

TROY RUTHS Presents

HIMMIS BACK TO THE FUTURE

PHYLOHMM: THE FUTURE

- Bring *back* the HMM *to the future*
- Updating the biological model with phylogenetics

Combining Phylogenetic and HMMs in
Biosequence Analysis

Adam Siepel & David Haussler (2004)

WHAT WE'LL DISCUSS

- Motivation for updating HMMs
- Design of Phylo-HMM
 - Tree Model
 - DNA Substitution
 - Evolutionary rate
 - ➔ Categories
 - ➔ Higher-order states
- Application to data & results

RECAP ON HMMs

- Dominant tool in biological sequence analysis
 - Gene prediction, homology searching, structure ...
- ➔ Balance simplicity and expressiveness

ANTIQUATED HMMs

Your HMMs disregard three decades of sequence evolution research.

...



“CS McFly”

“Biologist Biff”

ANTIQUATED HMMS

- Sites are independent
- Substitutions are homogeneous
- Evolutionary rates are consistent
- Functional categories are disregarded

What's the solution?



Unrealistic model!

ENTER PHYLOGENY

- Provides probabilistic models of evolution
- Based on
 - Topology of tree (relatedness)
 - Lengths of its branches (rates)
 - Pattern of substitution (categories)

➔ *Time-based*

➔ *Works across sequences*

HMM \neq PHYLOGENY

- Both are built on probabilistic models
- HMM operates along a sequence
- Phylogenetics operate between sequences

