



Joint Institute
of Engineering

SUN YAT-SEN UNIVERSITY

Carnegie Mellon University

Design and Implementation of Speech Recognition Systems

Fall 2014

Ming Li

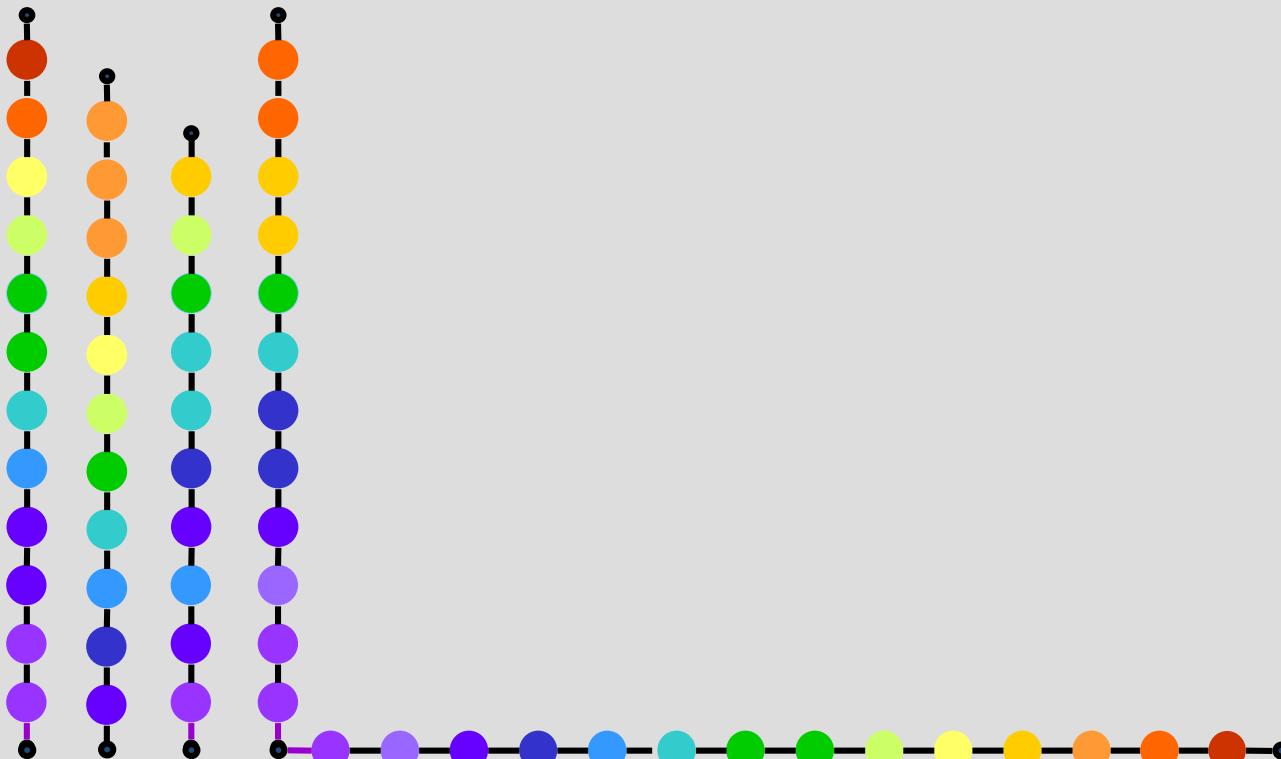
Class 7-8: HMMs

