Topic
PS1 out covers:L1-L5  due Mon 9/27	1	Thu, Sep 09	Introduction	Module I: Comparative Genomics	L1	Intro: Biology, Algorithms, Machine Learning
		Fri, Sep 10			R1	Recitation1: Probability, Statistics, Biology
	2	Tue, Sep 14		Foundations	L2	Global/local alignment/DynProg
		Thu, Sep 16			L3	StringSearch/Blast/DB Search
		Fri, Sep 17			R2	Recitation 2 - Multiple, Progressive, Phylogenetic, Whole-genome alignment
	3	Tue, Sep 21		Frontiers	L4	Comparative genomics, Selection, Evolutionary signatures of coding genes
		Thu, Sep 23			L5	Evolutionary Signatures of RNA structures, miRNAs, motifs, nucleotides
		Fri, Sep 24			R3	Recitation 3 - Evolutionary signatures and measures of selection
PS2 out covers:L6-R5  due Wed 10/13	4	Tue, Sep 28	Module II: Coding and non-coding genes	Foundations	L6	HMMs1 - Evaluation / Parsing
		Thu, Sep 30			L7	HMMs2 - PosteriorDecoding/Learning
		Fri, Oct 01			R4	Recitation4 - Posterior decoding review, Baum-Welch Learning
	5	Tue, Oct 05		Frontiers	L8	Gene finding in practice: CRFs, Discriminative training, GHMMs, RNAseq
		Thu, Oct 07			L9	Structural RNAs: Fold prediction and genome-wide annotation
		Fri, Oct 08			R5	Recitation 5 - SVMs and discriminative learning
PS3 out covers:L10-R8  due Mon 11/1	6	Thu, Oct 14	Module III: Networks and Gene Regulation	Foundations	L10	Expression Analysis: Clustering, Classification, Feature Selection
		Fri, Oct 15			R6	Recitation 6 - Microarrays/RNAseq expression analysis
	7	Tue, Oct 19		Foundations	L11	Regulatory Motif Discovery: Gibbs Sampling, Expectation Maximization
		Thu, Oct 21			L12	TF/miRNA target prediction and regulatory network inference
	8	Fri, Oct 22		Frontiers	R7	Recitation 7 - Entropy, Information, Background models
		Tue, Oct 26			L13	Regulatory Network Analysis, Bayesian Function Prediction
		Thu, Oct 28			L14	Epigenomics: Chromatin marks, chromatin states, and their dynamics
		Fri, Oct 29			R8	Recitation 8 - Regulatory genomics and epigenomics
PS4 out covers:L15-R10  due Mon 11/22	9	Tue, Nov 02	Module IV: Evolution and Phylo- genomics	Foundations	L15	Phylogenetics: Molecular Evolution, Tree Building, Phylogenetic inference
		Thu, Nov 04			L16	Phylogenomics: gene/species trees,reconciliation,coalescent,populations
		Fri, Nov 05			R9	Recitation 9 - Gene Trees, Species Trees, Reconciliation
	10	Tue, Nov 09		Foundations	L17	Population genetics: Statistical genetics and human disease mapping
		Fri, Nov 12			R10	Recitation 10 - Population genetics and genomics
	11	Tue, Nov 16		Frontiers	L18	Population genetics: Measuring natural selection in human populations
		Thu, Nov 18			L19	Population genetics: Learning population history from genetic data
		Fri, Nov 19			R10	Recitation 11 - Quiz Review
	12	Tue, Nov 23	In Class (Friendly) Quiz	L20	In Class Quiz (the only Quiz) - covers L1-R10	

- Foundations and frontiers. Practical/algorithmic problems

# Final project: mentoring, milestones

Week	Term Project	Psets	Week	Date	Lec
1	<b>I. Project Set-up.</b> Describe your previous research, areas of interest in computational biology, type of project that best fits your interests. Post these in a profile that lets your classmates know you and find potential partners. <b>Due Mon 9/27 with PS1</b>	PS1 out covers:L1-L5  due Mon 9/27	1	Thu, Sep 09	
2			2	Fri, Sep 10	
3			3	Tue, Sep 14	Module I: Comparative Genomics
4	<b>II. Brainstorming.</b> Identify previous project proposals, recent papers, and potential partners that match your areas of interest. List initial project ideas and partners. <b>Due Wed 10/13 with PS2</b>		4	Thu, Sep 16	
5			5	Fri, Sep 17	
6	<b>III. Proposal.</b> Form teams (2 or 1), specify project goals, division of work, milestones, datasets, challenges, aims, algorithms, in form of NIH proposal. <b>Due Mon 10/25.</b>		6	Tue, Sep 21	
7			7	Thu, Sep 23	
8	<b>b. Review.</b> Evaluate/discuss 3 peer proposals, NIH format, provide suggestions/advice. <b>Due Mon 11/1 with PS3.</b>		8	Fri, Sep 24	
9	<b>c. Response.</b> Bullet response to reviews, revise scope/aims as needed.		9	Tue, Sep 28	
10	<b>IV. Midterm progress report.</b> Continue making progress on proposed milestones. Adjust aims/deliverables as needed. Write outline of final report, with draft methods/figs/tables <b>Due Mon 11/22 with PS4</b>		10	Thu, Sep 30	Module II: Coding and non-coding genes
11			11	Fri, Oct 01	
12	<b>V. Final project report.</b> Finish methods, milestones, results, figures, write-up in conference paper format. In report, comment on lessons learned and overall project experience. <b>Due Mon 12/6.</b>		12	Tue, Oct 05	Module III: Networks and Gene Regulation
13			13	Thu, Oct 07	
14	<b>VI. Final class presentation</b> ~10min conference talk on your final project. <b>Thu 12/9</b>	Complete Final Project	14	Fri, Oct 08	
				Tue, Oct 12	
				Thu, Oct 14	
				Fri, Oct 15	
				Fri, Oct 15	Project Mentoring
				Tue, Oct 19	
				Thu, Oct 21	
				Fri, Oct 22	
				Tue, Oct 26	
				Thu, Oct 28	
				Fri, Oct 29	
				Fri, Oct 29	Project Mentoring
				Tue, Nov 02	
				Thu, Nov 04	
				Fri, Nov 05	Module IV: Evolution and Phylo-genomics
				Tue, Nov 09	
				Thu, Nov 11	
				Fri, Nov 12	
				Fri, Nov 12	Project Mentoring
				Tue, Nov 16	
				Thu, Nov 18	
				Fri, Nov 19	
				Tue, Nov 23	Final Project Mode
				Thu, Nov 25	
				Fri, Nov 26	
				Tue, Nov 30	
				Thu, Dec 02	
				Fri, Dec 03	
				Fri, Dec 03	Project Mentoring
				Tue, Dec 07	
				Thu, Dec 09	9:30am
				Thu, Dec 09	3pm

