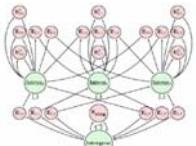# Computational Genomics

**10-810/02-710, Spring 2009**

## Model-based Comparative Genomics

**Eric Xing**

**Lecture 14, March 2, 2009**

**Reading:** class assignment

---

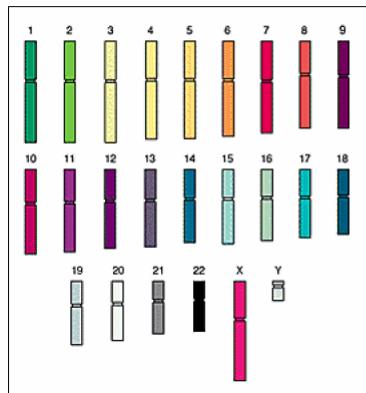# Uses of evolutionary theory

- Comparative genomics (this lecture)
  - Cladistics: figuring out closely related species, proteins, sequences
    - Drug design and testing
    - Building chimera : mixing genetic codes of species, genetic technology
    - Sequence prediction in related unsequenced species : use in sequencing, primer design, etc
  - Phylogenetic footprinting
    - Functional constraints on a genomic region inversely proportional to evolutionary rate, from neutral theory
    - Look at two concrete examples : transcription factor binding site (motif) and gene prediction

- Population genetics (module 4)
  - Population structure
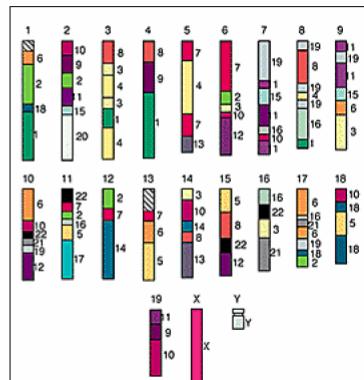  - Understanding evolutionary driving force underlying genome variation
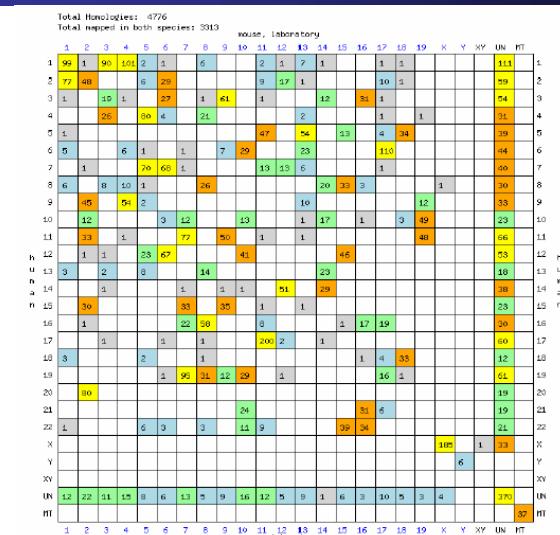
# Comparative Genomics

Human                                    Mouse

# A pairwise comparison between human and mouse genome

# Aligning One Locus



Exon 1  Exon 2  Exon 3  Exon 4
Intron 1  Intron 2  Intron 3
5'  3'

Splice site
GGTGAG

Splice site
CAG

Translation
Initiation
ATG

Branchpoint
CTG**A**C

Stop codon
TAG/TGA/TAA

# Example: a human/mouse ortholog



**Human Locus**

Alignment:

CDS

**Mouse Locus**

coding exons            intergenic regions
noncoding exons         strong alignment
introns                 weak alignment
intergenic regions

# Three Pairwise Alignments

---

# Paired HMM

Alignments correspond
1-to-1 with sequences
of states M, I, J



M
(+1, +1)

I
(+1, 0)

J
(0, +1)

```
-AGGCTATCACCTGACCTCCAGGCCGA--TGCCC---
TAG-CTATCAC--GACCGC-GGTCGATTTGCCCGACC
IMMJMMMMMMMJJMMMMMMMJMMMMMMMIIMMMMIII
```

# Let's score the transitions

Alignments correspond 1-to-1 with sequences of states M, I, J



```
-AGGCTATCACCTGACCTCCAGGCCGA--TGCCC---
TAG-CTATCAC--GACCGC-GGTCGATTTGCCCGACC
IMMJMMMMMMMJJMMMMMMJMMMMMMMIIMMMMMIII
```

# A Pair HMM for alignments

# Gene Finding



Exon 1, Exon 2, Exon 3, Exon 4, DNA, Intron 1, Intron 2, Intron 3, 5', 3'

Promoter TATA — Translation Initiation ATG — Splice site GGTGAG — Splice site CAG — Branchpoint CTG**A**C — Pyrimidine tract — Stop codon TAG/TGA/TAA — polyA signal