Oct 9, 13th

Thanks to Professor Bhiksha Raj for the contribution of the slides

DTW with multiple templates

TEMPLATES

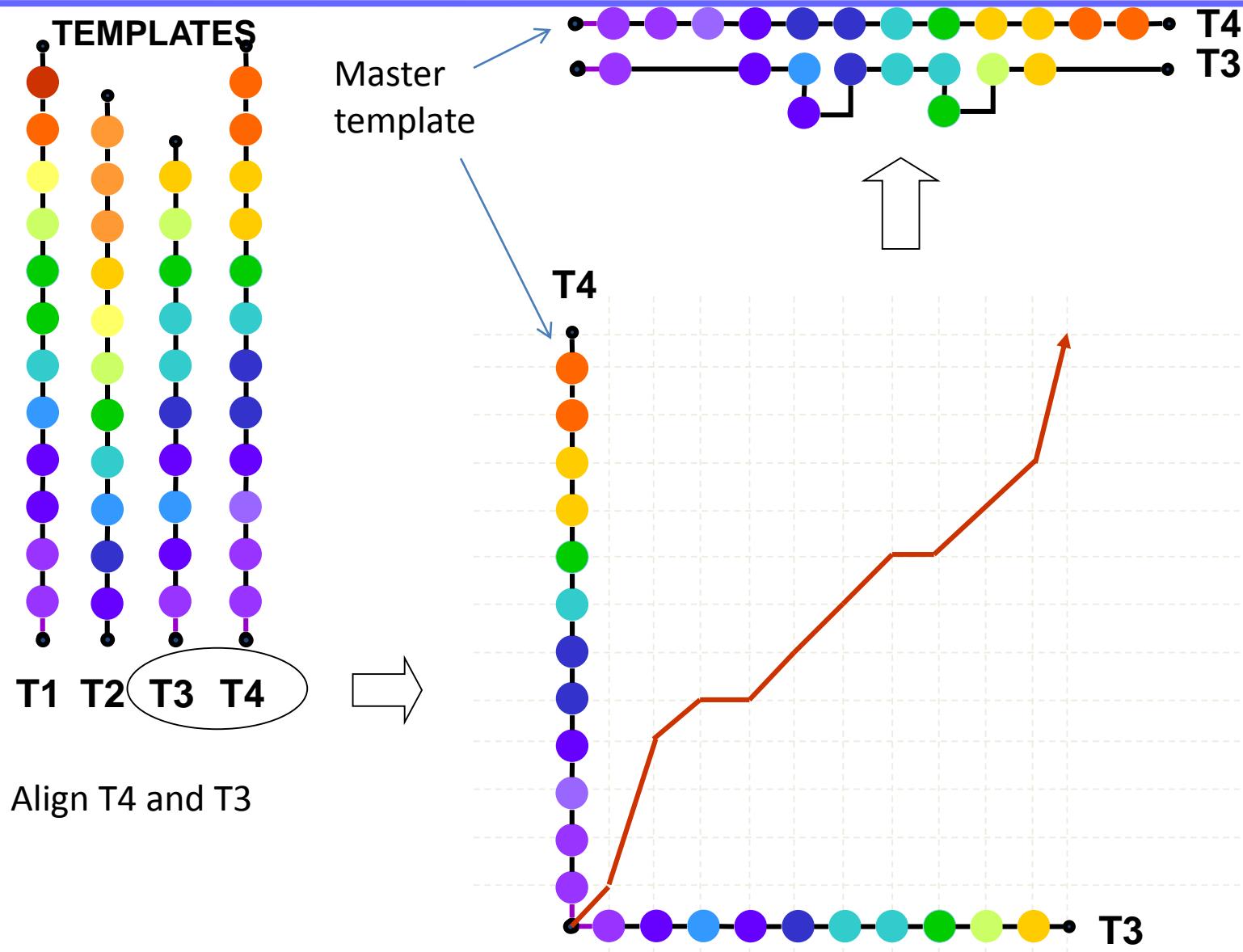


DATA

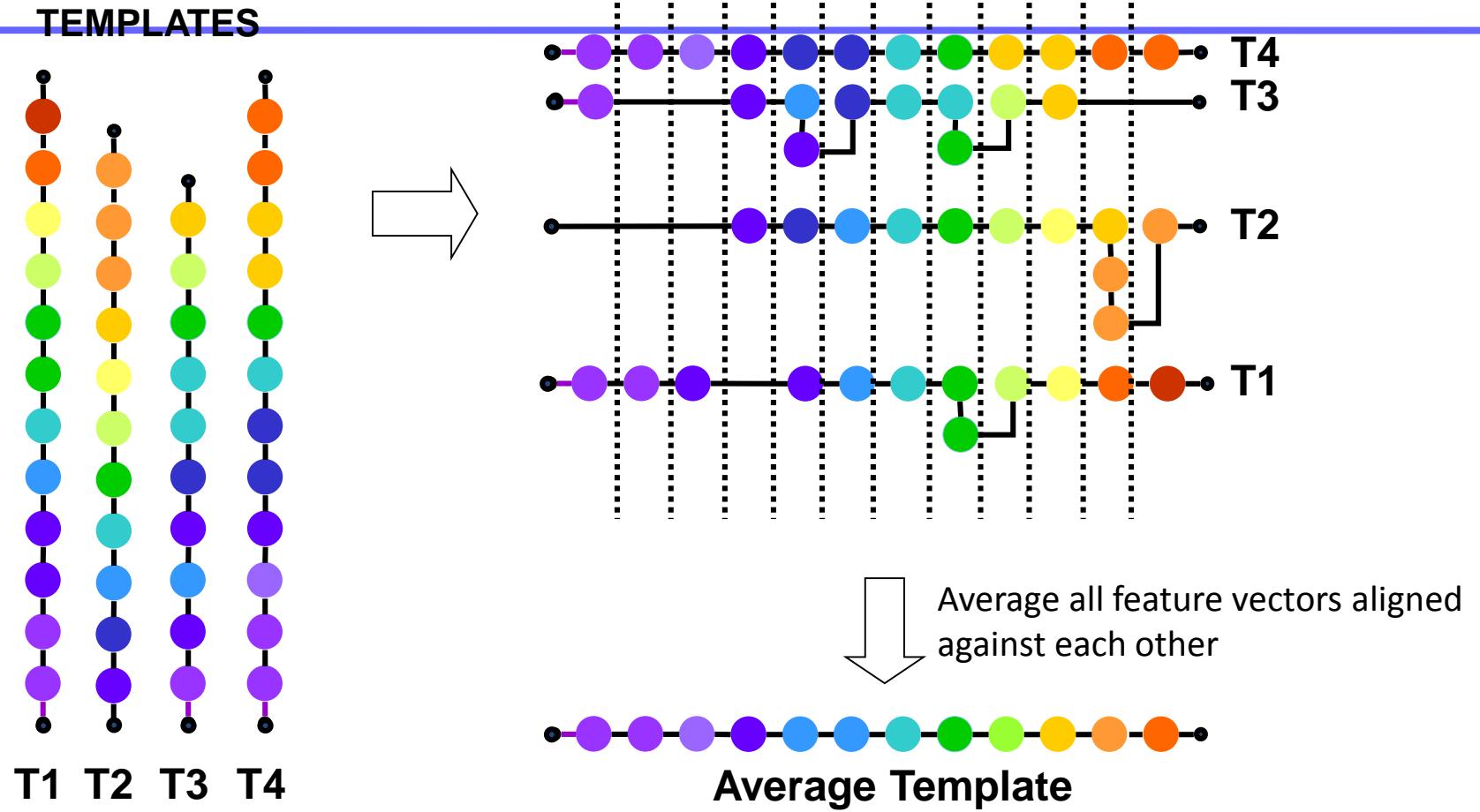
Improving the Templates

- Generalization by averaging the templates
- Generalization by reducing template length
- Accounting for variation within templates represented by the reduced model
- Accounting for varying segment lengths

