# Motifs

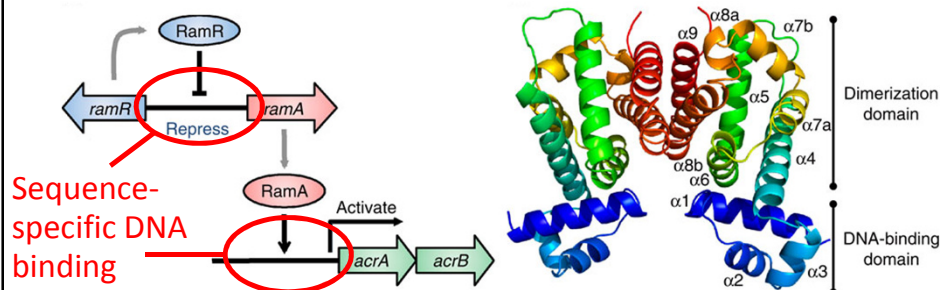**BCH394P/364C - Systems Biology / Bioinformatics**

**Edward Marcotte, Univ of Texas at Austin**
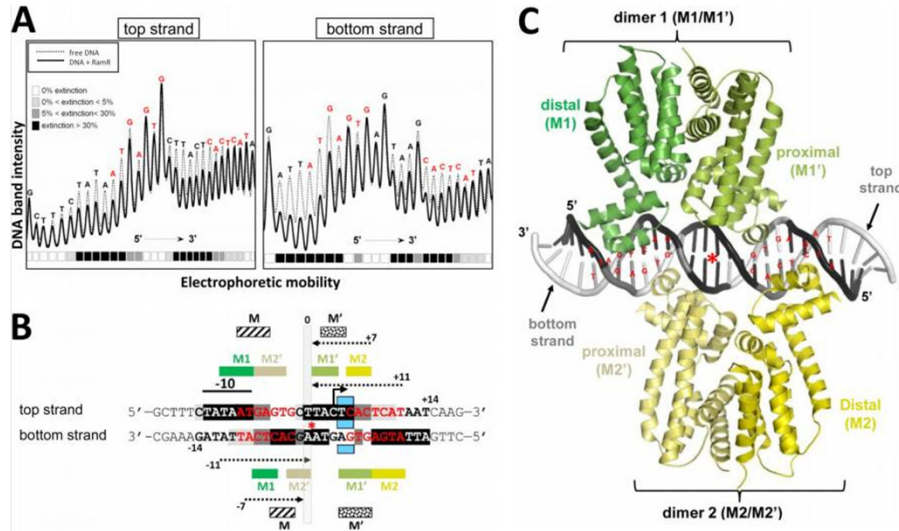
---

## An example transcriptional regulatory cascade
### Here, controlling *Salmonella* bacteria multidrug resistance



Sequence-specific DNA binding

RamR represses the *ramA* gene, which encodes the activator protein for the *acrAB* drug efflux pump genes.

