

# Operating System of the Cell

## Regulation of Gene Expression

Gold Lab Symposium

May 14-15, 2026

Gary D. Stormo

Erlanger Professor, emeritus

Department of Genetics

Edison Center for Genome Sciences

and Systems Biology



Washington University in St. Louis

# **Advice from Max Delbruck for speaking to a diverse audience with unknown level of background knowledge:**

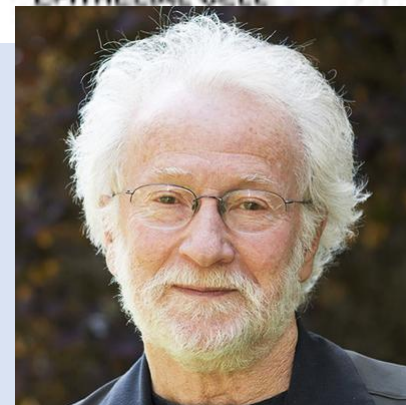
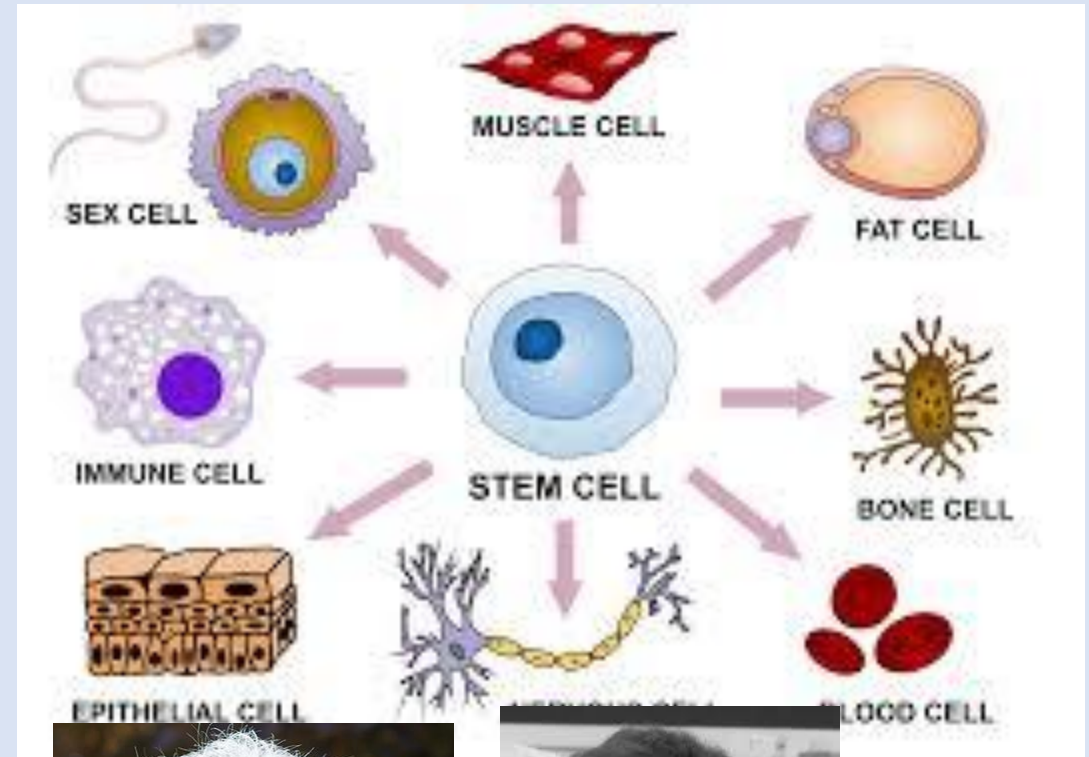
“Assume your audience has zero knowledge but infinite intelligence.”

# Operating System of the Cell

## Regulation of Gene Expression

Consider the various cell types within the human body. Although they all share identical DNA, their appearances and functions differ significantly. Each cell type activates only a particular subset of the entire gene repertoire.