11

---

# Recall generalized HMM gene finder

TAAT ATGTCCACGG GTATTGAG CATTGTACACGGG GTATTGAG CATGTAA TGAA

12

# Generalized Pair-HMM gene finder

# Hierarchical state transition in pHMM : knowledge of structure

## Allowing for inserted exons: knowledge of structure

## Motif finding: recap

**Recap:**

- **Functional regions in sequences often occur as small, noisy, repeating subsequences**
  - **In DNA : transcription factor binding sites, transcription start sites, splicing signals, etc**
  - **In proteins : transmembrane domains, phosphorylation sites, signal peptides. etc**

- **Subsequence similarity and functional significance go hand in hand**



TF binding locations in *S. cerevisiae*

Habib et al, PLoS CompBiol Feb 08

# Problem formulation

- Models : Consensus, RegEx, Weight Matrix



Paired

- Supervised motif detection : Given a set of sequences, and a PWM $w$ of length $k$, find the maximum likelihood set of $k$-mers which correspond to the WMM.
  - Given TF binding specificities, can we find the TF binding sites ?
- De novo (unsupervised) motif detection : Given a set of sequences, find the most overrepresented set of $k$-mers and the corresponding WMM $w$
  - Given a set of genes with similar expression (putatively co-regulated), can we find the TF binding sites common to them and the specificity of the corresponding TF?
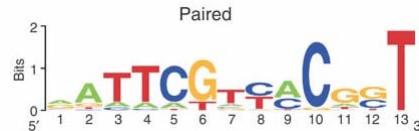
---

# Multi-species data pooling

- Simply pool together regulatory regions in related species
  - Bacterial DNA motifs, McGuire *et al* **Gen Res** 2000 & Gelfand *et al* **Nuc Ac Res** 2000
  - Hunchback TFBS in Drosophila species demonstration :

  CACCACTTTTTTATGCCGAGTTAAT          **D. melanogaster**

  GGTTTTTTCGATTCAATCGGTATA          **D. yakuba**

  AGTTAGCGTTTACCCTATTTTTTAC          **D. persimilis**

  GCATTTATCCTCTTTTTATAAGCTT          **D. mojavensis**

- What could be problematic with this approach ?

# Multi-species data pooling

- Biases analyses towards motifs in a bunch of closely related species – no explicit phylogenetic information used

- No distinction between paralogs and orthologs

- Variation in number of binding sites in orthologous CRMs much less than in CRMs of coregulated genes in same species
  - Signals in one species may be drowned out by cross species signals, or vice versa

# Orthologous sequence analysis

- What if the sequences are orthologous ?

| D. melanogaster | CTTTACGTATTTTAGTTATCGAGTTTATCTTCTGCTTGCTATCTCGCGC |
| D. yakuba | T--TACGTATTTTAGTTATCGAGTTTATCTTCTGCTTGCTATCTCGCGC |
| D. persimilis | GTTTACGTATTTTAGTTATCGAGTTTATCTTTCGCTT------TCTCGC |
| D. mojavensis | CTTTACGTATTTGAGTTATCAACTTTGT--TTTGCTT--TGCTTTTCGC |

- Functional regions like TFBSs are more conserved than background
- Phylogenetic dependencies between orthologs may be modelled to get more accurate scores for P( data | model )
- Paralogous sequence analysis also possible

# Chronology : from one to many

| Method | Single species | Multispecies |
|---|---|---|
| Combinatorial | Waterman (1985) | FootPrinter (2002) |
| LRT-like score threshold | Staden-PWM (1989) | rVista (2002) PhyloCon (2003) |
| Explicit mixture models | MEME (1994) | EMnEM (2004) PhyME (2004) CSMET (2008) |
| Gibbs Sampling | GMS (1993/1995) BioProspector (2001) AlignACE (2000) | Motif Sampler+ (2000) CompareProspector (2004) PhyloGibbs (2005) |
| HMM+ | Cister (2001) HMDM (2002) LOGOS (2003) BayCis (2008) | PhyloHMM (2004 –gene / 2008 - motif) PhyME (2004) MORPH (2008) CSMET (2008) |
| Ensemble models | EMD (2006) | - |

- First usage of term phylogenetic footprinting : Tagle *et al,* J Mol Biol 1988 : Regulatory regions of paralogous gamma and epsilon globin genes in Galago

- Google scholar hits for "phylogenetic footprinting"
  - 1988 – 1990 : 11 hits
  - 1991 – 2000 : 141 hits
  - 2001 – 2009 : 1850 hits
- Gibbs Sampling particularly easily adapted to incorporate phylogenetic footprinting

# FootPrinter: going combinatorial

# Footprinter

- For each node from the leaf up
  - Fill up a parsimony table W for each k-mer

- If the node n is a leaf
  - If the word is present in corr. string, $W[k]^{(n)} = 0$ else $W[k]^{(n)} = +\infty$
- Else
  - $W[k]^{(n)} = \sum\limits_{V \in child(n)} \min\limits_{t \in k\text{-mer}} (W[t]^{(v)} + d(t,k))$

- Choose most parsimonious k-mer(s) over the tree and alignment from table at root
- Time complexity = $O(n\, 4^{2k})$ for n species for fixed topology