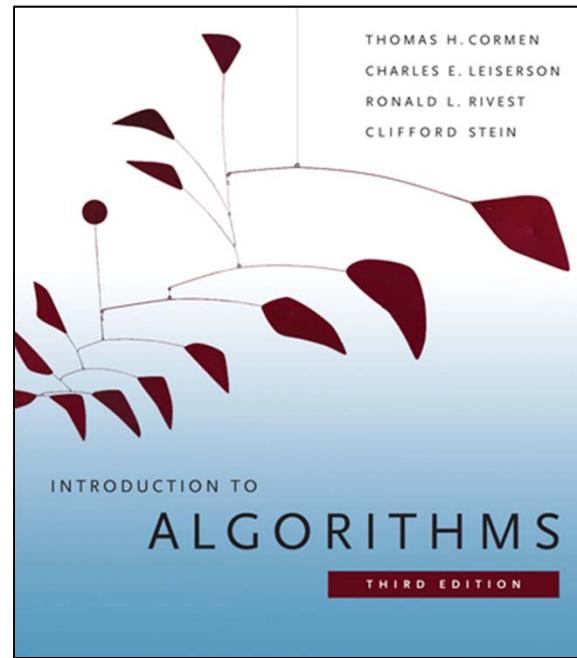


**For more info: [compbio.mit.edu](http://compbio.mit.edu)  
(or talk to any of us...)**



6.006

# *Introduction to Algorithms*



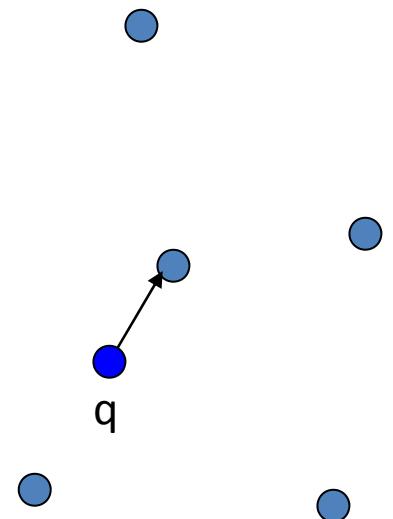
Lecture 25 2/3

(High-dimensional) geometry

Or: 6.850 “Geometric Computing” preview

# Nearest Neighbor

- Given: a set  $P$  of  $n$  points in  $\mathbb{R}^d$
- **Nearest Neighbor:** for any query  $q$ , returns a point  $p \in P$  minimizing the Euclidean distance  $\|p-q\|$



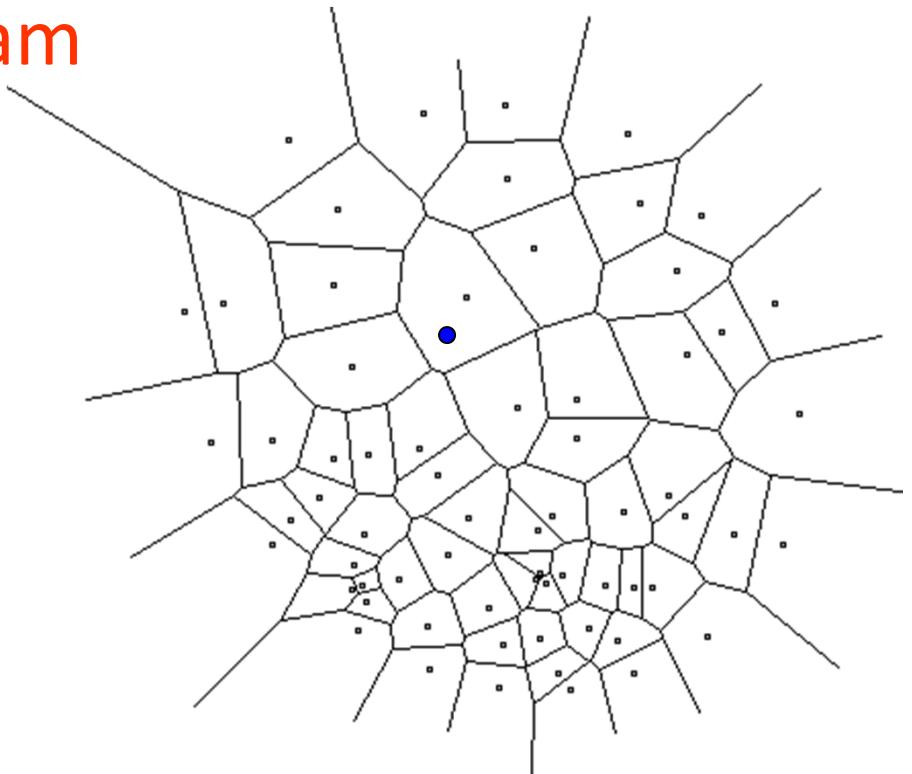
# The case of d=1

- Build BST
- Next-larger( $q$ ): finds the next element after element  $q$
- Next-smaller( $q$ ): analogous
- The closer of the two is the nearest neighbor
- Performance:
  - Space:  $O(n)$
  - Query time:  $O(\log n)$assuming balanced BST



# The case of d=2

- Compute **Voronoi diagram**
- Given  $q$ , perform **point location**
- Performance:
  - Space:  $O(n)$
  - Query time:  $O(\log n)$



# The case of $d > 2$

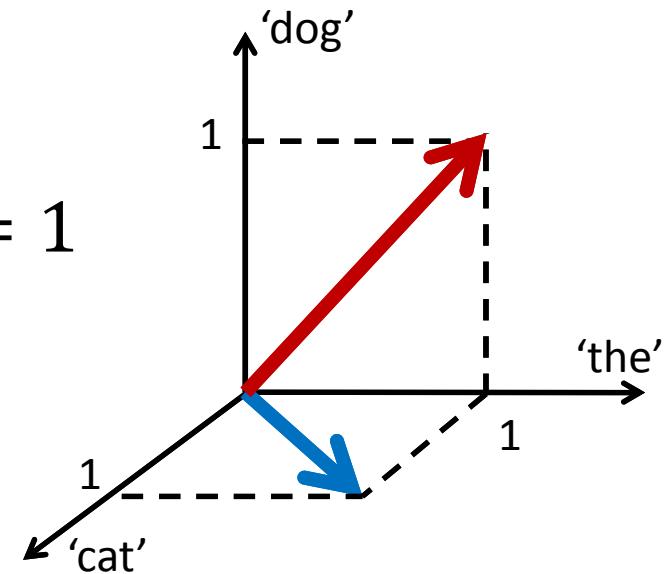
- Voronoi diagram has size  $n^{O(d)}$   
“Curse of dimensionality”
- We can also perform a linear scan:  $O(dn)$  time
- Both are pretty bad if  $d, n = \text{few million}$
- Why would  $d$  be a few million ?