RamR dimer

## Historically, DNA and RNA binding sites were defined biochemically (DNAse footprinting, gel shift assays, etc.)



Hydroxyl radical footprinting of *ramR-ramA* intergenic region with RamR

## Historically, DNA and RNA binding sites were defined biochemically (DNAse footprinting, gel shift assays, etc.)

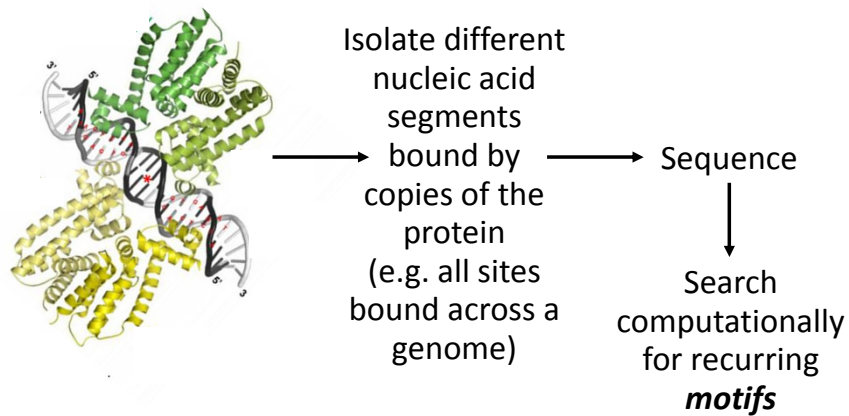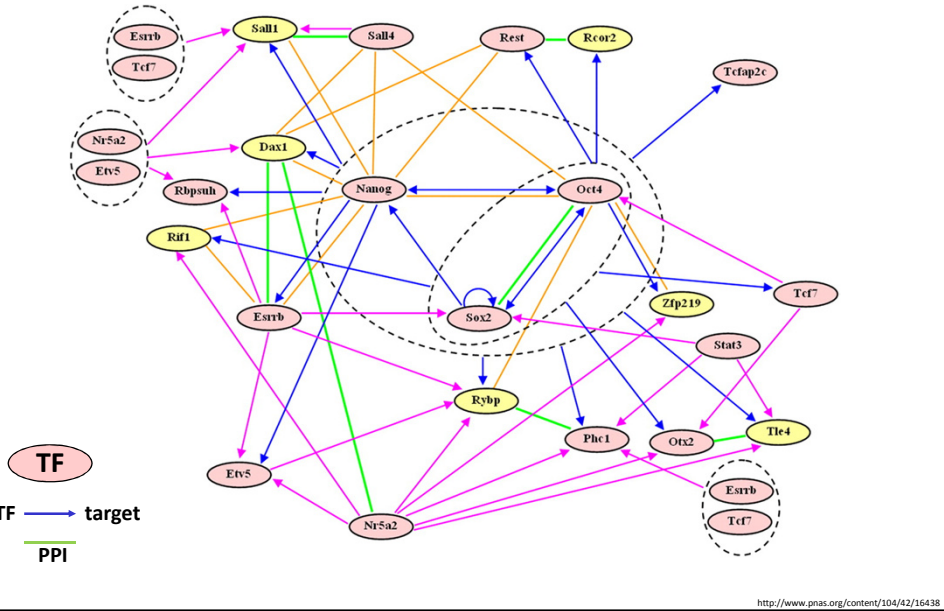Now, many binding motifs are discovered bioinformatically



Isolate different nucleic acid segments bound by copies of the protein (e.g. all sites bound across a genome) → Sequence → Search computationally for recurring *motifs*

# Transcription factor regulatory networks can be highly complex, e.g. as for embryonic stem cell regulators



TF

TF —→ target
PPI

---

## MOTIFS

| | |
|---|---|
| HEM13 | CCCATTGTTCTC |
| HEM13 | TTTCTGGTTCTC |
| HEM13 | TCAATTGTTTAG |
| ANB1 | CTCATTGTTGTC |
| ANB1 | TCCATTGTTCTC |
| ANB1 | CCTATTGTTCTC |
| ANB1 | TCCATTGTTCGT |
| ROX1 | CCAATTGTTTTG |

