

Lecture 26. Phylogeny methods, part 6 (Modeling rate variation among sites)

Joe Felsenstein

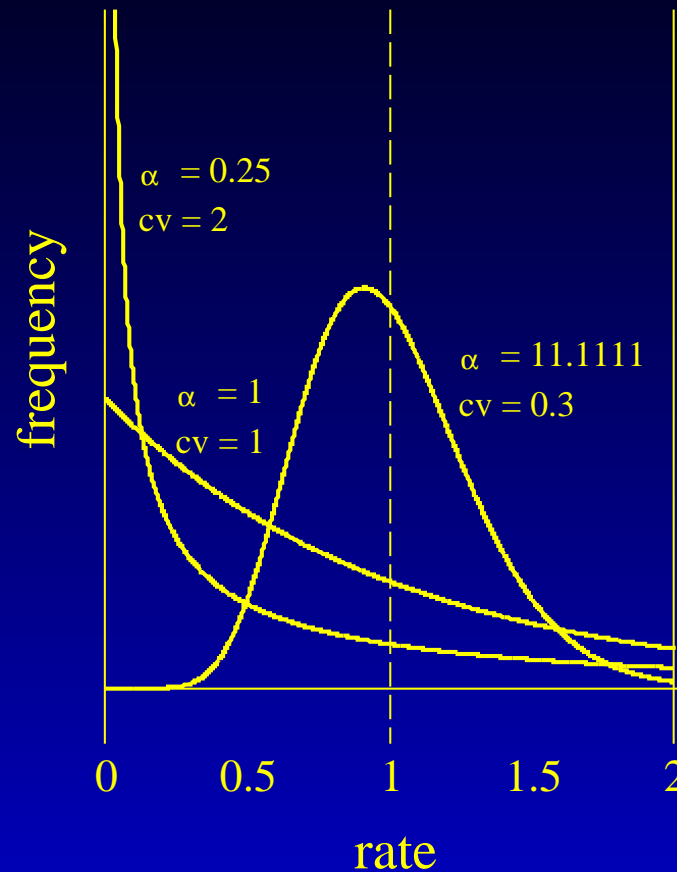
Department of Genome Sciences and Department of Biology

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 (if its mean is 1):

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

Gamma distributions



Gamma distributions with mean 1 and different coefficients of variation (standard deviation / mean). $\alpha = 1/CV^2$ is the “shape parameter” of the Gamma distribution

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 f_0 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 $\mathbf{L}^{(i)}(\mathbf{r}_i)$ is the likelihood of the tree for site i given that the rate of evolution at site i is \mathbf{r}_i , we can integrate this over a gamma density

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

so that the overall likelihood is

$$\mathbf{L} = \prod_{i=1}^m \left[\int_0^{\infty} \mathbf{f}(\mathbf{r}_i; \alpha) \mathbf{L}^{(i)}(\mathbf{r}_i) \mathbf{d}\mathbf{r}_i \right]$$

Unfortunately these integrals cannot be evaluated for trees with more than a few tips as the quantities $\mathbf{L}^{(i)}(\mathbf{r}_i)$ are complicated.

Hidden Markov Models

These are the most widely used models allowing rate variation to be correlated along the sequence.

We assume:

- There are a finite number of rates, m . Rate i is r_i .
- There are probabilities p_i of a site having rate i .
- 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 $\text{Prob}(j|i)$ of changing to rate j in the next site, given that it is at rate i in this site.
- 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.

