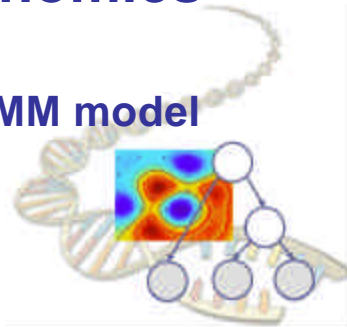


# Molecular Evolution and Comparative Genomics

--- the phylogenetic HMM model



10-810, CMB lecture 5---Eric Xing

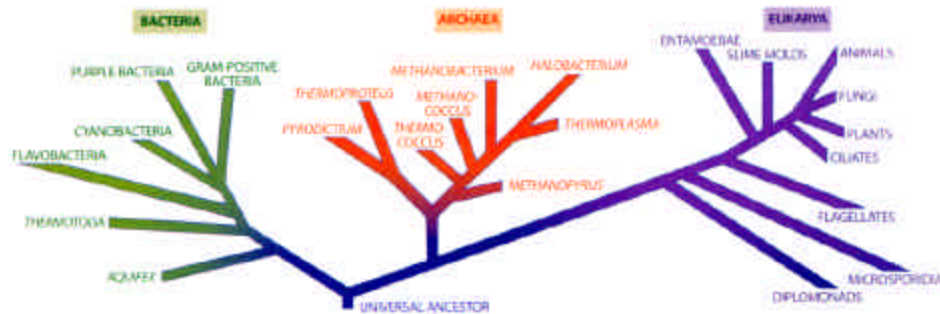
## Some important dates in history (billions of years ago)



Origin of the universe	15 ±4
Formation of the solar system	4.6
First self-replicating system	3.5 ±0.5
Prokaryotic-eukaryotic divergence	1.8 ±0.3
Plant-animal divergence	1.0
Invertebrate-vertebrate divergence	0.5
Mammalian radiation beginning	0.1

86 CSH Doolittle et al.

## The three kingdoms



*M. Madigan and E. Martz: 1997*

## Two important early observations

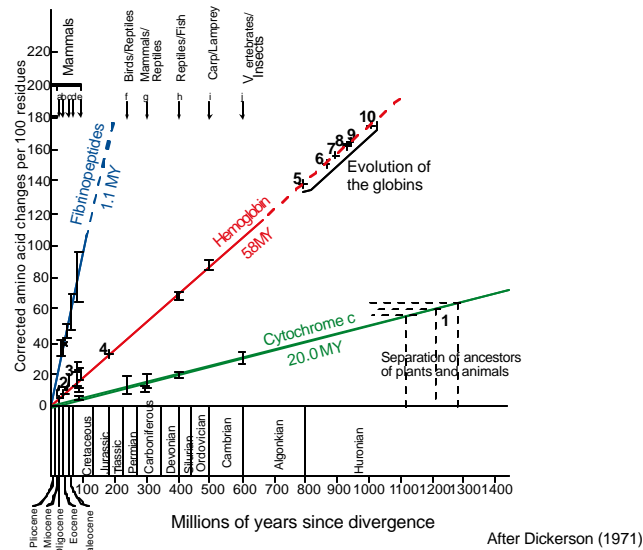


Different proteins evolve at different rates, and this seems more or less independent of the host organism, including its generation time.

It is necessary to adjust the observed percent difference between two homologous proteins to get a distance more or less linearly related to the time since their common ancestor. ( Later we offer a rational basis for doing this.)

An striking early version of these observations is next.

## Rates of macromolecular evolution



## 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



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.

Strictly speaking, only applicable to regions undergoing little selection.

## Some terminology



In evolution, homology (here of proteins), means similarity due to common ancestry.

A common mode of protein evolution is by duplication. Depending on the relations between duplication and speciation dates, we have two different types of homologous proteins. Loosely,

**Orthologues** the “same” gene in different organisms; common ancestry goes back to a speciation event.

**Paralogues** different genes in the same organism; common ancestry goes back to a gene duplication.

Lateral gene transfer gives another form of homology.



## Beta-globins: corrected pairwise distances



DISTANCES between protein sequences, calculated over 1 to 147.

