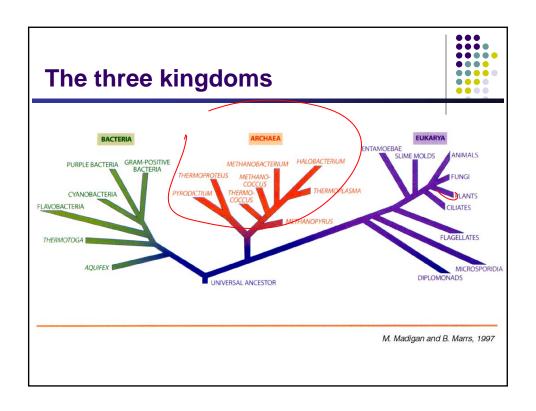


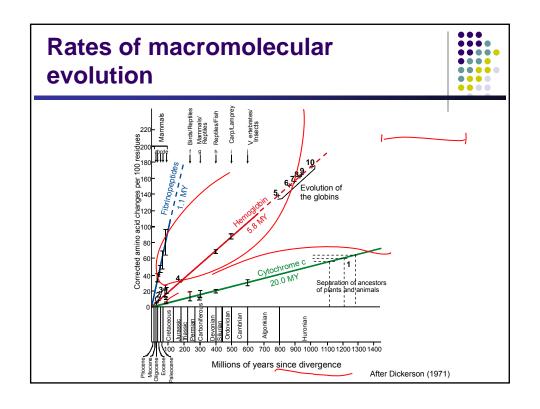
	ome important dates in illions of years ago)	history	
•	Origin of the universe	15 ±4	
	Formation of the solar system	4.6	
•	First self-replicating system	3.5 ± 0.5	
•	Prokaryotic-eukaryotic divergence	1.8 ±0.3	
•	Plant-animal divergence	1.0	
•	Invertebrate-vertebrate divergence	0.5	
•	Mammalian radiation beginning	0.1	
		(86 CSH Doolittle e	et al.)



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.)
- See nest slide ...



How does sequence variation arise?



IU

- 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)
 - Site independence.

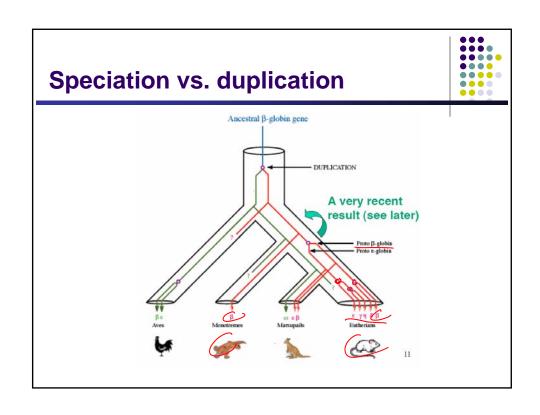


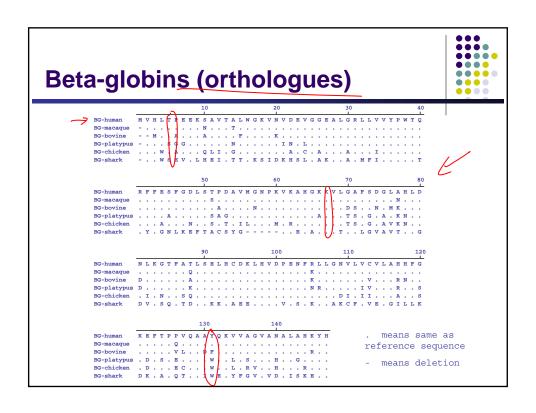
- Site homogeneity.
- Markovian: given current base, future substitutions independent of past.
- Temporal homogeneity: stationary Markov chain.
- Strictly speaking, only applicable to regions undergoing little selection.