My Primary Interest:  
Regulation of Gene Expression



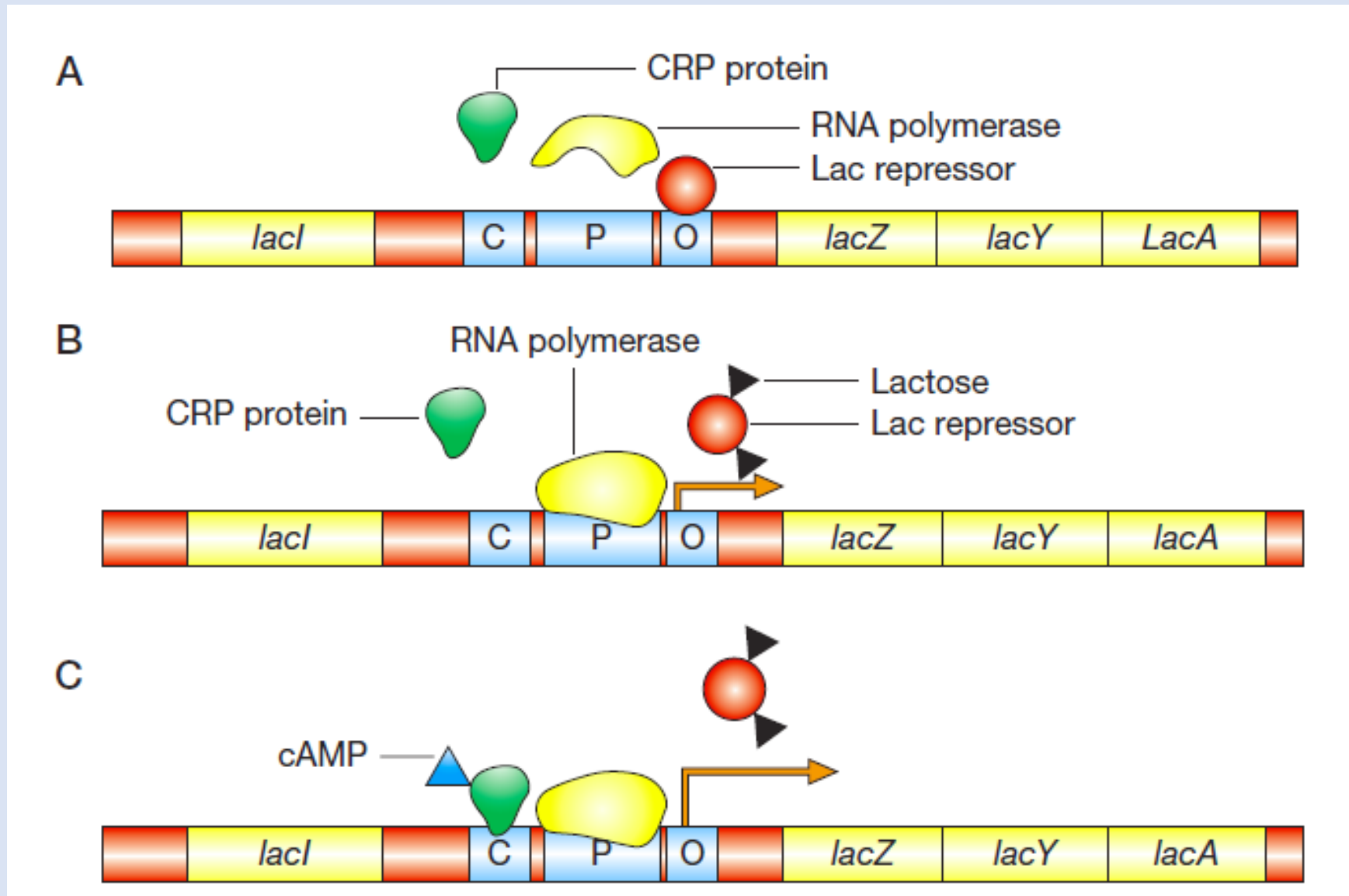
Larry Gold

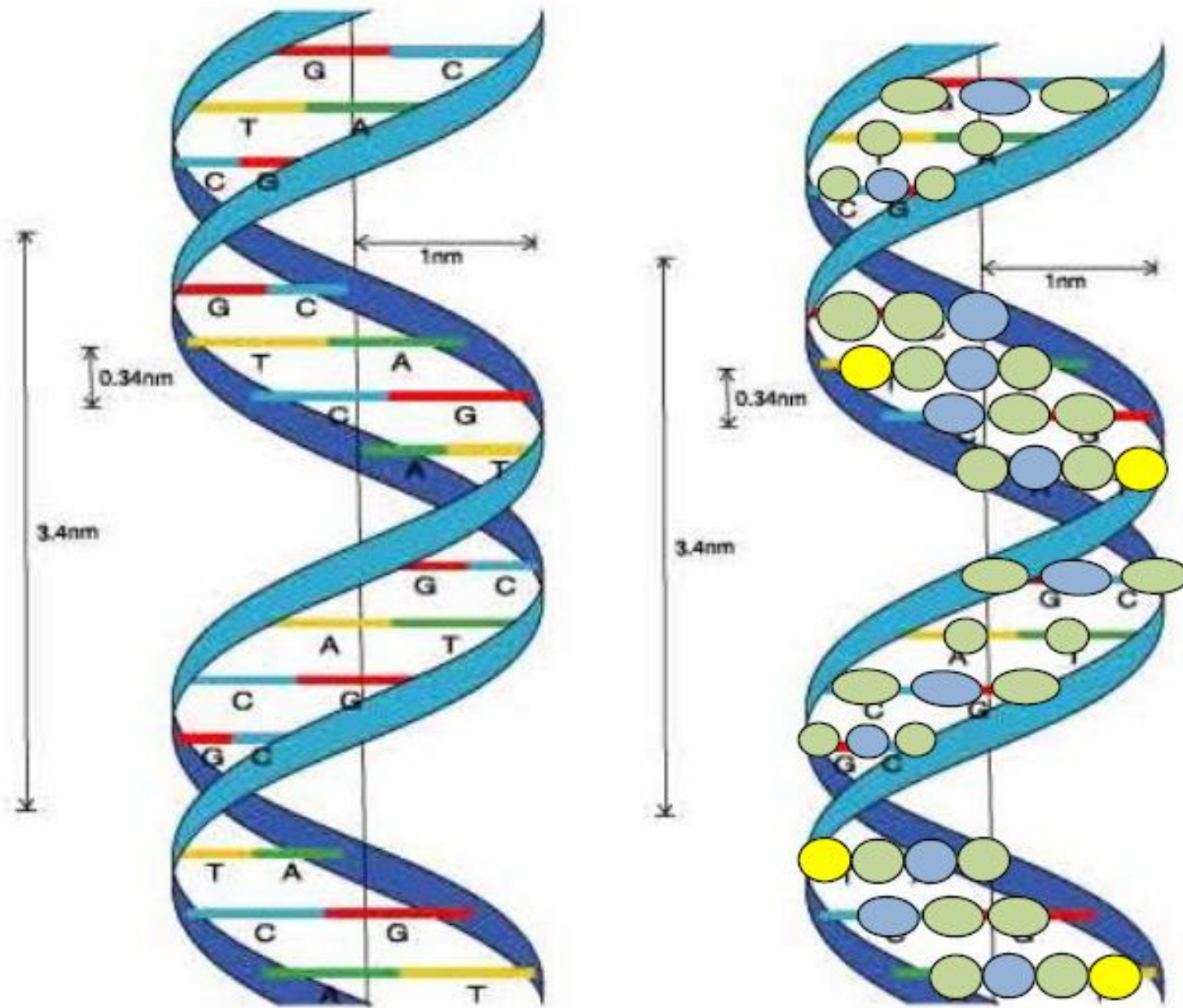


Tom Schneider

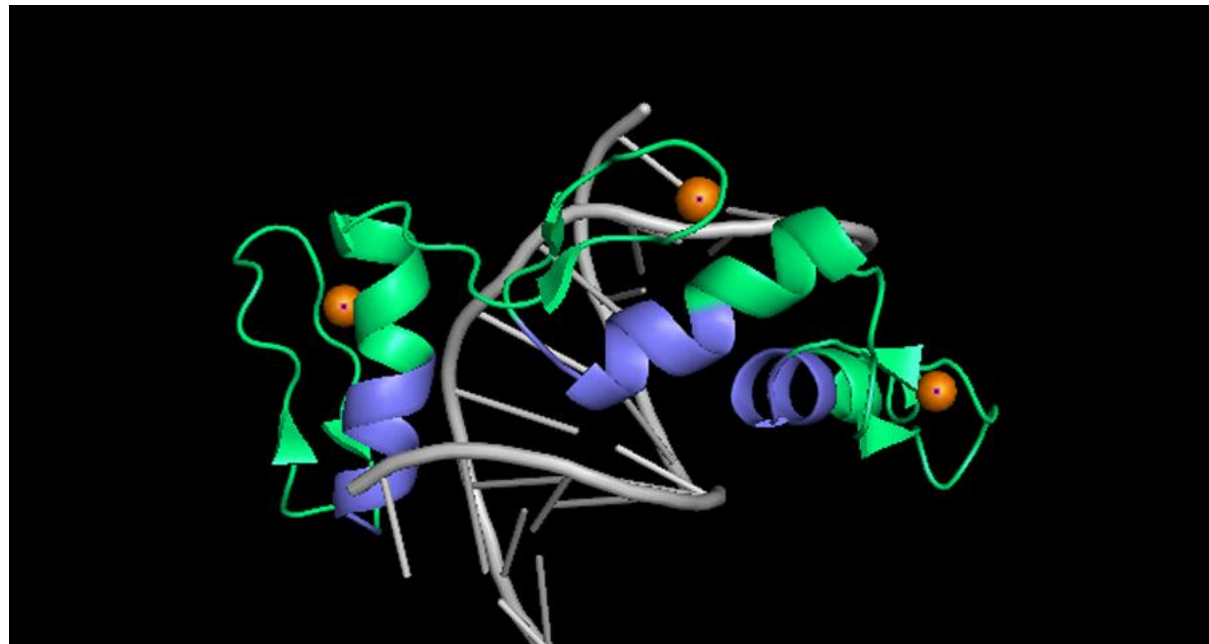
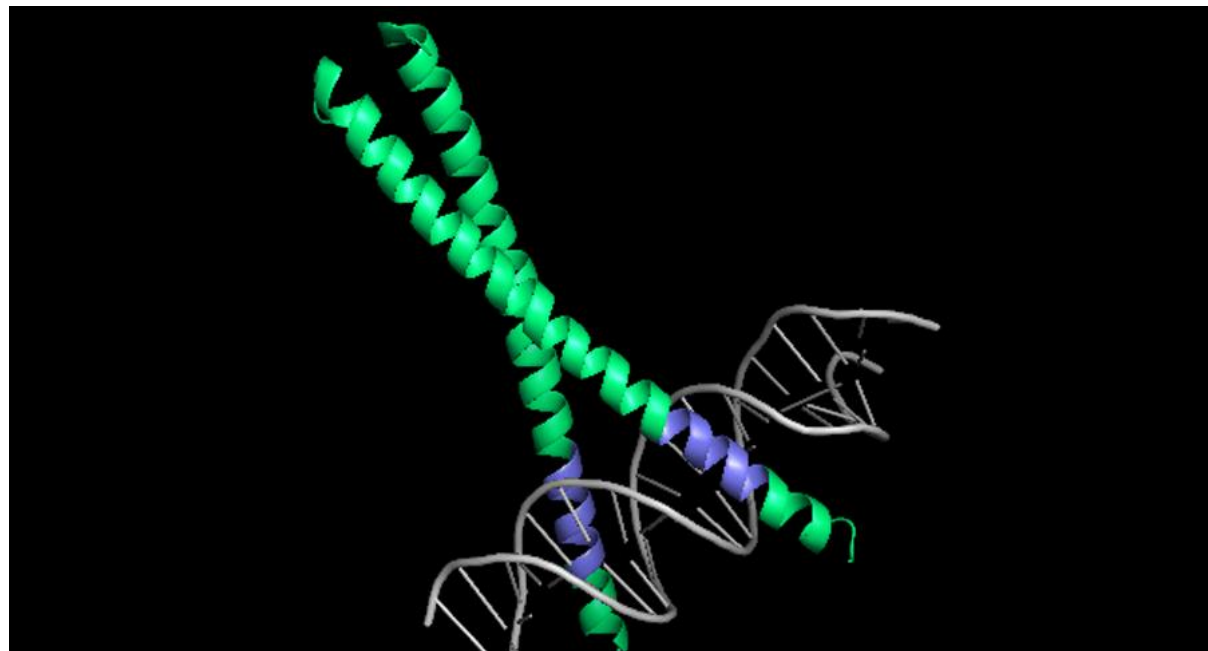
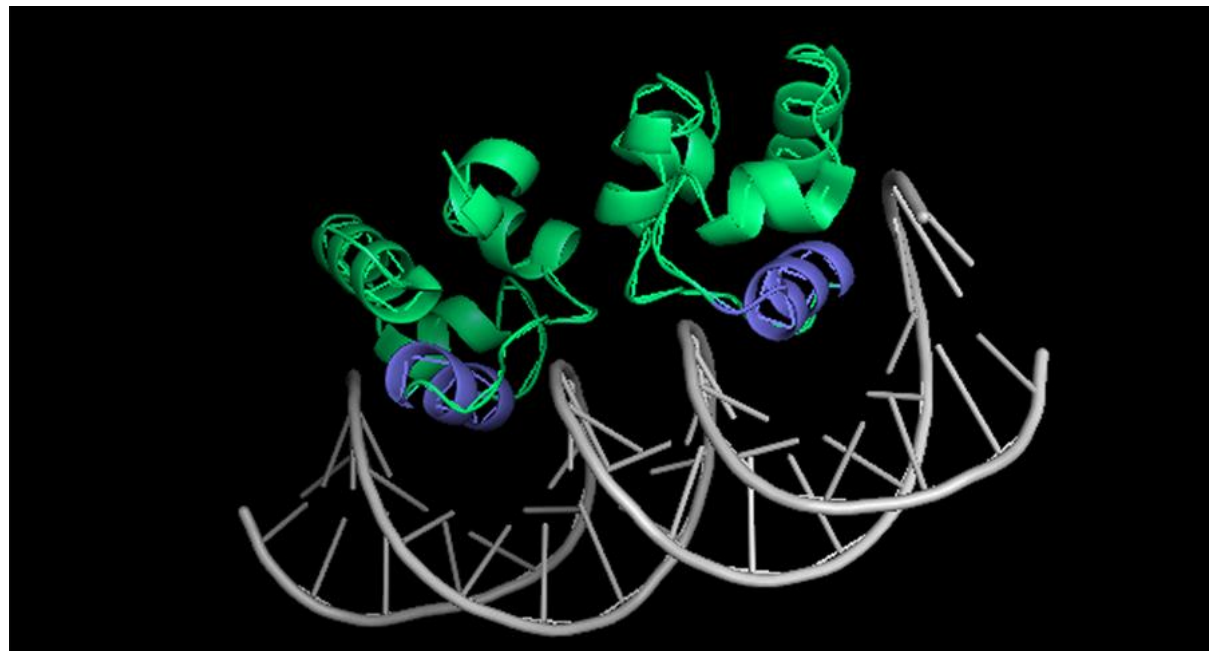
my primary interest:

# Lactose regulatory system, Jacob and Monod, 1961





The surface of DNA contains chemical features that can be “read” by the chemical surface of proteins



# Two Revolutions were occurring during my graduate school days

## 1. DNA sequencing revolution

DNA sequencing invented (1977)

Allowed new type of data for molecular biology research

## 2. Computer revolution

Symbolized by development of personal computers (IBM PC 1975,  
Microsoft 1975, Apple 1976)

## CRP binding site:

*“Under conditions of low glucose, turn on the expression of the adjacent gene.”*

**TAATGTGAGTTAGCTCACTCAT**

*“Under conditions of low glucose, turn on the expression of the adjacent gene.”*

**TAATGTGAGTTAGCTCACTCAT**  
**cgctGTGAccgtGgTCgCagtt**

Eventually no positions are completely conserved

**TAATGTGAGTTAGCTCACTCAT**  
**cgcTGTGAccgtGgTCgCagtT**  
**tttTtTGAtcgtttTCaCattT**  
**aacgTGAtagccgTCaaaca**

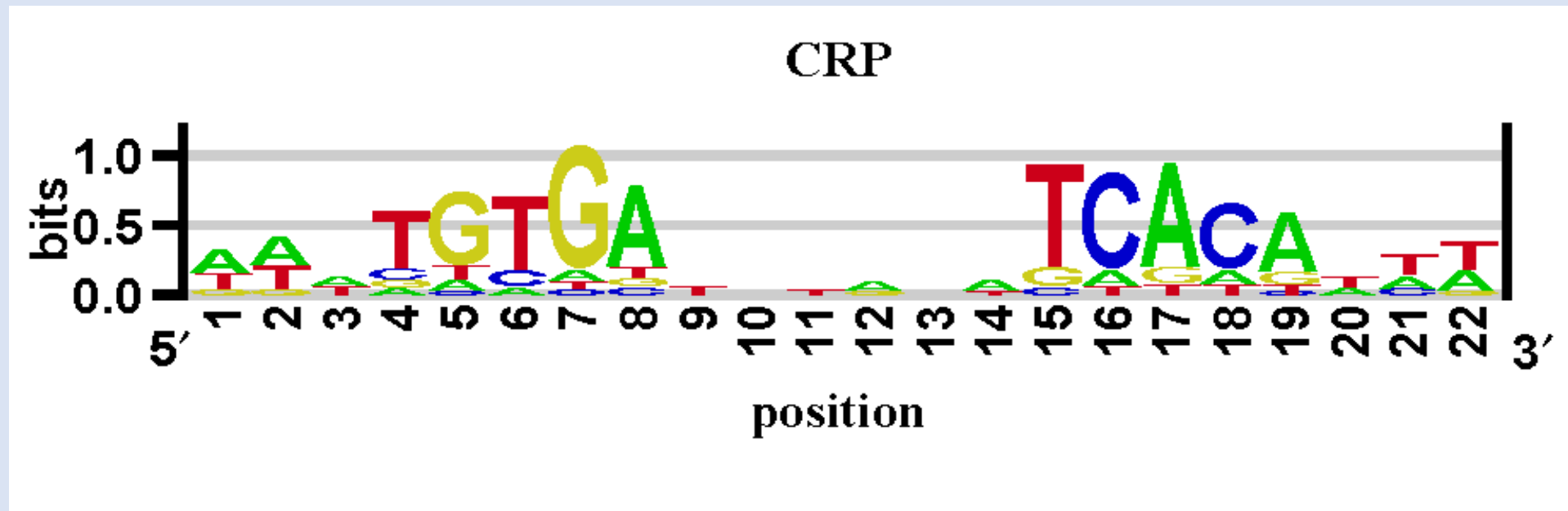
:

TGTGAnnnnnnTCACA

“consensus sequence” can be used  
to search allowing for mismatches

Eventually no positions are completely conserved

TAATGTGAGTTAGCTCACTCAT  
cgcTGTGAccgtGgTCgCagtT  
tttTtTGAtcgtttTCaCattT  
aaacgTGAtagccgTCaaacaa