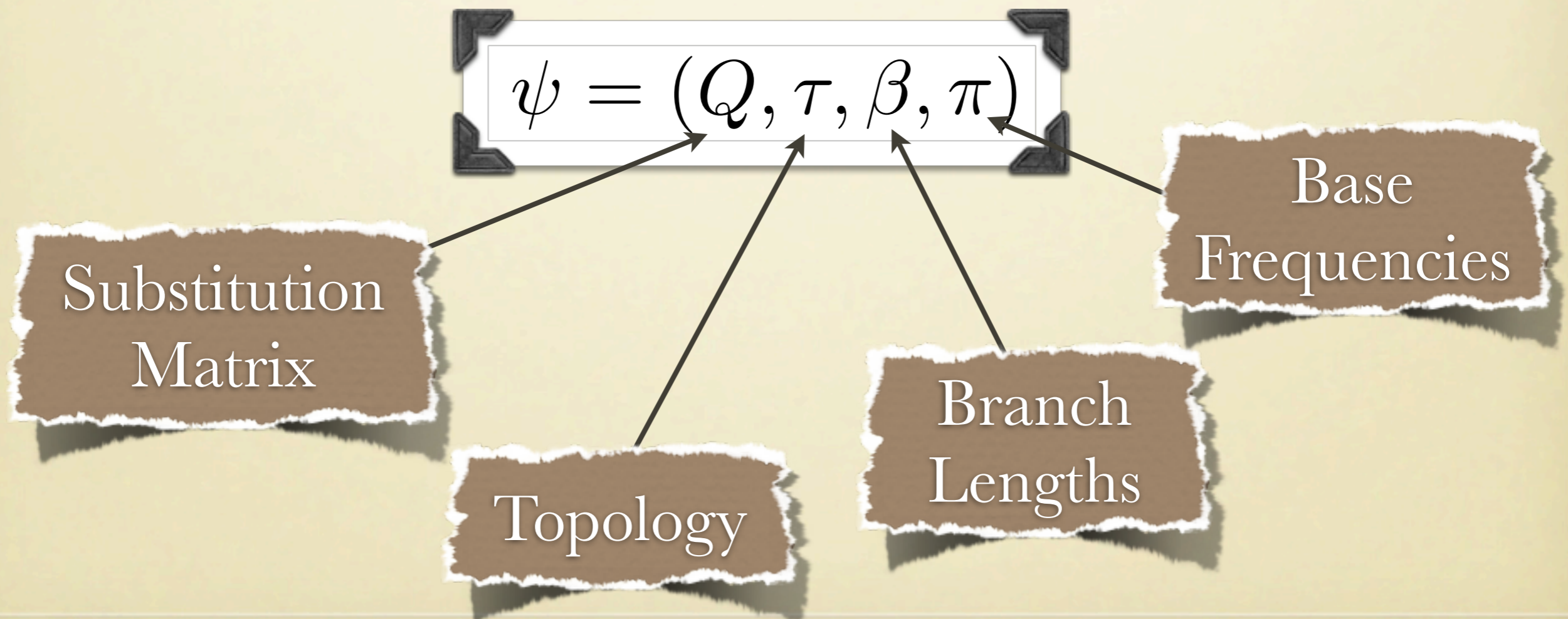


Space + Time = Phylo-HMM

THE METHOD

INPUT

- n aligned sequences of length L
- Phylogenetic tree relating the n taxa

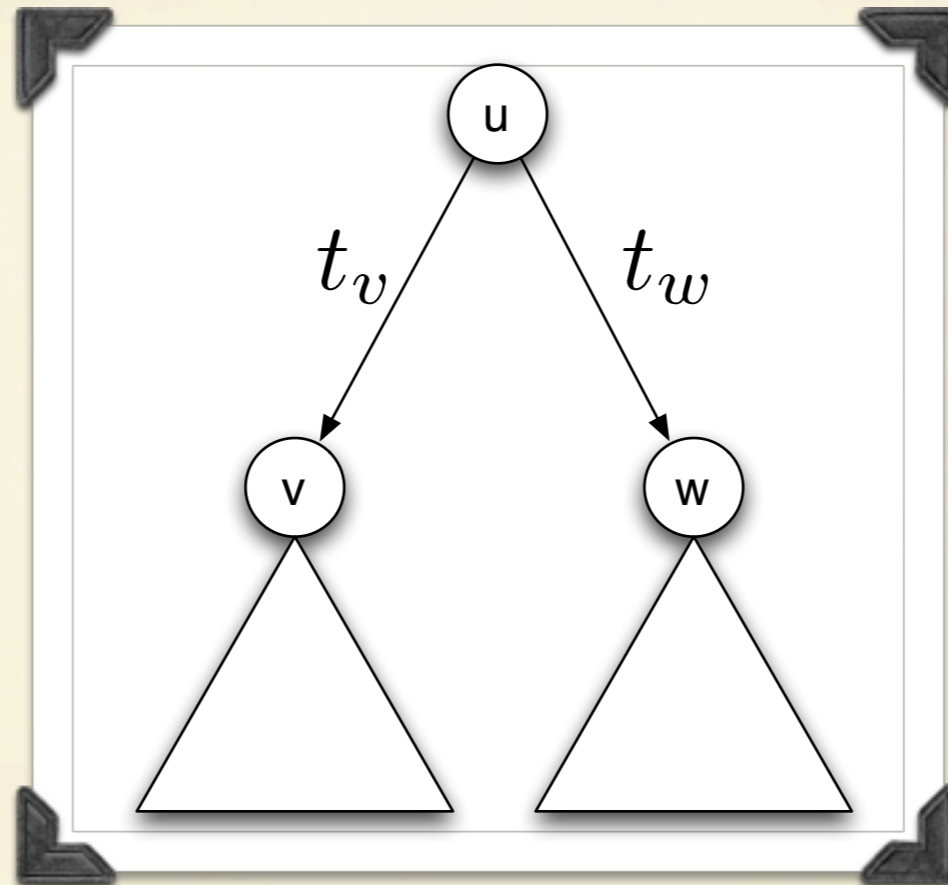


LIKELIHOOD OF A TREE

- Sites of the alignment are assumed independent
- Dynamic programming solution

$$P(X|\psi) = \prod_{i=1}^L P(X_i|\psi) = \sum_{\mathbb{L}} P(\mathbb{L}, X_i|\psi)$$

Labeling of
ancestral nodes



Recursion

$$P(L_u|a) = \begin{cases} I(a = x_u) & \text{if } u \text{ is a leaf} \\ \sum_b P(b|a, t_v) P(L_v|b) \sum_c P(c|a, t_w) P(L_w|c) & \text{otherwise} \end{cases}$$