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: uncorrected pairwise distances



- DISTANCES between protein sequences (calculated over: 1 to 147)
 - Below diagonal: observed number of differences
 - Above diagonal: number of differences per 100 amino acids

	hum	mac	bov	pla	chi	sha	0
hum		5	16	23	31	65	
mac	7	7	17	23	30	62	
bov	23	24		27	37	65	
pla	34	34	39	7	29	64	
chi	45	44	52	42	7	61	
sha	91	88	91	90	87		

Beta-globins: corrected pairwise distances



- DISTANCES between protein sequences (calculated over: 1 to 147)
- Below diagonal: observed number of differences
 - Above diagonal: number of differences per 100 amino acids
 - Correction method: Jukes-Cantor /

ווע	CCHOIT	memoc	i. Quices	Caritor				
		hum	mac	bov	pla	chi	sha	
	hum		5	17	27	37	108	1 /
	mac	7		18	27	36	102	
	bov	23	24		32	46	110	
	pla	34	34	39		34	106	
	chi	45	44	52	42		98	
	sha	91	88	91	90	87		

uman globins (paralogues)																																						
									10									20								3	0											1
alpha-human	Ξ	v	LS	3 P	A	D	K :	N	v	K	A	A	w	G	ĸ v	7 G	A	н	Α	G	E :	7 0	A	Е	A	ь	E R	м	F	L	s	F	P :	r 1	r			
beta-human	v	н	. 1	r.	E	E	. :	a A		т		L					-	N	v	D	. 1	7.	G			. (3.	L	L	v	v	Y	. 1	٠.				
delta-human	v	н	. 1	r.	E	E		. A		N		L					-	N	v	D.	A١	7.	G			. (3.	L	L	v	v	Y	. 1	٠.				
epsilon-human	v	н	F 1	C A	E	E	. 2	A		T	s	L		s	. 1	4 -	-	N	v	E	. 2	Α.	G			. (3.	L	L	v	V	Y	. 1	٠.				
gamma-human	G	н	F 1	E	E		. 2	Т	I	T	s	L					-	N	v	E	D A	Α.	G		T	. (з.	L	L	v	v	Y	. 1	٠.				
myo-human	-	G		. D	G	E	W (L	٠	L	N	v	٠	٠	-	. Е	٠.	D	Ι	P	G I	Ι.	Q	٠	V	. :	٠.	L	٠	K	G	H	. 1	Ξ.				
	40									50													60)								70						
alpha-human	K	т	Y E	7 P	н	F	- 1) L	s	н	G	s	Α	-			-	Q	v	K	G I	1 0	K	K	v.	A I) A	L	т	N	A	v .	A I	ı v	7			
beta-human	Q	R	F.	. E	s		G.			T	P	D		v	м	3 N	P	K			Α.				. :	L (з.	F	s	D	G	L		. I	4			
delta-human	Q	R	F.	. E	s		G.			s	P	D		v	м	3 N	P	K			Α.				. :	L (з.	F	s	D	G	L		. I	4			
epsilon-human	Q	R	F.	. D	s		G 1	ι.		s	P			Ι	L	3 N	I P	K			Α.				. :	L :	r s	F	G	D		I	K 1	1 1	1			
gamma-human	Q	R	F.	. D	s		G 1	Ι.		s	A			Ι	M (3 N	P	K			Α.				. :	L :	r s		G	D		I	κ.	. І	4			
myo-human	L	E	Κ.	. D	K	٠	K I	Ι.	K	s	E	D	E	M	K A	A S	E	D	L	•	K .		A	T	•	ь :	г.	٠	G	G	Ι	L	K I	C B	C			
						80									90								10	0							1	.10						
alpha-human	D	D :	мі	P N	A	L	s z	L	s	D	L	н	А	н	K I	R	v	D	P	v :	N I	7 B	L	L	s	н	2 1	L	v	т	L	A .	A I	1 1				
beta-human		N	LE	C	т	F	A :	٠.		E			C	D		. н	٠.			E					G :	N V	7.	v	C	v	÷	. 1	н.	. F				
delta-human		N	LE	C	т	F	. (E			C	D		. н	٠.			E		. F			G :	N V	7.	v	C	v		. 1	R 1	I F				
epsilon-human		N	LE	C P		F	A I	٠.		E			C	D		. н	٠.			E					G :	N V	7 M	v	I	I			т.	. F				
gamma-human			LE	C	т	F	A (2 .		E			C	D		. н	٠.			E					G :	N V	7.	v	т	v		. :	ī.	. F				
myo-human	G	H	н	S A	E	I	K 1		A	Q	s	٠	٠	T	. I	I K	I	P	v	K	Y 1	L E	F	Ι		Ε.	. 1	Ι	Q	v	٠	Q	S I	C E	I			
					1	L20								1	30								14	0														
alpha-human	P	A	E I	7 T	P	Α	V I	I A	s	L	D	K	F	L	A s	3 V	7 8	т	v	L	T S	3 B	Y	R	-	-		-	-									
beta-human	G	K				P	. (٠.	A	Y	Q		v	V	. (3.	A	N	A		A I	Ι.		Н														
delta-human	G	K				Q	M (٠.	A	Y	Q		v	V	. (3.	A	N	A		A I	Ι.		Н														
epsilon-human	G	K				Е	. (. 9	A	W	Q		L	V	s z	Α.	A	I	A		A I	Ι.		Н														
gamma-human	G					E																						٠	٠									
mvo-human		G :	D.	. G	A	D	A (2 G	Α	М	N		А		E 1	. F	R	K	D	M	Α.	. 19		K	E :	L (3 F	Q	G									