**Binding sites of the transcription factor ROX1**

YCHATTGTTCTC    **consensus**

| A | 002700000010 |
|---|---|
| C | 464100000505 |
| G | 000001800112 |
| T | 422087088261 |

**frequencies**



$$I_{seq}(i) = -\sum_b f_{b,i} \log_2 \frac{f_{b,i}}{P_b}$$

frequency of nuc b at position i

freq of nuc b in genome

**scaled by information content**
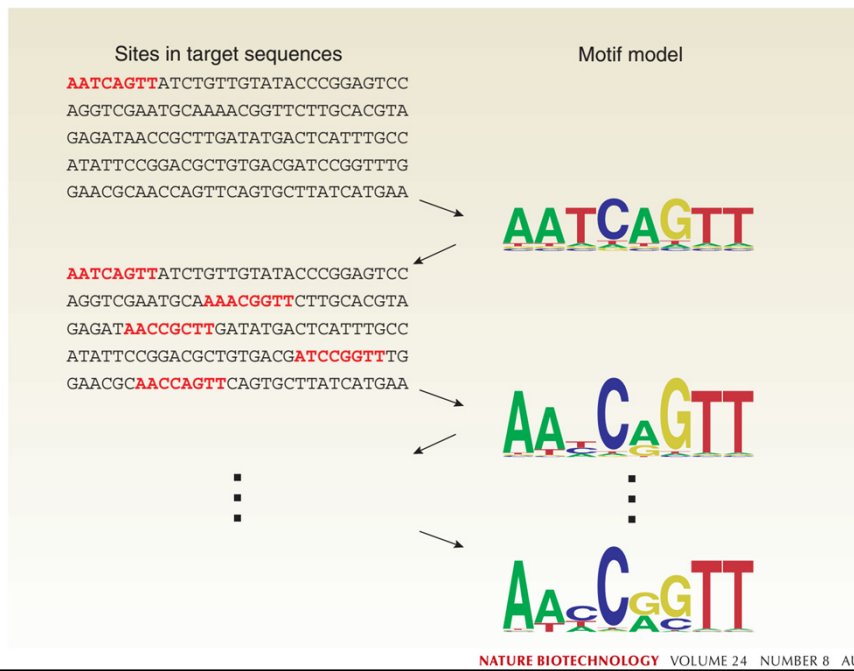
Bob Crimi

**So, here's the challenge:**

**Given a set of DNA sequences that contain a motif (e.g., promoters of co-expressed genes), how do we discover it computationally?**
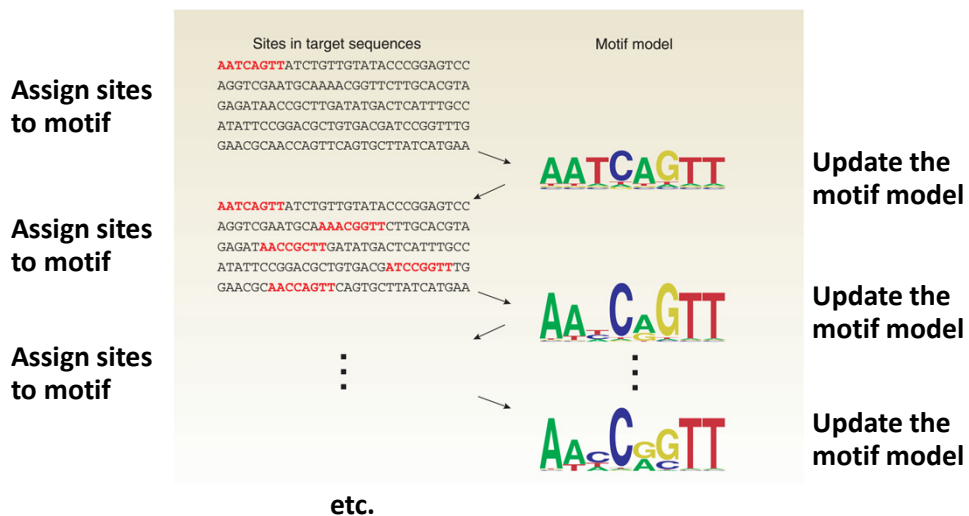
**Could we just count all instances of each *k*-mer?**

**Why or why not?**

**→ promoters and DNA binding sites are not well conserved**

# How does motif discovery work?



Sites in target sequences

**AATCAGTT**ATCTGTTGTATACCCGGAGTCC
AGGTCGAATGCAAAACGGTTCTTGCACGTA
GAGATAACCGCTTGATATGACTCATTTGCC
ATATTCCGGACGCTGTGACGATCCGGTTTG
GAACGCAACCAGTTCAGTGCTTATCATGAA