Root Call

$$P(X_i|\psi) = \sum_a \pi_a P(L_r|a)$$

DNA SUBSTITUTION

$$P(b|a, t)$$

- Probability that base b is substituted by base a over a branch of length t

$$Q_{\text{UNR}} = \begin{pmatrix} - & a & b & c \\ d & - & e & f \\ g & h & - & i \\ j & k & l & - \end{pmatrix} \quad Q_{\text{REV}} = \begin{pmatrix} - & a\pi_C & b\pi_G & c\pi_T \\ a\pi_A & - & d\pi_G & e\pi_T \\ b\pi_A & d\pi_C & - & f\pi_T \\ c\pi_A & e\pi_C & f\pi_G & - \end{pmatrix}$$

$$Q_{\text{HKY}} = \begin{pmatrix} - & \pi_C & \kappa\pi_G & \pi_T \\ \pi_A & - & \pi_G & \kappa\pi_T \\ \kappa\pi_A & \pi_C & - & \pi_T \\ \pi_A & \kappa\pi_C & \pi_G & - \end{pmatrix}$$

EVOLUTIONARY RATE

- Variate the rate of evolution by scaling the branches
- Discretize the gamma distribution into k rates

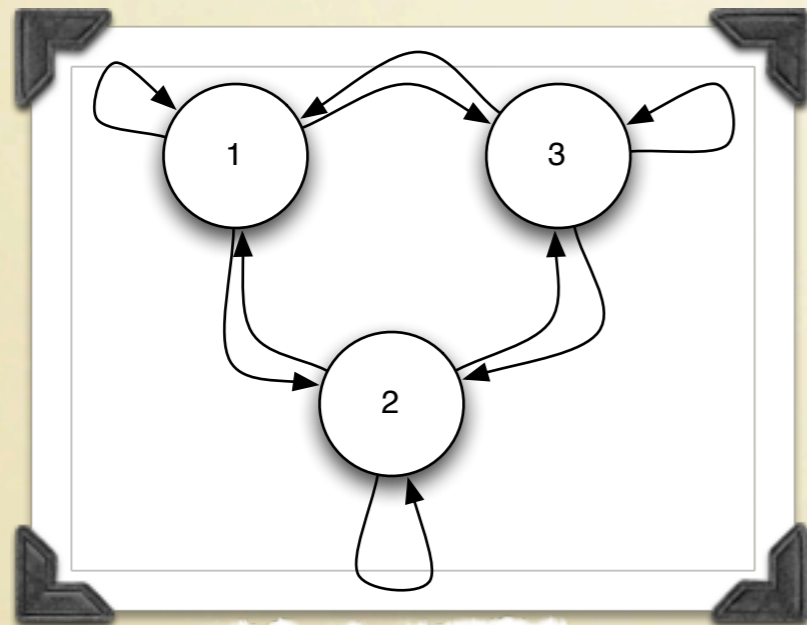
$$P(\mathbf{X}_i|\psi) = \sum_{j=1}^k \frac{1}{k} \cdot P(\mathbf{X}_i|(\mathbf{Q}, \tau, r_j\beta, \pi))$$

How do we assign rates?

Scaling the
branches

RATES HMM

- Autocorrelation (site i is the same as site $i+1$)
- Used in two step fitting process