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
alph	а	281	281	281	313	208
beta	82		7	30	31	1000
delta	a 82	10		34	33	470
epsi	l 89	35	39		21	402
gam	ma 85	39	42	29		470
myo	116	117	116	119	118	

Correcting distances between DNA and protein sequences

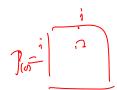


- Why it is necessary to adjust observed percent differences to get a distance measure which scales linearly with time?
- This is because we can have multiple and back substitutions at a given position along a lineage.
- All of the correction methods (with names like Jukes-Cantor, 2parameter Kimura, 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) = pr(next state S_i) | current state S_i)



Markov assumption:

p(next state S_j | current state S_i & 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



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

- = $\sum_{i} pr(\text{state } \underline{\text{after next}} \text{ is } S_{k}, \underline{\text{next state}} \text{ is } S_{i} | \text{ current state is } S_{i})$ [addition rule]
- = $\sum_i pr(\text{next state is } S_i)$ or sent state is S_i) x $pr(\text{state after next is } S_k \mid \text{corrent})$

state is S_i , next state is S_i)

[multiplication rule] -

 $\sum_{j} \overline{p}_{i,j} \times p_{j,k}$

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

P (North next [Markov assumption] here)

More generally,

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

Continuous-time version



- For any (s, t):
 - Let $p_{ij}(t) = pr(S_i \text{ at time } t+s \mid S_i \text{ at time } s)$ denote the stationary (time-homogeneous)
- Let $P(t) = (p_{ii}(t))$ denote the matrix of $p_{ii}(t)$'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.

Important approximation: when t is small,

 $P(t) \approx I + Qt$.



Interpretation of Q



- Roughly, q_{ij} is the **rate** of transitions of *i* to *j*, while q_{ij} = each row sum is 0 (Why?).
- Now we have the short-time approximation:

$$p_{i\neq j}(t+h)=q_{ij}h+o(h)$$

$$(p_{i=j}(t+h)) + 1 + q_{ii}h + o(h)$$

 $p_{i\neq j}(t+h) = q_{ij}h + o(h)$ where $p_{ij}(t+h)$ is the probability of transitioning from i at time t to j at time t+h

 Now consider the Chapman-Kolmogorov relation: (assuming we flave a continuous-time Markov chain, and let $p_i(t) = pr(S_i \text{ at time } t)$:

$$\frac{p_j(t+h) = \sum_i pr(S_i \text{ at } t, S_j \text{ at } t+h)}{pr(S_i \text{ at } t) pr(S_j \text{ at } t+h \mid S_i \text{ at } t)}$$

$$= p_j(t) \times (1 + q_{jj}h) + \sum_{i \neq j} p_i(t) \times hq_{ij}$$

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 $h\sqrt{0}$.