Below diagonal: observed number of differences

Above diagonal: **estimated number of substitutions per 100 amino acids**

Correction method: **Jukes-Cantor**

	hum	mac	bov	pla	chi	sha
hum	----	<b>5</b>	<b>17</b>	<b>27</b>	<b>37</b>	<b>108</b>
mac	7	----	<b>18</b>	<b>27</b>	<b>36</b>	<b>102</b>
bov	23	24	----	<b>32</b>	<b>46</b>	<b>110</b>
pla	34	34	39	----	<b>34</b>	<b>106</b>
chi	45	44	52	42	----	<b>98</b>
sha	91	88	91	90	87	----

## Human globins (paralogues)



	10	20	30
alpha-human	-VLSPADKTNVKAAWGKVGGAHAGEYGAERALERMFLSFPPTT		
beta-human	VH.T.EE.SA.T.L.....-NVD.V.G...G.LLVVY.W.		
delta-human	VH.T.EE..A.N.L.....-NVDV.V.G...G.LLVVY.W.		
epsilon-human	VHFTAE..AA.TSL.S.M--NVE.A.G...G.LLVVY.W.		
gamma-human	GHFTEE..ATITSL.....-NVEDA.G.T.G.LLVVY.W.		
myo-human	-G...DGEWQL.LNV.....E.DIPGH.Q.V.I.L.L.KGH.E.		

	40	50	60	70
alpha-human	KTYFPHF-DLSHGSA-----QVKGHGKKVADALTNVAHV			
beta-human	QRF.ES.G...TPD.VMGNPK..A...LG.FSDGL..L			
delta-human	QRF.ES.G...SPD.VMGNPK..A...LG.FSDGL..L			
epsilon-human	QRF.DS.GN..SP..ILGNPK..A...LTSFGD.IKNM			
gamma-human	QRF.DS.GN..SA..IMGNPK..A...LTS.GD.IK.L			
myo-human	LEK.DK.KH.KSEDEMKASEDL.K..AT.LT..GGILKKK			

	80	90	100	110
alpha-human	DDMPNALSALSDDLHAKLRVDPVNFKLLSHCLLVTLAAHL			
beta-human	.NLKGTFFAT..E..CD..H...E..R..GNV.VCV..H.F			
delta-human	.NLKGTFF.Q..E..CD..H...E..R..GNV.VCV..RNF			
epsilon-human	.NLKP.FAK..E..CD..H...E...GNVMVII..T.F			
gamma-human	.LKGTFQAQ..E..CD..H...E...GNV.VTV..I.F			
myo-human	GHHEAEIKP.AQS..T.HKIPVKYLEFI.E.IIQV.QSKH			

	120	130	140
alpha-human	FAEFTPAVHASLDKFLASVSTVLTLSKYR-----		
beta-human	GK...P.Q.AYO.VV.G.ANA.AH..H.....		
delta-human	GK...QMQ.AYO.VV.G.ANA.AH..H.....		
epsilon-human	GK...E.Q.AWQ.LVSA.AIA.AH..H.....		
gamma-human	GK...E.Q.WQ.MVTA.ASA.S.R.H.....		
myo-human	.GD.GADAQGANNA.ELFRKDMA.N.KELGFQG		

## Human globins: corrected pairwise distances



DISTANCES between protein sequences, calculated over 1 to 141.

Below diagonal: observed number of differences

Above diagonal: **estimated number of substitutions per 100 amino acids**

Correction method: **Jukes-Cantor**

	alpha	beta	delta	epsil	gamma	myo
alpha	----	<b>281</b>	<b>281</b>	<b>281</b>	<b>313</b>	<b>208</b>
beta	82	----	<b>7</b>	<b>30</b>	<b>31</b>	<b>1000</b>
delta	82	10	----	<b>34</b>	<b>33</b>	<b>470</b>
epsil	89	35	39	----	<b>21</b>	<b>402</b>
gamma	85	39	42	29	----	<b>470</b>
myo	116	117	116	119	118	----

## Correcting distances between DNA and protein sequences



Why it is necessary to **adjust** observed percent differences to get a distance measure which scales linearly with time?

This is because we can have **multiple** and **back substitutions** at a given position along a lineage.

All of the correction methods (with names like **Jukes-Cantor**, **2-parameter Kimura**, etc) are justified by simple probabilistic arguments involving Markov chains whose basis is worth mastering.

The same molecular evolutionary models can be used in **scoring** sequence alignments.

## Markov chain



State space = {A,C,G,T}.