# Example: Vector Space Model

[Salton, Wong, Yang 1975]

- Treat each document  $D$  as a vector of its words
  - One coordinate  $D(w)$  for every possible word  $w$
- Example:
  - $D_1$  = “the cat”
  - $D_2$  = “the dog”
- Similarity between vectors?
  - Dot product:

$$D_1 \cdot D_2 = 1$$



$$D_1 \cdot D_2 = \sum_w D_1(w) \cdot D_2(w)$$

# Vector Space Model ctd

[Salton, Wong, Yang 1975]

- We have

$$\|D_1 - D_2\|^2 = \|D_1\|^2 + \|D_2\|^2 - 2 D_1 \cdot D_2$$

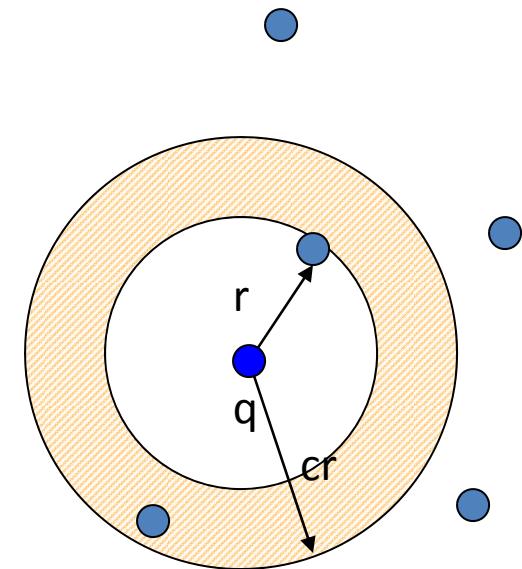
- If we normalize  $D_1, D_2$  then

$$\|D_1 - D_2\|^2 = 2 - 2 D_1 \cdot D_2$$

- Minimizing Euclidean distance = Maximizing Dot Product
- Many other applications: searching for similar bio sequences, similar images, etc

# Approximate Near Neighbor

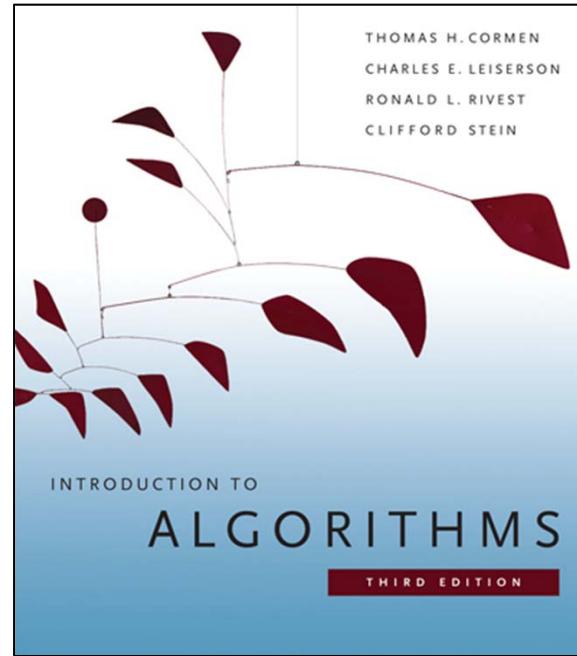
- $c$ -Approximate Nearest Neighbor: build data structure which, for any query  $q$ 
  - returns  $p' \in P$ ,  $\|p'-q\| \leq cr$ ,
  - where  $r$  is the distance to the nearest neighbor of  $q$
- Can beat the curse
- Example algorithm:
  - Space:  $O(dn+n^{1+1/c})$
  - Query time:  $O(dn^{1/c})$



- 6.850 “Geometric Computing”
- Next Spring
- Coming to lecture room near you!

**6.006**

# *Introduction to Algorithms*



**Lecture 26c: Beyond  
Prof. Erik Demaine**

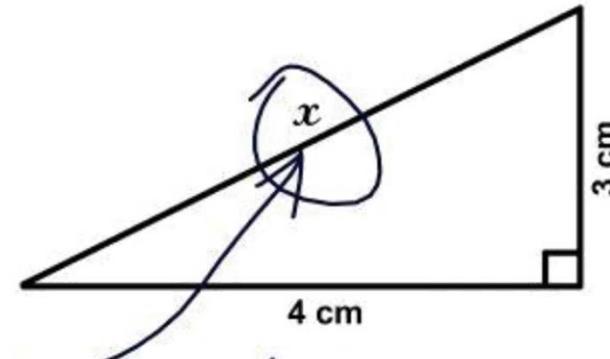
# Erik's Main Research Areas

- Data structures
- Graph algorithms
- Geometric folding algorithms
- Recreational algorithms

# 6.851: Advanced Data Structures

- Dynamic graphs
- Integer structures
- String structures
- Geometric DSS
- Reducing space
- Time travel
- Realistic models: cache, disk, GPUs, ...
- Best possible binary search trees

3. Find  $x$ .



Ocular Trauma - by Wade Clarke ©2005

## Searching for $x$ ?

Learn how with

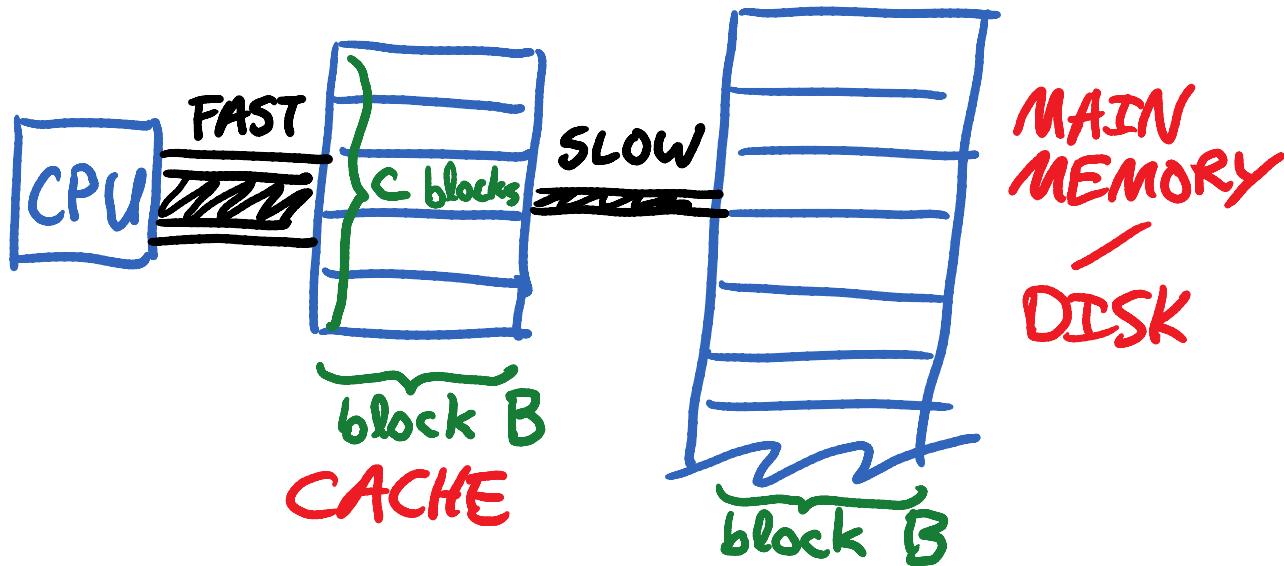
**Advanced Data Structures**

# Integer Data Structures

- Store  $n$  integers in range  $\{0, 1, \dots, u - 1\}$  subject to insert, delete, successor, and predecessor in
  - $O(\lg \lg u)$  time *[van Emde Boas]*
  - $O\left(\frac{\lg n}{\lg \lg u}\right)$  time *[fusion trees]*
  - $O\left(\sqrt{\frac{\lg n}{\lg \lg n}}\right)$  time *[combination]*