Probabilistic models for DNA changes

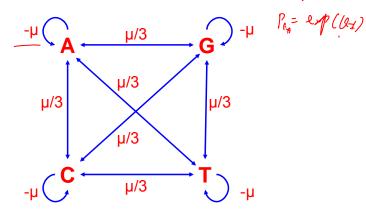
Orc: Elf: Dwarf:

Hobbit: Human: ACAGTGACGCCCCAAACGT ACAGTGACGCTACAAACGT CCTGTGACGTAACAAACGA CCTGTGACGTAGCAAACGA CCTGFGACGTAGCAAACGA

The Jukes-Cantor model (1969)



• Substitution rate:



Transition probabilities under the Jukes-Cantor model

the simplest symmetrical model for DNA evolution



- IID assumption:
 - All sites change independently
 - All sites have the same stochastic process working at them
- Equiprobablity assumption:
 - Make up a fictional kind of event, such that when it happens the site changes to one of the 4 bases chosen at random equiprobably
- Equilibrium condition:



 No matter how many of these fictional events occur, provided it is not zero, the chance of ending up at a particular base is 1/4.



• Solving differentially equation system P' = QP

Transition probabilities under the Jukes-Cantor model (cont.)

• Prob transition matrix:

A C G T

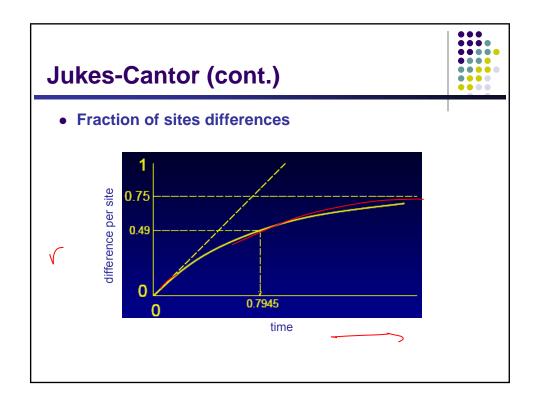
P(t) = C S(t) S(t) S(t) S(t)

G S(t) S(t) S(t) S(t)

T S(t) S(t) S(t) T(t)

Where we can derive:

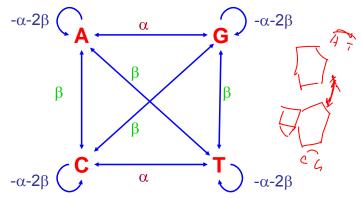
$$r(t) = \frac{1}{4}(1 - 3e^{-\frac{4}{3}\mu t})$$
Homework!



Kimura's K2P model (1980)



• Substitution rate:



- which allows for different rates of transition and transversions.
- Transitions (rate α) are much more likely than transversions (rate β).

Kimura (cont.)



• Prob transition matrix:

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

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

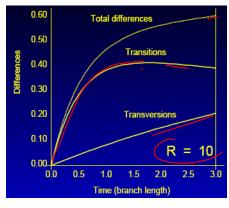
 $u(t) = \frac{1}{4} (1 + e^{-4\beta t} - e^{-2(\alpha + \beta)t})$
 $r(t) = 1 - 2s(t) - u(t)$

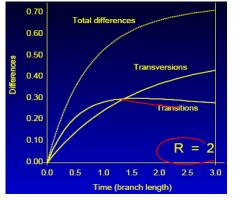
 By proper choice of α and β one can achieve the overall rate of change and Ts=Tn ratio R you want (warning: terminological tangle).