# Next Challenge

## Outline of motif discovery problem

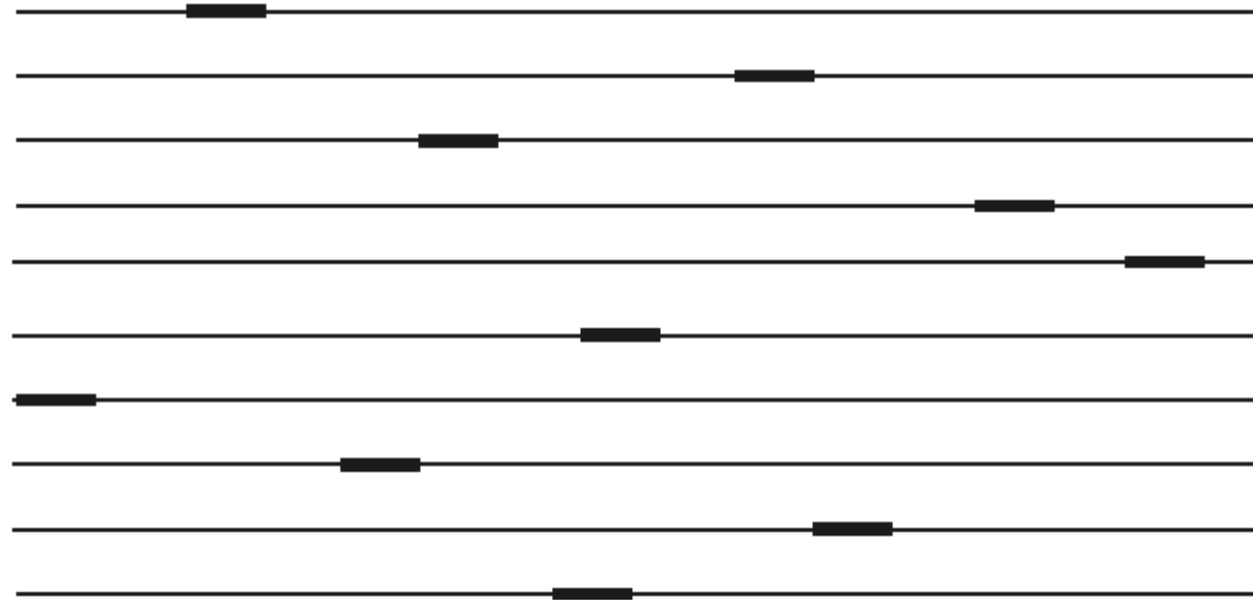


Fig. 6.1. A general schematic of the motif finding problem. Each *long thin line* represents a single DNA sequence. The *dark segments* within each line represent the binding sites whose positions are unknown in advance and we are trying to discover.

## Motif Discovery, only need “regulatory regions” of co-regulated genes

Input: Genes regulated by CRP in *E. coli*

CE1CG

\TAATGTTTGTGCTGGTTTTTGTGGCATCGGGCGAGAATAGCGCGTGGTGTGAAAGACTGTTTTTTTGTATCGTTTTTCACAAAAATGGAAGTCCACAGTCTTGACAG\  
ECOARABOP

\GACAAAAACGCGTAACAAAAGTGTCTATAATCACGGCAGAAAAGTCCACATTGATTATTTGCACGGCGTCACACTTTGCTATGCCATAGCATTATTTTATCCATAAG\  
ECOBGLR1

\ACAAATCCCAATAAATACTTAATTATTGGGATTTGTTATATATAAATTTATAAAATTCCTAAAATTACACAAAAGTTAATAACTGTGAGCATGGTCATATTTTTTATCAAT\  
ECOCR

