Computational Genomics

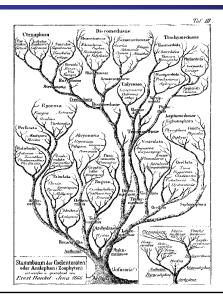
Molecular Evolution: Phylogenetic trees

Eric Xing
Lecture 14, March 6, 2007



Reading: DTW book, Chap 12 DEKM book, Chap 7, 8

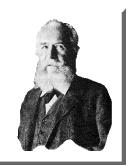
Phylogeny





trees for most known groups of living

organisms.



Ernst Haeckel (1834-1919)

Phylogenetic Inference



• Given a multiple alignment, how do we construct the tree?

A - GCTTGTCCGTTACGAT

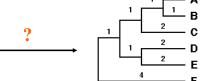
B - ACTTGTCTGTTACGAT

- ACTTGTCCGAAACGAT

- ACTTGACCGTTTCCTT

- AGATGACCGTTTCGAT

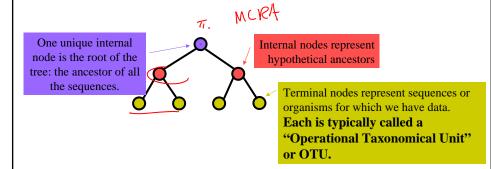
- ACTACACCCTTATGAG

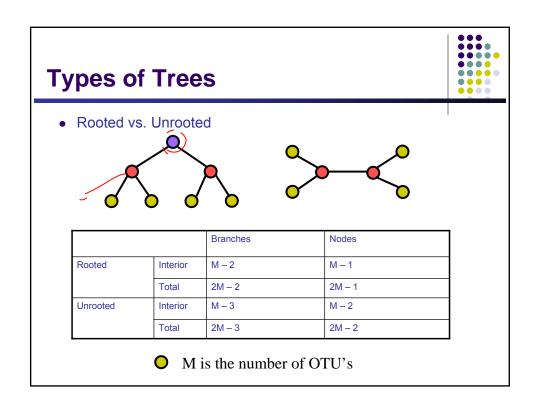


What Is a Tree?



- A tree is a mathematical structure which represents a model of an actual evolutionary history of a group of sequences or organisms. T= {7.7.5 M}
 - In other words, it is an evolutionary hypothesis.
- A tree consists of nodes connected by branches.



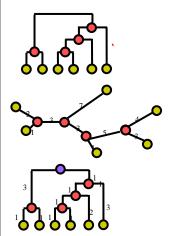


The number of rooted and unrooted trees: **Possible Number of** Number Rooted Unrooted of OTU's trees trees

More different kinds of trees



 Different kinds of trees can be used to depict different aspects of evolutionary history



- 1. Cladogram: simply shows relative order of common ancestry
- 2. Additive trees:
 a cladogram with branch lengths,
 also called phylograms and metric trees
- 3. Ultrametric trees:
 (dendograms) special kind of additive tree in which
 the tips of the trees are all equidistant from the root

Phylogeny methods



Basic principles:

- Degree of sequence difference is proportional to length of independent sequence evolution
- Only use positions where alignment is pretty certain avoid areas with (too many) gaps

Major methods:

Clustering methods	UPGMA Neighbor-joining WPGMA			
†	Single linkage			
l .	Complete linkage			
Objective	Least-squares distance			
criterion-	Maximum parsimony			
based	Minimum evolution			
methods	Maximum likelihood			

UPGMA



- Construction of a distance tree using clustering with the Unweighted Pair Group Method with Arithmatic Mean (UPGMA)
- First, construct a distance matrix:

A - GCTTGTCCGTTACGAT

B - ACTTGTCTGTTACGAT

C - ACTTGTCCGAAACGAT

D - ACTTGACCGTTTCCTT

E - AGATGACCGTTTCGAT

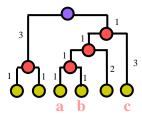
F - ACTACACCCTTATGAG

	Α	В	С	D	Е
В	2				
С	4	4			
D	6	6	6		
Е	6	6	6	4	
F	8	8	8	8	8

From http://www.icp.ucl.ac.be/~opperd/private/upgma.html

Ultrametric Trees







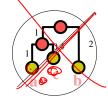
• Non-negativity: $d(a,b) \ge 0$

• Distinctness: d(a,b) = 0 if and only if a = b

• Symmetry: d(a,b) = d(b,a)

