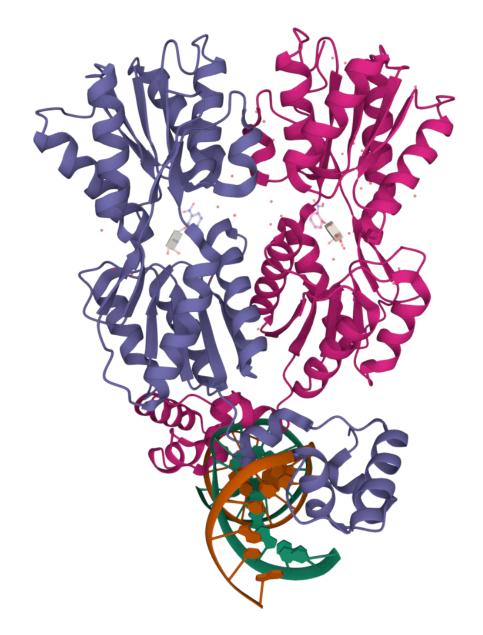
### Protein-DNA recognition



https://www.rcsb.org/3d-view/1EFA

Whole LacI https://www.rcsb.org/3d-view/1EFA

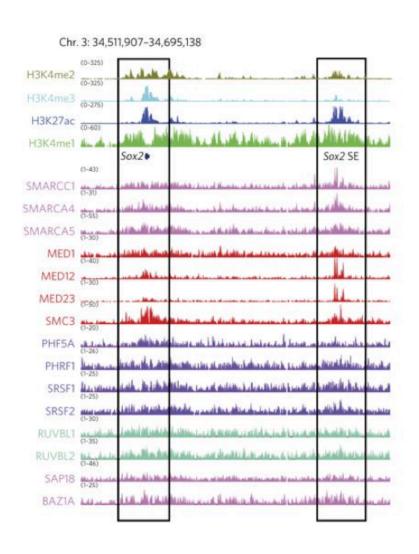
DNA-binding domain https://www.rcsb.org/3d-view/1L1M

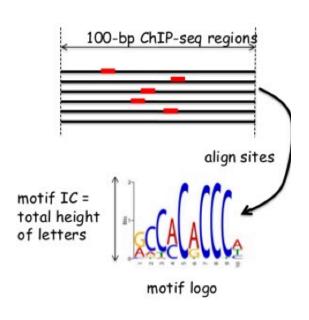
On non-spec DNA https://www.rcsb.org/3d-view/1OSL

Model of transition

https://proteopedia.org/wiki/index.php/Lac\_repressor

#### Methods to characterize protein-DNA specificity





Binding localization, occupancy Motif identification

## Nobel Prize in Physiology or Medicine 1965



Photo from the Nobel Foundation archive.

François Jacob

Prize share: 1/3



Photo from the Nobel Foundation archive.

André Lwoff

Prize share: 1/3



Photo from the Nobel Foundation archive.