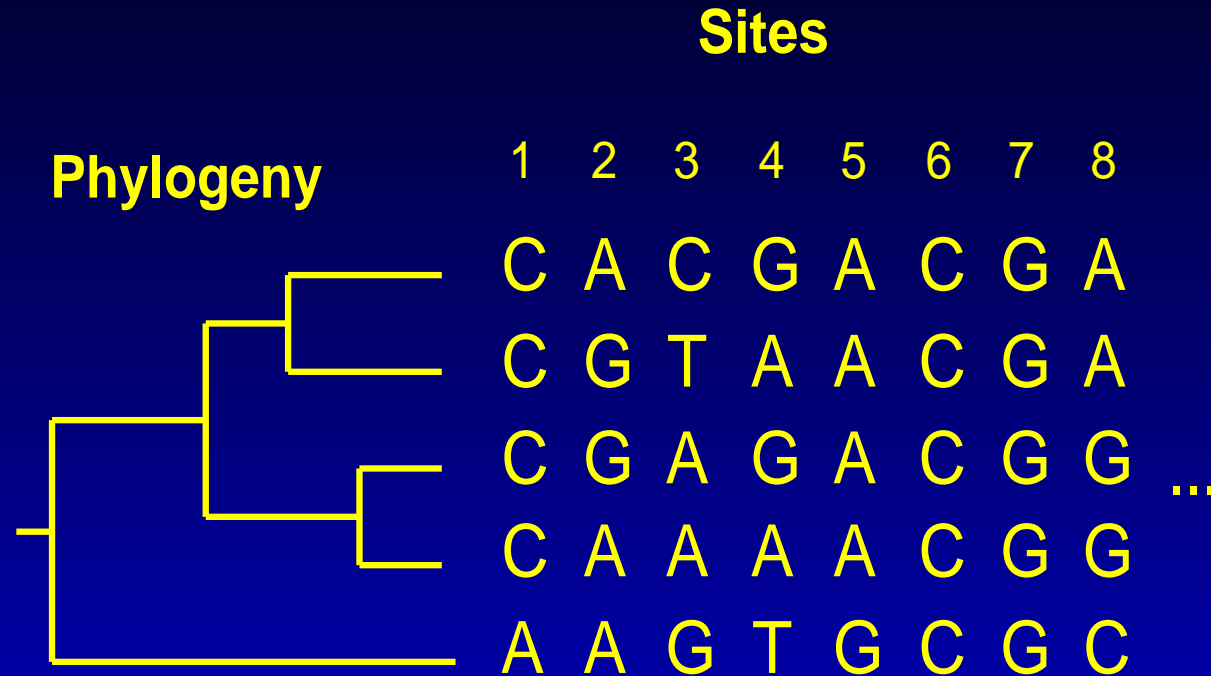
Likelihood with a[n] HMM

Suppose that we have a way of calculating, for each possible rate at each possible site, the probability of the data at that site given that rate. This is

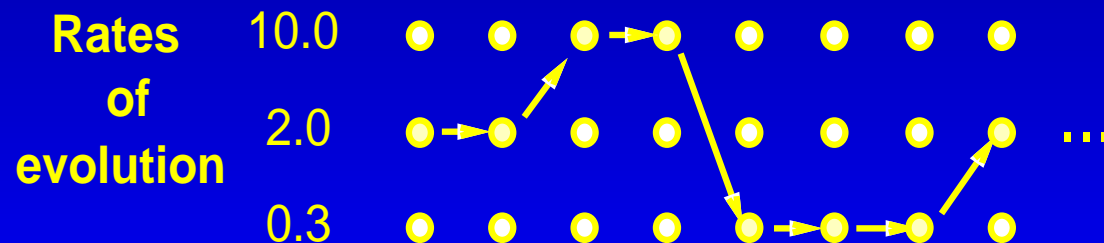
$$\text{Prob} \left(\mathbf{D}^{(i)} \mid \mathbf{r}_j \right)$$

To get the overall probability of all data, sum over all possible paths through the array of sites \times rates, weighting each combination of rates by its probability:

A Hidden Markov Model for rates in a phylogeny



Hidden Markov chain:



Hidden Markov Models

If there are a number of hidden rate states, with state i having rate r_i

$$\begin{aligned} \text{Prob}(\mathbf{D} \mid \mathbf{T}) &= \sum_{\mathbf{i}_1} \sum_{\mathbf{i}_2} \dots \sum_{\mathbf{i}_p} \text{Prob}(\mathbf{r}_{\mathbf{i}_1}, \mathbf{r}_{\mathbf{i}_2}, \dots, \mathbf{r}_{\mathbf{i}_p}) \\ &\quad \times \text{Prob}(\mathbf{D} \mid \mathbf{T}, \mathbf{r}_{\mathbf{i}_1}, \mathbf{r}_{\mathbf{i}_2}, \dots, \mathbf{r}_{\mathbf{i}_m}) \end{aligned}$$

Evolution is independent once each site has had its rate specified

$$\begin{aligned} \text{Prob}(\mathbf{D} \mid \mathbf{T}, \mathbf{r}_1, \mathbf{r}_2, \dots, \mathbf{r}_p) &= \\ \prod_{\mathbf{i}=1}^p \text{Prob}(\mathbf{D}^{(\mathbf{i})} \mid \mathbf{T}, \mathbf{r}_{\mathbf{i}}). \end{aligned}$$

Seems impossible ...