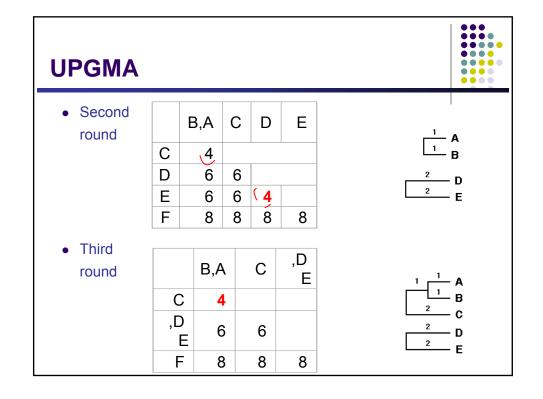
• Triangle Inequality: $d(a,c) \le d(a,b) + d(b,c)$

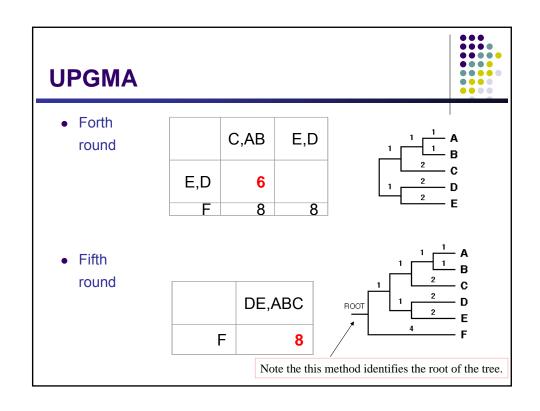


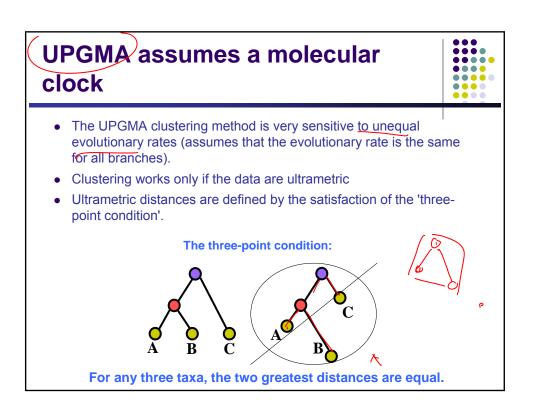


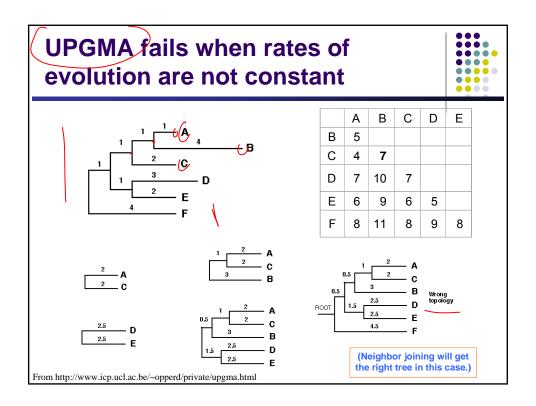
- Ultrametric must obey one additional rule:
 - Three point condition: $d(a,b) \le max(d(a,c), d(b,c))$

UPGMA First round С Ε D Α В dist(A,B),C = (distAC + distBC) / 2 = 4dist(A,B),D = (distAD + distBD) / 2 = 62 В dist(A,B),E = (distAE + distBE) / 2 = 6dist(A,B),F = (distAF + distBF) / 2 = 8C 4 4 6 6 D 6 Ε 6 6 6 4 B,A) С D Ε F 8 8 8 8 8 С 4 D 6 6 Choose the most similar pair, Ε 6 6 4 cluster them together and calculate the new distance matrix. F 8 8 8 8









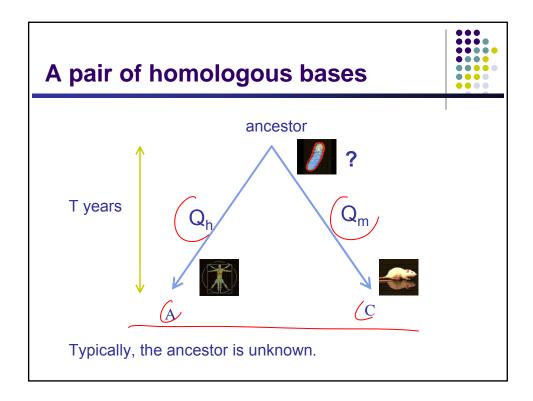
Phylogeny based upon the molecular clock



- Evidence for a human mitochondrial origin in Africa: African sequence diversity is twice as large as that of non-African
- Gyllensten and colleagues estimate that the divergence of Africans and non-Africans occurred 52,000 to 28,000 years ago.