$p(i,j) = \text{pr}(\text{next state } S_j \mid \text{current state } S_i)$

**Markov assumption:**

$\text{pr}(\text{next state } S_j \mid \text{current state } S_i \text{ \& any configuration of states before this}) = p(i,j)$

Only the *present* state, not previous states, affects the probs of moving to next states.

## The multiplication rule



$\text{pr}(\text{state after next is } S_k \mid \text{current state is } S_i)$

$= \sum_j \text{pr}(\text{state after next is } S_k, \text{next state is } S_j \mid \text{current state is } S_i)$  [addition rule]

$= \sum_j \text{pr}(\text{next state is } S_j \mid \text{current state is } S_i) \times \text{pr}(\text{state after next is } S_k \mid \text{current state is } S_j, \text{next state is } S_j)$  [multiplication rule]

$= \sum_j p_{i,j} \times p_{j,k}$  [Markov assumption]

$= (i,k)\text{-element of } P^2, \text{ where } P=(p_{i,j}).$

More generally,

$\text{pr}(\text{state } t \text{ steps from now is } S_k \mid \text{current state is } S_i) = (i,k) \text{ element of } P^t$



## Continuous-time version



For any  $(s, t)$ ,

let  $p_{ij}(t) = pr(S_j \text{ at time } t+s \mid S_i \text{ at time } s)$  denote the stationary (time-homogeneous) **transition probabilities**.

Let  $P(t) = (p_{ij}(t))$  denote the matrix of  $p_{ij}(t)$ 's.

Then for any  $(t, u)$ :  $P(t+u) = P(t) P(u)$ .

It follows that  $P'(t) = \exp(Qt)$ , where  $Q = P'(0)$  ( the derivative of  $P(t)$  at  $t = 0$  ).

$Q$  is called the **infinitesimal matrix (transition rate matrix)** of  $P(t)$ , and satisfies

$$P'(t) = QP(t) = P(t)Q.$$

## Interpretation of Q



Roughly,  $q_{ij}$  is the **rate** of transitions of  $i$  to  $j$ , while  $q_{ii} = -\sum_{j \neq i} q_{ij}$ , so each row sum is 0.

Now we have the short-time approximation:

$$p_{i \neq j}(t+h) = q_{ij} h + o(h), \quad \text{and} \quad p_{ii}(t+h) = 1 + q_{ii} h + o(h),$$

Now consider the Chapman-Kolmogorov relation:

(assuming we have a continuous -time Markov chain, and let  $p_j(t) = pr(S_j \text{ at time } t)$ )

$$\begin{aligned} p_j(t+h) &= \sum_i pr(S_i \text{ at } t, S_j \text{ at } t+h) \\ &= \sum_i pr(S_i \text{ at } t) pr(S_j \text{ at } t+h \mid S_i \text{ at } t) \\ &= p_j(t) \times (1 + h q_{jj}) + \sum_{i \neq j} p_i(t) \times h q_{ij} \end{aligned}$$

i.e.,  $h^{-1}[p_j(t+h) - p_j(t)] = p_j(t)q_{jj} + \sum_{i \neq j} p_i(t)q_{ij}$ , which becomes  $P' = QP$  as  $h \rightarrow 0$ .

Important approximation: when  $t$  is small,  $P(t) \approx I + Qt$ .



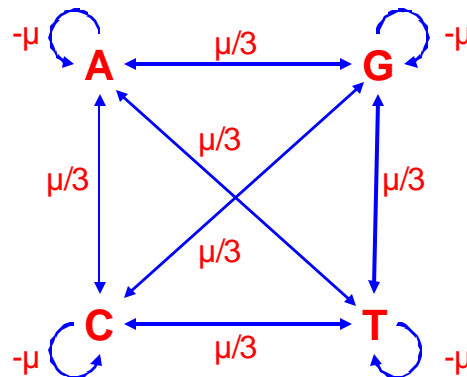
## Probabilistic models for DNA changes

Orc:           ACAGTGACGCCCAAACGT  
Elf:           ACAGTGACGCTACAAACGT  
Dwarf:        CCTGTGACGTAACAAACGA  
Hobbit:       CCTGTGACGTAGCAAACGA  
Human:        CCTGTGACGTAGCAAACGA