© Eric Xing @ CMU, 2005-2009

---

# Motif Sampler + footprinting

- A simple idea that works reasonably well
- Wasserman et al, 2000
  - Given an input multiple alignment **A**
  - Compute a score for conservation across the alignment
  - Filter out all regions of the alignment with a score below a threshold *t*
  - Perform Gibbs Motif Sampling on the remaining alignment **A'**

| | |
|---|---|
| D. melanogaster | CTTTACGTATTTTAGTTATCGATTTTATTTTCTGCTTGCTATCTCGCGC |
| D. yakuba | T--TACGTATTTTAGTTATCGATTTTATTTTCTGCTTGCTATCTCGCGC |
| D. persimilis | GTTTACGTATTTTAGTTATCGATTTTAGTTTTCGCTT------TCTCGC |
| D. mojavensis | CTTTACGTATTTGAGTTATCAATTTTGGTTTTTGCTT--TGCTTTTCGC |

# Motif Sampler + footprinting

- A simple idea that works reasonably well
- Wasserman et al, 2000
  - Given an input multiple alignment **A**
  - Compute a score for conservation across the alignment
  - Filter out all regions of the alignment with a score below a threshold *t*
  - Perform Gibbs Motif Sampling on the remaining alignment **A'**

| | |
|---|---|
| D. melanogaster | CTTTACGTATTTTAGTTATCGATTTTATTTTCTGCTTGCTATCTCGCGC |
| D. yakuba | T--TACGTATTTTAGTTATCGATTTTATTTTCTGCTTGCTATCTCGCGC |
| D. persimilis | GTTTACGTATTTTAGTTATCGATTTTAGTTTTCGCTT------TCTCGC |
| D. mojavensis | CTTTACGTATTTGAGTTATCAATTTTGGTTTTTGCTT--TGCTTTTCGC |

---

# Motif Sampler + footprinting

- A simple idea that works reasonably well
- Wasserman et al, 2000
  - Given an input multiple alignment **A**
  - Compute a score for conservation across the alignment
  - Filter out all regions of the alignment with a score below a threshold *t*
  - Perform Gibbs Motif Sampling on the remaining alignment **A'**

| | |
|---|---|
| D. melanogaster | CTTTACGTATTTTAGTTATCGATTTTATTTTCTGCTTGCTATCTCGCGC |
| D. yakuba | T--TACGTATTTTAGTTATCGATTTTATTTTCTGCTTGCTATCTCGCGC |
| D. persimilis | GTTTACGTATTTTAGTTATCGATTTTAGTTTTCGCTT------TCTCGC |
| D. mojavensis | CTTTACGTATTTGAGTTATCAATTTTGGTTTTTGCTT--TGCTTTTCGC |

- Motifs sampled according to:

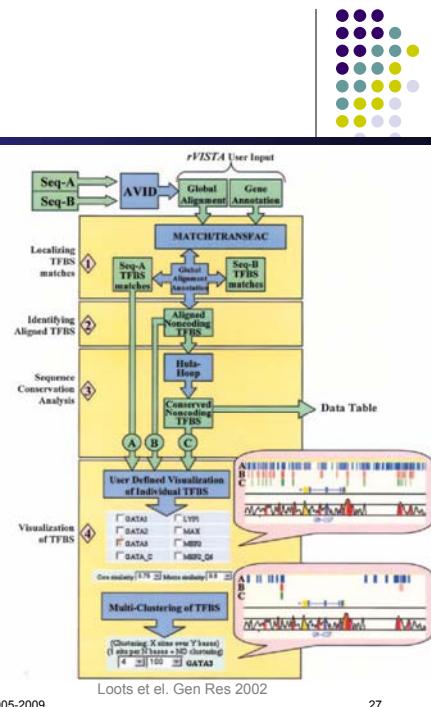$$\frac{\prod_{k=a}^{a+w} pm_{k-a,R_k}}{\prod_{k=a}^{a+w} p0_{k,R_k}}$$

**Score (likelihood) under PWM (motif) scenario**
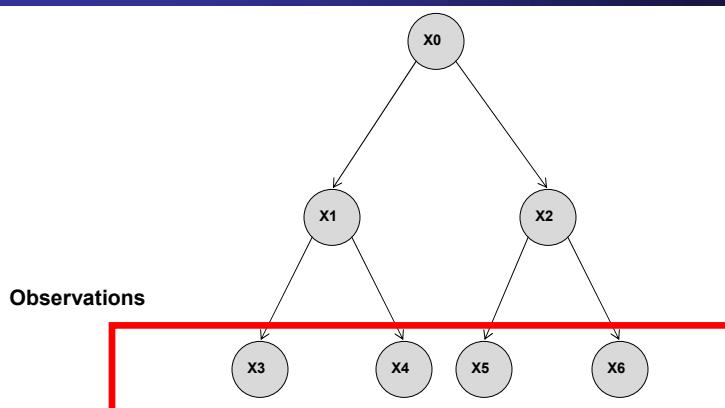
**Score (likelihood) under background scenario**