**AATCAGTT**ATCTGTTGTATACCCGGAGTCC
AGGTCGAATGCA**AAACGGTT**CTTGCACGTA
GAGAT**AACCGCTT**GATATGACTCATTTGCC
ATATTCCGGACGCTGTGACG**ATCCGGTTT**TG
GAACGC**AACCAGTT**CAGTGCTTATCATGAA

Motif model

NATURE BIOTECHNOLOGY VOLUME 24 NUMBER 8 AUGUST 2006

# How does motif discovery work?

**Assign sites to motif**

**Assign sites to motif**

**Assign sites to motif**

**Update the motif model**

**Update the motif model**

**Update the motif model**

etc.

**How does motif discovery work?**

**Motif finding often uses _expectation-maximization_ _i.e._ alternating between building/updating a motif model and assigning sequences to that motif model.**

**Searches the space of possible motifs for optimal solutions without testing everything.**

**Most common approach = _Gibbs sampling_**

---

**Detecting Subtle Sequence Signals: A Gibbs Sampling Strategy for Multiple Alignment**

Charles E. Lawrence, Stephen F. Altschul, Mark S. Boguski, Jun S. Liu, Andrew F. Neuwald, John C. Wootton

**We will consider N sequences, each with a motif of length w:**



$A_k$ = position in seq k of motif

N seqs

k

w

$q_{ij}$ = probability of finding nucleotide (or aa) j at position i in motif
   i ranges from 1 to w
   j ranges across the nucleotides (or aa)
$p_j$ = background probability of finding nucleotide (or aa) j
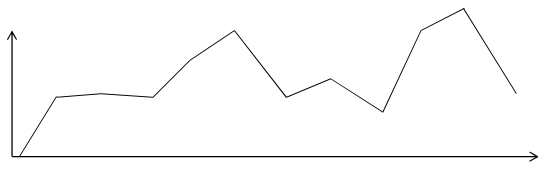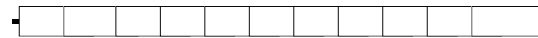
**Detecting Subtle Sequence Signals: A Gibbs Sampling Strategy for Multiple Alignment**

Charles E. Lawrence, Stephen F. Altschul, Mark S. Boguski, Jun S. Liu, Andrew F. Neuwald, John C. Wootton

**NOTE: You won't give any information at all about what or where the motif should be!**

## Start by choosing w and randomly positioning each motif:

$A_k$ = position in seq k of motif

N seqs

k

**Completely randomly positioned!**

$q_{ij}$ = probability of finding nucleotide (or aa) j at position i in motif
   i ranges from 1 to w
   j ranges across the nucleotides (or aa)
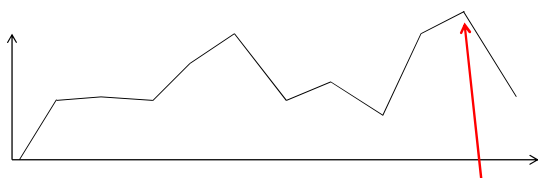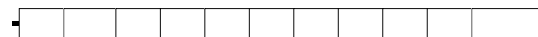$p_j$ = background probability of finding nucleotide (or aa) j

---

**Detecting Subtle Sequence Signals: A Gibbs Sampling Strategy for Multiple Alignment**

Charles E. Lawrence, Stephen F. Altschul, Mark S. Boguski, Jun S. Liu, Andrew F. Neuwald, John C. Wootton

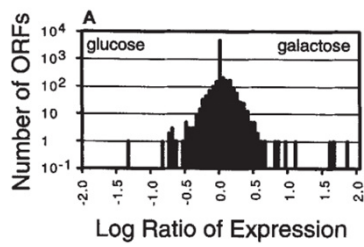## Predictive update step: Randomly choose one sequence, calculate $q_{ij}$ and $p_j$ from N-1 remaining sequences