## The Jukes-Cantor model (1969)



Substitution rate:



the simplest symmetrical model for DNA evolution

## Jukes-Cantor (cont.)



Prob transition matrix:

$$S(t) = \begin{matrix} & \begin{matrix} A & C & G & T \end{matrix} \\ \begin{matrix} A \\ C \\ G \\ T \end{matrix} & \begin{pmatrix} r(t) & s(t) & s(t) & s(t) \\ s(t) & r(t) & s(t) & s(t) \\ s(t) & s(t) & r(t) & s(t) \\ s(t) & s(t) & s(t) & r(t) \end{pmatrix} \end{matrix}$$

Where we can derive:

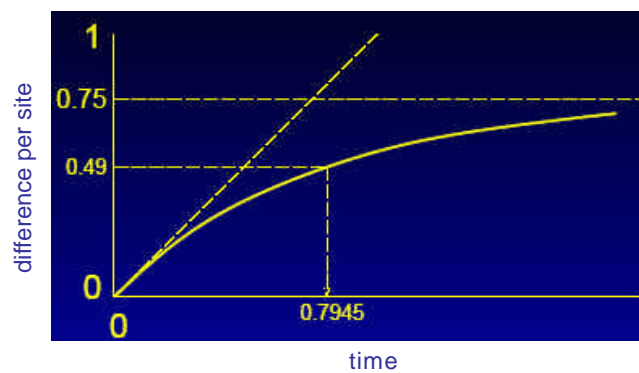
$$r(t) = \frac{1}{4} (1 + 3 e^{-4\alpha t})$$

$$s(t) = \frac{1}{4} (1 - e^{-4\alpha t})$$

## Jukes-Cantor (cont.)



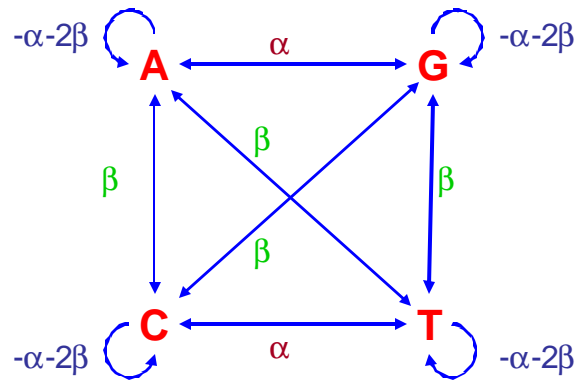
Fraction of sites differences



## Kimura's K2P model (1980)



Substitution rate:



which allows for different rates of transition and transversions.

**Transitions (rate  $\alpha$ )** are much more likely than **transversions (rate  $\beta$ )**.

## Kimura (cont.)



Prob transition matrix:

$$S(t) = \begin{pmatrix} r(t) & s(t) & u(t) & s(t) \\ s(t) & r(t) & s(t) & u(t) \\ u(t) & s(t) & r(t) & s(t) \\ s(t) & u(t) & s(t) & r(t) \end{pmatrix}$$

Where

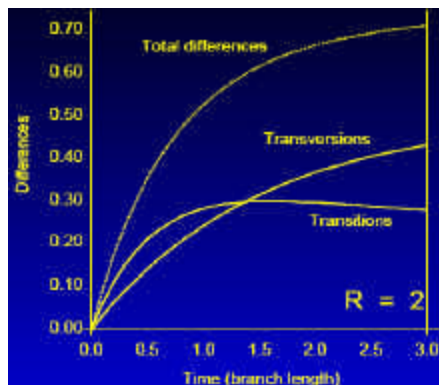
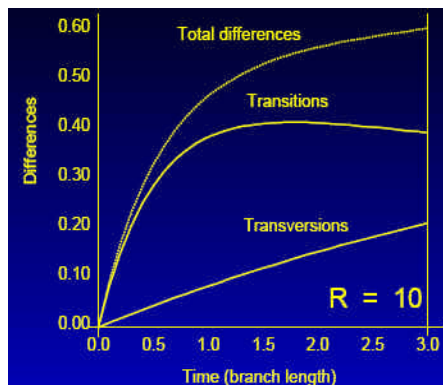
$$\begin{aligned} s(t) &= \frac{1}{4} (1 - e^{-4\beta t}) \\ u(t) &= \frac{1}{4} (1 + e^{-4\beta t} - e^{-2(\alpha+\beta)t}) \\ r(t) &= 1 - 2s(t) - u(t) \end{aligned}$$

By proper choice of  $\alpha$  and  $\beta$  one can achieve the overall rate of change and  $T_s=T_v$  ratio  $R$  you want (*warning: terminological tangle*).