# rVISTA

- Select motifs with a PWM score greater than a threshold
- Screens for motifs above a certain threshold for nucleotide conservation
- Two step screening a common way to capture both overrepresentation & conservation
  - Loots et al, Gen Res 2002 [rVISTA]
  - Kellis et al, Nature 2003
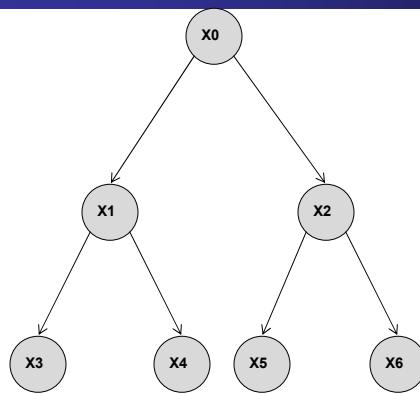  - Wang & Stormo, Bioinf 2003 [PhyloCon]



Loots et el. Gen Res 2002

---

# Phylogenetic model: recap



**Observations**

$P(D|M) = P(X0, X1, X2, X3, X4, X5, X6 \mid \tau, \beta, \theta, \pi)$

$= \sum_{X0, X1, X2} P(X3|X1; tree)\ P(X4|X1; tree)\ P(X5|X2; tree)\ P(X6|X2; tree)\ P(X2|X0;$

$tree)\ P(X1|X0; tree)\ P(X0| tree)$

# Phylogenetic model: recap

X0

X1

X2

X3    X4    X5    X6

•**Topology** – how the observations are "tied together": τ

•**Branch lengths** – the length for which the CTMP runs: β

•**Parameters of CTMP** – characterizing the substitution model: θ

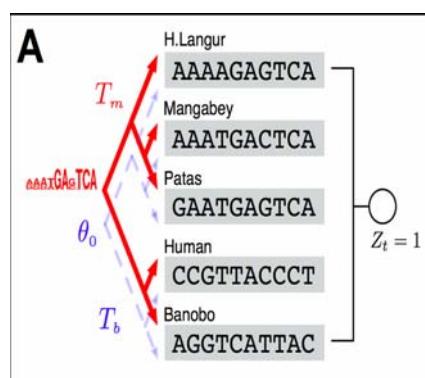•**Distribution at root** - maybe stationary dist of CTMP : π

P(D|M) = P(X0, X1, X2, X3, X4, X5, X6 | τ, β, θ, π)

= Σ$_{X0, X1, X2}$ P(X3|X1; tree) P(X4|X1; tree) P(X5|X2; tree) P(X6|X2; tree) P(X2|X0;

tree) P(X1|X0; tree) P(X0| tree)

---

# EMnEM: model based approaches

- Mixture model
- Each block of k-mer could be generated from a background model with probability **1-π$_m$** or from a motif model with probability π$_m$
- Bernoulli draw for the mixture indicator
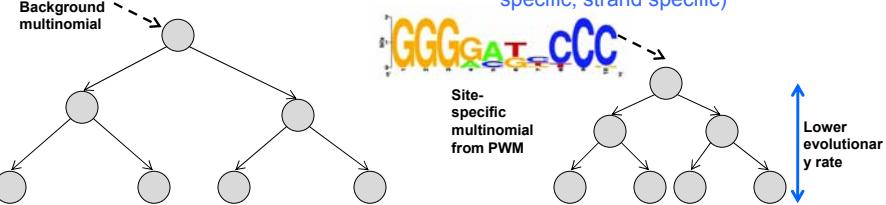- In the spirit of MEME

A

$T_m$

H.Langur
AAAAGAGTCA

Mangabey
AAATGACTCA

Patas
GAATGAGTCA

AAAxGAgTCA

$θ_0$

Human
CCGTTACCCT

$Z_t = 1$

$T_b$

Banobo
AGGTCATTAC

# Function specific phylogenetic models

- Background model $T_b$
  - Topology invariant unless evidence otherwise
  - Substitution matrix invariant unless evidence otherwise
  - Branch lengths longer than functional sites
  - Root distribution : background frequency

  **Background multinomial**

- Motif site-specific, strand-specific model $T_{m, k, +/-}$
  - Topology invariant unless evidence otherwise
  - Substitution matrix invariant unless evidence otherwise
  - Branch lengths shorter than background sites
  - Root distribution : from PWM (site specific, strand specific)

  **Site-specific multinomial from PWM**

  GGGGAT CCC

  **Lower evolutionary rate**

© Eric Xing @ CMU, 2005-2009

---

# EMnEM: Expectation maximization on mixtures of phylogenies

A
X   Y

$$L = \prod_{i=0}^{N-w} \sum_{m_i} \boxed{p(m_i)} \prod_{k=i}^{i+w-1} \sum_{b=0}^{3} p(X_k, Y_k \mid \boxed{A_{kb}}, m_i) \boxed{p(A_{kb} \mid m_i)}$$