**Randomly choose →**

**Update model w/ these**

**background frequency of symbol j**

**count of symbol j at position i**

**$\Sigma b_j$**

$$q_{i,j} = \frac{c_{i,j} + b_j}{N - 1 + B}$$

**$p_j$ is calculated similarly from the counts outside the motifs**

$q_{ij}$ = probability of finding nucleotide (or aa) j
   i ranges from 1 to w
   j ranges across the nucleotides (or aa)
$p_j$ = background probability of finding nucleotide (or aa) j

7

**Over many iterations, this magically converges to the most enriched motifs. Note, it's stochastic:**
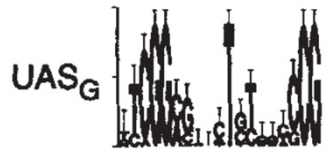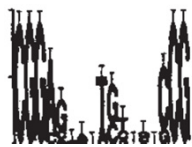


**3 runs on the same data**

Finding DNA regulatory motifs within
unaligned noncoding sequences clustered
by whole-genome mRNA quantitation

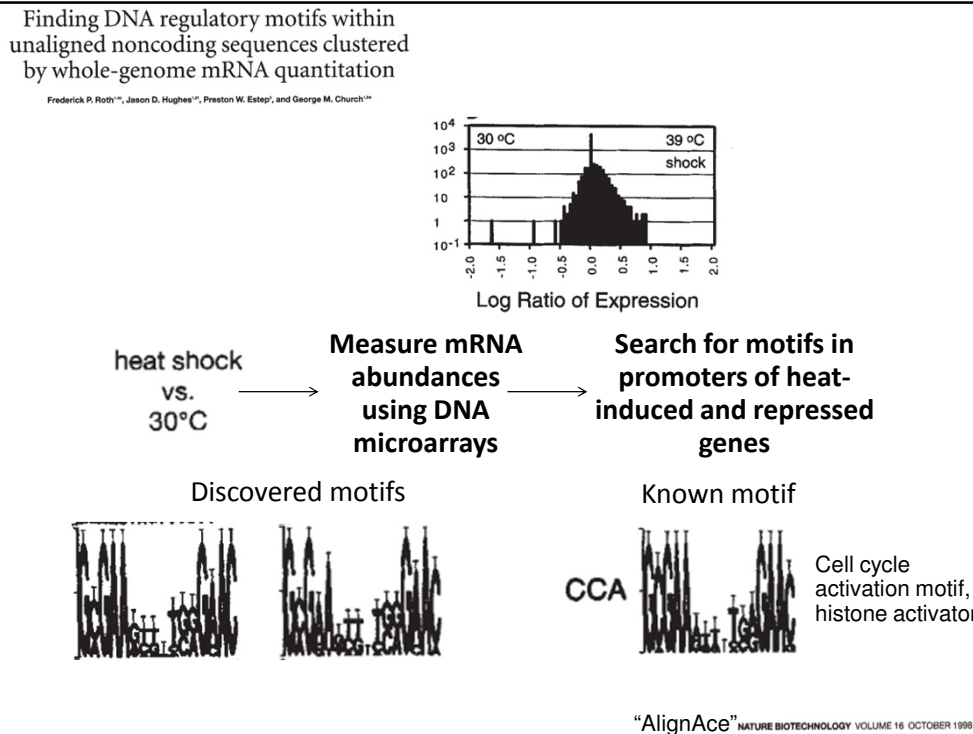Frederick P. Roth, Jason D. Hughes, Preston W. Estep, and George M. Church



galactose
vs.
glucose

→

**Measure mRNA abundances using DNA microarrays**

→

**Search for motifs in promoters of glucose vs galactose controlled genes**

Discovered motifs

Known motif



UAS_G

Galactose
upstream
activation
sequence

Finding DNA regulatory motifs within unaligned noncoding sequences clustered by whole-genome mRNA quantitation