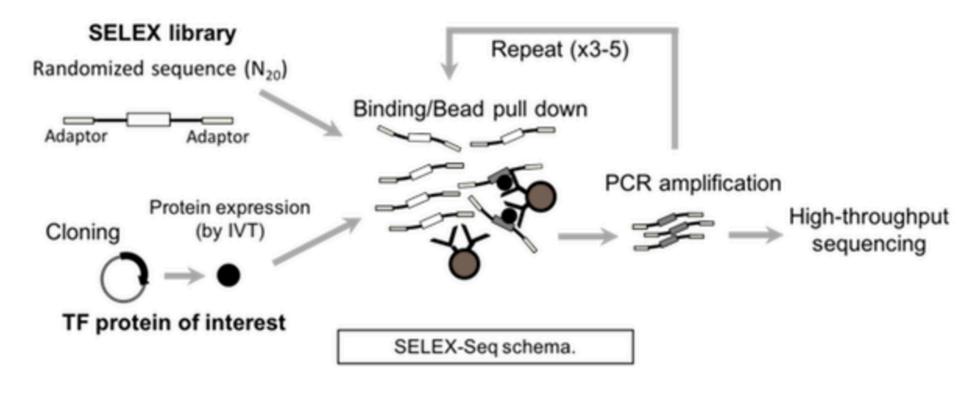
Jacques Monod

Prize share: 1/3

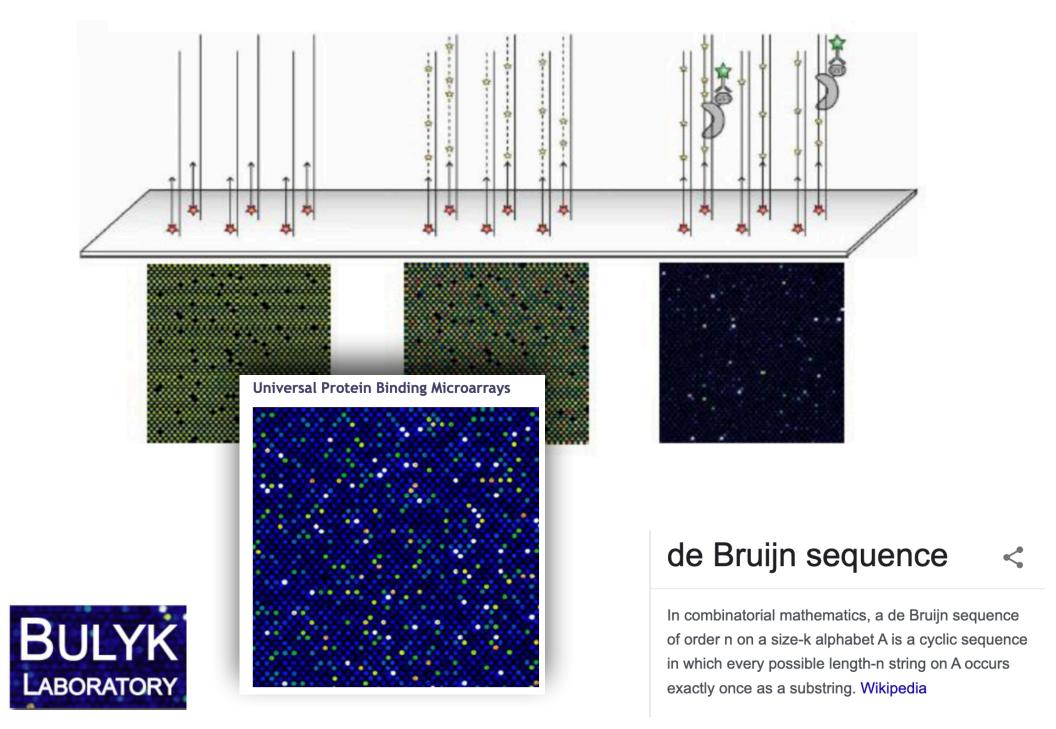
The Nobel Prize in Physiology or Medicine 1965 was awarded jointly to François Jacob, André Lwoff and Jacques Monod "for their discoveries concerning genetic control of enzyme and virus synthesis"

#### Methods to characterize protein-DNA specificity

#### **SELEX-Seq**



#### Methods to characterize protein-DNA specificity

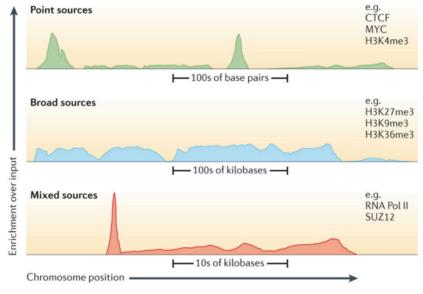


## ChIP-seq

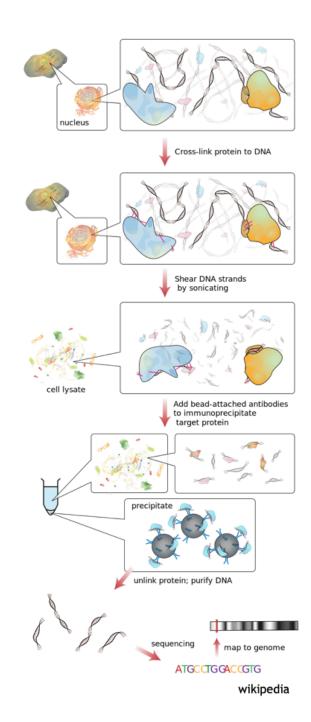
Use antibodies to pull down target of interest!

#### Related techniques:

- ChIP-exo
- CUT&RUN (no X-linking)



https://www.nature.com/articles/nrg3642



#### DNase I hypersensitive sites (DHS)

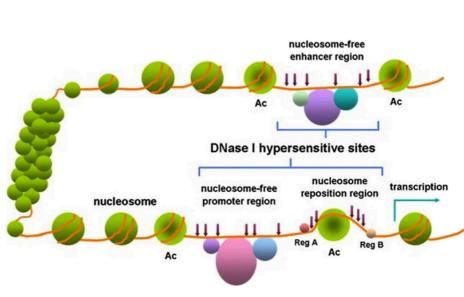
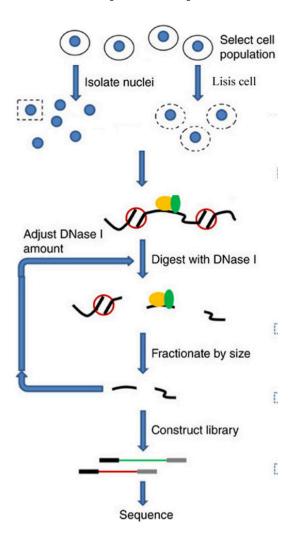


Fig. 1. DHSs within chromatin (Wang et al., 2012).