Mixture parameter    Ancestral base    PWM

- E-step :
  - Mixture parameter:
  $$\langle m_i \rangle = p(m_i | X, Y) = \frac{p(m_i)p(X,Y|m_i)}{p(X,Y)}$$

  $$p(X,Y|m_i) = \prod_{k=i}^{i+w-1} \sum_{b=0}^{3} p(X_k, Y_k | A_{kb}, m_i) p(A_{kb}|m_i)$$

  $$p(X,Y) = \sum_{m_i} p(X,Y|m_i)p(m_i)$$

  - Ancestral nucleotide:
  $$\langle A_{ib} \rangle = p(A_{ib}|X_i, Y_i) = \sum_{m_i} p(A_{ib}|X_i, Y_i, m_i)p(m_i) = \sum_{m_i} \frac{p(A_{ib})p(X_i, Y_i|A_{ib}, m_i)}{p(X_i, Y_i|m_i)} p(m_i)$$

- M-step :
  $$\langle \ln L_c \rangle = \sum_{i=0}^{N-w} \sum_{m_i} \langle m_i \rangle \left[ \ln \pi_m + \sum_{k=i}^{i+w-1} \sum_{b=0}^{3} \langle A_{kb} \rangle \left( \ln p(X_k, Y_k \mid A_{kb}, m_i) + \ln f_{mkb} \right) \right]$$
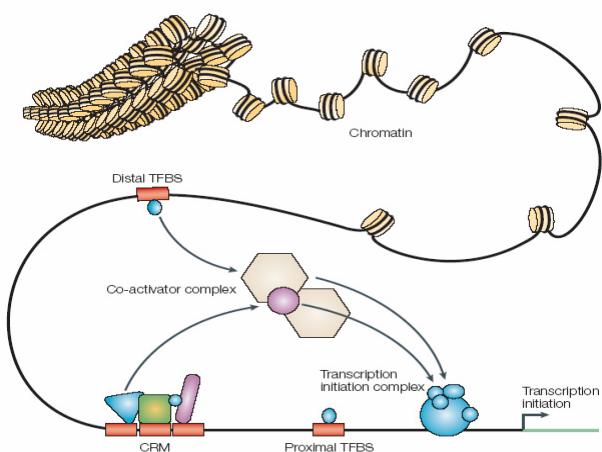
  $$\frac{\partial \langle \ln L_c \rangle}{\partial \pi_m} = 0, \quad \frac{\partial \langle \ln L_c \rangle}{\partial f_{mkb}} = 0 \text{ and } \frac{\partial \langle \ln L_c \rangle}{\partial x_{mk}} = 0$$

# CRM: putting the pieces together

- HMMs !



Courtesy: Simone Scalabria
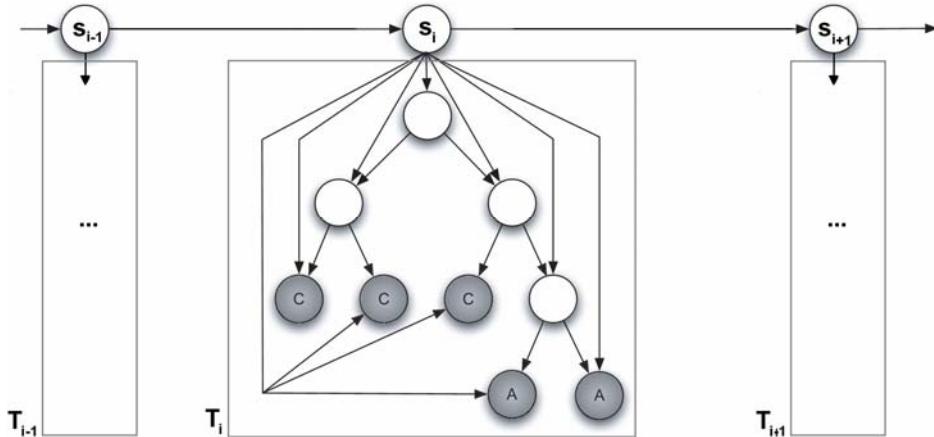
# A vanilla HMM …



ACATTGCCATACCATAATCCTTAATT…

- Emits a symbol at every discrete step
- A run of the HMM outputs a sequence
- PhyloHMM outputs a vector of characters
- A run of the PhyloHMM outputs a multiple sequence alignment