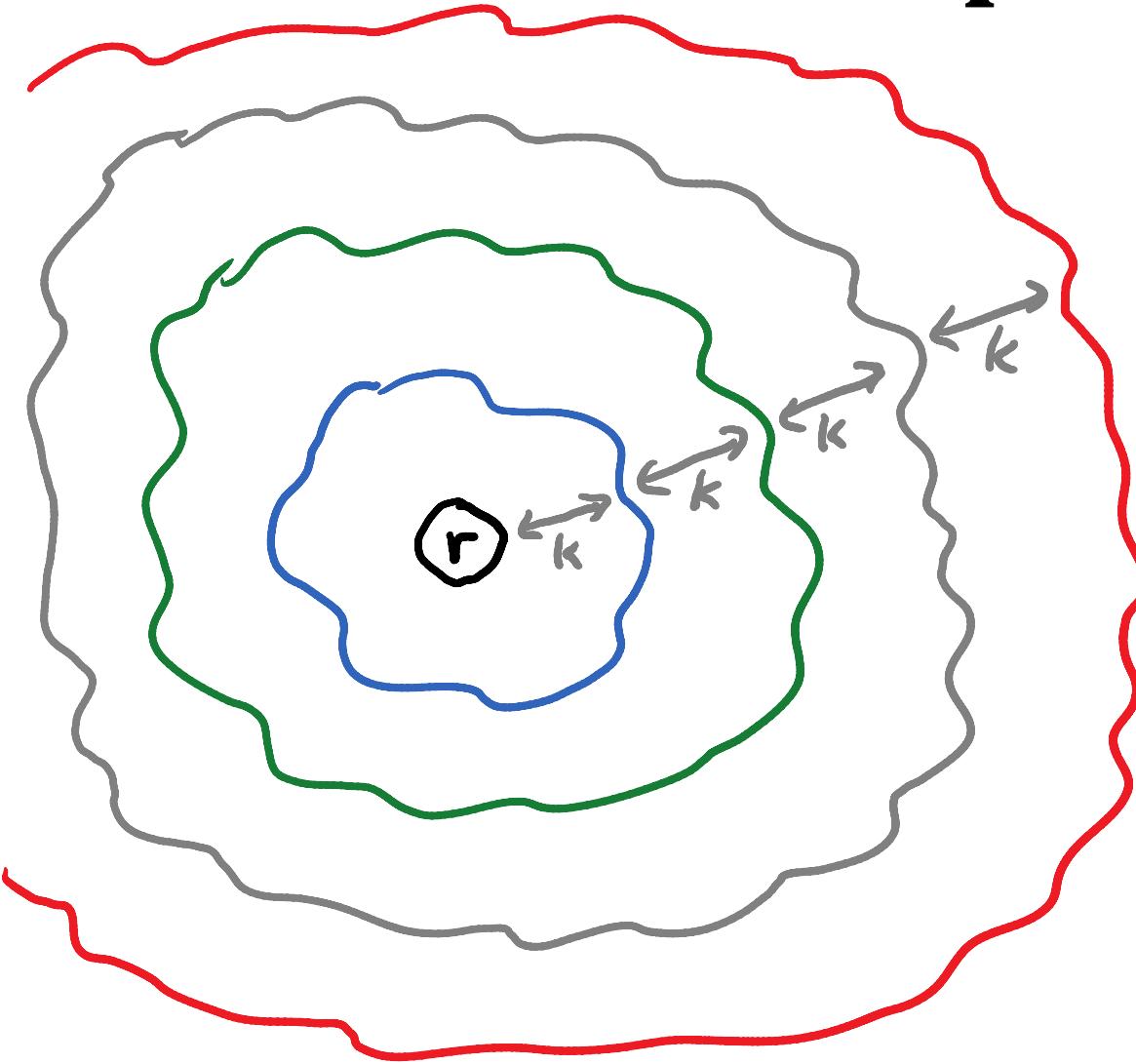
# Cache-Efficient Data Structures

- Memory transfers happen in blocks:



- Searching takes  $O(\log_B n)$  memory transfers
- Sorting takes  $O\left(\frac{n}{B} \log_C \frac{n}{B}\right)$  memory transfers
- ...even if you don't know  $B$  or  $C$ !

# Approximation Algorithms in Planar Graphs (& more)



- Delete every  $k$ th BFS layer
- Rings are width- $O(k)$  “trees”
  - Most problems solved in  $n^{O(k)}$
- Patch up at a factor  $\approx 1 + 1/k$
- Approximate better as  $k \rightarrow \infty$

## 6.849: Geometric Folding Algorithms: Linkages, Origami, Polyhedra (Fall 2010)

Prof. Erik Demaine

[Home] [Problem Sets] [Project] [Lectures] [Problem Session Notes]

### Lecture 5 Video [\[previous\]](#)

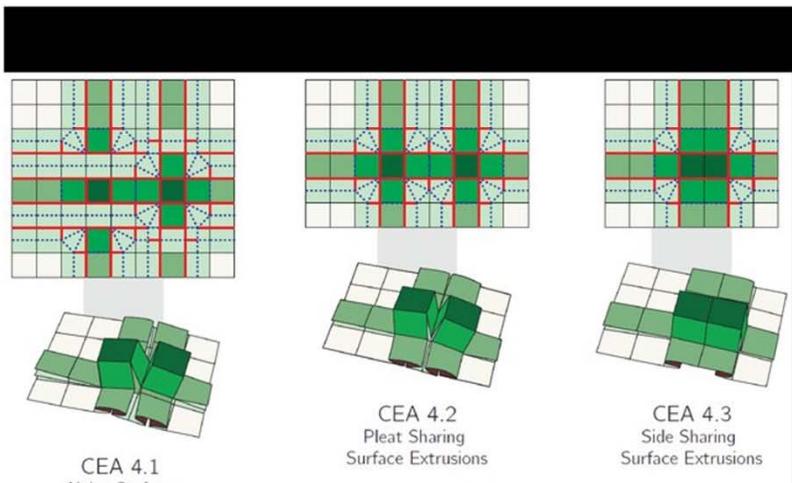
[\[+\]](#) Universal hinge patterns: box pleating, polycubes; orthogonal maze folding. NP-hardness: introduction, reductions; simple foldability; crease pattern flat foldability; disk packing (for tree

Handwritten notes, page 1/7 • [\[previous page\]](#) • [\[next page\]](#) • [\[PDF\]](#)



Download Video: [360p](#)

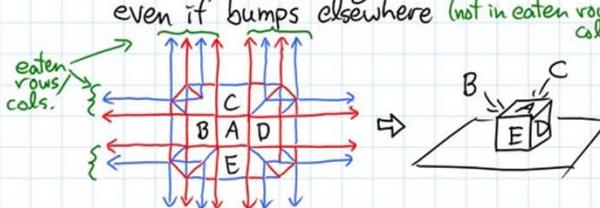
Slides, page 5/20 • [\[previous page\]](#) • [\[next page\]](#) • [\[PDF\]](#)