Frederick P. Roth, Jason D. Hughes, Preston W. Estep, and George M. Church
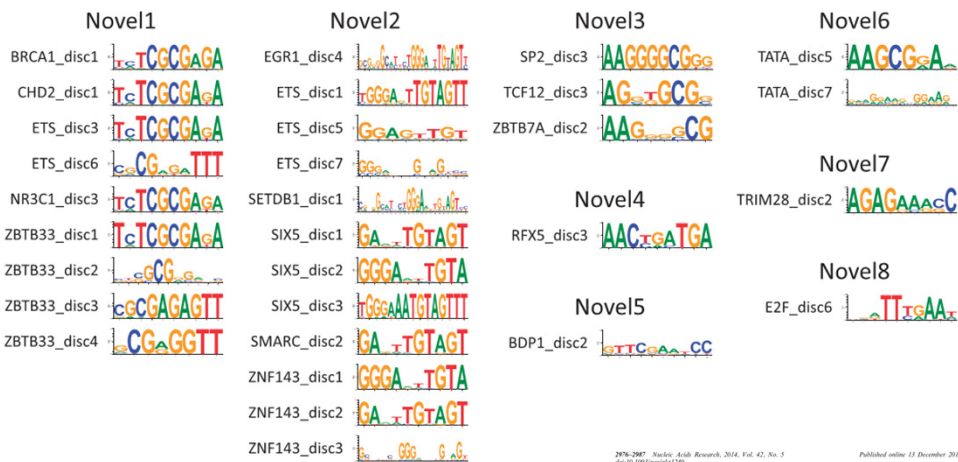


"AlignAce" NATURE BIOTECHNOLOGY VOLUME 16 OCTOBER 1998

---

**If you need them, we now know the binding motifs for 100's of transcription factors at 1000's of distinct sites in the human genome, including many new motifs.**
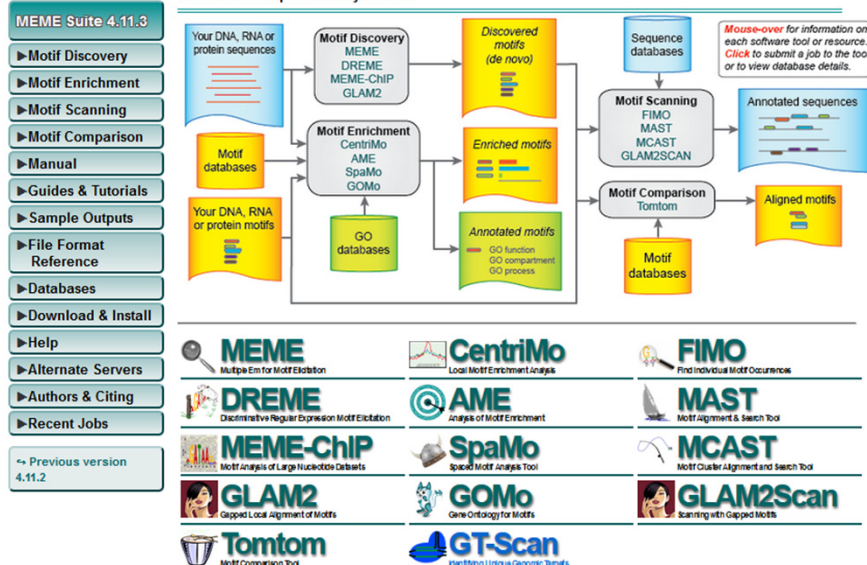e.g., http://compbio.mit.edu/encode-motifs/



2976–2987 Nucleic Acids Research, 2014, Vol. 42, No. 5
doi:10.1093/nar/gkt1249
Published online 13 December 2013

Systematic discovery and characterization of regulatory motifs in ENCODE TF binding experiments

**Here's a good place to start if you want to do this practically:** http://meme-suite.org/



**Note: online MEME suite can sometimes be quite laggy. GibbsCluster is a good alternative for peptide motifs:** http://www.cbs.dtu.dk/services/GibbsCluster/



**Both can also be installed on your own computer**