\CACAAAGCGAAAGCTATGCTAAAACAGTCAGGATGCTACAGTAATACATTGATGTACTGCATGTATGCAAAGGACGTCACATTACCGTGCAGTACAGTTGATAGC\  
ECOCYA

\ACGGTGCTACACTTGTATGTAGCGCATCTTTCTTTACGGTCAATCAGCAAGGTGTTAAATTGATCACGTTTTAGACCATTTTTTCGTGCGTAAAATAAAAAACC\  
ECODEOP2

\AGTGAATTATTTGAACCAGATCGCATTACAGTGATGCAAACCTGTAAGTAGATTTCTTAATTGTGATGTGTATCGAAGTGTGTTGCGGAGTAGATGTTAGAATA\  
ECOGALE

\GCGCATAAAAAACGGCTAAATTCTTGTGTAAACGATTCCACTAATTTATTCCATGTACACTTTTTCGCATCTTTGTTATGCTATGGTTATTTTCATACCATAAGCC\  
ECOILVBR

\GCTCCGGCGGGGTTTTTTTGTATATCTGCAATTCAGTACAAAAAGGTGATCAACCCCTCAATTTTTCCCTTTGCTGAAAAATTTTCCATTGTCTCCCTGTAAAAGCTGT\  
ECOCYB

\GCTCCGGCGGGGTTTTTTTGTATATCTGCAATTCAGTACAAAAAGGTGATCAACCCCTCAATTTTTCCCTTTGCTGAAAAATTTTCCATTGTCTCCCTGTAAAAGCTGT\  
ECOCYB

\GCTCCGGCGGGGTTTTTTTGTATATCTGCAATTCAGTACAAAAAGGTGATCAACCCCTCAATTTTTCCCTTTGCTGAAAAATTTTCCATTGTCTCCCTGTAAAAGCTGT\  
ECOCYB

\GCTCCGGCGGGGTTTTTTTGTATATCTGCAATTCAGTACAAAAAGGTGATCAACCCCTCAATTTTTCCCTTTGCTGAAAAATTTTCCATTGTCTCCCTGTAAAAGCTGT\  
ECOCYB

\GCTCCGGCGGGGTTTTTTTGTATATCTGCAATTCAGTACAAAAAGGTGATCAACCCCTCAATTTTTCCCTTTGCTGAAAAATTTTCCATTGTCTCCCTGTAAAAGCTGT\  
ECOCYB

\GCTCCGGCGGGGTTTTTTTGTATATCTGCAATTCAGTACAAAAAGGTGATCAACCCCTCAATTTTTCCCTTTGCTGAAAAATTTTCCATTGTCTCCCTGTAAAAGCTGT\  
ECOCYB

\GCTCCGGCGGGGTTTTTTTGTATATCTGCAATTCAGTACAAAAAGGTGATCAACCCCTCAATTTTTCCCTTTGCTGAAAAATTTTCCATTGTCTCCCTGTAAAAGCTGT\  
ECOCYB

\GCTCCGGCGGGGTTTTTTTGTATATCTGCAATTCAGTACAAAAAGGTGATCAACCCCTCAATTTTTCCCTTTGCTGAAAAATTTTCCATTGTCTCCCTGTAAAAGCTGT\  
ECOCYB

\GCTCCGGCGGGGTTTTTTTGTATATCTGCAATTCAGTACAAAAAGGTGATCAACCCCTCAATTTTTCCCTTTGCTGAAAAATTTTCCATTGTCTCCCTGTAAAAGCTGT\  
ECOCYB

\GCTCCGGCGGGGTTTTTTTGTATATCTGCAATTCAGTACAAAAAGGTGATCAACCCCTCAATTTTTCCCTTTGCTGAAAAATTTTCCATTGTCTCCCTGTAAAAGCTGT\  
ECOCYB

\GCTCCGGCGGGGTTTTTTTGTATATCTGCAATTCAGTACAAAAAGGTGATCAACCCCTCAATTTTTCCCTTTGCTGAAAAATTTTCCATTGTCTCCCTGTAAAAGCTGT\  
ECOCYB

\GCTCCGGCGGGGTTTTTTTGTATATCTGCAATTCAGTACAAAAAGGTGATCAACCCCTCAATTTTTCCCTTTGCTGAAAAATTTTCCATTGTCTCCCTGTAAAAGCTGT\  
ECOCYB

\GCTCCGGCGGGGTTTTTTTGTATATCTGCAATTCAGTACAAAAAGGTGATCAACCCCTCAATTTTTCCCTTTGCTGAAAAATTTTCCATTGTCTCCCTGTAAAAGCTGT\  
ECOCYB

\GCTCCGGCGGGGTTTTTTTGTATATCTGCAATTCAGTACAAAAAGGTGATCAACCCCTCAATTTTTCCCTTTGCTGAAAAATTTTCCATTGTCTCCCTGTAAAAGCTGT\  
ECOCYB

\GCTCCGGCGGGGTTTTTTTGTATATCTGCAATTCAGTACAAAAAGGTGATCAACCCCTCAATTTTTCCCTTTGCTGAAAAATTTTCCATTGTCTCCCTGTAAAAGCTGT\  
ECOCYB

\GCTCCGGCGGGGTTTTTTTGTATATCTGCAATTCAGTACAAAAAGGTGATCAACCCCTCAATTTTTCCCTTTGCTGAAAAATTTTCCATTGTCTCCCTGTAAAAGCTGT\  
ECOCYB

\GCTCCGGCGGGGTTTTTTTGTATATCTGCAATTCAGTACAAAAAGGTGATCAACCCCTCAATTTTTCCCTTTGCTGAAAAATTTTCCATTGTCTCCCTGTAAAAGCTGT\  
ECOCYB

\GCTCCGGCGGGGTTTTTTTGTATATCTGCAATTCAGTACAAAAAGGTGATCAACCCCTCAATTTTTCCCTTTGCTGAAAAATTTTCCATTGTCTCCCTGTAAAAGCTGT\  
ECOCYB

\GCTCCGGCGGGGTTTTTTTGTATATCTGCAATTCAGTACAAAAAGGTGATCAACCCCTCAATTTTTCCCTTTGCTGAAAAATTTTCCATTGTCTCCCTGTAAAAGCTGT\  
ECOCYB

\GCTCCGGCGGGGTTTTTTTGTATATCTGCAATTCAGTACAAAAAGGTGATCAACCCCTCAATTTTTCCCTTTGCTGAAAAATTTTCCATTGTCTCCCTGTAAAAGCTGT\  
ECOCYB

\GCTCCGGCGGGGTTTTTTTGTATATCTGCAATTCAGTACAAAAAGGTGATCAACCCCTCAATTTTTCCCTTTGCTGAAAAATTTTCCATTGTCTCCCTGTAAAAGCTGT\  
ECOCYB

\GCTCCGGCGGGGTTTTTTTGTATATCTGCAATTCAGTACAAAAAGGTGATCAACCCCTCAATTTTTCCCTTTGCTGAAAAATTTTCCATTGTCTCCCTGTAAAAGCTGT\  
ECOCYB

\GCTCCGGCGGGGTTTTTTTGTATATCTGCAATTCAGTACAAAAAGGTGATCAACCCCTCAATTTTTCCCTTTGCTGAAAAATTTTCCATTGTCTCCCTGTAAAAGCTGT\  
ECOCYB

\GCTCCGGCGGGGTTTTTTTGTATATCTGCAATTCAGTACAAAAAGGTGATCAACCCCTCAATTTTTCCCTTTGCTGAAAAATTTTCCATTGTCTCCCTGTAAAAGCTGT\  
ECOCYB

\GCTCCGGCGGGGTTTTTTTGTATATCTGCAATTCAGTACAAAAAGGTGATCAACCCCTCAATTTTTCCCTTTGCTGAAAAATTTTCCATTGTCTCCCTGTAAAAGCTGT\  
ECOCYB

\GCTCCGGCGGGGTTTTTTTGTATATCTGCAATTCAGTACAAAAAGGTGATCAACCCCTCAATTTTTCCCTTTGCTGAAAAATTTTCCATTGTCTCCCTGTAAAAGCTGT\  
ECOCYB

\GCTCCGGCGGGGTTTTTTTGTATATCTGCAATTCAGTACAAAAAGGTGATCAACCCCTCAATTTTTCCCTTTGCTGAAAAATTTTCCATTGTCTCCCTGTAAAAGCTGT\  
ECOCYB

\GCTCCGGCGGGGTTTTTTTGTATATCTGCAATTCAGTACAAAAAGGTGATCAACCCCTCAATTTTTCCCTTTGCTGAAAAATTTTCCATTGTCTCCCTGTAAAAGCTGT\  
ECOCYB

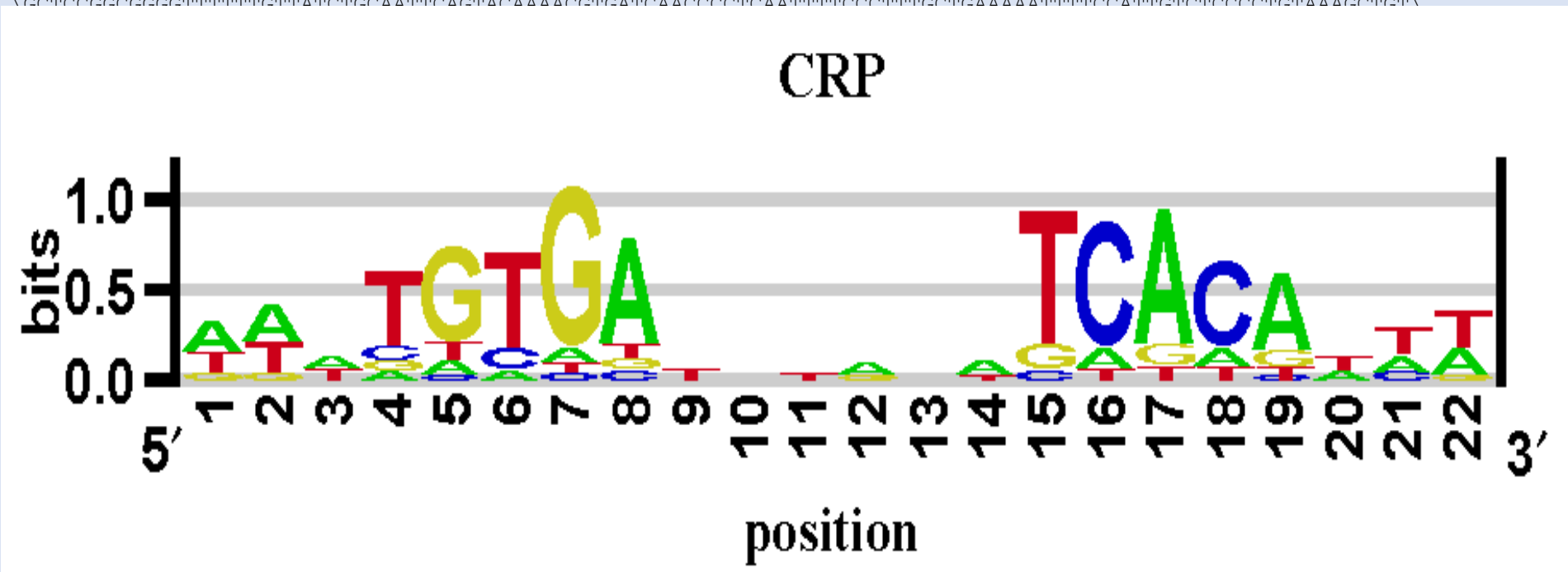
\GCTCCGGCGGGGTTTTTTTGTATATCTGCAATTCAGTACAAAAAGGTGATCAACCCCTCAATTTTTCCCTTTGCTGAAAAATTTTCCATTGTCTCCCTGTAAAAGCTGT\  
ECOCYB

\GCTCCGGCGGGGTTTTTTTGTATATCTGCAATTCAGTACAAAAAGGTGATCAACCCCTCAATTTTTCCCTTTGCTGAAAAATTTTCCATTGTCTCCCTGTAAAAGCTGT\  
ECOCYB

\GCTCCGGCGGGGTTTTTTTGTATATCTGCAATTCAGTACAAAAAGGTGATCAACCCCTCAATTTTTCCCTTTGCTGAAAAATTTTCCATTGTCTCCCTGTAAAAGCTGT\  
ECOCYB

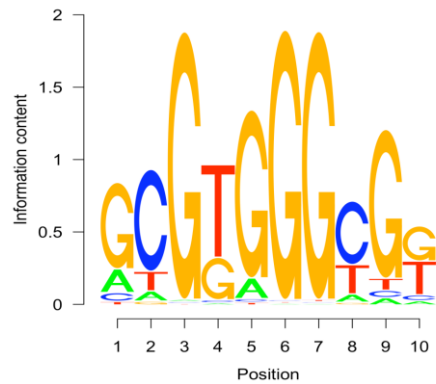
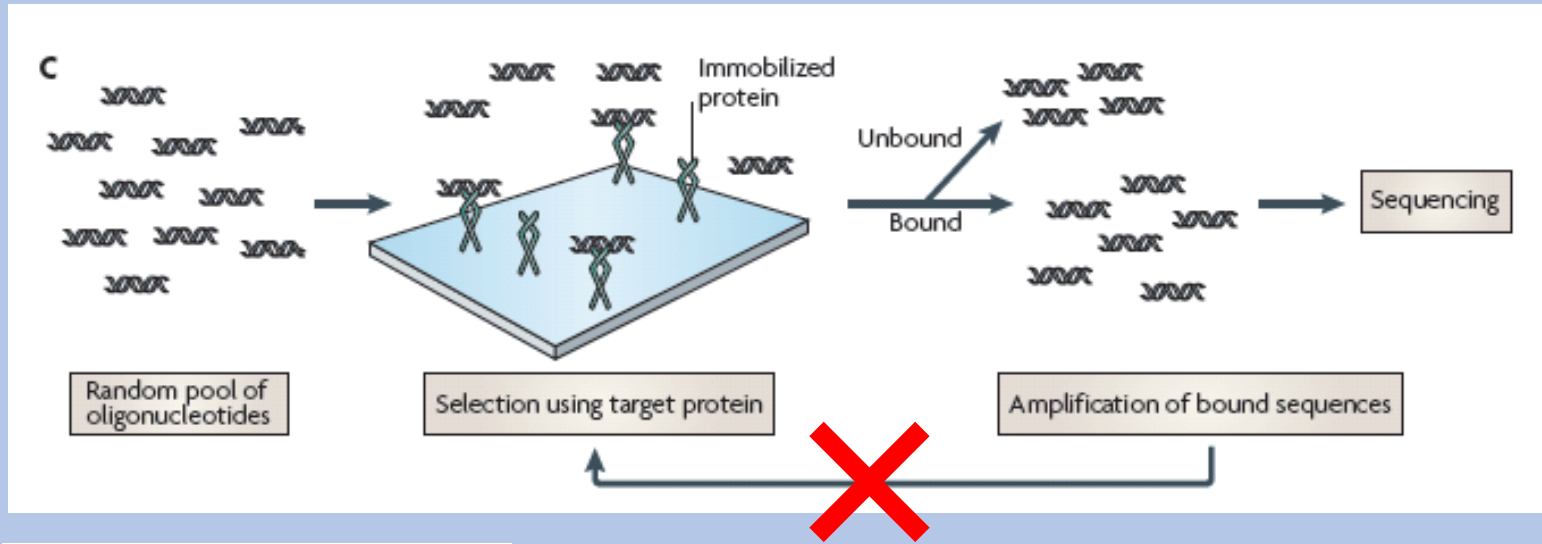
\GCTCCGGCGGGGTTTTTTTGTATATCTGCAATTCAGTACAAAAAGGTGATCAACCCCTCAATTTTTCCCTTTGCTGAAAAATTTTCCATTGTCTCCCTGTAAAAGCTGT\  
ECOCYB

\GCTCCGGCGGGGTTTTTTTGTATATCTGCAATTCAGTACAAAAAGGTGATCAACCCCTCAATTTTTCCCTTTGCTGAAAAATTTTCCATTGTCTCCCTGTAAAAGCTGT\  
ECOCYB



# Inferring Binding Energies from Selected Binding Sites

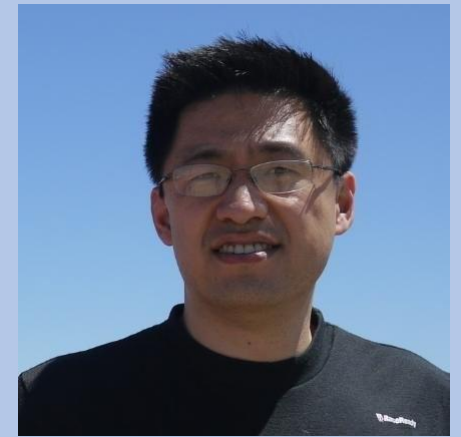
Yue Zhao, David Granas, Gary D. Stormo\*



Further development and generalization of the algorithm

**BEESEM: estimation of binding energy models using HT-SELEX data**

Shuxiang Ruan<sup>1</sup>, S. Joshua Swamidass<sup>2</sup> and Gary D. Stormo<sup>1,\*</sup>



Many groups working to determine motifs for human TFs: Wolfe, Bulyk, Hughes,...  
Jussi Taipale automated the process using robots

## Multiplexed massively parallel SELEX for characterization of human transcription factor binding specificities

Arttu Jolma,<sup>1,2</sup> Teemu Kivioja,<sup>1,3</sup> Jarkko Toivonen,<sup>3</sup> Lu Cheng,<sup>3</sup> Gonghong Wei,<sup>1</sup>  
Martin Enge,<sup>2</sup> Mikko Taipale,<sup>1</sup> Juan M. Vaquerizas,<sup>4</sup> Jian Yan,<sup>1</sup> Mikko J. Sillanpää,<sup>5</sup>  
Martin Bonke,<sup>1</sup> Kimmo Palin,<sup>3</sup> Shaheynoor Talukder,<sup>6</sup> Timothy R. Hughes,<sup>6</sup>  
Nicholas M. Luscombe,<sup>4</sup> Esko Ukkonen,<sup>3</sup> and Jussi Taipale<sup>1,2,7</sup>

*Genome Res.* 2010 20: 861-873

~20 TFs

## DNA-Binding Specificities of Human Transcription Factors

Arttu Jolma,<sup>1,2,8</sup> Jian Yan,<sup>1,8</sup> Thomas Whittington,<sup>1</sup> Jarkko Toivonen,<sup>3</sup> Kazuhiro R. Nitta,<sup>1</sup> Pasi Rastas,<sup>3</sup>  
Ekaterina Morgunova,<sup>1</sup> Martin Enge,<sup>1</sup> Mikko Taipale,<sup>2</sup> Gonghong Wei,<sup>2</sup> Kimmo Palin,<sup>2</sup> Juan M. Vaquerizas,<sup>4</sup>  
Renaud Vincentelli,<sup>5</sup> Nicholas M. Luscombe,<sup>4</sup> Timothy R. Hughes,<sup>6</sup> Patrick Lemaire,<sup>7</sup> Esko Ukkonen,<sup>3</sup> Teemu Kivioja,<sup>1,2,3</sup>  
and Jussi Taipale<sup>1,2,\*</sup>