DTW with multiple templates



DTW with multiple templates

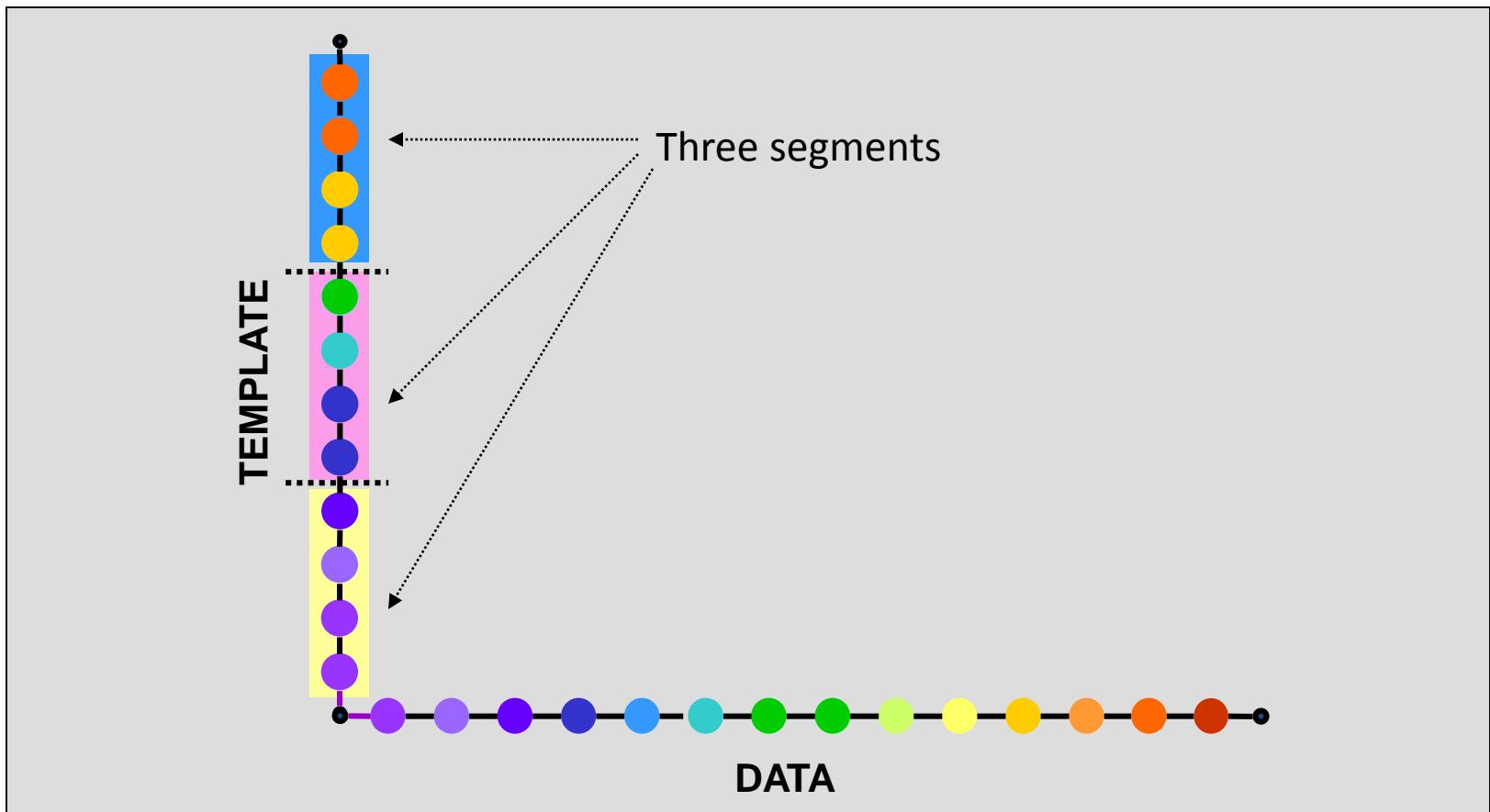


Align T4/T2 and T4/T1, similarly; then average all of them

Improving the Templates

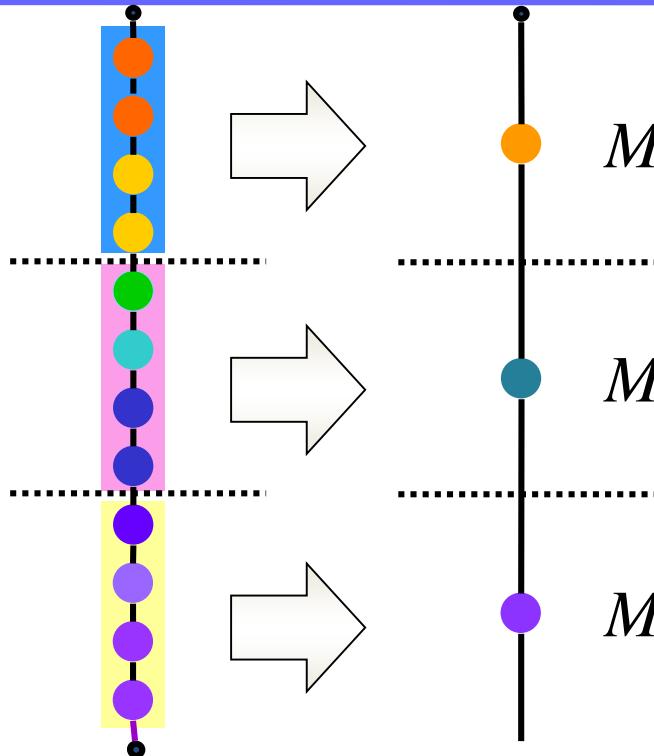
- Generalization by averaging the templates
- Generalization by reducing template length
- Accounting for variation within templates represented by the reduced model
- Accounting for varying segment lengths

Example: Single Templates With Three Segments



The feature vectors within each segment are assumed to be similar to each other

Averaging Each Template Segment



$$\text{Model Vector} = \frac{1}{N} \sum_{i \in \text{segment}} \text{vector}(i)$$

$$\text{Model Vector} = \frac{1}{N} \sum_{i \in \text{segment}} \text{vector}(i)$$

$$\text{Model Vector} = \frac{1}{N} \sum_{i \in \text{segment}} \text{vector}(i)$$