# Inference of e(i,a) using expectation maximization

## Precise physical models of protein—DNA interaction from high-throughput data

Justin B. Kinney, Gašper Tkačik, and Curtis G. Callan, Jr.\*

Physics Department and Lewis Sigler Institute for Integrative Genomics, Princeton University, Princeton, NJ 08544

Contributed by Curtis G. Callan, Jr., November 8, 2006 (sent for review September 30, 2006)

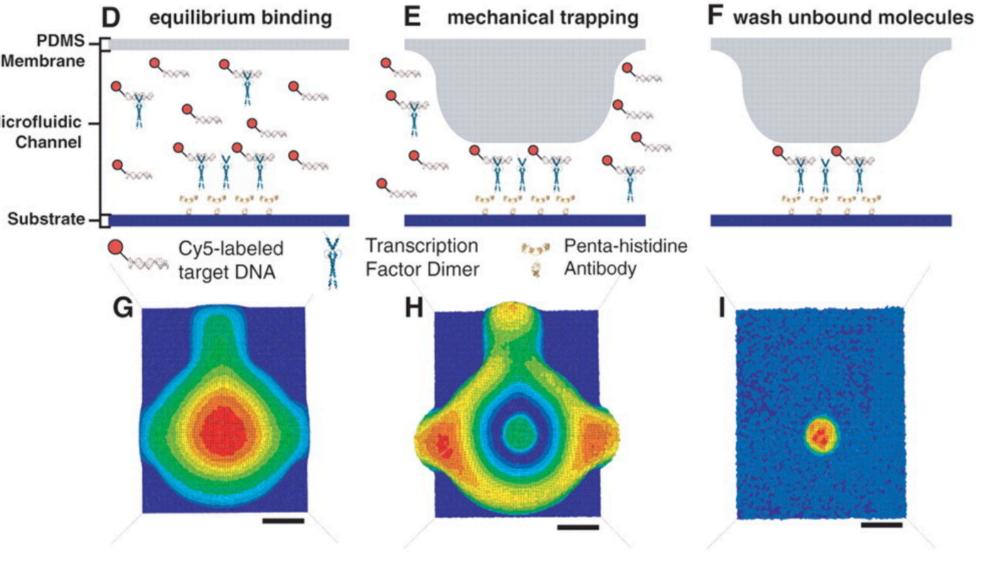
BIOINFORMATICS

Vol. 22 no. 14 2006, pages e141–e149 doi:10.1093/bioinformatics/btl223

### Statistical mechanical modeling of genome-wide transcription factor occupancy data by MatrixREDUCE

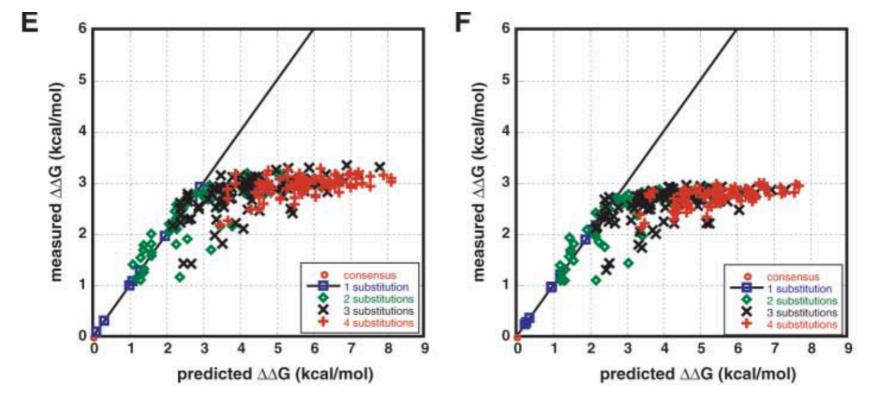
Barrett C. Foat<sup>1</sup>, Alexandre V. Morozov<sup>2</sup> and Harmen J. Bussemaker<sup>1,3,\*</sup>

<sup>1</sup>Department of Biological Sciences, Columbia University, New York, NY 10027, USA, <sup>2</sup>Center for Studies in Physics and Biology, The Rockefeller University, New York, NY 10021, USA and <sup>3</sup>Center for Computational Biology and Bioinformatics, Columbia University, New York, NY 10032, USA



#### A Systems Approach to Measuring the Binding Energy Landscapes of Transcription Factors

Sebastian J. Maerkl<sup>1,2</sup> and Stephen R. Quake<sup>2</sup>\*



**Fig. 2.** Binding affinities of C-terminally tagged TFs MAX iso A (**A**), MAX iso B (**B**), Pho4p (**C**), and Cbf1p (**D**) to all sequence permutations of  $N_{-4}$  to  $N_{-1}$ . Sequences  $N_{-3}$  to  $N_{-1}$  are plotted on the category axis, with the fourth base,  $N_{-4}$ , displayed as clusters of four columns per category. (**E** and **F**) Comparisons of predicted changes in the Gibbs free energy ( $\Delta\Delta G$ ) against measured values for MAX isoforms A and B are shown, respectively. All predicted values were calculated from PWMs assuming base independence.