*Cell* 152, 327–339, January 17, 2013

~240 TFs

## Determination and Inference of Eukaryotic Transcription Factor Sequence Specificity

*Cell* 158, 1431–1443, September 11, 2014 Weirauch et al

>1000 TFs

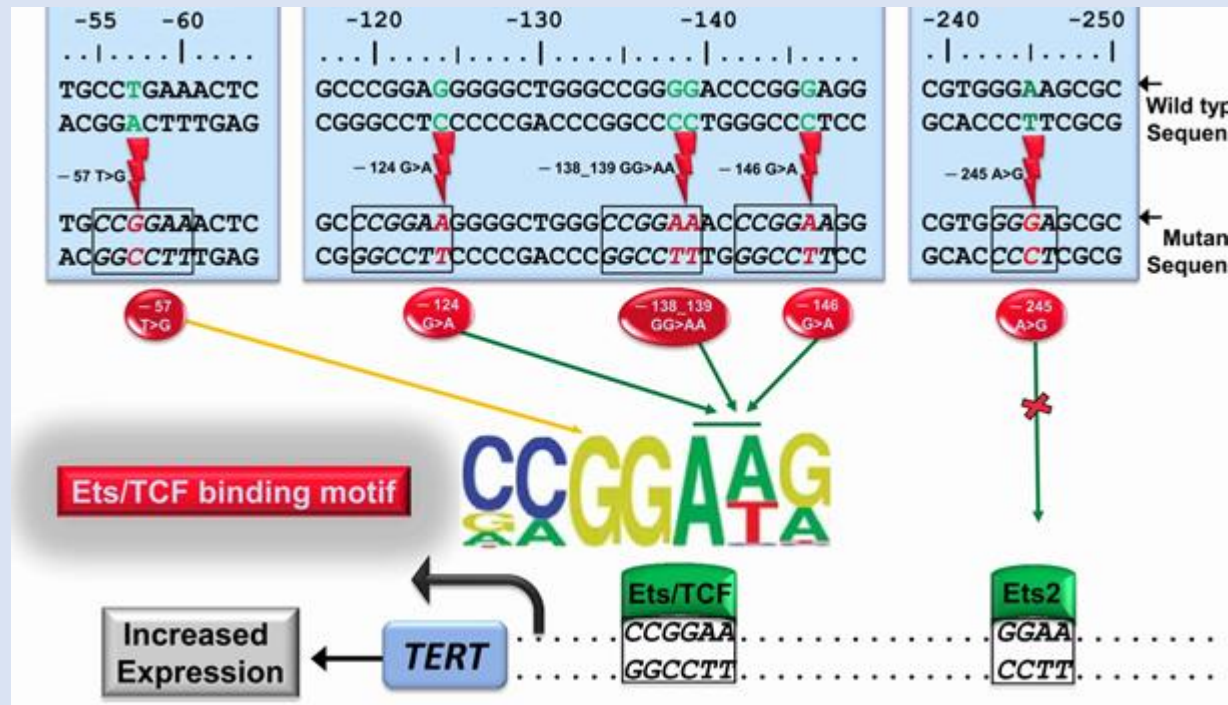
A Biomedical applications of TF motifs:

Genome-Wide Association Studies (GWAS) identify DNA variants associated (correlated) with phenotypes:

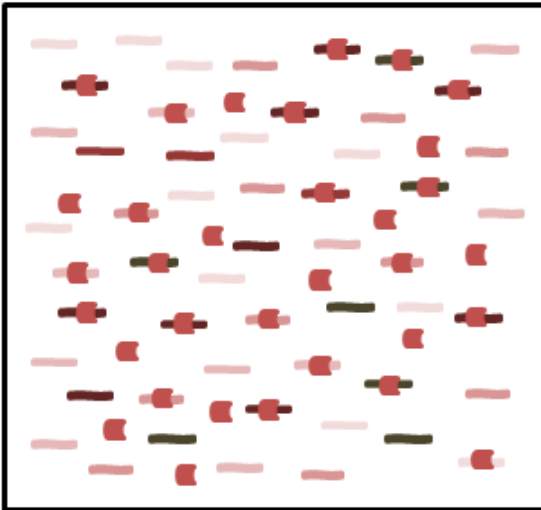
- normal human variations
- disease causing

The vast majority of GWAS hits (>90%) are in non-coding DNA. Suggests alterations in gene expression rather than gene function cause phenotypes

Telomerase Reverse Transcriptase (TERT) promoter mutations are the most frequent non-coding mutations in cancer, occur in over 50 types of cancer and are very common in some (80% melanoma, 70% glioblastoma...) Create an ETS binding site that leads to high expression of the TERT protein, facilitating cancer growth.



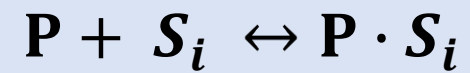
A.



B.



## Spec-seq: Specificity by sequencing



$$K_A(S_i) = \frac{[P \cdot S_i]}{[P][S_i]}$$

$$K_A(S_1) : K_A(S_2) : \dots : K_A(S_n)$$

$$= \frac{[P \cdot S_1]}{[S_1]} : \frac{[P \cdot S_2]}{[S_2]} : \dots : \frac{[P \cdot S_n]}{[S_n]}$$

# Specificity of the Lac repressor

Wild-type sequences

O1: **AATTGTGAG CGG ATAACAATT**

Randomized libraries

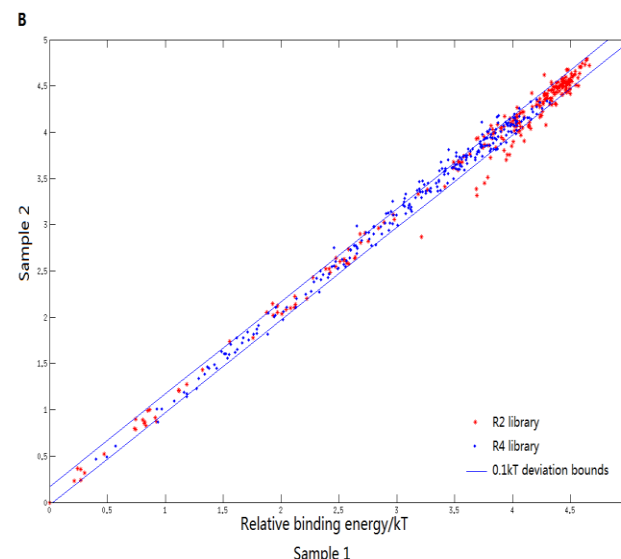
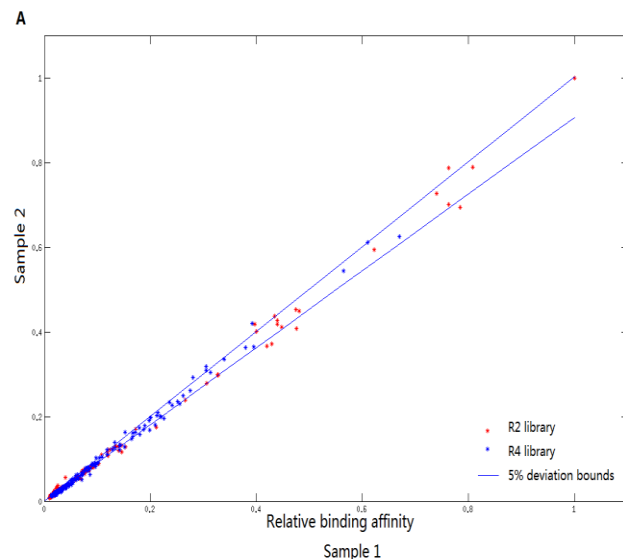
R2: **AATTGTNAN C G NTNACAATT**