## Kimura (cont.)



Transitions, transversions expected under different R:



## The general time-reversible model

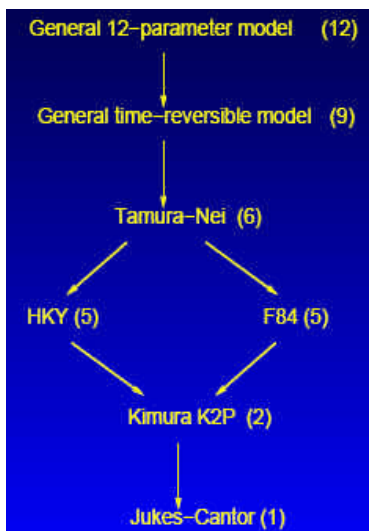


It maintains "detailed balance" so that the probability of starting at (say) A and ending at (say) T in evolution is the same as the probability of starting at T and ending at A:

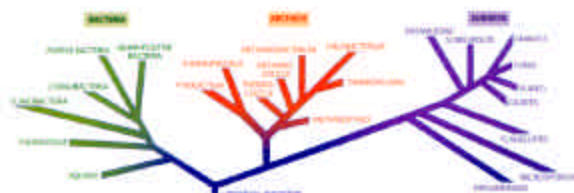
	A	C	G	T
A	?	$ap_C$	$\beta p_G$	$?p_T$
C	$ap_A$	?	$dp_G$	$ep_T$
G	$\beta p_A$	$dp_C$	?	$?p_T$
T	$?p_A$	$ep_C$	$?p_G$	?

And there is of course the **general 12-parameter model** which has arbitrary rates for each of the 12 possible changes (from each of the 4 nucleotides to each of the 3 others). (Neither of these has formulas for the transition probabilities, but those can be done numerically.)

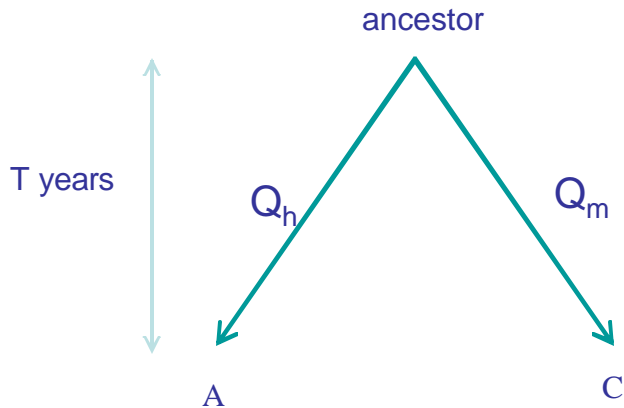
## Relation between models



## Phylogenetic trees



## A pair of homologous bases



Typically, the ancestor is unknown.

## More assumptions



- $Q_h = s_h Q$  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  $\mathbf{p}$  of  $Q$ .
- $Q$  is **reversible**: for  $a, b, t \geq 0$   
$$\mathbf{p}(a)P(t, a, b) = P(t, b, a)\mathbf{p}(b),$$
  
(detailed balance).

## The stationary distribution



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

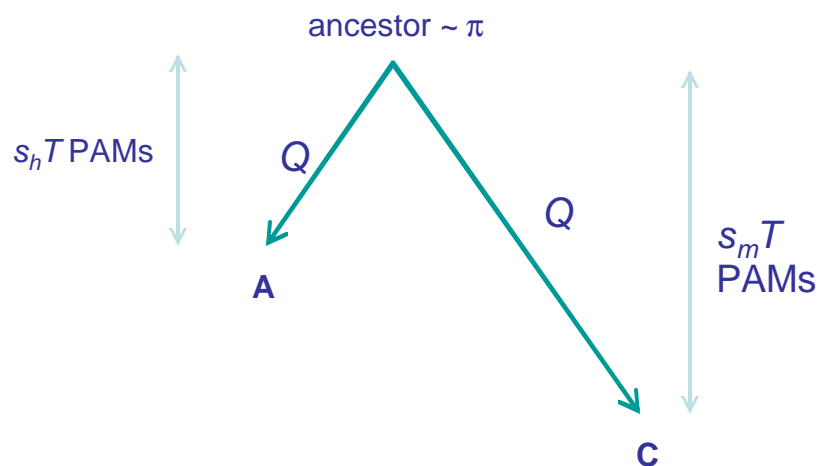
$$\sum_i p(i) P(i, j) = p(j).$$

**Exercise.** Given any initial distribution, the distribution at time  $t$  of a chain with transition matrix  $P$  converges to  $\mathbf{p}$  as  $t \rightarrow \infty$ . Thus,  $\mathbf{p}$  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.)