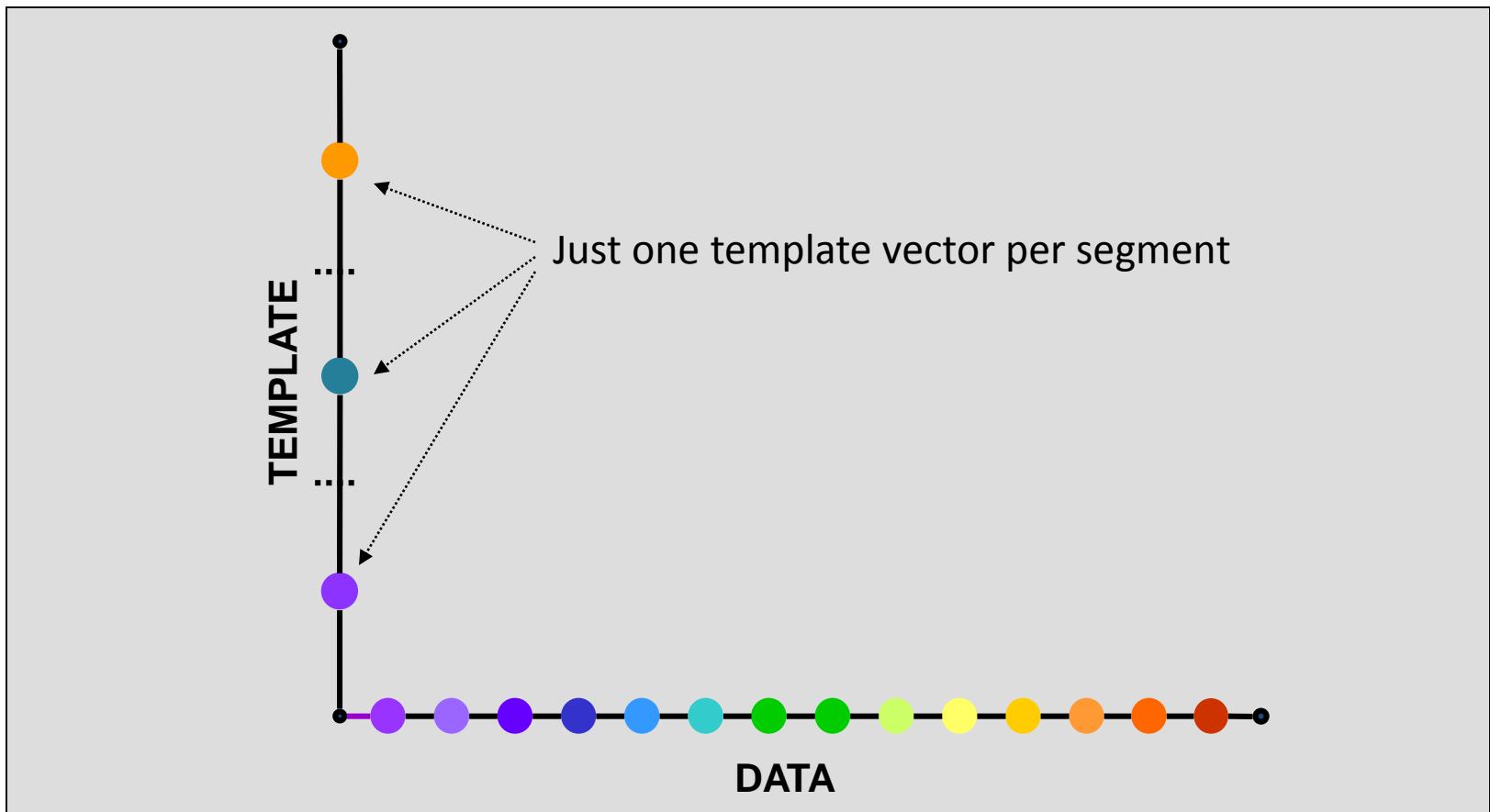
$$m_j = \frac{1}{N_j} \sum_{i \in \text{segment}(j)} x(i)$$