6.849 | Lecture 5 | Sept. 22, 2010

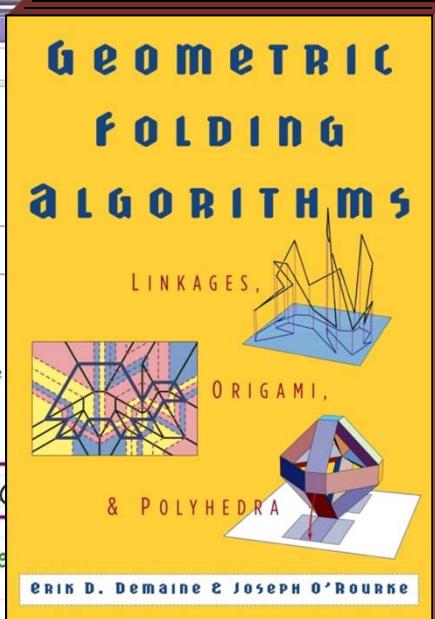
Universal hinge patterns: (for origami transforms)  
[Benednou, Demaine, Demaine, Ovadry 2010]

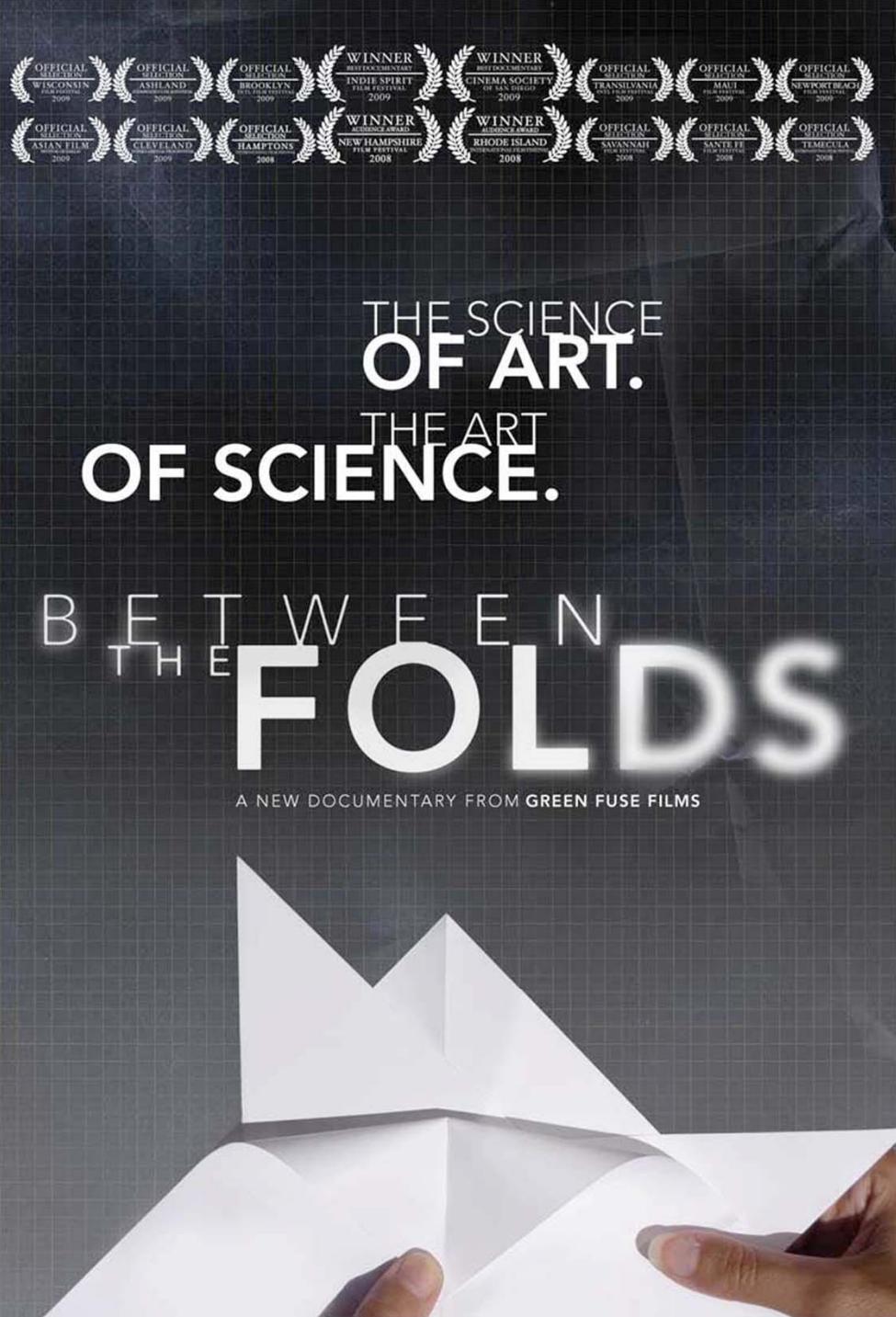
- suppose crease pattern required to be subset of fixed "hinge pattern"  
(e.g. Origamizer uses completely different creases for every model)
- $n \times n$  box-pleat pattern can make any polycube of  $O(n)$  cubes, seamless:
  - cube gadget turns  $O(1)$  rows & columns into a cube sticking out of sheet  $\sim$  even if bumps elsewhere (not in eaten rows/cols.)



- to make a tree of cubes: (=any polycube)
  - make a leaf
  - conceptually remove it { "postorder traversal"
  - repeat
- actually need to reserve space ahead of time for all the cube gadgets

Handwritten notes, page 1/7 • [\[previous page\]](#) • [\[next page\]](#) • [\[PDF\]](#)