Kimura (cont.)



• Transitions, transversions expected under different R:





Other commonly used models



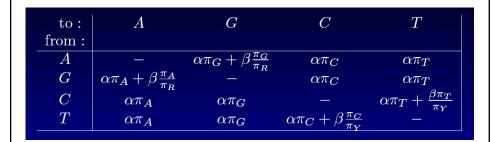
- Two models that specify the equilibrium base frequencies
 (you provide the frequencies A; C; G; T and they are set up to
 have an equilibrium which achieves them), and also let you
 control the transition/transversion ratio:
- The Hasegawa-Kishino-Yano (1985) model:

to:	A	G	C	T
${ m from}:$				
\overline{A}	_	$\alpha \pi_G + \beta \pi_G$	$\alpha\pi_C$	$lpha\pi_T$
G	$\alpha \pi_A + \beta \pi_A$		$lpha\pi_C$	$lpha\pi_T$
C	$lpha\pi_A$	$lpha\pi_G$	_	$\alpha \pi_T + \beta \pi_T$
T	$lpha\pi_A$	$lpha\pi_G$	$\alpha \pi_C + \beta \pi_C$	_

Other commonly used models



• The **F84 model** (Felsenstein)



• where $\pi_R = \pi_A + \pi_G$ and $\pi_Y = \pi_C + \pi_T$ (The equilibrium frequencies of purines and pyrimidines)

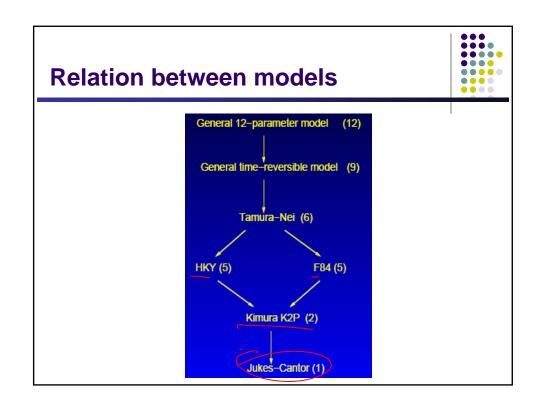
The general time-reversible model

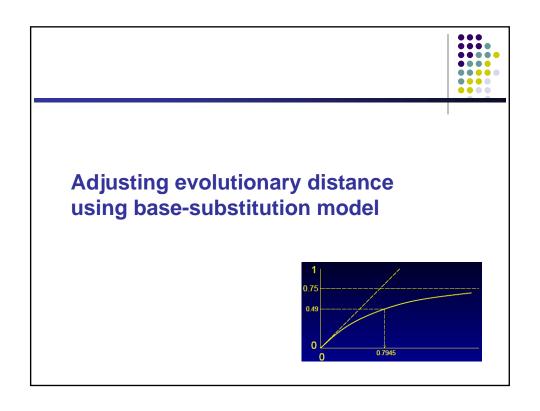


 It maintains "detailed balance" so that the probability of starting at (say) A and ending at (say) T in evolution is the same as the probability of starting at T and ending at A:

	Α	C	G	Т
Α	_	$a\pi_{C}$	$oldsymbol{eta}\pi_{\!\scriptscriptstyle G}$	$\gamma \pi_{T}$
С	$\alpha \pi_{\!\scriptscriptstyle A}$	_	$\delta\pi_{\!\scriptscriptstyle G}$	$\varepsilon\pi_{T}$
G	$\beta\pi_{\!\scriptscriptstyle A}$	$\delta \pi_{\!\scriptscriptstyle C}$	_	$V\Pi_T$
Т	γπ	$arepsilon\pi_{\!\scriptscriptstyle C}$	$V\pi_{\!\scriptscriptstyle G}$	_

- 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.)