Evolution is independent once each site has had its rate specified

$$\text{Prob}(\mathbf{D}|\mathbf{T}, \mathbf{r}_1, \mathbf{r}_2, \dots, \mathbf{r}_m) = \prod_{i=1}^m \text{Prob}(\mathbf{D}^{(i)}|\mathbf{T}, \mathbf{r}_i).$$

To compute the likelihood we sum over all ways rate states could be assigned to sites:

$$\begin{aligned} \mathbf{L} &= \text{Prob}(\mathbf{D} | \mathbf{T}) \\ &= \sum_{i_1=1}^m \sum_{i_2=1}^m \dots \sum_{i_p=1}^m \text{Prob}(\mathbf{r}_{i_1}, \mathbf{r}_{i_2}, \dots, \mathbf{r}_{i_p}) \\ &\quad \times \text{Prob}(\mathbf{D}^{(1)} | \mathbf{r}_{i_1}) \text{Prob}(\mathbf{D}^{(2)} | \mathbf{r}_{i_2}) \dots \text{Prob}(\mathbf{D}^{(n)} | \mathbf{r}_{i_p}) \end{aligned}$$

Problem: The number of rate combinations is very large. With 100 sites and 3 rates at each, it is $3^{100} \simeq 5 \times 10^{47}$. This makes the summation impractical.

Factorization and the algorithm

Fortunately, the terms can be reordered:

$$\begin{aligned} \mathbf{L} &= \text{Prob}(\mathbf{D} \mid \mathbf{T}) \\ &= \sum_{\mathbf{i}_1=1}^m \sum_{\mathbf{i}_2=1}^m \dots \sum_{\mathbf{i}_p=1}^m \text{Prob}(\mathbf{i}_1) \text{Prob}(\mathbf{D}^{(1)} \mid \mathbf{r}_{\mathbf{i}_1}) \\ &\quad \times \text{Prob}(\mathbf{i}_2 \mid \mathbf{i}_1) \text{Prob}(\mathbf{D}^{(2)} \mid \mathbf{r}_{\mathbf{i}_2}) \\ &\quad \times \text{Prob}(\mathbf{i}_3 \mid \mathbf{i}_2) \text{Prob}(\mathbf{D}^{(3)} \mid \mathbf{r}_{\mathbf{i}_3}) \\ &\quad \vdots \\ &\quad \times \text{Prob}(\mathbf{i}_p \mid \mathbf{i}_{p-1}) \text{Prob}(\mathbf{D}^{(p)} \mid \mathbf{r}_{\mathbf{i}_p}) \end{aligned}$$

Using Horner's Rule

and the summations can be moved each as far rightwards as it can go:

$$\begin{aligned} \mathbf{L} = & \sum_{\mathbf{i}_1=1}^{\mathbf{m}} \text{Prob}(\mathbf{i}_1) \text{Prob}(\mathbf{D}^{(1)} | \mathbf{r}_{\mathbf{i}_1}) \\ & \sum_{\mathbf{i}_2=1}^{\mathbf{m}} \text{Prob}(\mathbf{i}_2 | \mathbf{i}_1) \text{Prob}(\mathbf{D}^{(2)} | \mathbf{r}_{\mathbf{i}_2}) \\ & \sum_{\mathbf{i}_3=1}^{\mathbf{m}} \text{Prob}(\mathbf{i}_3 | \mathbf{i}_2) \text{Prob}(\mathbf{D}^{(3)} | \mathbf{r}_{\mathbf{i}_3}) \\ & \vdots \\ & \sum_{\mathbf{i}_p=1}^{\mathbf{m}} \text{Prob}(\mathbf{i}_p | \mathbf{i}_{p-1}) \text{Prob}(\mathbf{D}^{(p)} | \mathbf{r}_{\mathbf{i}_p}) \end{aligned}$$

Recursive calculation of HMM likelihoods

The summations can be evaluated innermost-outwards. The same summations appear in multiple terms. We can then evaluate them only once. A huge saving results. The result is this algorithm:

Define $\mathcal{P}_i(j)$ as the probability of everything at or to the right of site i , given that site i has the j -th rate.