m_j is the *model* vector for the j^{th} segment

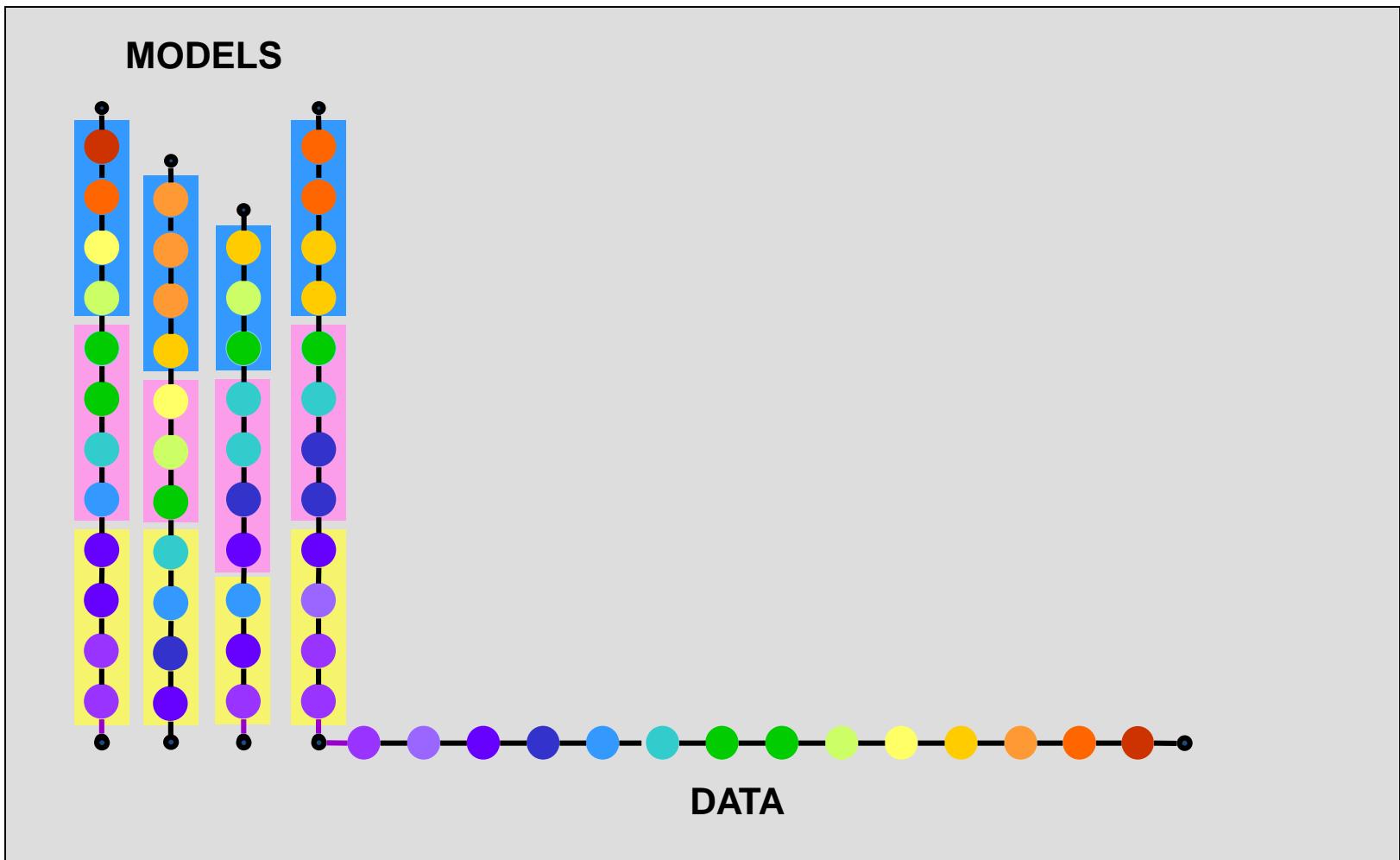
N_j is the number of vectors in the j^{th} segment

$x(i)$ is the i^{th} feature vector

Template With One Model Vector Per Segment



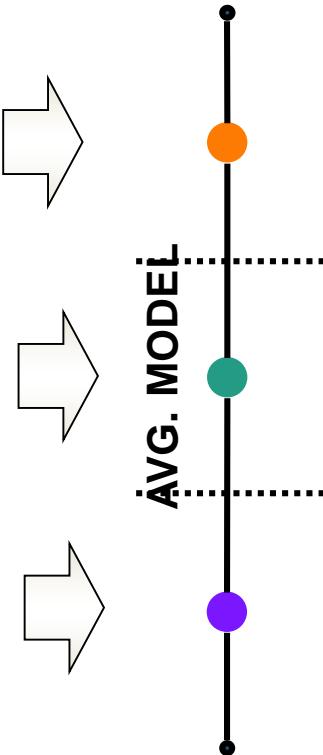
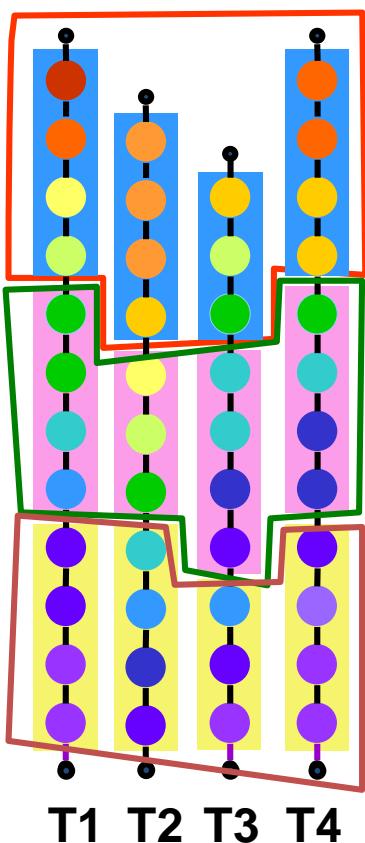
DTW with multiple models