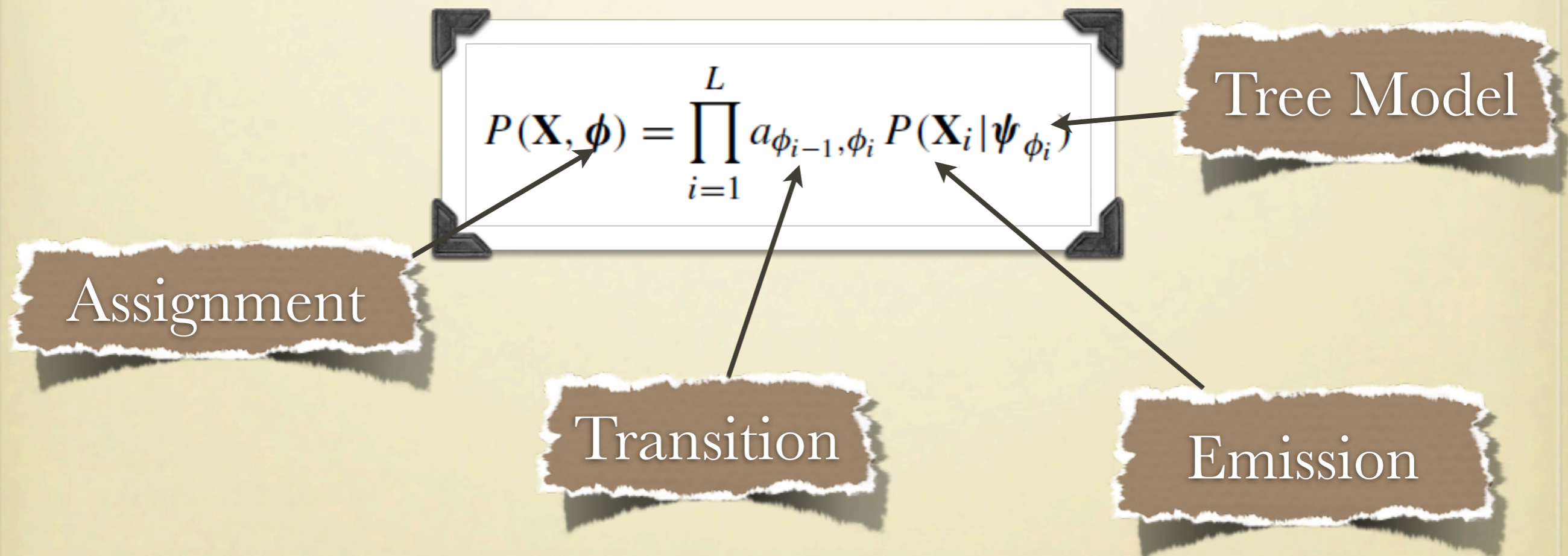
$$k = 3$$

$$c_{j,l} = \frac{1 - \lambda}{k}$$
$$c_{j,j} = \lambda + \frac{1 - \lambda}{k}$$

Transitions

CATEGORY HMMs

- Use tree models for “functional categories”
- Topologies may vary, but are usually the same



CATEGORY X RATES HMM

- Rate and function are orthogonal
- Create HMM that incorporates both
- Take the cross product of states, transitions
 - ➔ scale the tree models

$$\psi'_{i,i'} = r_{i,i'} \psi_i = (\mathbf{Q}_i, \tau_i, r_{i,i'} \beta_i, \pi_i, \alpha_i).$$



What about
slow evolving
coding regions?

HIGHER-ORDER STATES

- Emissions are context-dependent
- Adjust alphabet size to $|\Sigma|^{N+1}$
- Increases complexity
- In practice, $N = 2$ or 3

Complexity

$$O(nL|\Sigma|^{N+1})$$

$$P(L_u|a_1a_2) = \begin{cases} I(a_1a_2 \text{ matches } x_{u,1}x_{u,2}) & \text{if } u \text{ is a leaf} \\ \sum_{b_1b_2} P(b_1b_2|a_1a_2, t_v) P(L_v|b_1b_2) \sum_{c_1c_2} P(c_1c_2|a_1a_2, t_w) P(L_w|c_1c_2) & \text{otherwise} \end{cases}$$