The Jukes-Cantor model



Common ancestor of human and orang

$$Q = \begin{bmatrix} -3\alpha & \alpha & \alpha & \alpha \\ \alpha & -3\alpha & \alpha & \alpha \\ \alpha & \alpha & -3\alpha & \alpha \\ \alpha & \alpha & \alpha & -3\alpha \end{bmatrix}$$

t time unit

Human (now)

Consider e.g. the 2nd

position in a-globin2 Alu1.
$$r = (1+3e^{-4\alpha t})/4$$
, $s = (1-e^{-4\alpha t})/4$.

Definition of PAM



• Let P(t) = exp(Qt). Then the A, G element of P(t) is

$$pr(G \text{ now} \mid A \text{ then}) = (1 - e^{-4\alpha t})/4.$$

- Same for all pairs of different nucleotides.
- Overall rate of change $k = 3\alpha t$.



- PAM = accepted point mutation
 - When k ≠ .01 described as 1 PAM
 - Put $t = .01/3\alpha$ (1/300 α .) Then the resulting $P = P(1/300\alpha)$ is called the PAM(1) matrix. P(t)
- Why use PAMs?

Evolutionary time, PAM



- Since sequences evolve at different rates, it is convenient to rescale time so that *1 PAM* of evolutionary time corresponds to *1%* expected substitutions.
- For Jukes-Cantor, $k = 3\alpha t$ is the expected number of substitutions in [0,t], so is a distance. (Show this.)
 - Set $3\alpha t = 1/100$, or $t = 1/300\alpha$, so $1 PAM = 1/300\alpha$ years.

Distance adjustment



- For a pair of sequences, $k = 3\alpha t$ is the desired metric, but not observable. Instead, pr(different) is observed. So we use a model to convert pr(different) to k.
- This is completely analogous to the conversion of $\theta = pr(recombination)$



to genetic (map) distance (= expected number of crossovers) using the Haldane map function

$$\theta = 1/2 \times (1 - e^{-2d}),$$

assuming the no-interference (Poisson) model.

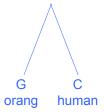
Towards Jukes-Cantor adjustment



• E.g., 2nd position in a-globin Alu 1

common ancestor

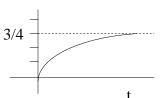
 Assume that the common ancestor has A, G, C or T with probability 1/4.



• Then the chance of the nt differing

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

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



Jukes-Cantor adjustment



 If we suppose all nucleotide positions behave identically and independently, and n_≠ differ out of n, we can invert this, obtaining

$$\widehat{k} = -\frac{3}{4} \times \log \left(1 - \frac{4}{3} n_{\pm} / n\right)$$

- This is the corrected or adjusted fraction of differences (under this simple model). × 100 to get PAMs
- The analogous simple model for amino acid sequences has

$$\widehat{k} = -\frac{19}{20} \times \log \left(1 - \frac{20}{19} n_{\neq} / n \right)$$

 \times 100 for PAM.

Illustration



1. Human and bovine beta-globins are aligned with no deletions at 145 out of 147 sites. They differ at 23 of these sites. Thus $n_{\neq}/n = 23/145$, and the corrected distance using the Jukes-Cantor formula is (natural logs)

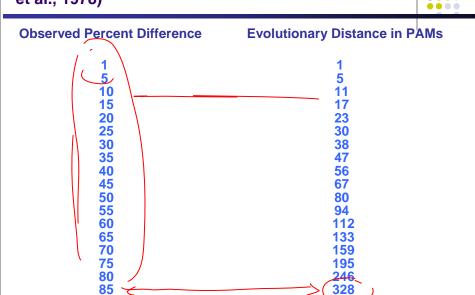
```
-19/20 \times \log(1-20/19 \times 23/145) = 17.3 \times 10^{-2}.
```

2. The human and gorilla sequences are aligned without gaps across all 300 bp, and differ at 14 sites. Thus $n_{\neq}/n = 14/300$, and the corrected distance using the Jukes-Cantor formula is

$$-3/4 \times \log(1-4/3 \times 14/300) = 4.8 \times 10^{-2}$$
.