Ingman, M., Kaessmann, H., Pääbo, S. & Gyllensten, U. (2000) Nature 408: 708-713.



How does sequence variation arise?



- Mutation:
 - (a) Inherent: DNA replication errors are not always corrected.
 - (b) External: exposure to chemicals and radiation.
- Selection: Deleterious mutations are removed quickly.
 Neutral and rarely, advantageous mutations, are tolerated and stick around.
- **Fixation**: It takes time for a new variant to be established (having a stable frequency) in a population.

Modeling DNA base substitution



- Strictly speaking, only applicable to regions undergoing little selection.
- Standard assumptions (sometimes weakened)
 - 1. Site independence.
 - 2. Site homogeneity.
 - 3. Markovian: given current base, future substitutions independent of past.
 - 4. Temporal homogeneity: stationary Markov chain.

More assumptions



- $Q_h \neq s_h Q$ and $Q_m = s_m Q$, for some positive s_h , s_m , and a rate matrix Q.
- The ancestor is sampled from the stationary distribution π of Q.
- Q is **reversible**: for $a, b, t \ge 0$ $\pi(a)P(t,a,b) = P(t,b,a)\pi(b),$ (detailed balance).





The stationary distribution

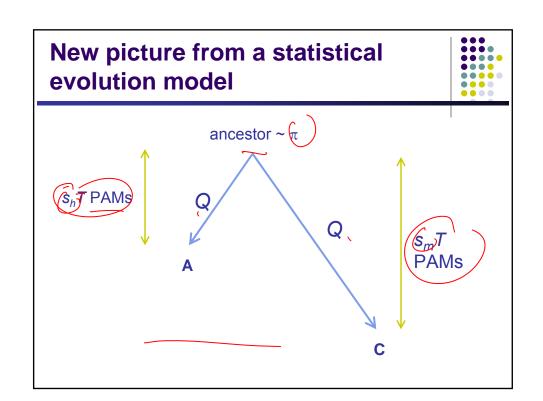


A probability distribution π on {A,C,G,T} is a stationary distribution of the Markov chain with transition probability matrix P = P(i,j), if for all j,

$$\sum_i \, \pi(i) \; P(i,j) = \pi(j).$$

- **Exercise**. Given any initial distribution, the distribution at time t of a chain with transition matrix P converges to π as $t \to \infty$. Thus, π is also called an **equilibrium** distribution.
- **Exercise**. For the Jukes-Cantor and Kimura models, the uniform distribution is stationary. (Hint: diagonalize their infinitesimal rate matrices.)

We often assume that the ancestor sequence is i.i.d π .



The Jukes-Cantor adjustment



 Assume that the common ancestor has A, G, C or T with probability 1/4.

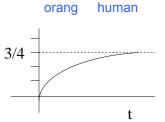


• Then the chance of the nt differing

$$(p_{\neq} = 3/4 \times (1 - e^{(8d)})$$

= 3/4 \times (1 - e^{-4k/3}), since $k = 2 \times 3\alpha t$

When k = .01, described as 1 PAM



Joint probability of A and C



• Under the model in the previous slides, the joint probability is

$$(p(A,C) \Rightarrow \sum_{a} \pi(a) p(A \mid S_{h}T, Q, a) p(C \mid S_{m}T, Q, a)$$

$$= \sum_{a} \pi(A) p(a \mid S_{h}T, Q, A) p(C \mid S_{m}T, Q, a)$$

$$= \pi(A) p(C \mid S_{h}T + S_{m}T, Q, A)$$

$$= (t, A, C)$$



- where t = s_hT+ s_mT is the (evolutionary) distance between A and C.
 Note that s_h, s_m and T are not identifiable.
- The matrix *F*(t) is symmetric. It is equally valid to view A as the ancestor of C or vice versa.