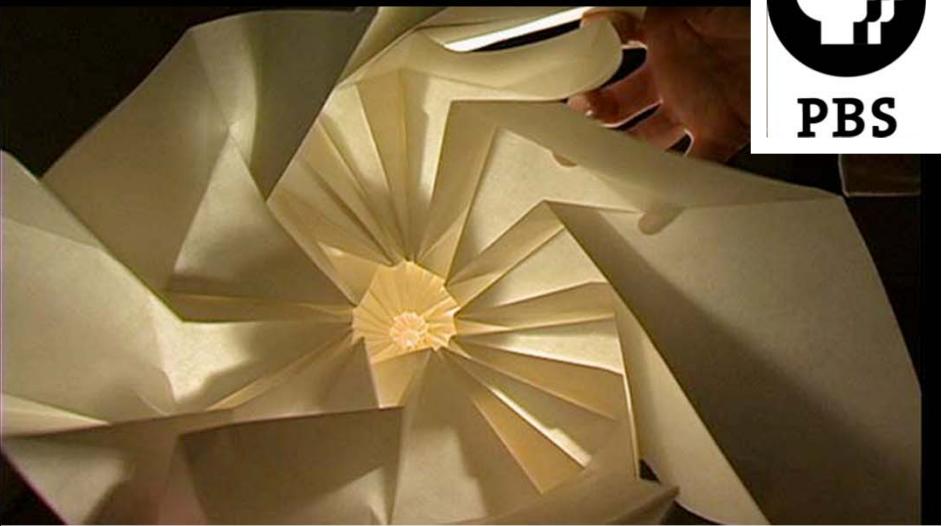




# THE SCIENCE OF ART. THE ART OF SCIENCE.

# BETWEEN FOLDS

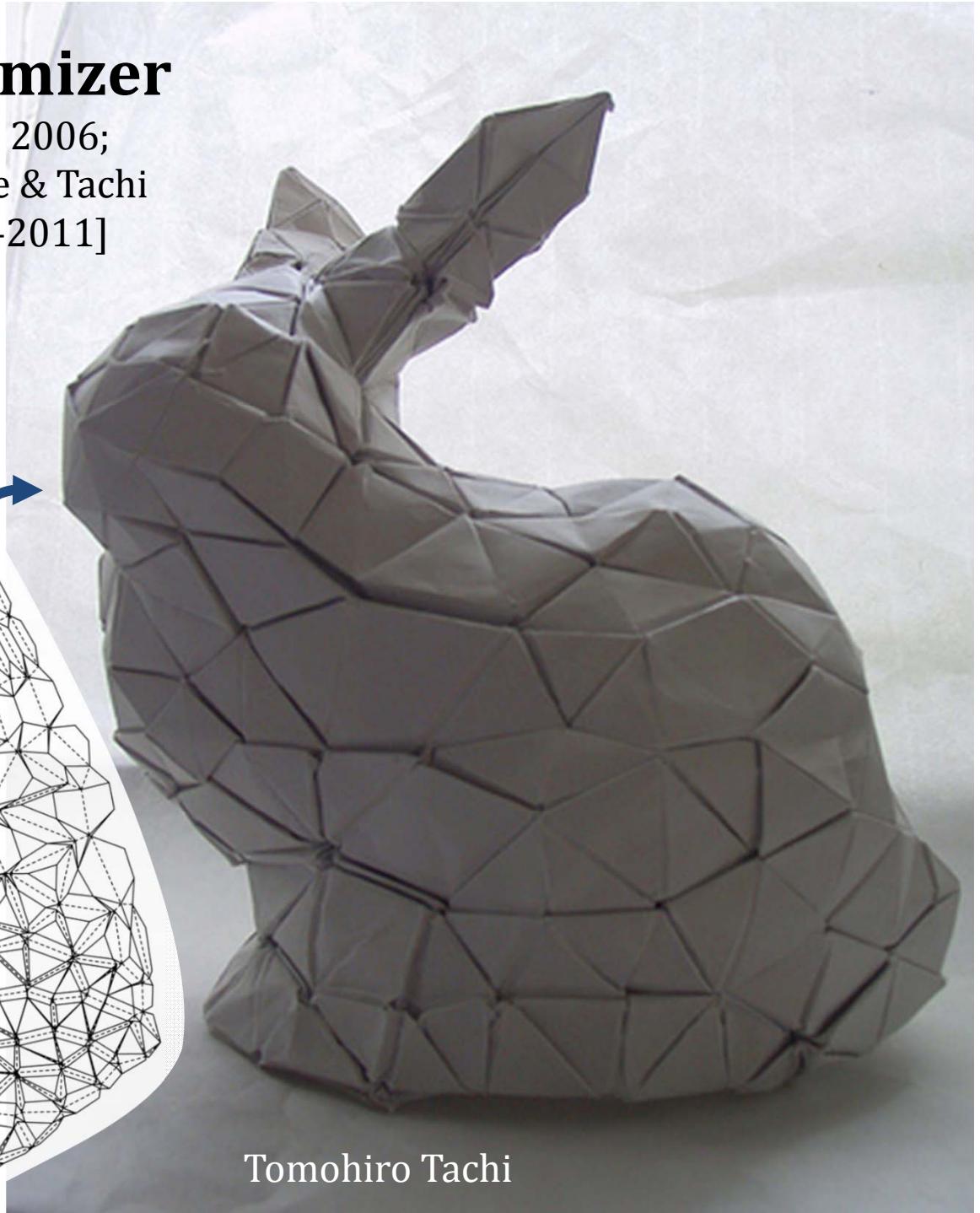
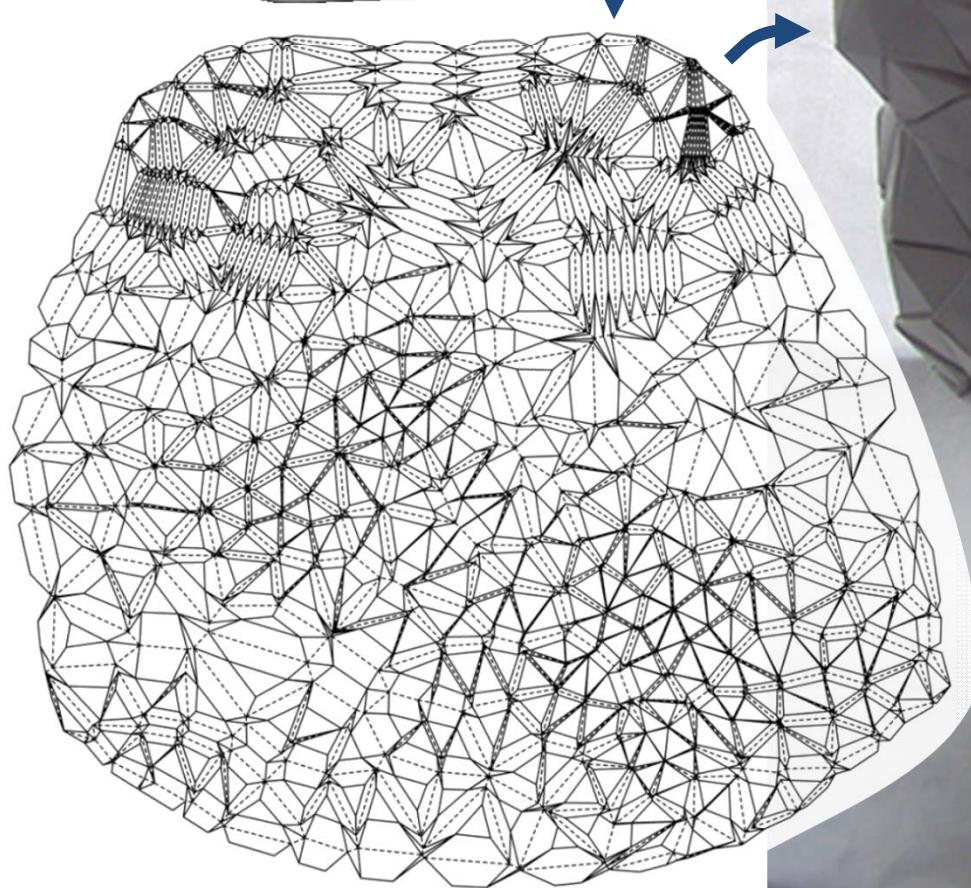
A NEW DOCUMENTARY FROM GREEN FUSE FILMS





# Origamizer

[Tachi 2006;  
Demaine & Tachi  
2009–2011]