# Phylo-HMM

- The emission vector $O_i$ is shaded in gray



Courtesy: McAuliffe

35

# Phylo-HMM

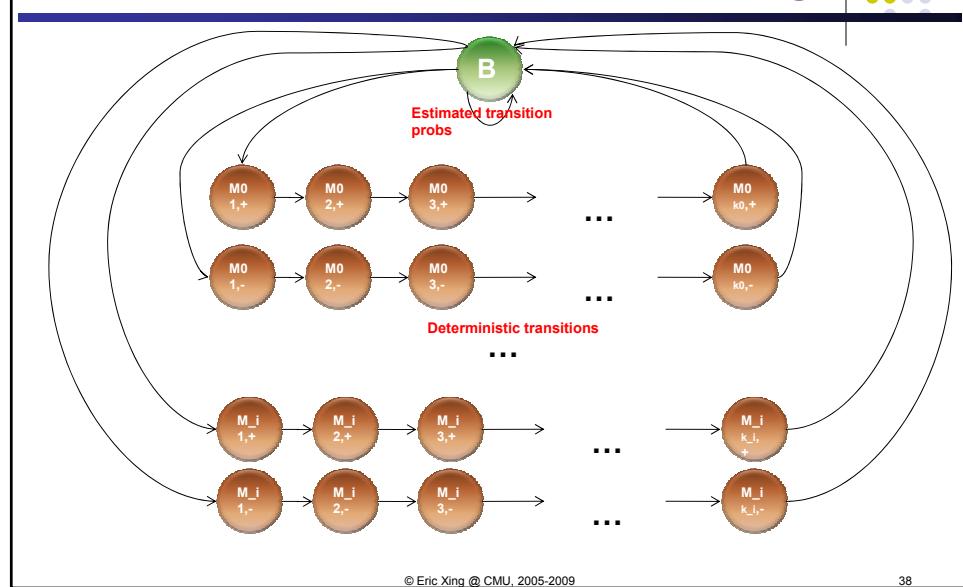- Emits a vector at each step, generates alignment in a run



```
CATTGCAT...
CAATGAAT...
CATTGTTT...
AATTATTT...
AGTTAGTT...
```

36

# State space comparisons



Standard HMM

$X = \text{TAACGGCAGA}\ldots$

Phylo HMM

Evolutionary parameters
Associated with each
state to help calculate
emission probability

$X = \begin{array}{l} \text{TAACGGCAGA}\ldots \\ \text{TTAGGCAAGG}\ldots \\ \text{AAGGCGCCGA}\ldots \end{array}$

Courtesy: Siepel

---

# HMM state space for motif finding



**B**

Estimated transition
probs

M0 1,+    M0 2,+    M0 3,+    ...    M0 k0,+

M0 1,-    M0 2,-    M0 3,-    ...    M0 k0,-

Deterministic transitions
...

M_i 1,+    M_i 2,+    M_i 3,+    ...    M_i k_l, +

M_i 1,-    M_i 2,-    M_i 3,-    ...    M_i k_l,-

# More realism, more parameters

---

# Phylo-HMM

- A normal HMM, except the emission probabilities are a multinomial distribution over the space of [**ATGC**]$^n$, n being the number of sequences in the alignment

- $4^n$ emission probabilities can be pre computed

- But usually calculated on the fly using Felsenstein's Pruning Algorithm - a special case of the GM Belief Propagation Algorithm on trees
  - Siepel & Haussler, RECOMB 2004, for gene finding
  - Ray et al. PLoS CompBio 2008, adapted for motif finding

# Analogy with HMM

**Set of states** ⟶ $(S, \psi, N, n, T, b)$ ⟵ **Initial probs**

**Transition probs**

**Set of associated phylogenetic models**  **Set of alphabets**  **No of species**

- Emission probability $= P(O_i \mid S_i = s)$

  $= P(O_i \mid \text{phylogenetic model}_s)$

  $= P(\text{Alignment column}_i \mid \psi_s)$

  => Calculate using the Pruning Algorithm

- Apply standard Viterbi (maximize joint) or posterior decoding on the Forward-Backward matrix

- Baum-Welch algorithm (E-M) for unsupervised settings

- Exactly analogous to single species motif finding
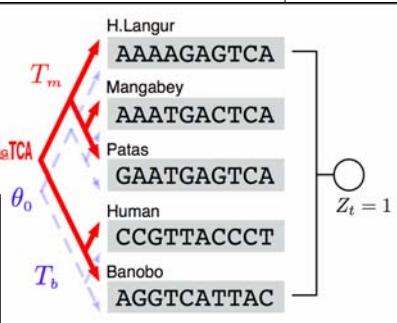
# Missing motifs

# Functional turnover

# CSMET : Phylogenetic mixtures of phylogenies

- Mixture models for evolutionary model selection (EMnEM) →
- Bernoulli draw for mixture variable



- What if the mixture variables are phylogenetically related ?
- Output of a "functional" phylogeny