Estimating the evolutionary distance between two sequences



- Suppose two aligned protein sequences $a_1...a_n$ and $b_1...b_n$ are separated by t PAMs.
- Under a reversible substitution model that is IID across sites, the likelihood of t is

$$L(t) = p(a_1 \dots a_n, b_1 \dots b_n \mid \text{model})$$

$$= \prod_k (t, a_k, b_k)$$

$$= \prod_{a,b} (t, a, b)^{(a,b)}$$

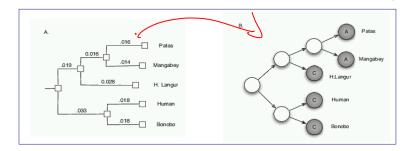


- where $c(a,b) = \# \{k : a_k = a, b_k = b\}.$
- Maximizing this quantity gives the maximum likelihood estimate of *t*. This generalizes the distance correction with Jukes-Cantor.

Phylogeny





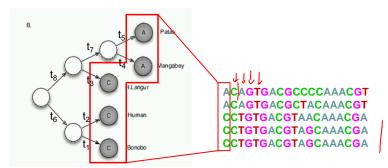


- The shaded nodes represent the observed nucleotides at a given site for a set of organisms
- The unshaded nodes represent putative ancestral nucleotides
- Transitions between nodes capture the dynamic of evolution

Likelihood methods



• A tree, with branch lengths, and the data at a single site.



• Since the sites evolve independently on the same tree,

$$L = P(D | \overline{\mathcal{I}}) = \prod_{i=1}^{m} P(D^{(i)} | T)$$

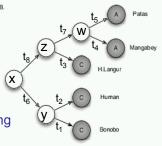
Likelihood at one site on a tree



 We can compute this by summing over all assignments of states x, y, z and w to the interior nodes:

$$P(D^{(i)} | T) = \sum_{x} \sum_{y} \sum_{z} \sum_{z} P(A, A, C, C, C, x, y, z, w | T)$$

 Due to the Markov property of the tree, we can factorize the complete likelihood according to the tree topology:



$$P(A, A, C, C, C, x, y, z, w | T) = P(x) P(y | x, t_6) P(C | y, t_1) P(C | y, t_2)$$

$$P(z | x, t_8) P(C | y, t_3)$$

$$P(w | z, t_7) P(A | y, t_4) P(A | y, t_5)$$

• Summing this up, there are 256 terms in this case!

Getting a recursive algorithm



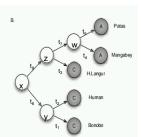
• when we move the summation signs as far right as possible:

$$P(D^{(i)} | T) = \sum_{x} \sum_{y} \sum_{z} \sum_{w} P(A, A, C, C, C, x, y, z, w | T) = \sum_{x} P(x)$$

$$\left(\sum_{y} P(y | x, t_{6}) \quad P(C | y, t_{1}) P(C | y, t_{2}) \right)$$

$$\left(\sum_{z} P(z | x, t_{8}) \quad P(C | z, t_{3}) \right)$$

$$\left(\sum_{w} P(w | z, t_{7}) P(A | w, t_{4}) P(A | w, t_{5}) \right)$$



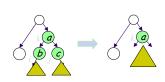
Felsenstein's Pruning Algorithm



• To calculate $P(x_1, x_2, ..., x_N | T, t)$

Initialization:

Set
$$k = 2N - 1$$



UL

Recursion: Compute $P(L_k \mid a)$ for all $a \in \Sigma$

If k is a leaf node:

Set
$$P(L_k | a) = 1(a = x_k)$$



If k is not a leaf node:

- 1. Compute $P(L_i \mid b)$, $P(L_j \mid c)$ for all b and c, for daughter nodes i, j
- 2. Set $P(L_k \mid a) = \sum_{b, c} P(b \mid a, t_i) P(L_i \mid b) P(c \mid a, t_j) P(L_j \mid c)$

Termination:

Likelihood at this column =
$$P(x_1, x_2, ..., x_N | T, t) = \sum_a P(L_{2N-1} | a)P(a)$$

• This algorithm can easily handle Ambiguity and error in the sequences (how?)

Finding the ML tree

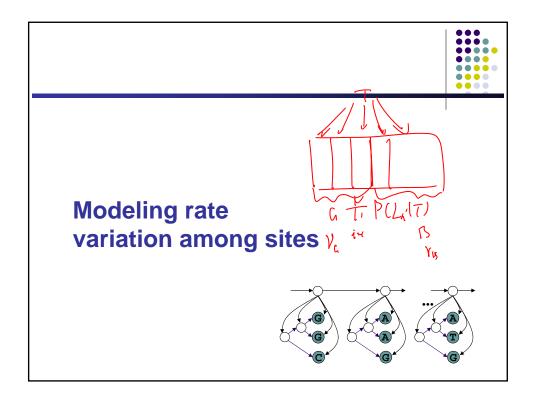


- So far I have just talked about the computation of the likelihood for one tree with branch lengths known.
- To find a ML tree, we must search the space of tree topologies, and for each one examined, we need to optimize the branch lengths to maximize the likelihood.