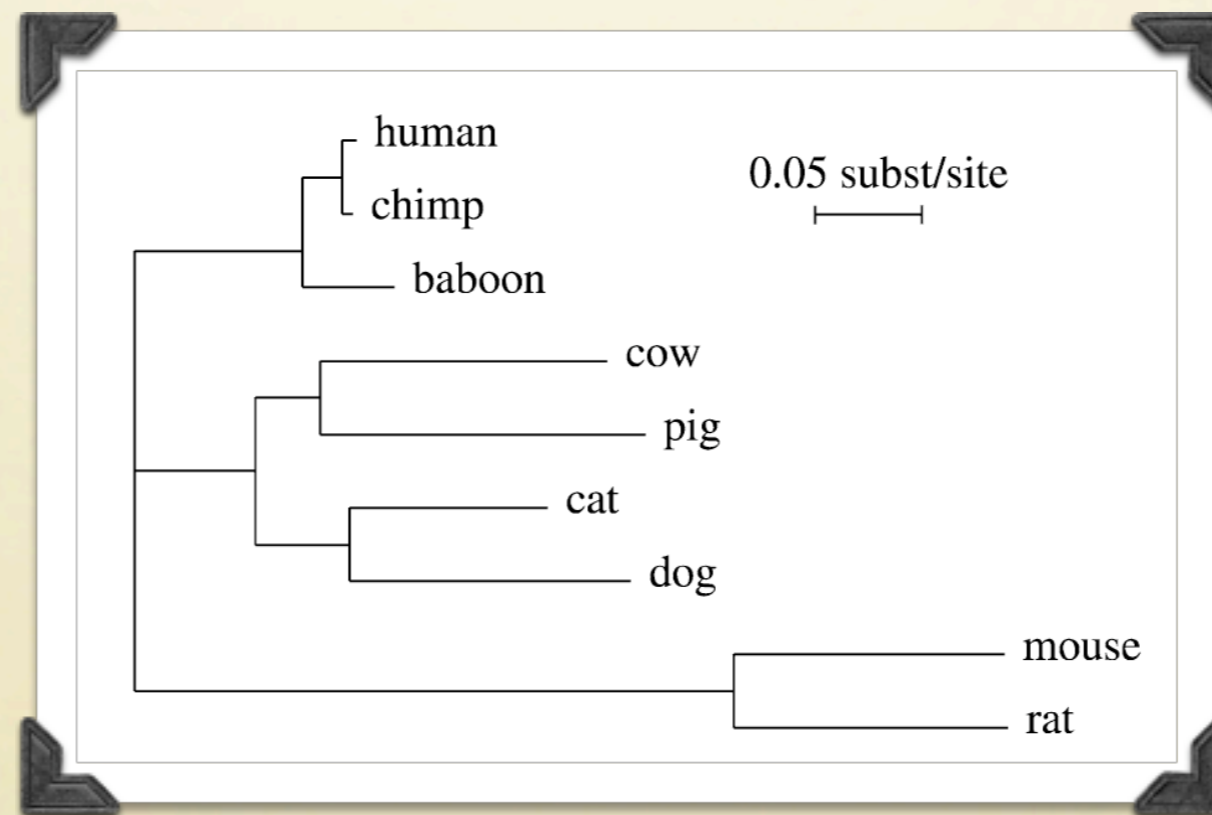
SUMMARY

- Assume k rate, q functional categories
- HMM of order \mathcal{N}
- Estimate transition probabilities of categories
 - ▶ Compute $kq \times L$ emission probabilities
 - ▶ Train autocorrelation
 - ▶ Run Viterbi

THE RESULTS

DATA

- Used portions of huge multiple alignment
- Trained using counting and annotations