Now we can immediately see for the last site that for each possible rate category i_p

$$\mathcal{P}_p(i_p) = \text{Prob} \left(\mathbf{D}^{(p)} \mid \mathbf{r}_{i_p} \right)$$

(as “at or to the right of” simply means “at” for that site).

Recursive calculation

More generally, for site $\ell < p$ and its rates \mathbf{i}_ℓ

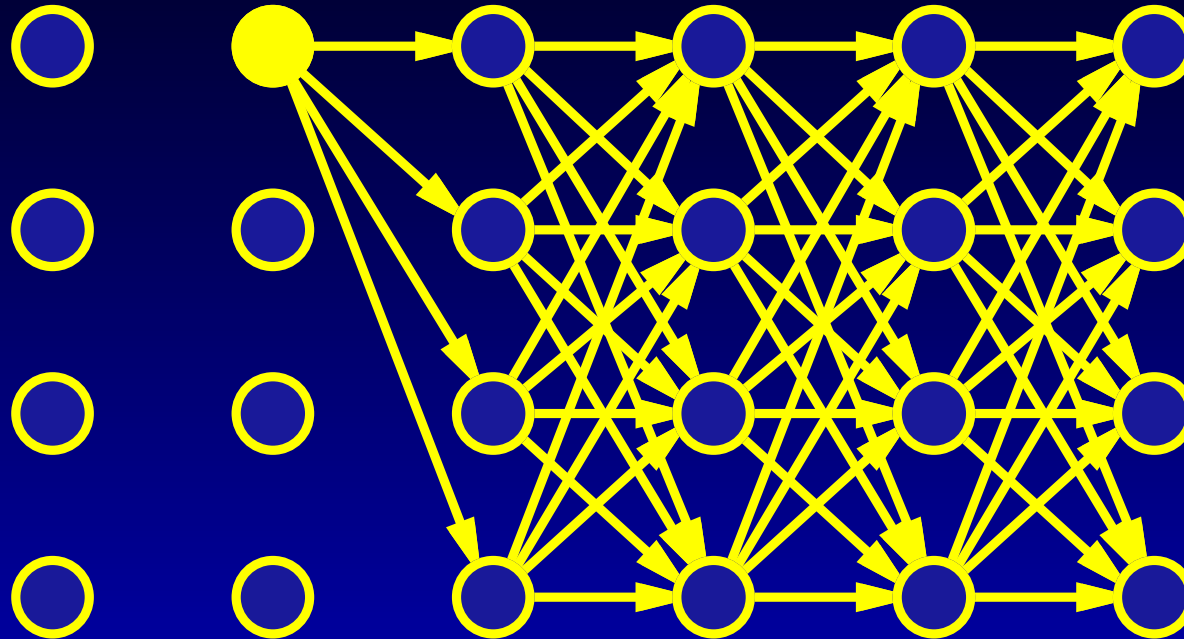
$$\mathcal{P}_\ell(\mathbf{i}_\ell) = \text{Prob} \left(\mathbf{D}^{(\ell)} \mid \mathbf{r}_{\mathbf{i}_\ell} \right) \sum_{\mathbf{i}_{\ell+1}=1}^m \text{Prob} (\mathbf{i}_{\ell+1} \mid \mathbf{i}_\ell) \mathcal{P}_{\ell+1}(\mathbf{i}_{\ell+1})$$

We can compute the \mathcal{P} 's recursively using this, starting with the last site and moving leftwards down the sequence. Finally we have the $\mathcal{P}_1(\mathbf{i}_1)$ for all m states. These are simply weighted by the equilibrium probabilities of the Markov chain of rate categories:

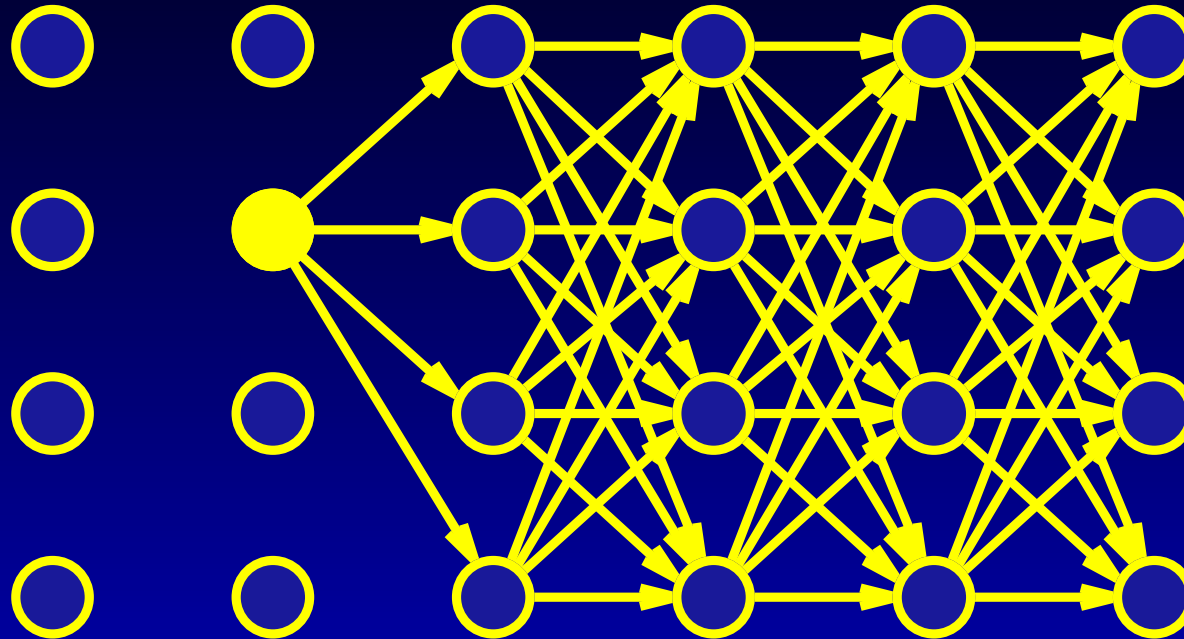
$$\mathbf{L} = \text{Prob} (\mathbf{D} \mid \mathbf{T}) = \sum_{\mathbf{i}_1=1}^m \text{Prob} (\mathbf{i}_1) \mathcal{P}_1(\mathbf{i}_1)$$

An entirely similar calculation can be done from left to right, remembering that the transition probabilities $\text{Prob} (\mathbf{i}_k \mid \mathbf{i}_{k+1})$ would be different in that case.

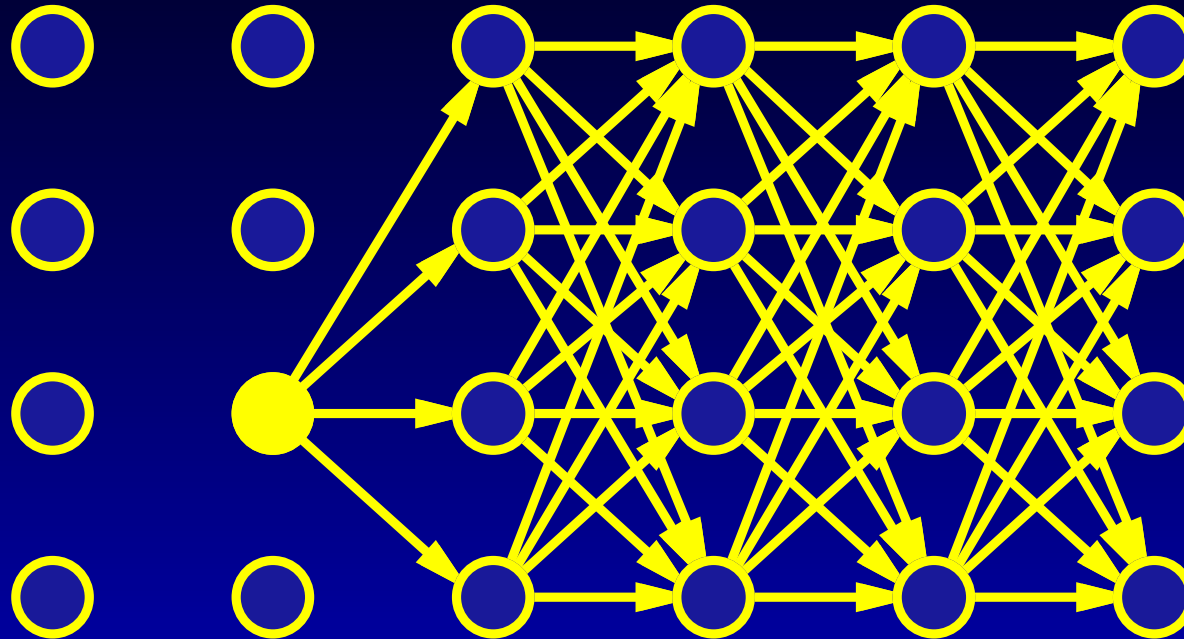
Computing all possible paths from one node.