Bayesian phylogeny methods



- Bayesian inference has been applied to inferring phylogenies (Rannala and Yang, 1996; Mau and Larget, 1997; Li, Pearl and Doss, 2000).
 - All use a prior distribution on trees. The prior has enough influence on the result that its reasonableness should be a major concern. In particular, the depth of the tree may be seriously affected by the distribution of depths in the prior.
 - All use Markov Chain Monte Carlo (MCMC) methods. They sample from the posterior distribution.
 - When these methods make sense they not only get you a point estimate of the phylogeny, they get you a distribution of possible phylogenies.



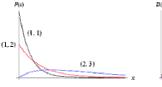
A model of variation in evolutionary rates among sites

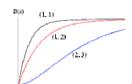


 The basic idea is that the rate at each site is drawn independently from a distribution of rates. The most widely used choice is the Gamma distribution, which has density function:

$$f(r) = \frac{\lambda^{\alpha} r^{\alpha - 1} e^{-\lambda r}}{\Gamma(\alpha)} = \frac{r^{\alpha - 1} e^{-r/\theta}}{\Gamma(\alpha) \theta^{\alpha}}$$

• Gamma distributions (α, θ)





Unrealistic aspects of the model:



- There is no reason, aside from mathematical convenience, to assume that the Gamma is the right distribution.
- A common variation is to assume there is a separate probability f₀ of having rate 0.
- Rates at different sites appear to be correlated, which this model does not allow.
- Rates are not constant throughout evolution, they change with time.

Rates varying among sites



If L⁽ⁱ⁾(r_i) is the likelihood of the tree for site i given that the rate of evolution at site i is r_i, we can integrate this over a gamma density:

$$\underline{\underline{L}}^{(i)} = \int_0^\infty (f(\mathbf{r}_i; \alpha) \underline{\underline{L}}^{(i)}(\mathbf{r}_i) d\mathbf{r}_i$$

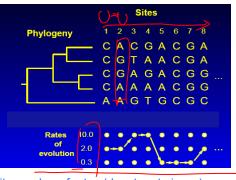
• so that the overall likelihood is

$$L = \prod_{i=1}^{m} \left[\int_{0}^{\infty} f(\mathbf{r}_{i}; \alpha) L^{(i)}(\mathbf{r}_{i}) d\mathbf{r}_{i} \right]$$

• Unfortunately these integrals cannot be evaluated for trees with more than a few tips as the quantities $L^{(i)}(r_i)$ becomes complicated.

Modeling rate variation among sites

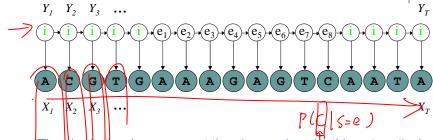




- There are a finite number of rates (denote rate i as r_i).
- There are probabilities p_i of a site having rate i.
- A process not visible to us ("hidden") assigns rates to sites.
- The probability of our seeing some data are to be obtained by summing over all possible combinations of rates, weighting appropriately by their probabilities of occurrence.

Rocall the HMM





- The shaded nodes represent the observed nucleotides at particular sites of an-organism's genome
- For discrete Y_i, widely used in computational biology to represent segments of sequences
 - gene finders and motif finders
 - profile models of protein domains
 - models of secondary structure

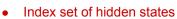
Definition (of HMM)



Observation space

Alphabetic set: $C = \{c_1, c_2\}$ Euclidean space: R^{σ}

 $C = \{c_1, c_2, \dots, c_K\} \qquad (y_1) \longrightarrow (y_2) \longrightarrow$



$$I = \{1, 2, \cdots, M\}$$



Transition probabilities between any two states

$$p(y_t^j = 1 | y_{t-1}^i = 1) = a_{i,j},$$
or
$$p(y_t | y_{t-1}^i = 1) \sim \text{Multinomial}(a_{i,1}, a_{i,2}, \dots, a_{i,M}), \forall i \in I.$$

Start probabilities

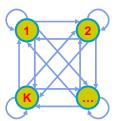
$$p(y_1) \sim \text{Multinomial}(\pi_1, \pi_2, ..., \pi_M)$$
.