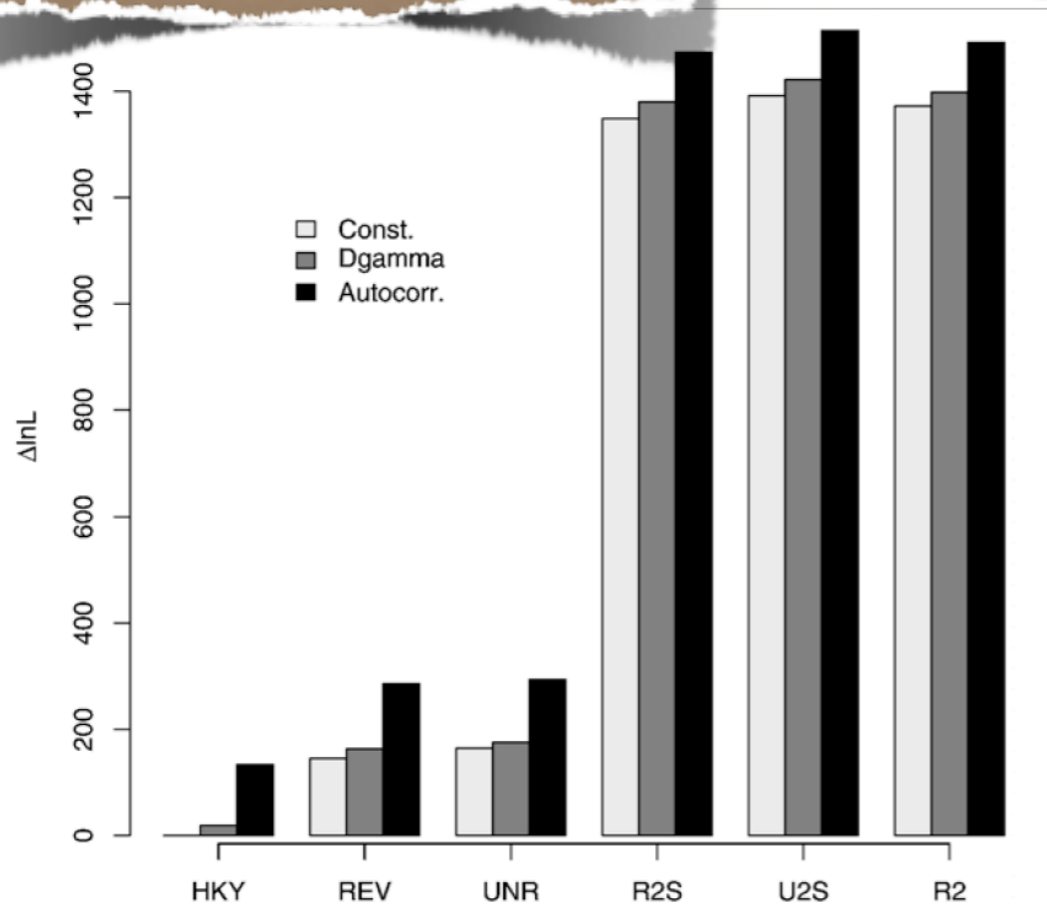
TEST

- Compared models with likelihood ratio test (LRT)
- 5 substitution models (includes higher order)
 - REV, HKY, UNR, R2, R2S, U2S
- 3 rate variations -constant, gamma, autocorrelation

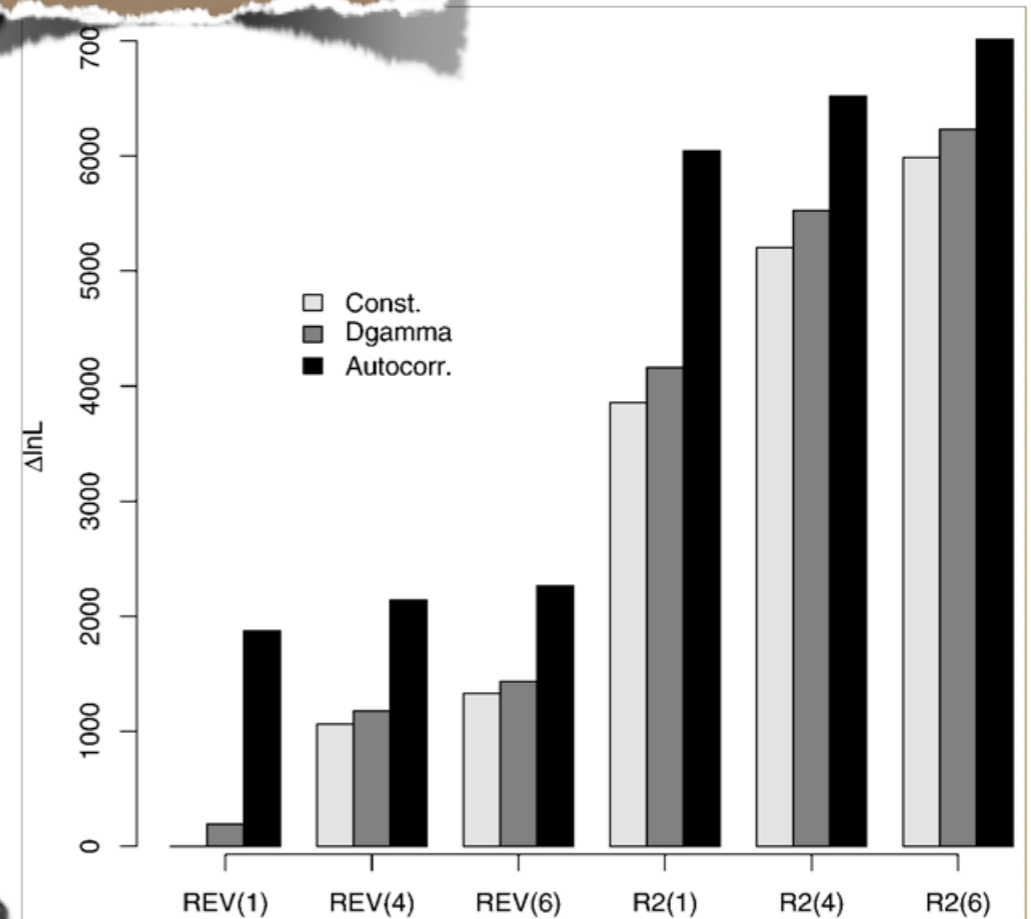
RESULTS

- Higher states give largest boost

Ancestral Repeat



WNT2





QUESTIONS?