# CSMET

- CSMET-HMM :
  - An HMM with emission vector [A,T,G,C]*
  - Each vector is the output of a generative process involving a mixture of trees
  - Mixture indicator variables themselves generated by a phylogeny
  - Similar scheme to PhyloHMM, except for calculating emission probs
- Schematic of generation and ML inference
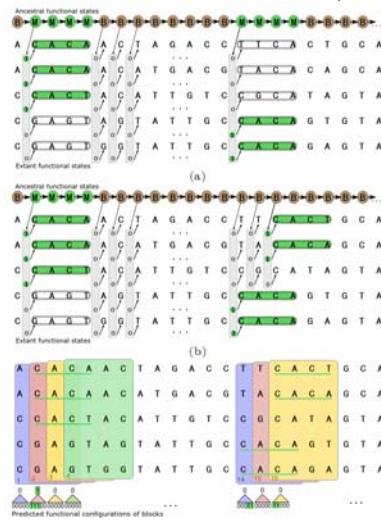- CSMET-HMM : **Corr character sets**

$$(S, TN, TF, N, F, n, T, b)$$

**Set of nucleotide phylogenetic models corr to each annotation**

**Set of functional phylogenetic models corr to each state in S**
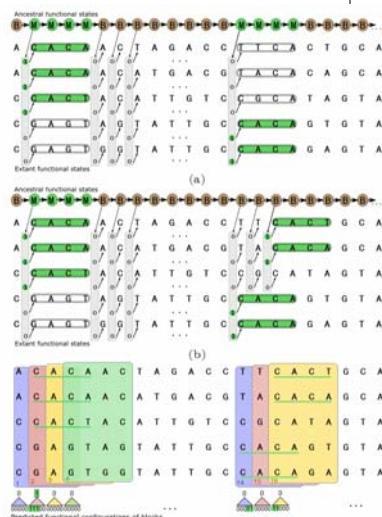
---

# CSMET

- To calculate emission probabilities:
  - Calculate likelihoods of nucleotide data for each subtree of the nucleotide phylogeny
  - Calculate likelihood of functional indicators for the functional phylogeny
  - Putting the likelihoods together using conditional independences
  - Marginalize out hidden variables
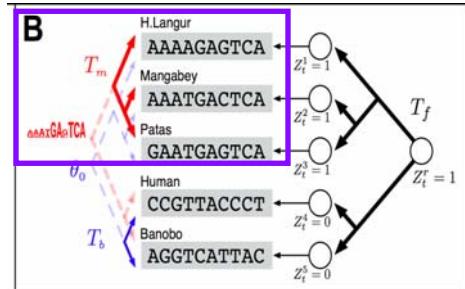- The rest would be analogous to an HMM !

# Likelihoods on partial phylogenies

- Marginalize out observed nucleotides present in parts of the phylogeny we are not interested in
- Turns out to be equivalent to calculating the likelihood of the data on the subtree !



$$P(A_l' | T'^{(l)}) = \sum_{A_l''} P(A_l', A_l'' | T^{(l)}) =$$

$$\sum_{A_l''} \sum_{\mathbf{v}_{1:K'}} P(\mathbf{V}_{1:K'} = \mathbf{v}_{1:K'}, \mathbf{V} = A_l'', \mathbf{V}' = A_l')$$
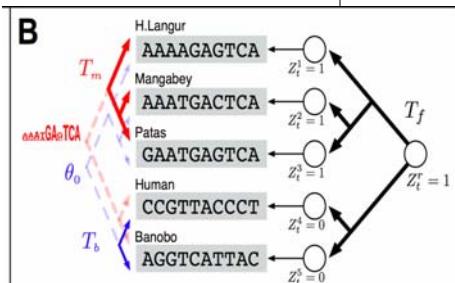
---

# CSMET: toolkit for calculations

- Calculating likelihoods on the nucleotide phylogeny and functional phylogeny

**Nucleotide phylogeny :**
**F84 model – simplest arbitrary stationary distribution**



$$Q_N = \begin{pmatrix} * & (1+\kappa/\pi_Y)\iota\pi_C & \iota\pi_A & \iota\pi_G \\ (1+\kappa/\pi_Y)\iota\pi_T & * & \iota\pi_A & \iota\pi_G \\ \iota\pi_T & \iota\pi_C & * & (1+\kappa/\pi_R)\iota\pi_G \\ \iota\pi_T & \iota\pi_C & (1+\kappa/\pi_R)\iota\pi_A & * \end{pmatrix}$$

**Functional phylogeny**
**Jukes Cantor model**

$$P_F = \begin{pmatrix} \frac{1}{2} + \frac{1}{2}e^{-2\beta} & \frac{1}{2} - \frac{1}{2}e^{-2\beta} \\ \frac{1}{2} - \frac{1}{2}e^{-2\beta} & \frac{1}{2} + \frac{1}{2}e^{-2\beta} \end{pmatrix}$$