R3.1: **AATTGTGAG C N G NNNNCAATT**

R3.2: **AATTGNNNN C N G ATAACAATT**

R4: **AATTGTNAN CC GG NTNACAATT**

-10 -6 -1 -0 0 +0 +1 6 10



WT operator is asymmetric

4 libraries: vary both sequence and spacing

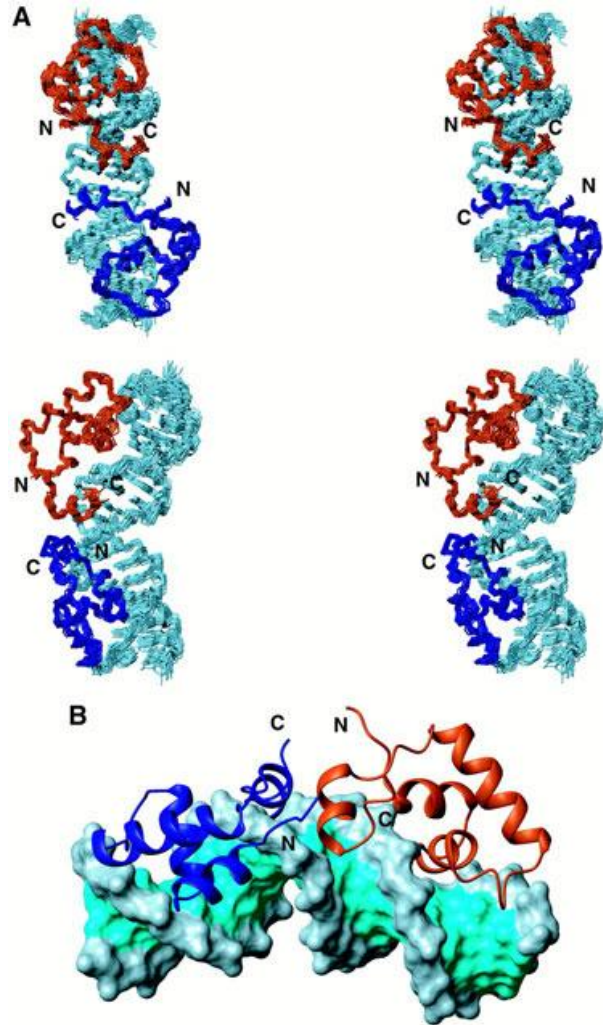
2560 different binding sites

Highly reproducible:  
~5% variance in affinity  
~0.1kT variance in energy

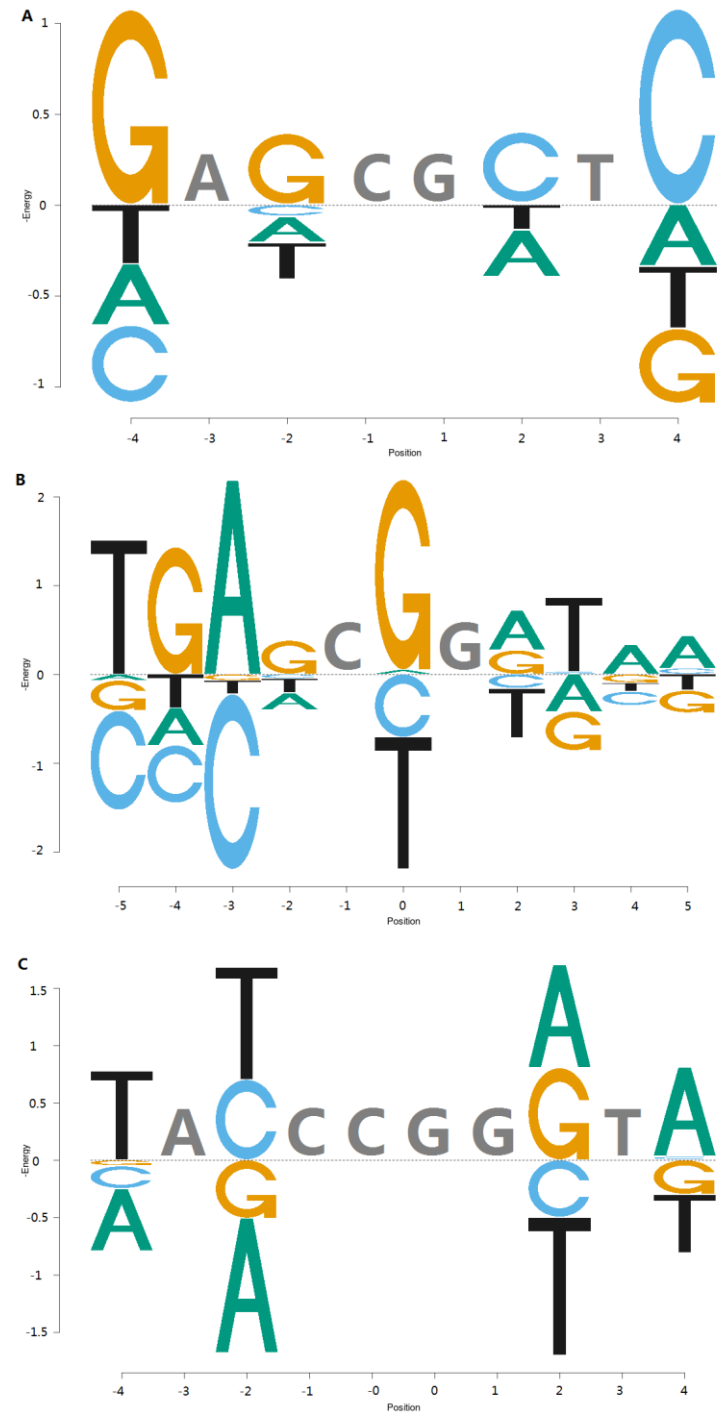


Zuo, Stormo (2014) Genetics

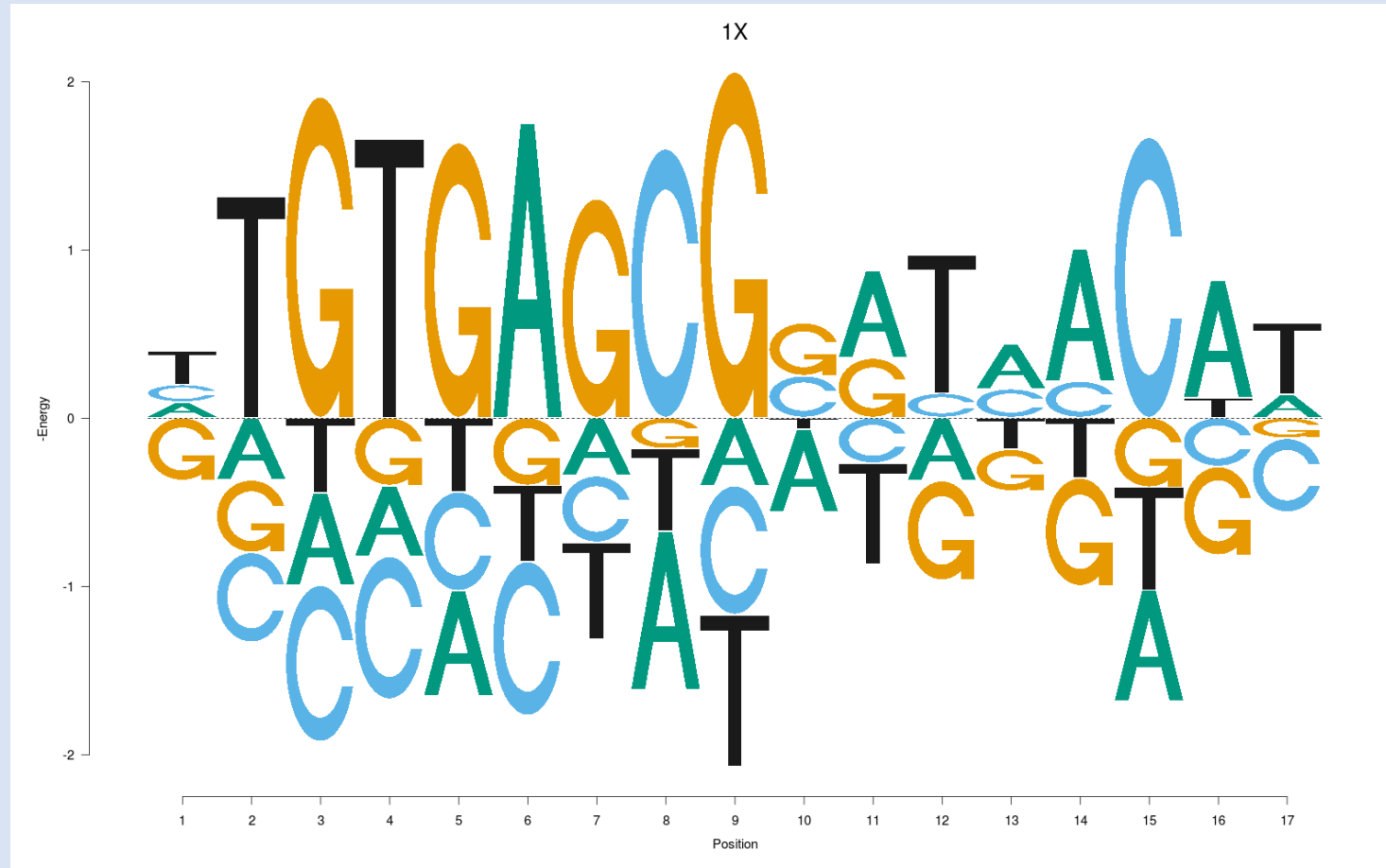
## Three-dimensional structure of the dimeric lac HP62–O1 operator complex.



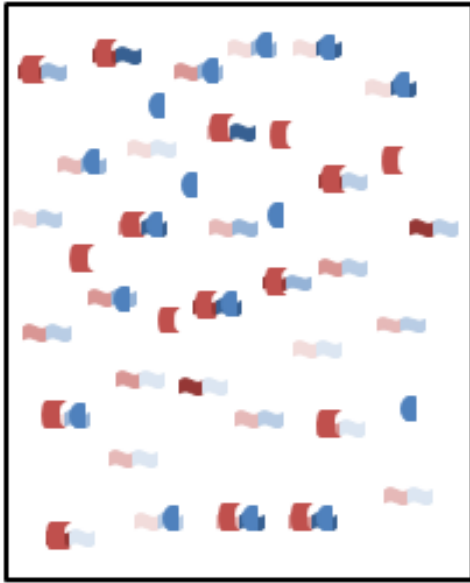
Kalodimos C G et al. *EMBO J.* 2002;21:2866-2876



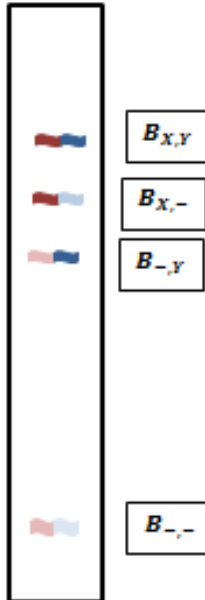
# Specificity of complete “3 spacer” binding sites



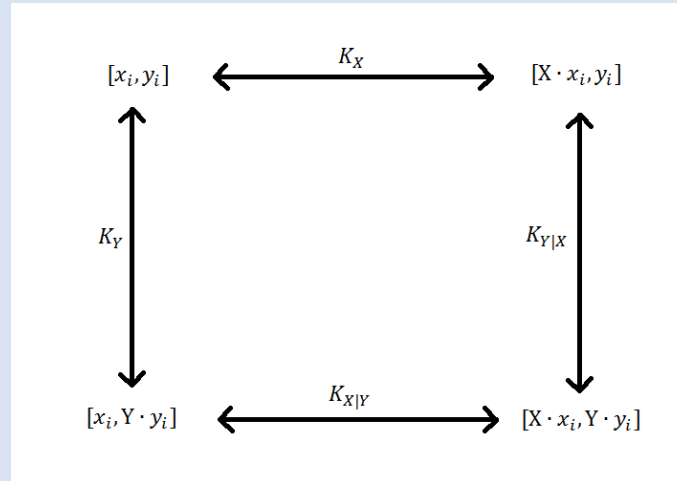
A.



B.



**Coop-seq:** get all of the important parameters in one experiment, including cooperativity



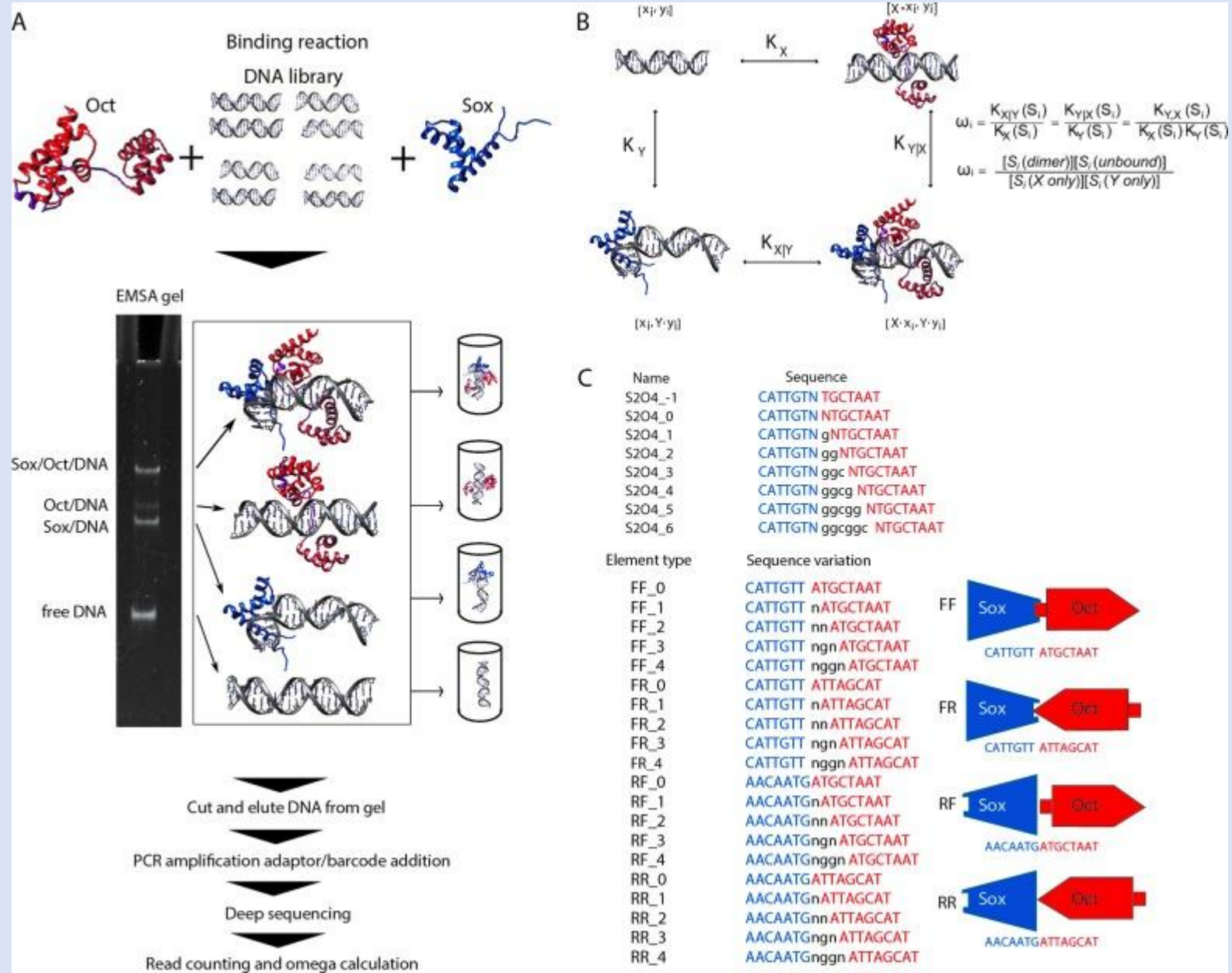
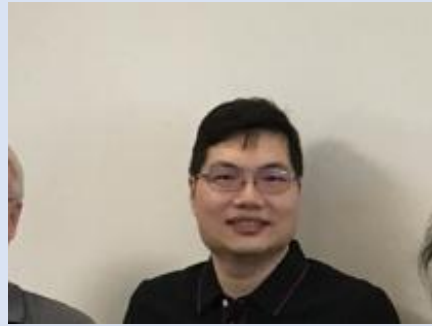
$$K_X(x_1): K_X(x_2): \dots : K_X(x_n) = \frac{N(x_1|B_{X,-})}{N(x_1|B_{-,-})} \cdot \frac{N(x_2|B_{X,-})}{N(x_2|B_{-,-})} \cdot \dots \cdot \frac{N(x_n|B_{X,-})}{N(x_n|B_{-,-})}$$

$$K_{X|Y}(x_1): K_{X|Y}(x_2): \dots : K_{X|Y}(x_n) = \frac{N(x_1|B_{X,Y})}{N(x_1|B_{-,Y})} \cdot \frac{N(x_2|B_{X,Y})}{N(x_2|B_{-,Y})} \cdot \dots \cdot \frac{N(x_n|B_{X,Y})}{N(x_n|B_{-,Y})}$$