#### A Systems Approach to Measuring the Binding Energy Landscapes of Transcription Factors

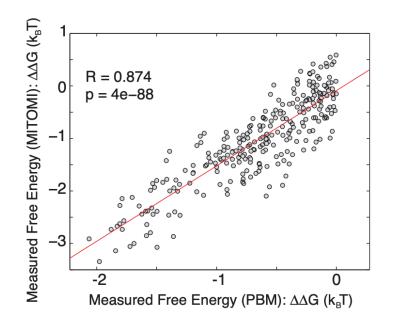
Sebastian J. Maerkl<sup>1,2</sup> and Stephen R. Quake<sup>2</sup>\*

### Non-specific binding

## Protein—DNA binding in the absence of specific base-pair recognition

Ariel Afek<sup>a</sup>, Joshua L. Schipper<sup>b</sup>, John Horton<sup>b</sup>, Raluca Gordân<sup>b,1</sup>, and David B. Lukatsky<sup>a,1</sup>

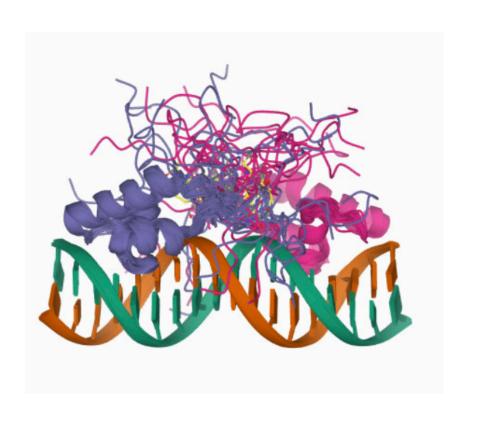
analysis. First, the presence of the specific motif leads statistically, on average, to at most  $\sim -2~k_BT$  free-energy difference compared with the negative control. Second, the magnitude of the identified



~1 kT per bp

<sup>&</sup>lt;sup>a</sup>Department of Chemistry, Ben-Gurion University of the Negev, Be'er Sheva 8410501 Israel; and <sup>b</sup>Center for Genomic and Computational Biology, Department of Biostatistics and Bioinformatics, Duke University, Durham, NC 27708

## Non-specific binding





https://www.rcsb.org/3d-view/1OSL

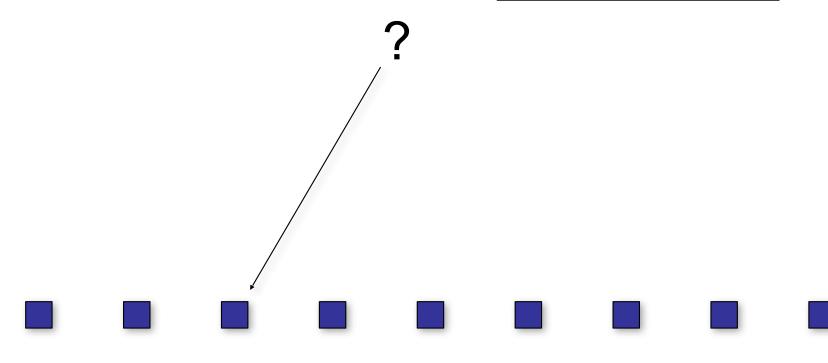
https://www.rcsb.org/3d-view/1L1M

# Information theory

#### How much information do you need?

To find an object among N decoys

Need at least: log<sub>2</sub>N bits



### How much information do you need?

To find an object among N decoys

Need at least:  $I_{min} = log_2 N$  bits

Microbes N= $10^{6-7}$   $I_{min}$ = 20-23 bits Multicell euk N= $10^{8-10}$   $I_{min}$ =27-33 bits

#### How much information is contained in a motifs

- Conserved position = 2 bits
- T or C are equally likely = 1bit
- more generally:

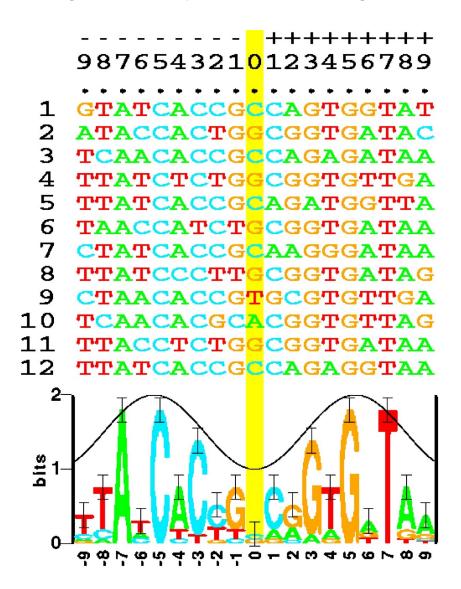
$$I_k = \sum_{x=A,T,G,C} p_k(x) \log_2 \left[ p_k(x) / g(x) \right]$$