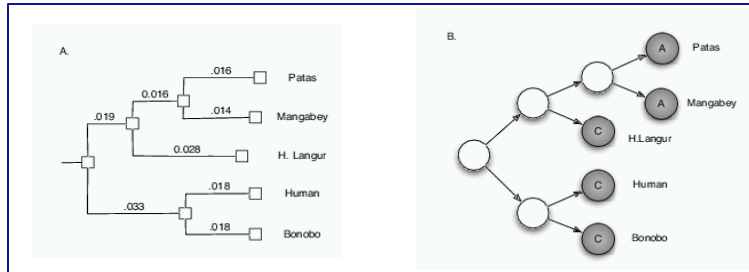
We often assume that the ancestor sequence is i.i.d  $\pi$ .

## New picture





# 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

# 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

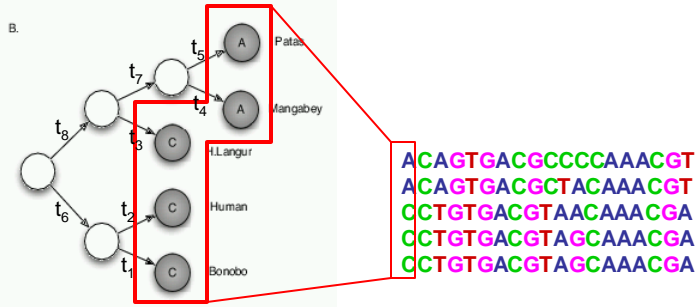
## Major methods:

- Parsimony phylogeny methods
- Likelihood methods

# Likelihood methods



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



Since the sites evolve independently on the same tree,

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

# Likelihood at one site on a tree



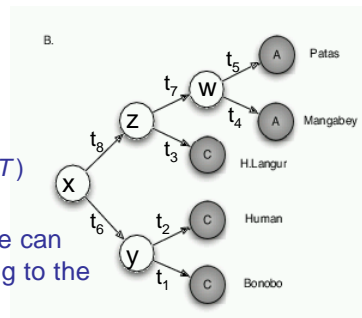
We can compute this by summing over all assignments of states x, y, z and w to the interior nodes:

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

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

$$P(A,A,C,C,C,x,y,z,w | T) = P(x) P(y | x, t_6) P(A | y, t_1) P(C | y, t_2) P(z | x, t_6) P(C | y, t_3) P(w | z, t_7) P(C | y, t_4) P(C | y, t_5)$$

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:

$$P(D^{(t)} | T) = \sum_x \sum_y \sum_z \sum_w P(A, A, C, C, C, x, y, z, w | T) =$$

$$\sum_x P(x)$$

$$\left( \sum_y P(y | x, t_6) P(A | y, t_4) P(C | y, t_2) \right)$$

$$\left( \sum_z P(z | x, t_8) P(C | z, t_3) \right)$$

$$\left( \sum_w P(w | z, t_7) P(C | w, t_4) P(C | w, t_5) \right)$$

## Felsenstein's Pruning Algorithm



To calculate  $P(x_1, x_2, \dots, x_N | T, t)$

### Initialization:

Set  $k = 2N - 1$

### Recursion: Compute $P(L_k | a)$ 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  $P(L_i | b)$ ,  $P(L_j | b)$  for all  $b$ , for daughter nodes  $i, j$

2. Set  $P(L_k | a) = \sum_{b, c} P(b | a, t_i) P(L_i | b) P(c | a, t_j) P(L_j | c)$

### Termination:

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

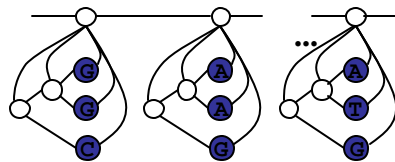
## Finding the ML tree



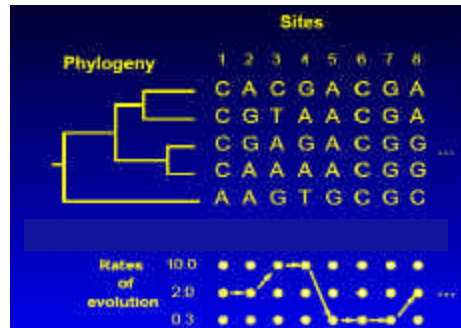
So far I have just talked about the computation of the likelihood for one tree with branch lengths known.