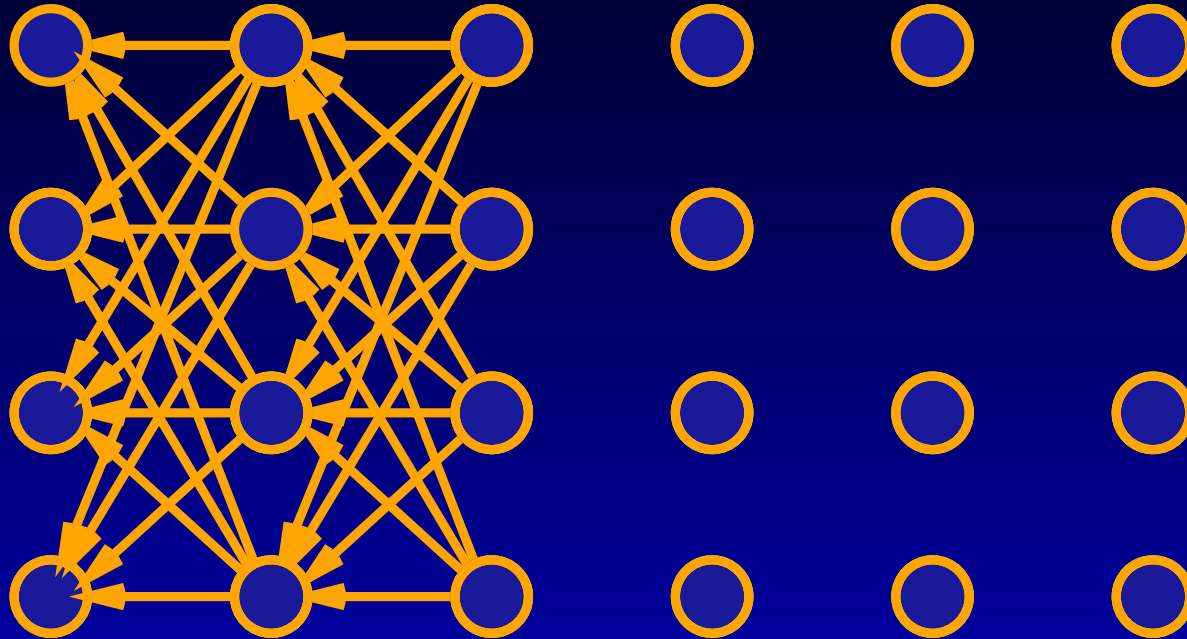
Computing all possible paths starting at another node