where p(x) is the frequency of x at positions k

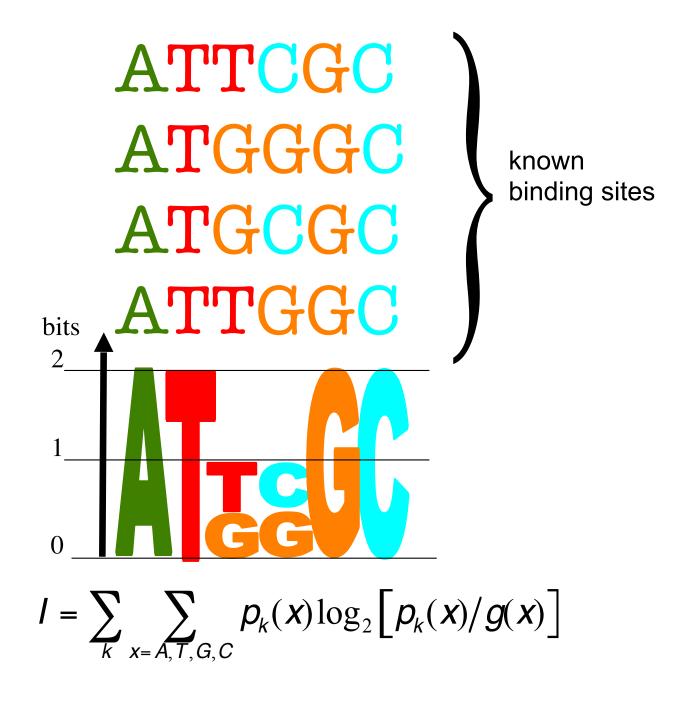
#### **The Information Content**

of the whole motif

$$I = \sum_{k=1}^{L} I_k$$



#### How much information is contained in a motifs



#### Results

#### We analyzed 969 TF motifs

- calculate their Information Contents I

- compare to  $I_{min}$ 

Bacteria

 $N=10^{6-7} I_{min} = 20-23 bits$ 

Yeast

 $N=10^7$   $I_{min}=24 \ bits$ 

Multicell euk N= $10^{8-10}$   $I_{min}$ = 27-33 bits

Bacteria: 20-25 bits

14 bits Yeast:

Multicell eukaryotes: 12 bits

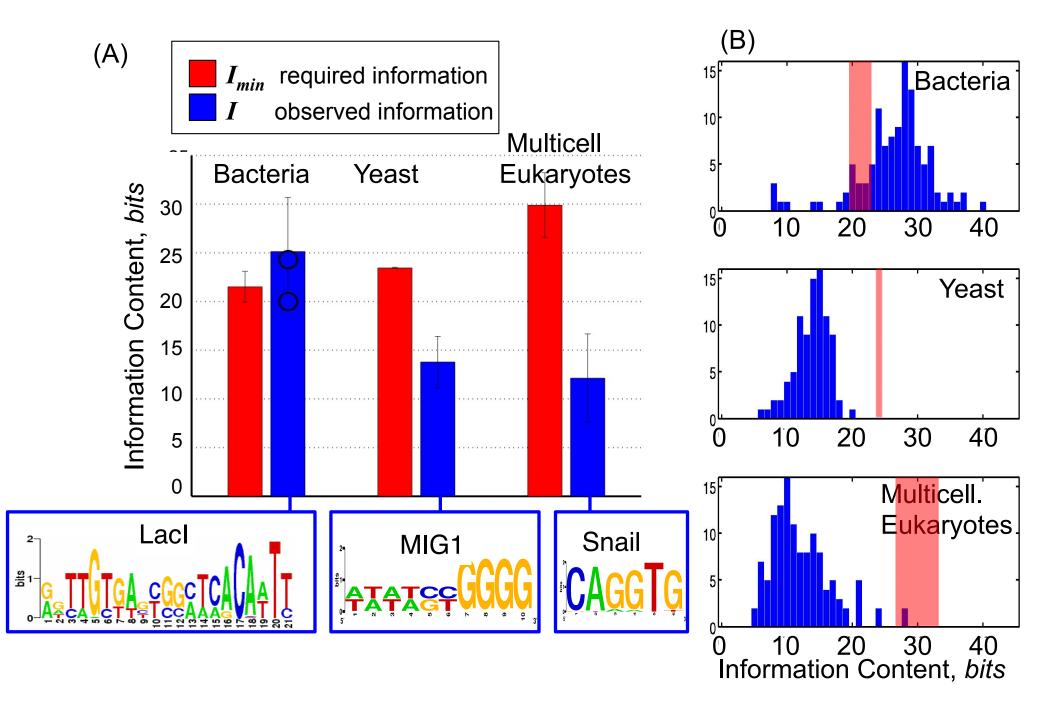
<u>Trends Genet.</u> 2009 Oct;25(10):434-40. doi: 10.1016/j.tig.2009.08.003. Epub 2009 Oct 6. Paperpile



Different gene regulation strategies revealed by analysis of binding motifs.

Wunderlich Z<sup>1</sup>, Mirny LA.