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.

## Phylogenetic HMM

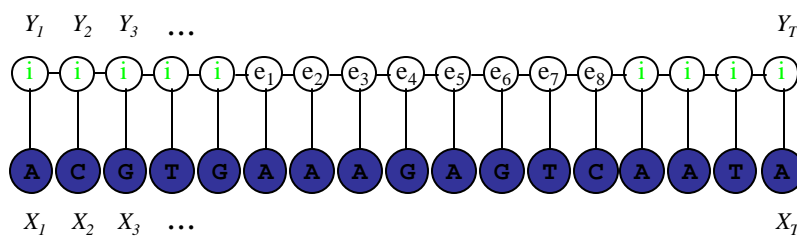


## Modeling rate variation among sites



- 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.

## Recall 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)



**Definition:** A hidden Markov model (HMM)

- **Observation alphabet**  $\Sigma = \{ b_1, b_2, \dots, b_M \}$
- **Set of hidden states**  $Q = \{ 1, \dots, K \}$
- **Transition probabilities** between any two states

$$a_{ij} = P(y=j|y=i)$$

$$a_{i1} + \dots + a_{iK} = 1, \text{ for all states } i = 1 \dots K$$

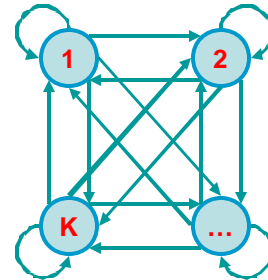
- **Start probabilities**  $a_{0i}$

$$a_{01} + \dots + a_{0K} = 1$$

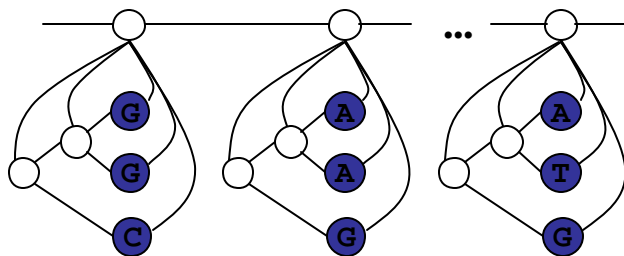
- **Emission probabilities** associated with each state

$$e_{ib} = P(x_i = b | y_i = k)$$

$$e_{i,b1} + \dots + e_{i,bM} = 1, \text{ for all states } i = 1 \dots K$$



## Hidden Markov Phylogeny



- This yields a gene finder that exploits evolutionary constraints
- Based on sequence data from 12-15 primate species, McAuliffe et al (2003) obtained sensitivity of 100%, with a specificity of 89%.
  - Genscan (state-of-the-art gene finder) yield a sensitivity of 45%, with a specificity of 34%.

## The Forward Algorithm



We can compute  $f_k(t)$  for all  $k, t$ , using dynamic programming!

### Initialization:

$$\begin{aligned} f_0(0) &= 1 \\ f_k(0) &= 0, \text{ for all } k > 0 \end{aligned}$$

### Iteration:

$$f_i(t) = e_i(\mathbf{X}_t) \sum_k f_k(t-1) a_{ki} \quad (a_{0k} \text{ is a vector of initial probability})$$

### Termination:

$$P(x) = \sum_k f_k(T)$$

## The Backward Algorithm



We can compute  $b_k(t)$  for all  $k, t$ , using dynamic programming

### Initialization:

$$b_k(T) = 1, \text{ for all } k$$

### Iteration:

$$b_k(t) = \sum_l e_l(\mathbf{X}_{t+1}) a_{kl} b_l(t+1)$$

### Termination:

$$P(x) = \sum_l a_{0l} e_l(x_1) b_l(1)$$

## 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 sites use the same phylogeny? Function-specific phylogeny?
- Other evolutionary events: duplication, rearrangement, lateral transfer, etc.