Segment all templates
Average each region into a single point

DTW with multiple models

MODELS



$$m_j = \frac{1}{\sum_k N_{k,j}} \sum_{i \in segment_k(j)} x_k(i)$$

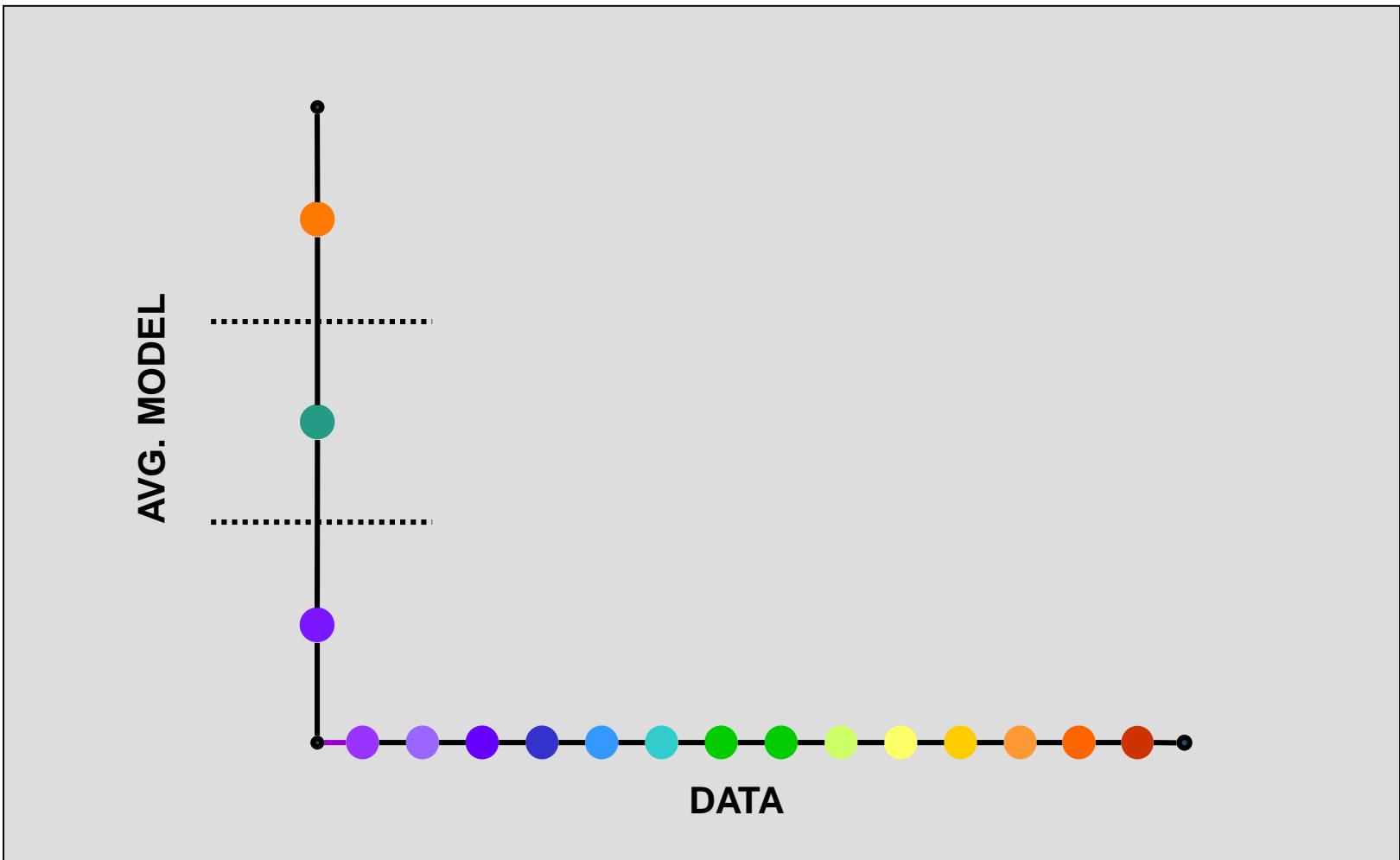
$segment_k(j)$ is the j^{th} segment of the k^{th} training sequence