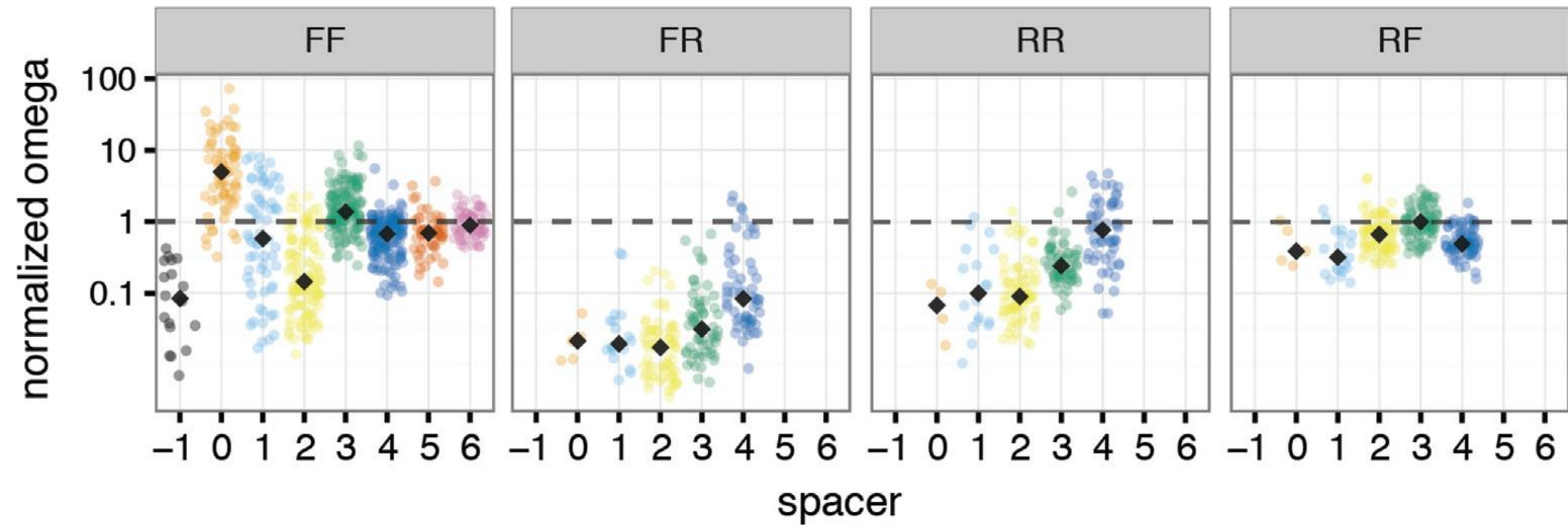
$$\omega_i = \frac{K_{X|Y}(S_i)}{K_X(S_i)} = \frac{K_{Y|X}(S_i)}{K_Y(S_i)} = \frac{K_{X,Y}(S_i)}{K_X(S_i)K_Y(S_i)}$$

# Quantitative profiling of selective Sox/POU pairing on hundreds of sequences in parallel by Coop-seq

(NAR, 2017) Chang et al, collaboration with [Ralf Jauch lab](#)



# Sox2/Oct4



Aberrant homeodomain–DNA cooperative dimerization underlies distinct developmental defects in two dominant CRX retinopathy models

Yiqiao Zheng, Gary D. Stormo, and Shiming Chen

*Genome Res.* 2025 35(2):242-256.

Two mutations in the CRX DNA-binding domain of the protein

E80A changes Glu at position 80 to Ala

K88N changes Lys at position 88 to Asn



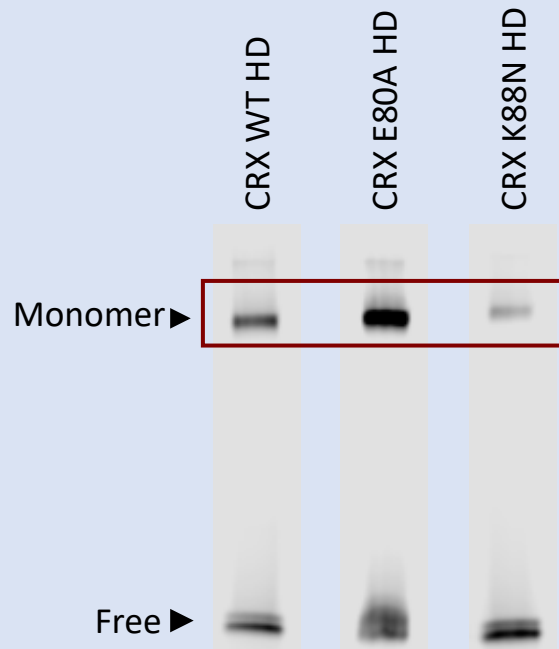
# Spec-seq captures CRX DNA binding specificity changes at monomer motif

CRX WT consensus

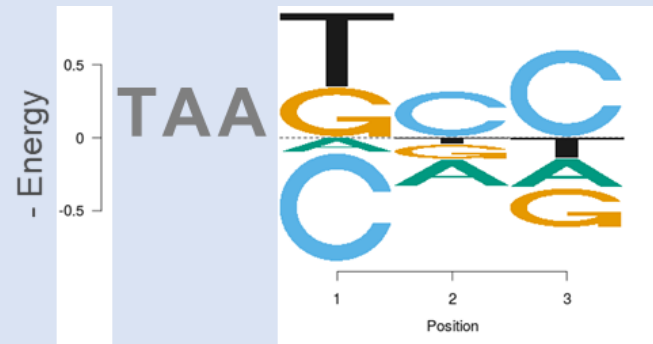
TAATCC

Monomer library

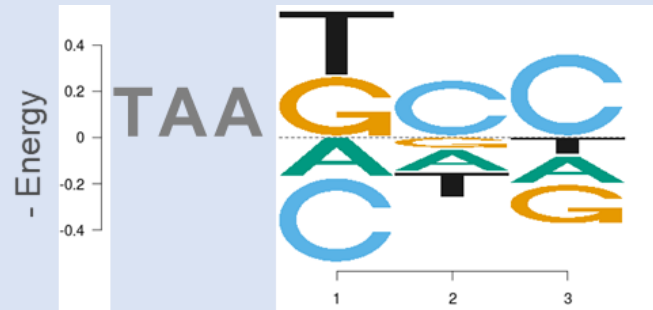
... TAA NNN ...



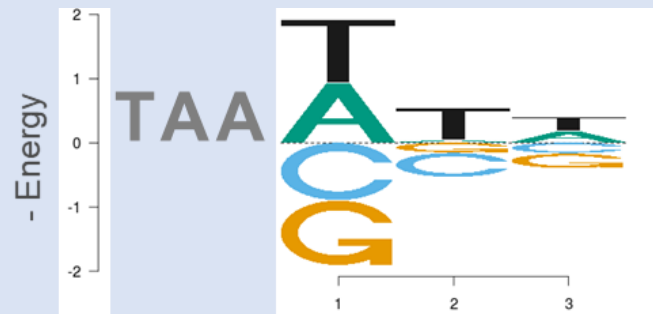
WT HD



E80A HD: **better tolerate nucleotide variants**



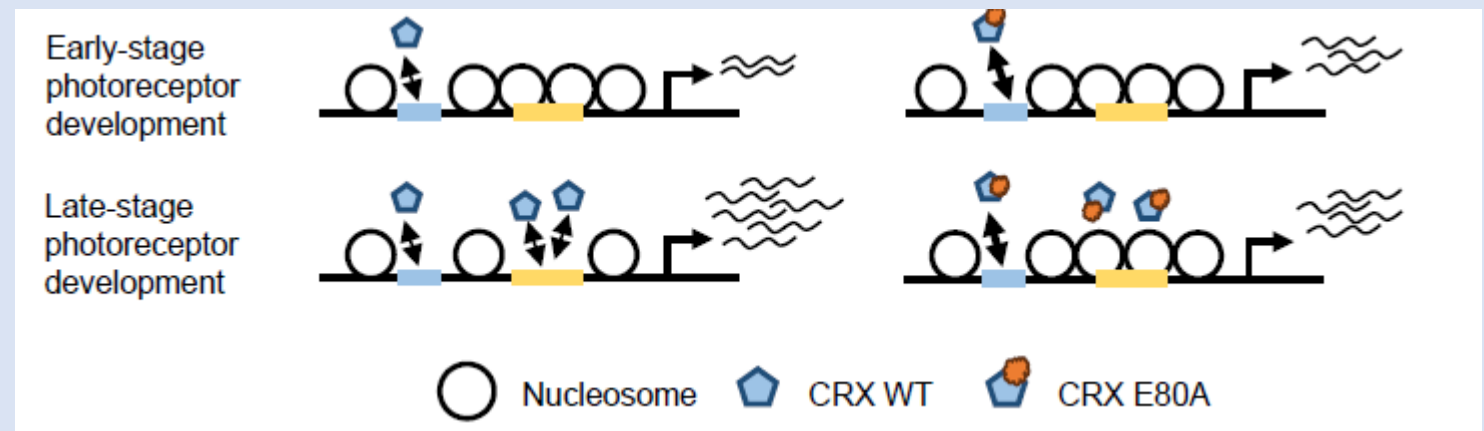
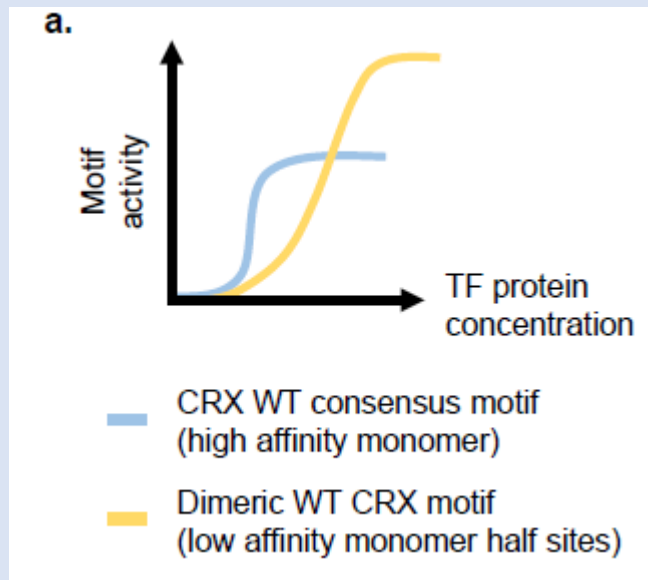
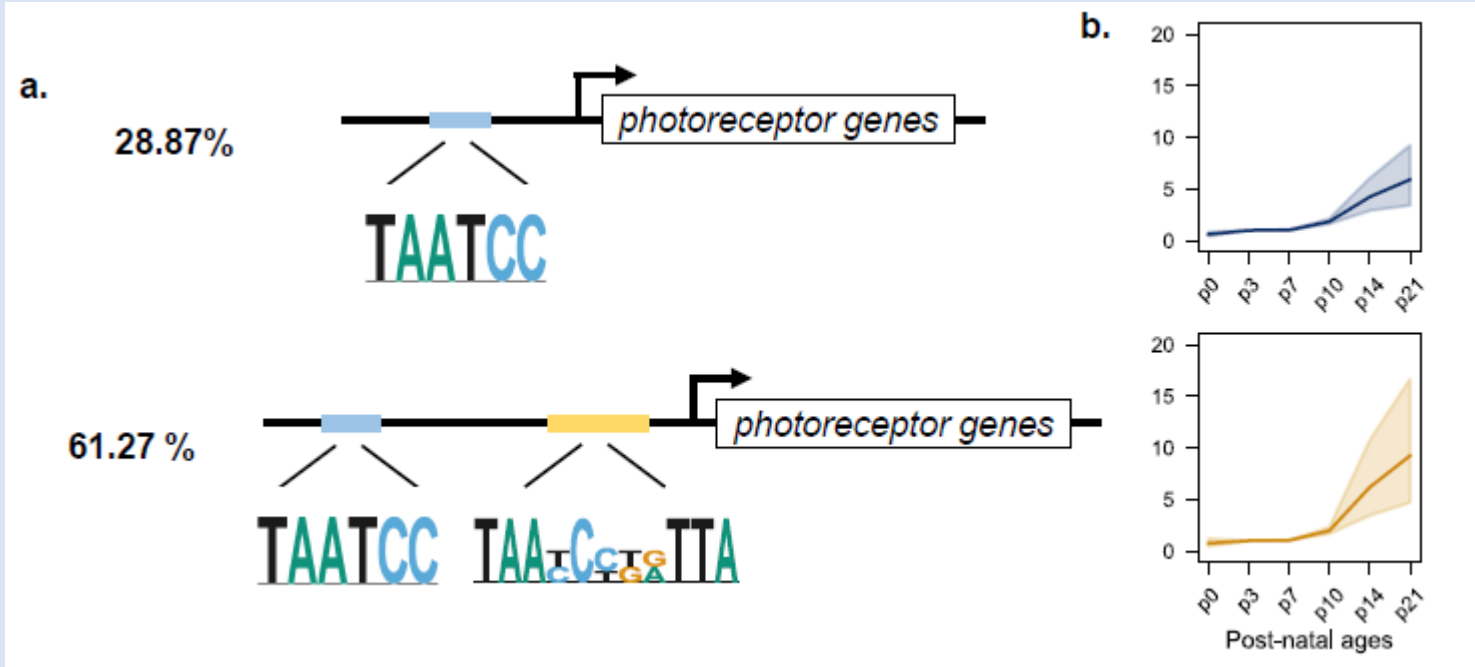
K88N HD: **altered specificity**