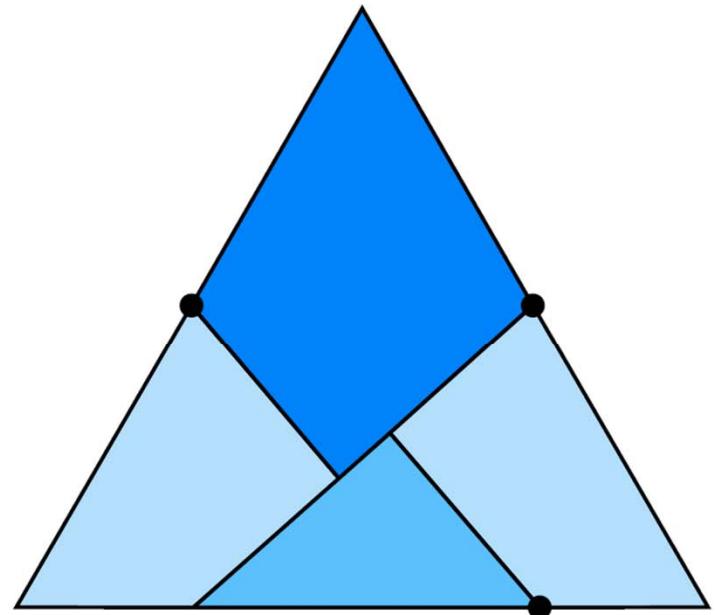
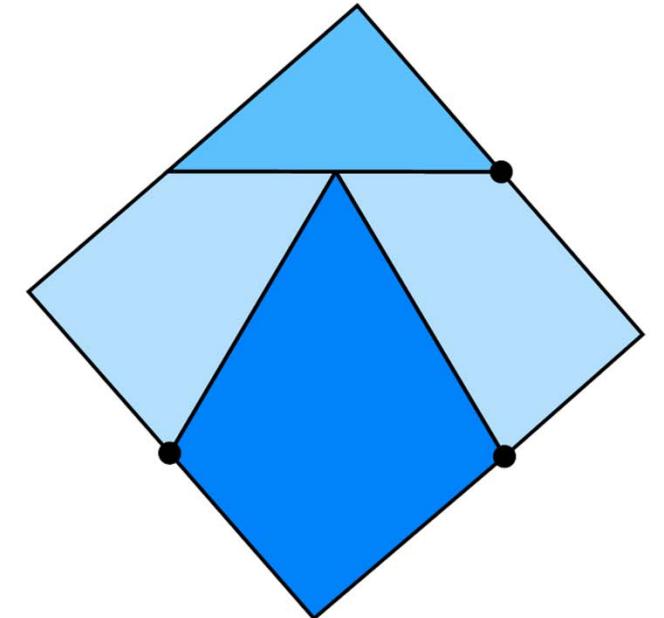
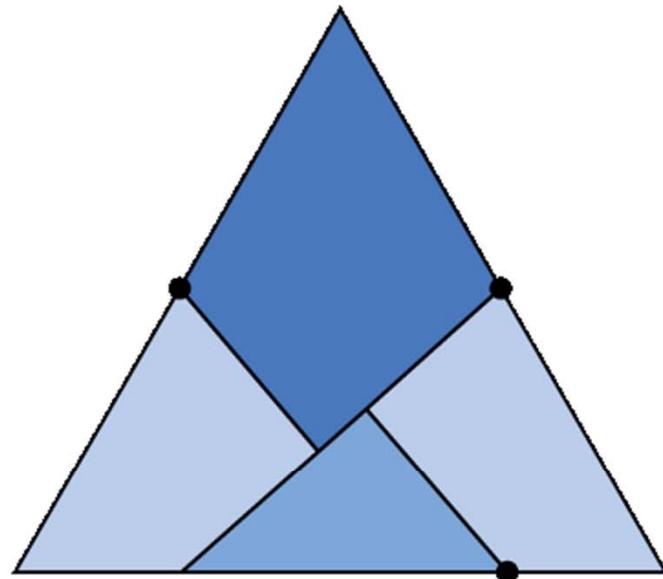
## “Natural Cycles”

Erik & Martin Demaine

Renwick Gallery,  
Smithsonian American  
Art Museum, 2012



# Hinged Dissection

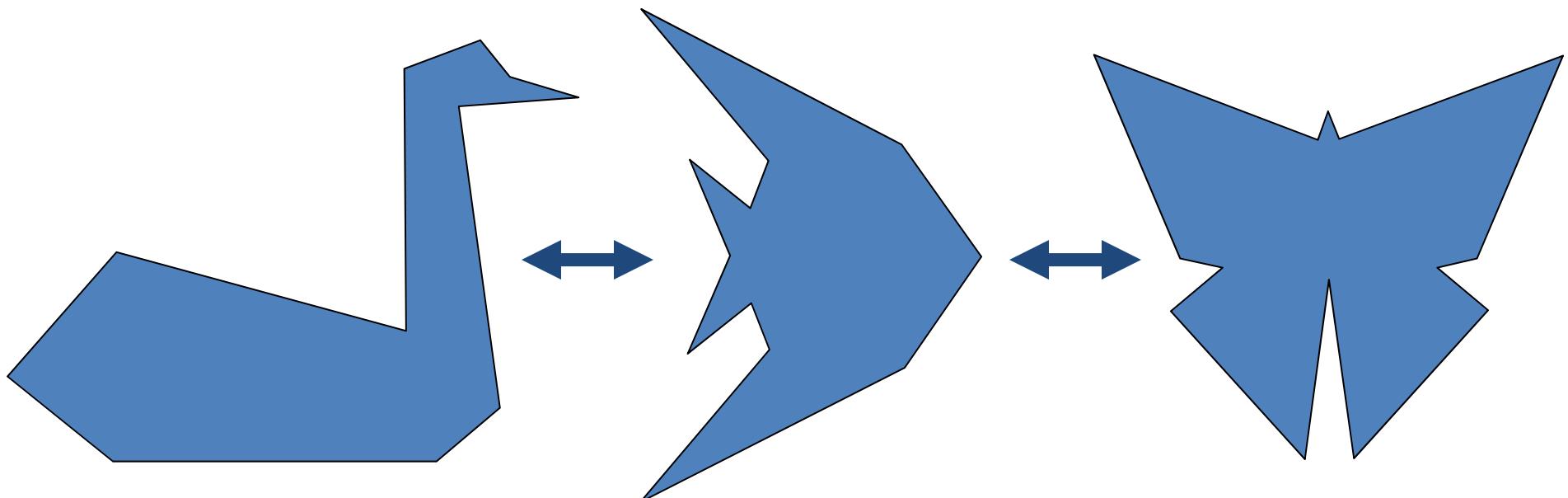


[Dudeney 1902]

# Hinged Dissection Universality

[Abbott, Abel, Charlton, Demaine, Demaine, Kominers 2008]

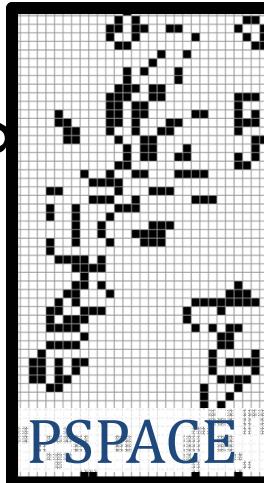
- Theorem: For *any finite set* of polygons of equal area, there is a hinged dissection that can fold into any of the polygons, *continuously without self-intersection*



# Complexity of Games & Puzzles

[Demaine, Hearn & many others]

unbounded  
bounded



PSPACE



PSPACE



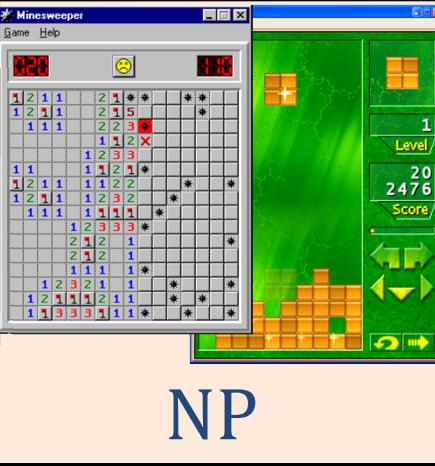
EXPTIME



Undecidable



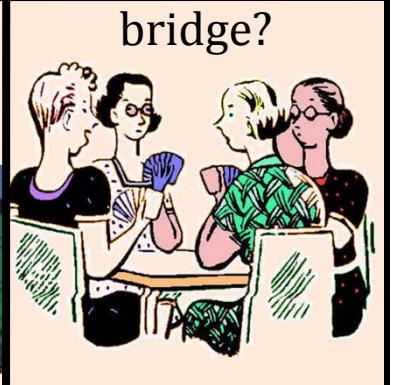
P



NP



PSPACE



bridge?