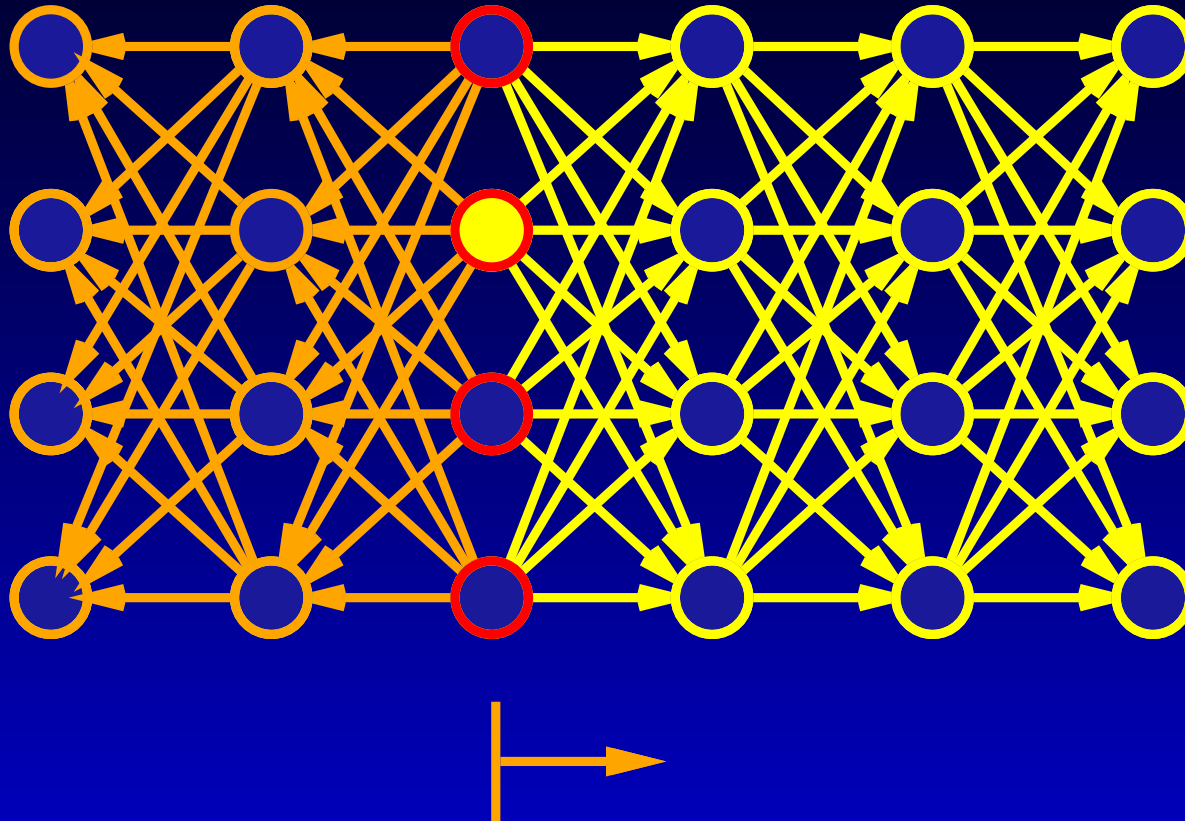
Note that they re-use the same terms.



The same algorithm can work forwards.



The forwards-backwards algorithm.



References

- Churchill, G.A. 1989. Stochastic models for heterogeneous DNA sequences. *Bulletin of Mathematical Biology* **51**: 79-94. [First paper using HMMs on sequences, without trees]
- Felsenstein, J. and G. A. Churchill. 1996. A hidden Markov model approach to variation among sites in rate of evolution *Molecular Biology and Evolution* **13**: 93-104. [HMMs for rates in ML trees]
- Jin, L. and M. Nei. 1990. Limitations of the evolutionary parsimony method of phylogenetic analysis. *Molecular Biology and Evolution* **7**: 82-102. [Gamma distributed rates in a distance]
- Olsen, G. J. 1987. Earliest phylogenetic branchings: comparing rRNA-based evolutionary trees inferred with various techniques. *Cold Spring Harbor Symposia on Quantitative Biology* **52**: 825-837. [Lognormal rates in a distance]

references, cont'd

- Waddell, P. J. and M. A. Steel. 1997. General time-reversible distances with unequal rates across sites: mixing Γ and inverse Gaussian distributions with invariant sites. *Molecular Phylogenies and Evolution* **8**: 398-414. [**More generalized rate distributions in a more generalized distance**]
- Yang, Z. 1994. Maximum-likelihood estimation of phylogeny from DNA sequences when substitution rates differ over sites. *Molecular Biology and Evolution* **10**: 1396-1401. [**Rates varying in gamma distribution in an ML tree method for few species**]
- Yang, Z. 1994. Maximum likelihood phylogenetic estimation from DNA sequences with variable rates over sites: approximate methods. *Journal of Molecular Evolution* **39**: 306-314. [**First paper on HMMs for ML trees, used to approximate Gamma distributions for more species**]
- Yang, Z. 1995. A space-time process model for the evolution of DNA sequences. *Genetics* **139**: 993-1005. [**Also allowing for correlated rates along the molecule**]

How it was done

This projection produced as a PDF, not a PowerPoint file, and viewed using the Full Screen mode (in the View menu of Adobe Acrobat Reader):

- using the `prospect` style in LaTeX,
- using LaTeX to make a `.dvi` file,
- using `dvi2ps` to turn this into a Postscript file,
- using `ps2pdf` to mill it into a PDF file, and
- displaying the slides in Adobe Acrobat Reader.

Result: nice slides using freeware.