E80A  
K88N

Homeodomain

# E80A has lost cooperativity



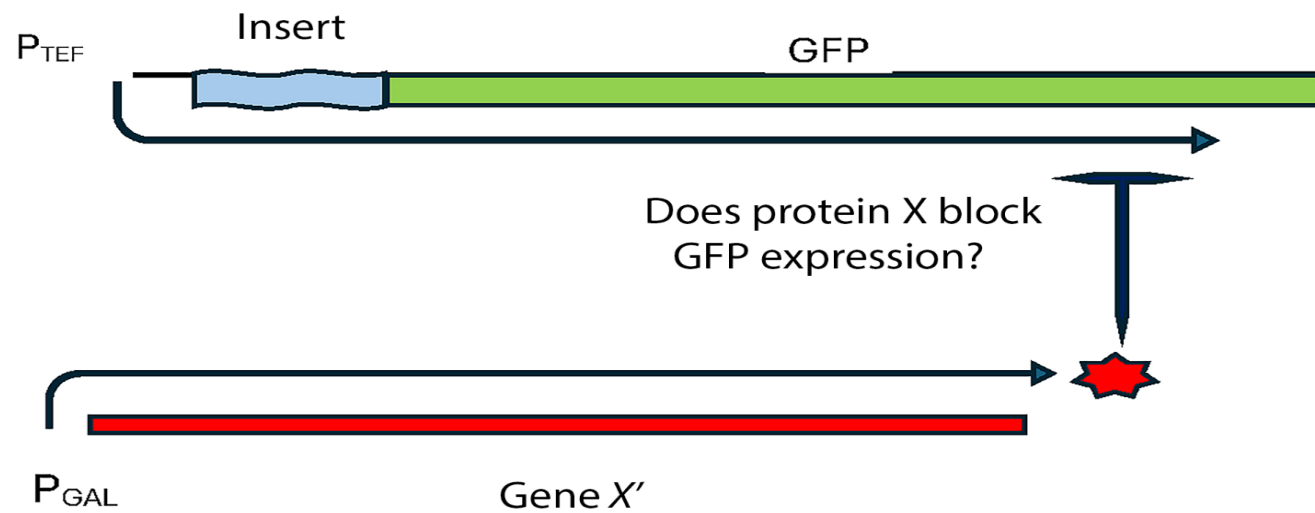
Regulation of Transcription Initiation is Ubiquitous  
but post-transcriptional initiation is also common

Regulation can occur at  
Transcription termination  
RNA splicing  
RNA translation  
RNA degradation  
Protein modification  
Protein degradation

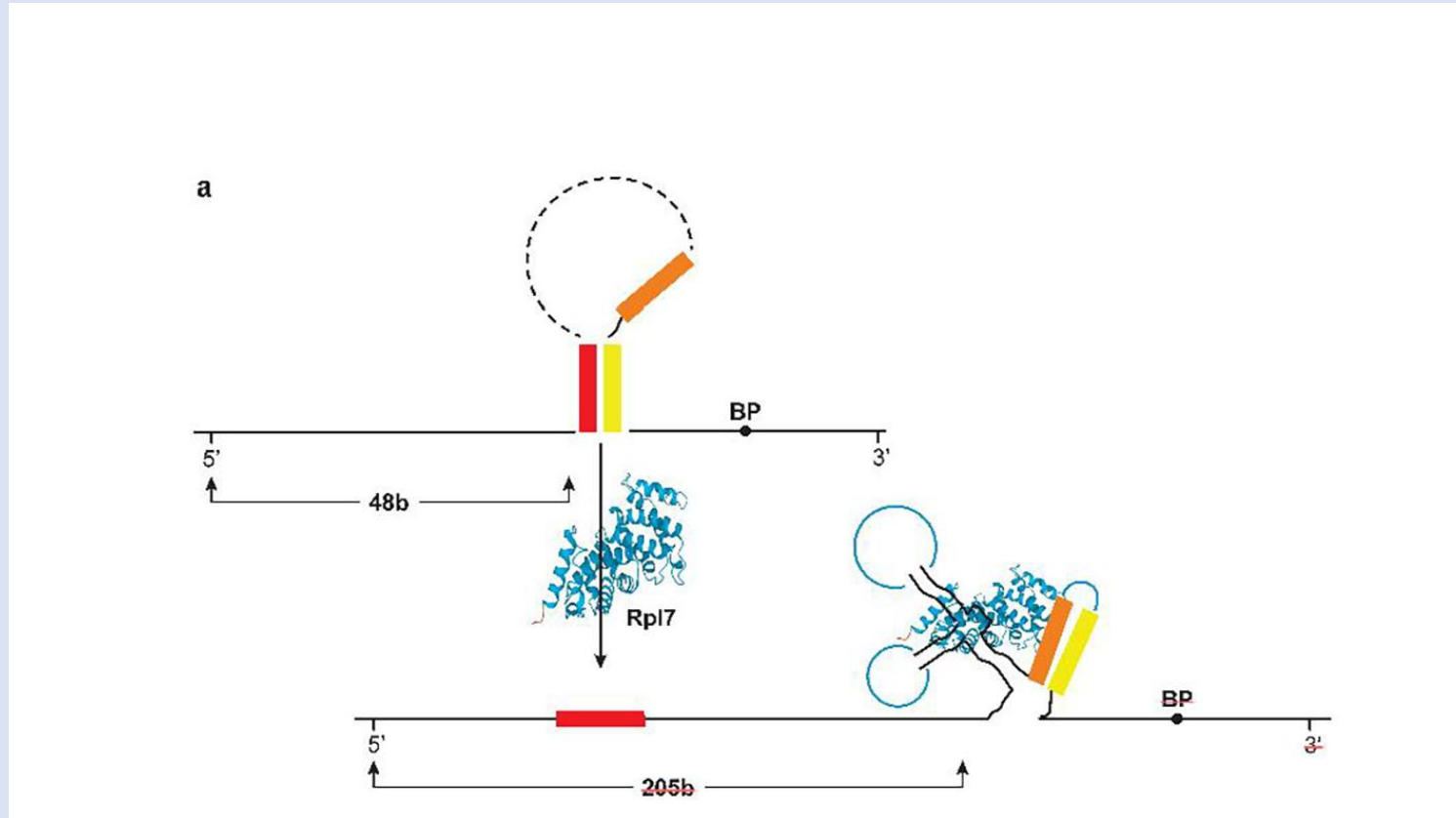
Auto-Regulation is of particular interest to me

A gene can control its own expression, setting its  
appropriate level (think cruise control)

# Autoregulation of yeast ribosomal protein genes



**Two alternative RNA structures. One serves as an enhancer to facilitate splicing and protein expression. The protein binds to the alternative structure to inhibit splicing and expression.**



# Acknowledgements

## From the beginning:

Larry Gold  
Andrzej Ehrenfeucht  
Tom Schneider

Gerry Hertz, George Hartzell  
Charles Lawrence, John Heumann  
Alan Lapedes, Chris Workman  
Sam Levy, Dana Fields, Eric Snyder



## Wash U postdocs, students, research assistants, collaborators:

Fugen Li	Alan Kwan	Paul Taghert, Russ van Gelder
Takis Benos	Kai Tan	Jeff Milbrandt, Jeff Magee
Debraj GuhaThakurta	Yongmei Ji	Perren Cobb, Richard Hotchkiss, Mark Watson
Tao Zhao	Robin Matlib	Rafi Kopan, Buddy Brownstein, Tim Buchman
Yong Yin	Ritesh Agrawal	Jeff Gordon
Jiajian Liu	Liwei Chan	Charles Gu, Mike Province,
Guoyan Zhao	Xing Xu	Susan Dutcher
Burr Fontaine	Billy Li	Howard McLeod
Chris Man	Yiing Lin	Bob Waterston
Dana Homsy	Ting Wang	Alan Permutt, Jean Shaffer
Basab Roy	Yue Zhao	Mark Johnston, Barak Cohen
David Granas	Aaron Spivak	Jeremy Buhler, Weixong Zhang
Larry Schreifer	Ryan Christensen	Steve Beverley
Nnamdi Ihuegbu	Kenny Chang	Shin Imai, Christina Strong
Shuxiang Ruan	Zheng Zuo	Tim Ley, Jacqueline Payton,
Manishi Pandey	Shane Chu	Jim Havranek
Lala Motlhabi	Ronak Patel	Linda Sandell
Barrett Foat	Gurmukh Sahota	Ken Murphy, Xiaowei Wang



Funding: NIH, DOE

## Outside collaborators:

Scot Wolfe, Chris Link, Javier Irazoqui, Ralf Jauch, Tony Gerber, Jan Gorodkin, Petko Petkov, Shiming Chen