Emission probabilities associated with each state

$$p(x_t \mid y_t^i = 1) \sim \text{Multinomial}(b_{i,1}, b_{i,2}, \dots, b_{i,K}), \forall i \in \mathbb{I}.$$

or in general:

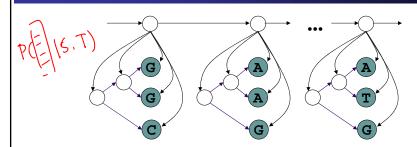
$$p(\mathbf{x}_t | \mathbf{y}_t^i = 1) \sim f(\cdot | \theta_i), \forall i \in \mathbb{I}.$$



State automata

Hidden Markov Phylogeny





- Replacing the standard emission model with a tree
 - A process not visible to us (.hidden") assigns rates to sites. It is a Markov process working along the sequence.
 - For example it might have transition probability Prob () of changing to rate / in the next site, given that it is at rate / in this site.
- These are the most widely used models allowing rate variation to be correlated along the sequence.

The Forward Algorithm



 $\alpha_1^k = P(x_1, y_1^k = 1)$

• We can compute α_t^k for all k, t, using dynamic programming!

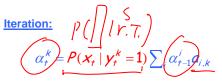
Initialization:

$$\alpha_{1}^{k} = P(x_{1} | y_{1}^{k} = 1)P(y_{1}^{k} = 1)$$

$$= P(x_{1} | y_{1}^{k} = 1)P(y_{1}^{k} = 1)$$

$$= P(x_{1} | y_{1}^{k} = 1)\pi_{k}$$

$$= P(x_{1} | y_{1}^{k} = 1)\pi_{k}$$



Termination:

$$P(\mathbf{x}) = \sum_{k} \alpha_{T}^{k}$$

The Backward Algorithm



• We can compute β_t^k for all k, t, using dynamic programming!

Initialization:

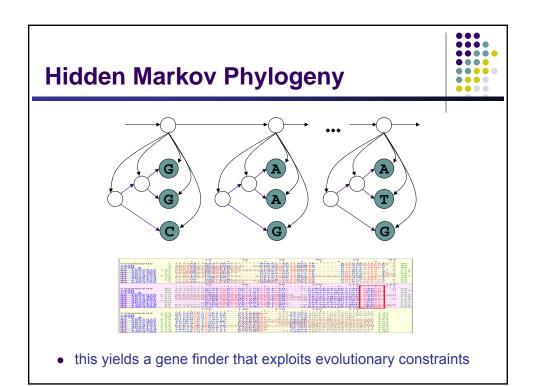
$$\beta_T^k = 1, \ \forall k$$

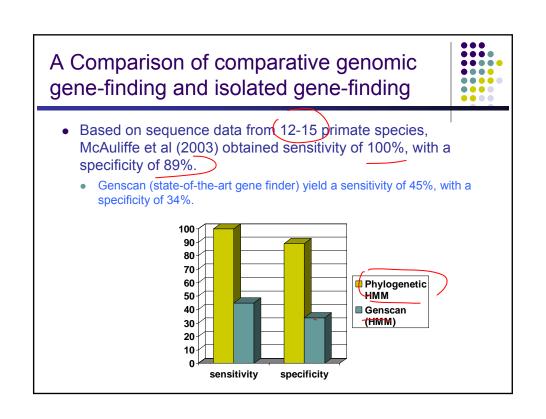
Iteration:

$$\beta_t^k = \sum_i a_{k,i} P(\mathbf{x}_{t+1} | \mathbf{y}_{t+1}^k = 1) \beta_{t+1}^i$$

Termination:

$$P(\mathbf{x}) = \sum_{k} \alpha_1^k \beta_1^k$$





Open questions (philosophical)



Observation:

- Finding a good phylogeny will help in finding the genes.
- Finding the genes will help to find biologically meaningful phylogenetic trees

Which came first, the chicken or the egg?

Open questions (technical)



- How to learn a phylogeny (topology and transition prob.)?
- Should different site use the same phylogeny? Functionspecific phylogeny?
- Other evolutionary events: duplication, rearrangement, lateral transfer, etc.

Acknowledgments



- Terry Speed: for some of the slides modified from his lectures at UC Berkeley
- **Phil Green** and **Joe Felsenstein**: for some of the slides modified from his lectures at Univ. of Washington
- Itai Yanai: for some of the slides modified from his lectures at Weizmann Institute of Science