Correspondence between observed a.a. differences and the evolutionary distance (Dayhoff et al., 1978)

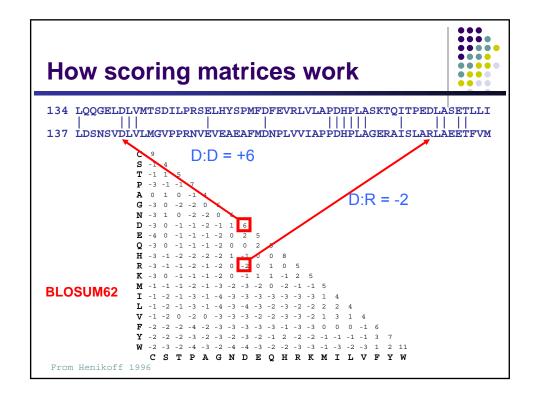






Scoring matrices for alignment





Statistical motivation for alignment scores



```
Alignment: AGCTGATCA...
AACCGGTTA... Hypotheses: H = homologous (indep. sites, Jukes-Cantor)
R = random (indep. sites, equal freq.)
```

$$pr(data \mid \mathcal{H}) = pr(AA \mid \mathcal{H})pr(GA \mid \mathcal{H})pr(CC \mid \mathcal{H})...$$

$$= (1-p)^{a} p^{d}, \text{ where } a = \text{\#agreements, } d = \text{\#disagreements, } p = \frac{3}{4}(1-e^{-8at}).$$

$$pr(data \mid R) = pr(AA \mid R)pr(GA \mid R)pr(CC \mid R)...$$

$$= (\frac{1}{4})^{a}(\frac{3}{4})^{d}$$

$$\Rightarrow \log\{\frac{pr(data \mid \mathcal{H})}{pr(data \mid R)}\} = a\log\frac{1-p}{1/4} + d\log\frac{p}{3/4} = a \times \sigma + d \times (-\mu).$$

- Since p < 3/4, $\sigma = log((1-p)/(1/4)) > 0$, while $-\mu = log(p/(3/4)) < 0$.
- Thus the alignment score = $a \times \sigma + d \times (-\mu)$, where the match score $\sigma > 0$, and the mismatch penalty is $-\mu < 0$.

Large and small evolutionary distances

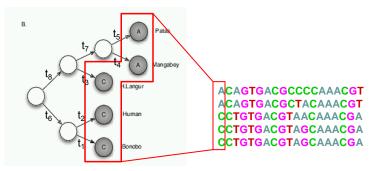


- Recall that
 - $p = (3/4)(1-e^{-8\alpha t}),$
 - $\sigma = \log((1-p)/(1/4))$,
 - $-\mu = log(p/(3/4))$.
- Now note that if αt ≈ 0,
 - then p ≈ 6αt, and 1-p ≈ 1, and so σ ≈ log4, while -μ ≈ log8αt is large and negative.
 - That is, we see a big difference in the two values of σ and μ for small distances.
- Conversely, if αt is large,
 - $p = (3/4)(1-\varepsilon)$, hence $p/(3/4) = 1-\varepsilon$, giving $\mu = -\log(1-\varepsilon) \approx \varepsilon$, while $1-p = (1+3\varepsilon)/4$, $(1-p)/(1/4) = 1+3\varepsilon$, and so $\sigma = \log(1+3\varepsilon) \approx 3\varepsilon$.
 - Thus the scores are about 3 (for a match) to 1 (for a mismatch) for large distances. This makes sense, as mismatches will on average be about 3 times more frequent than matches.
- the matrix which performs best will be the matrix that reflects the evolutionary separation of the sequences being aligned.

What about multiple alignment



• Phylogenetic methods: a tree, with branch lengths, and the data at a single site.



 See next lecture for how to compute likelihood under this hypothesis

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