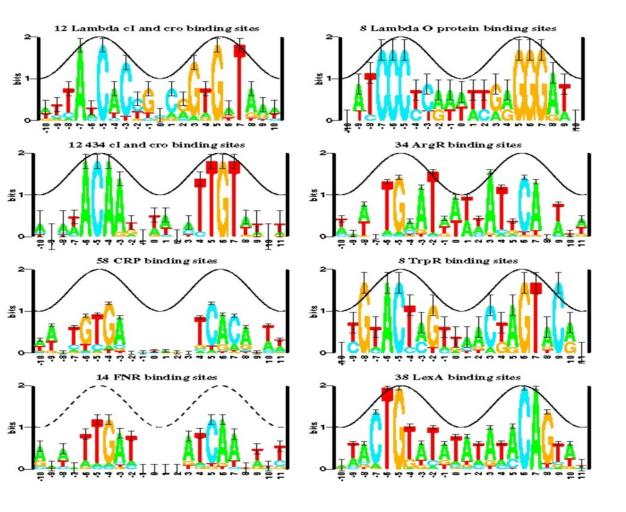
Figure1



## Examples

Bacteria: 26 bits

Eukaryotes: 12-14 bits



Cbf1 (S. cerevisiae bHLH)



Zif268 (M. musculus C<sub>2</sub>H<sub>2</sub> zinc finger)



Literature consensus: GCGKGGGCG

Ceh-22 (C. elegans NK homeodomain)



Literature consensus: CACTNNA

Oct-1 (H. sapiens POU homeodomain)



Literature consensus: ATGCAAAT

### Information deficiency => binding decoys

Prob of a **hits** (mistake) ≥

Number of **hits** per genome ≥

$$2^{-1}$$
 $N2^{-1} = 2^{I_{\min} - I}$ 

 $I_{min}$ -I

~0 bits Bacteria:

Multicell eukaryotes: 18 bits Number of hits

~1 per genome

~104-106 per genome

assuming 90% chromatinized DNA

<u>Trends Genet.</u> 2009 Oct;25(10):434-40. doi: 10.1016/j.tig.2009.08.003. Epub 2009 Oct 6. Paperpile



Different gene regulation strategies revealed by analysis of binding motifs.

Wunderlich Z<sup>1</sup>, Mirny LA.

## Cell Biology

- Widespread non-functional binding
- Binding ≠ gene expression
- Carefully interpret experimental data

Expected spurious binding ~20% of promoter fro any TF.

ChIP-on-chip significance analysis reveals large-scale binding and regulation by human transcription factor oncogenes

Adam A. Margolin<sup>a,b,c,1</sup>, Teresa Palomero<sup>d,e</sup>, Pavel Sumazin<sup>b</sup>, Andrea Califano<sup>a,b,d,2,3</sup>, Adolfo A. Ferrando<sup>d,e,f,2,3</sup>, and Gustavo Stolovitzky<sup>b,c,2,3</sup>

MYC binds 48% of all promoters, (!!!) NOTCH1 19%, HES1 18%.

## Cell Biology

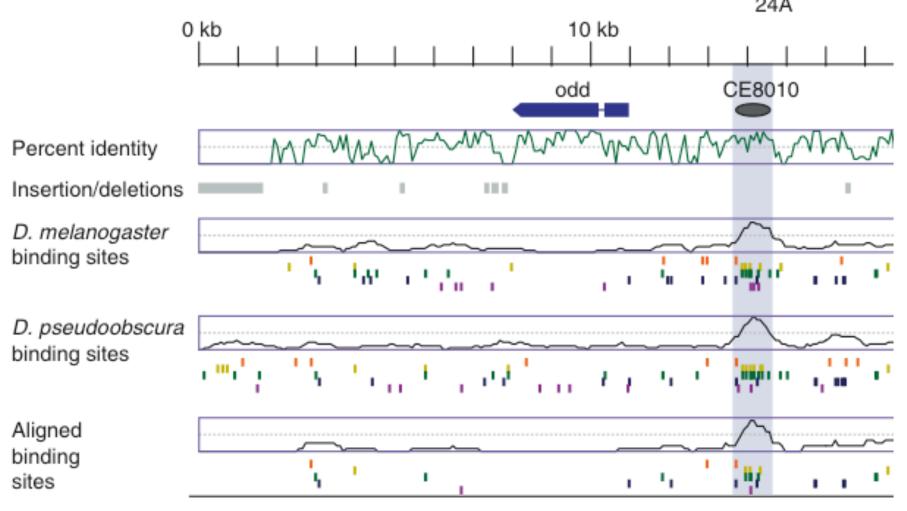
- Widespread non-functional binding
- Binding ≠ gene expression
- Carefully interpret experimental data

#### Need cluster of site to specify a genomic region

In a region of 1 kb composed of the sites of 3–10 different TFs, we calculate  $n_{cluster} = 10$ –20 sites. This lower limit on the number of required binding sites is remarkably consistent with the mean of  $\sim$ 20 sites per 1 kb observed in fly developmental enhancers [28].