NEXPTIME

0 players  
(simulation)

1 player  
(puzzle)

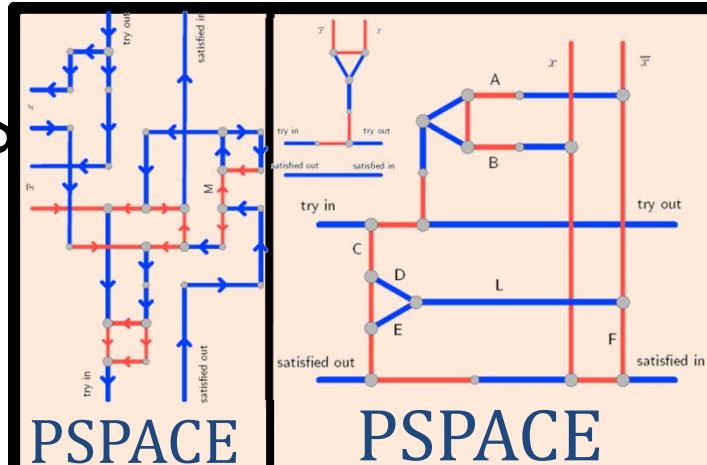
2 players  
(game)

team,  
imperfect info

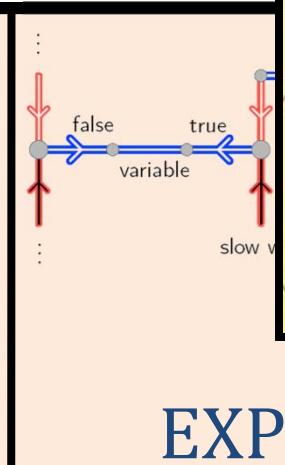
# Constraint Logic

[Hearn & Demaine 2009]

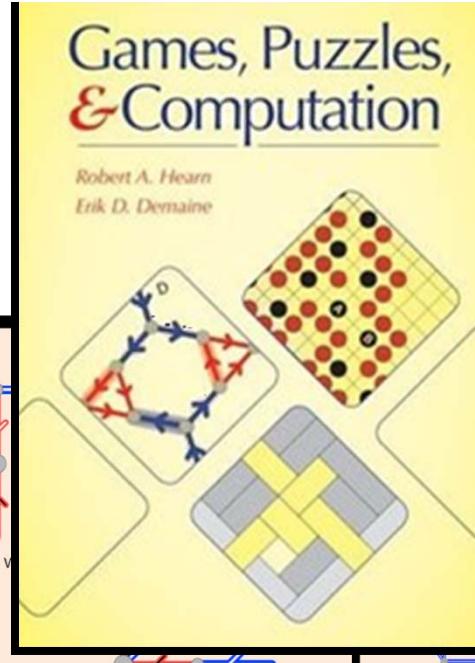
unbounded



PSPACE

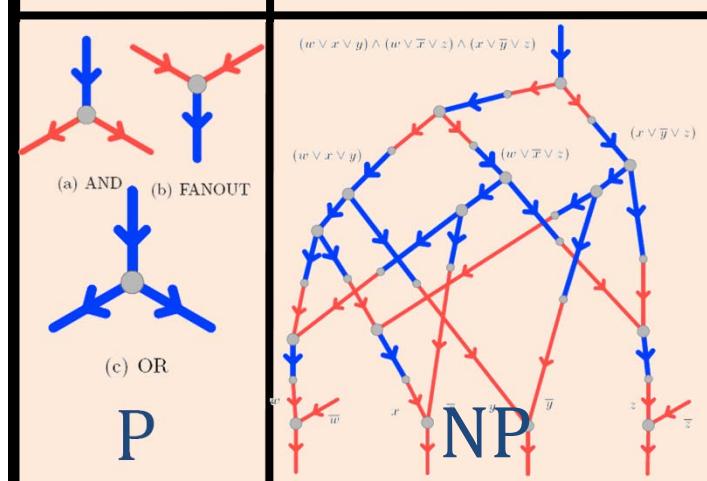


EXPTIME

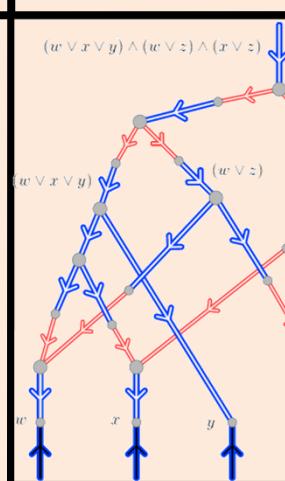


Undecidable

bounded



P



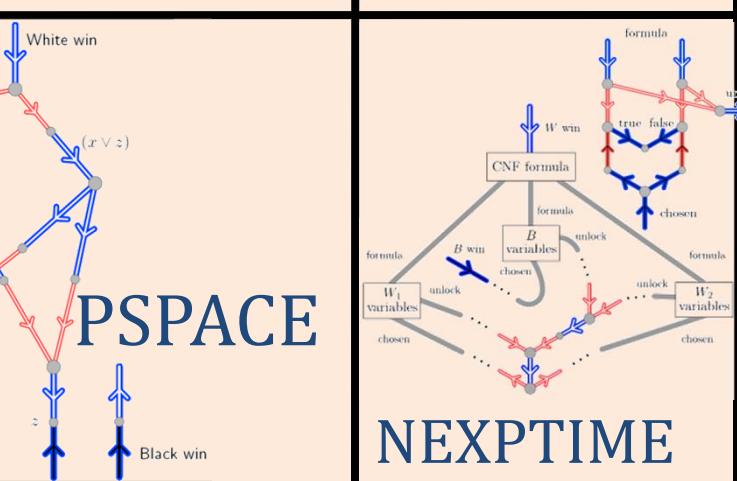
PSPACE

0 players  
(simulation)

1 player  
(puzzle)

2 players  
(game)

team,  
imperfect info



NEXPTIME

# Follow-On Algorithms Classes

- 6.046: Intermediate Algorithms
- 6.047: Computational Biology
- 6.854: Advanced Algorithms
- 6.849: Geometric Folding Algorithms
- 6.850: Geometric Computing
- 6.851: Advanced Data Structures
- 6.852: Distributed Algorithms
- 6.853: Algorithmic Game Theory
- 6.855: Network Optimization
- 6.856: Randomized Algorithms
- 6.857: Network and Computer Security

# Follow-On Theory Classes

- 6.045: Automata, Computability, Complexity
- 6.840: Theory of Computing
- 6.841: Advanced Complexity Theory
- 6.842: Randomness & Computation
- 6.845: Quantum Complexity Theory