**Likelihoods on partial phylogenies**

$$P(A_l' | T'^{(l)}) = \sum_{A_l''} P(A_l', A_l'' | T^{(l)}) =$$

$$\sum_{A_l''} \sum_{\mathbf{v}_{1:K'}} P(\mathbf{V}_{1:K'} = \mathbf{v}_{1:K'}, \mathbf{V} = A_l'', \mathbf{V}' = A_l')$$
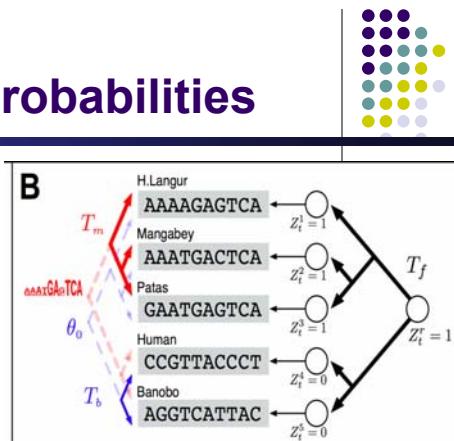
# CSMET : emission probabilities

- Emission prob : Prob of block surrounding particular aligned site

- Again, analogous to an HMM, with one twist : $Z_i$ s not observed

**Joint Probability for an instantiated block**

$$P\left(\mathbf{A}_t, \mathbf{z}_t, z_t^r\right) = P\left(\mathbf{A}_t \mid Z_t = \mathbf{z}_t, T_m, T_b\right) P\left(Z_t = \mathbf{z}_t \mid Z_t^r = z_t^r, T_f\right)$$

$$P\left(Z_t^r = z_t^r\right) = P\left(\mathbf{A}'_t \mid T'_m\right) P\left(\mathbf{A}''_t \mid T'_b\right) P\left(\mathbf{z}_t \mid z_t^r, T_a\right) P\left(z_t^r\right).$$

**Conditional probability for the block**

$$P\left(\mathbf{A}_t \mid Z_t = \mathbf{z}_t, T_m, T_b\right) = P\left(\mathbf{A}'_t \mid T'_m\right) P\left(\mathbf{A}''_t \mid T'_b\right) =$$

$$\prod_{l=1}^{L} P\left(A'_l(t) \mid T'_m(l)\right) P\left(A''_l(t) \mid T'_b\right).$$



**Emission probability for the block (marginalized)**

$$P\left(\mathbf{A}_t \mid z_t^r\right) = \sum_{z_t} P\left(\mathbf{A}_t, \mathbf{z}_t \mid z_t^r\right) = \sum_{z_t} P\left(\mathbf{A}'_t(\mathbf{z}_t) \mid T'_m(\mathbf{z}_t)\right)$$

$$P\left(\mathbf{A}''_t(\mathbf{z}_t) \mid T'_b(\mathbf{z}_t)\right) P\left(\mathbf{z}_t \mid T_a, z_t^r\right),$$

---

# Chronology : aspects of footprinting

- Footprinting + Gibbs Sampling
  - **Motif Sampler+ : 2 species alignment**
  - **CompareProspector : Pairwise alignment**
  - **PhyloGibbs : Multiple alignment**

- Footprinting + HMM
  - **PhyloHMM : Emission of HMM generated by a CTMP phylogenetic tree, no tolerance for functional turnover**
  - **PhyME**
  - **CSMET : Emission of HMM generated by a mixture of CTMP phylogenetic trees, explicit tolerance for functional turnover**

- **Footprinting + alignment**
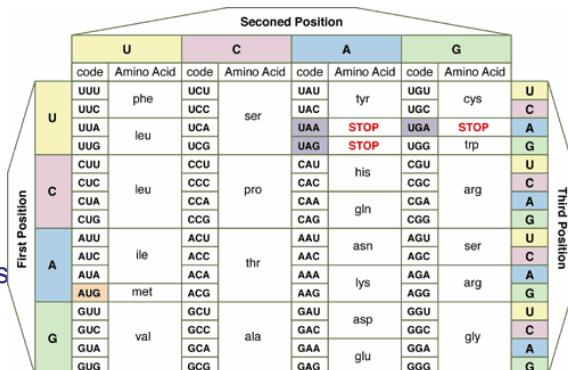  - **OrthoMEME**
  - **MORPH**

# Can we do even better ?

- Footprinting improves with
  - More knowledge about the functional component we are searching for : what to look for in a single species
  - More knowledge about how it evolves : what to look for in related species
- We know a lot about both aspects for protein coding regions, or genes
- Initial footprinting algorithms on genes and proteins

# Evolution of codons

- Genes evolve at a level of higher granularity
  - Nucleotide
  - Codon
- HMM states corresponding to codons
- How to choose priors for transition probabilities ?



Courtesy: Bioephemera.com

# Incorporating evolutionary processes

- Selection
- Transition probabilities can reflect
  - Synonymous transitions more frequent than non synonymous ones
  - How much more frequent ?
  - Selection parameters estimated from data



Courtesy: Bioephemera.com

---

# Summary

- Use genomic representation of functional component
- Use evolutionary models of functional component
- Can be used for non-sequence data too :
  - Gene regulatory network
  - Expression levels : microarray data

# Acknowledgments

- **Serafim Batzoglou**: for some of the slides adapted or modified from his lecture slides at Stanford University
- **Lior Pachter'**: for some of the slides modified from his lectures at UC Berkeley
- Acknowledgements for some graphics on corresponding slides