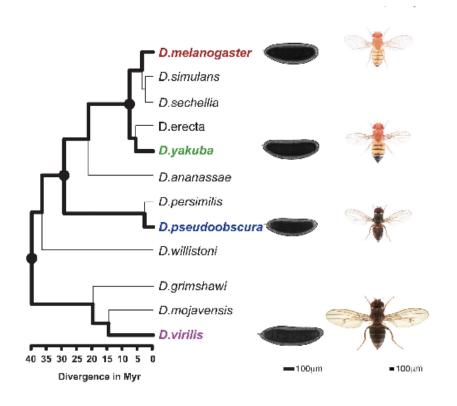


## Significant evolutionary flexibility



Genome Biology 2004, 5:R61

Computational identification of developmental enhancers: conservation and function of transcription factor binding-site clusters in Drosophila melanogaster and Drosophila pseudoobscura Benjamin P Berman\*\*, Barret D Pfeiffer\*\*, Todd R Laverty\*, Steven L Salzberg§, Gerald M Rubin\*\*\*, Michael B Eisen\*\*¶¥ and Susan E Celniker\*\*



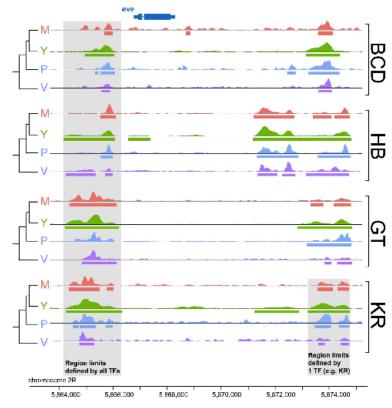
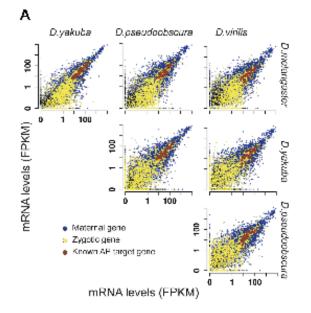


Figure 2. Comparison of binding profiles of BCD, GT, HB and KR at the even-skipped locus in the four species *D.melanogaster*. *D.yakuba*, *D.pseudoobscura* and *D.vitilis*. An illustration of the two types of comparisons made in this study are highlighted in grey; trans-species comparison for each single TF (right) or trans-TFs comparisons (left). For simplicity, the species names were shortened using their initial: *D.melanogaster* (M), *D.yakuba* (Y), *D* 



OPEN ACCESS Freely available online



#### Extensive Divergence of Transcription Factor Binding in Drosophila Embryos with Highly Conserved Gene Expression

Mathilde Paris<sup>1</sup>\*, Tommy Kaplan<sup>1,2</sup>, Xiao Yong Li<sup>1,3</sup>, Jacqueline E. Villalta<sup>2</sup>, Susan E. Lott<sup>1,4</sup>, Michael B. Eisen<sup>1,2,3</sup>\*

September 2013 | Volume 9 | Issue 9 | e1003748

## Significant evolutionary flexibility

#### Abstract

To better characterize how variation in regulatory sequences drives divergence in gene expression, we undertook a systematic study of transcription factor binding and gene expression in blastoderm embryos of four species, which sample much of the diversity in the 40 million-year old genus Drosophila: D. melanogaster, D. yakuba, D. pseudoobscura and D. virilis. We compared gene expression, measured by mRNA-seq, to the genome-wide binding, measured by ChIP-seq, of four transcription factors involved in early anterior-posterior patterning. We found that mRNA levels are much better conserved than individual transcription factor binding events, and that changes in a gene's expression were poorly explained by changes in adjacent transcription factor binding. However, highly bound sites, sites in regions bound by multiple factors and sites near genes are conserved more frequently than other binding, suggesting that a considerable amount of transcription factor binding is weakly or non-functional and not subject to purifying selection.



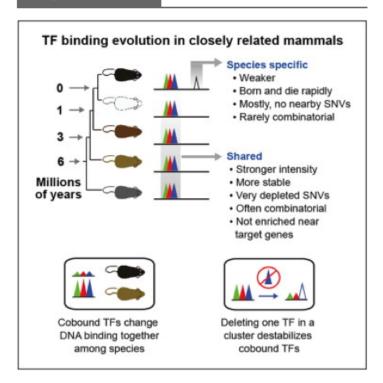


#### Extensive Divergence of Transcription Factor Binding in Drosophila Embryos with Highly Conserved Gene Expression

Mathilde Paris<sup>1</sup>\*, Tommy Kaplan<sup>1,2</sup>, Xiao Yong Li<sup>1,3</sup>, Jacqueline E. Villalta<sup>2</sup>, Susan E. Lott<sup>1,4</sup>, Michael B. Eisen<sup>1,2,3</sup>\*

## Significant evolutionary flexibility

#### Graphical Abstract

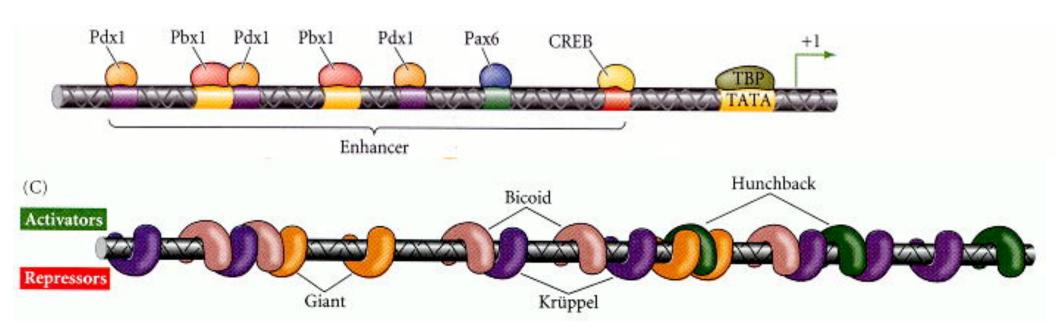


# Cooperativity and Rapid Evolution of Cobound Transcription Factors in Closely Related Mammals

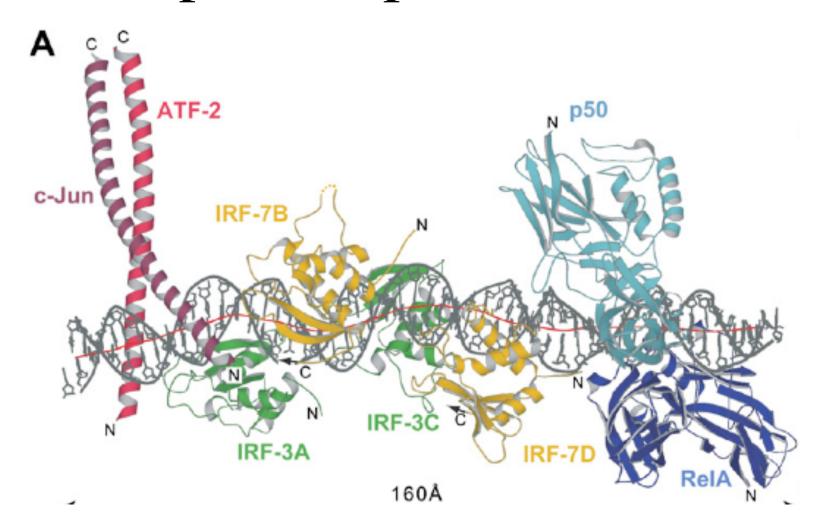
Klara Stefflova, 1,8 David Thybert, 2,8 Michael D. Wilson, 3 Ian Streeter, 2 Jelena Aleksic, 4,5 Panagiota Karagianni, 6 Alvis Brazma, 2 David J. Adams, 7 Iannis Talianidis, 6 John C. Marioni, 2 Paul Flicek, 2,7,\* and Duncan T. Odom<sup>1,7,\*</sup>

# Transcription factors and promoters

#### **Eukaryotes**



## Few protein-protein interactions

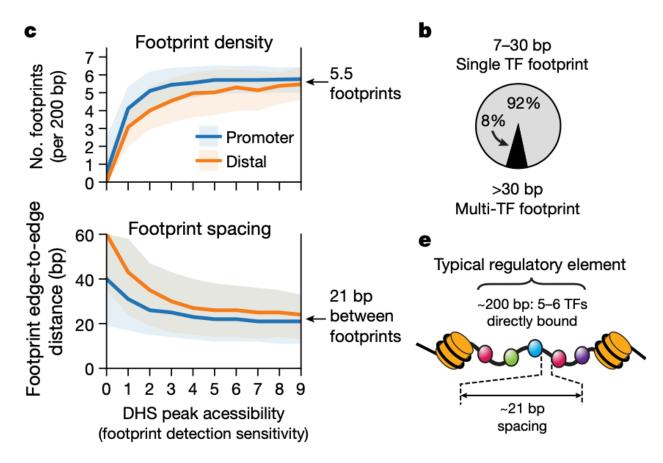




## An Atomic Model of the Interferon-β Enhanceosome

Daniel Panne, 1 Tom Maniatis, 2 and Stephen C. Harrison 1,\*

## Clusters of regulatory motifs seen by DHS footprinting



## Global reference mapping of human transcription factor footprints

https://doi.org/10.1038/s41586-020-2528-x

Received: 30 January 2020

Accepted: 25 June 2020

#### **Article**

## Low overlap of transcription factor DNA binding and regulatory targets

https://doi.org/10.1038/s41586-025-08916-0

Lakshmi Mahendrawada¹, Linda Warfield¹, Rafal Donczew¹,2 & Steven Hahn¹⊠

Received: 12 July 2023

Accepted: 19 March 2025

Published online: 16 April 2025

Check for updates

DNA sequence-specific transcription factors (TFs) modulate transcription and chromatin architecture, acting from regulatory sites in enhancers and promoters of eukaryotic genes<sup>1,2</sup>. How multiple TFs cooperate to regulate individual genes is still unclear. In yeast, most TFs are thought to regulate transcription via binding to upstream activating sequences, which are situated within a few hundred base pairs upstream of the regulated gene<sup>3</sup>. Although this model has been validated for individual TFs and specific genes, it has not been tested in a systematic way. Here we integrated information on the binding and expression targets for the near-complete set of yeast TFs and show that, contrary to expectations, there are few TFs with dedicated activator or repressor roles, and that most TFs have a dual function. Although nearly all protein-coding genes are regulated by one or more TFs, our analysis revealed limited overlap between TF binding and gene regulation. Rapid depletion of many TFs also revealed many regulatory targets that were distant from detectable TF binding sites, suggesting unexpected regulatory mechanisms. Our study provides a comprehensive survey of TF functions and offers insights into interactions between the set of TFs expressed in a single cell type and how they contribute to the complex programme of gene regulation.