m_j is the model vector for the j^{th} segment

$N_{k,j}$ is the number of training vectors in the j^{th} segment of the k^{th} training sequence

$x_k(i)$ is the i^{th} vector of the k^{th} training sequence

DTW with multiple models



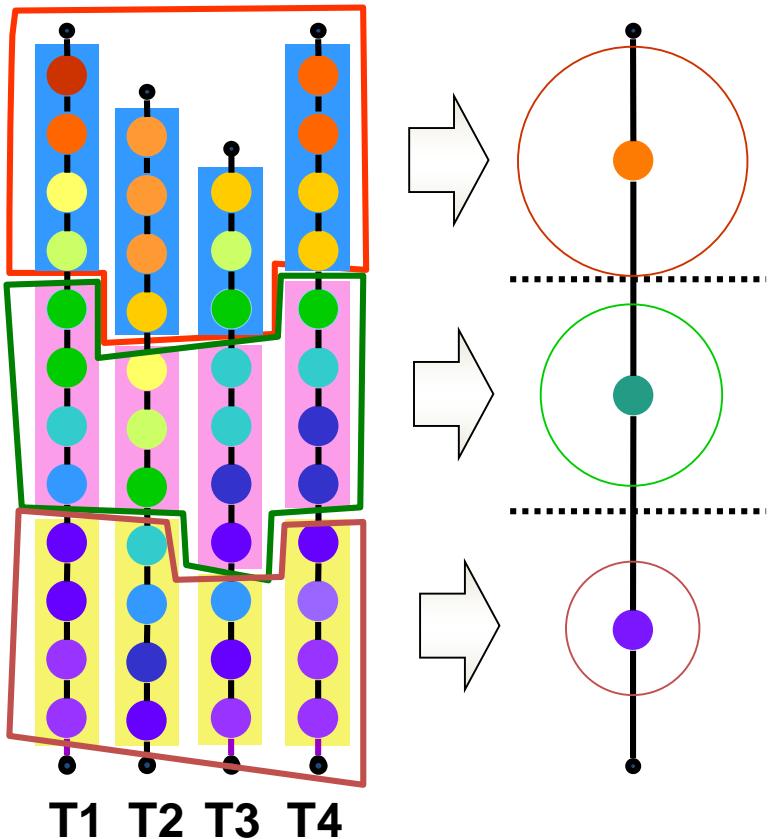
Segment all templates, average each region into a single point
To get a simple average model, which is used for recognition

Improving the Templates

- Generalization by averaging the templates
- Generalization by reducing template length
- Accounting for variation within templates represented by the reduced model
- Accounting for varying segment lengths

DTW with multiple models

MODELS



We can define the covariance for each segment using the standard formula for covariance

$$C_j = \frac{1}{\sum_k N_{k,j}} \sum_{i \in segment_k(j)} (x_k(i) - m_j) (x_k(i) - m_j)^T$$

m_j is the *model* vector for the j^{th} segment

C_j is the covariance of the vectors in the j^{th} segment

DTW with multiple models

- The distance function must be modified to account for the covariance
- Mahalanobis distance:
 - Normalizes contribution of all dimensions of the data
$$d(x, m_j) = (x - m_j)^T C_j^{-1} (x - m_j)$$
 - x is a data vector, m_j is the mean of a segment, C_j is the covariance matrix for the segment
- Negative Gaussian log likelihood:
 - Assumes a Gaussian distribution for the segment and computes the probability of the vector on this distribution

$$\text{Gaussian}(x; m_j, C_j) = \frac{1}{\sqrt{(2\pi)^D |C_j|}} e^{-0.5(x - m_j)^T C_j^{-1} (x - m_j)}$$

$$d(x, m_j) = -\log(\text{Gaussian}(x; m_j, C_j)) = 0.5 \log((2\pi)^D |C_j|) + 0.5(x - m_j)^T C_j^{-1} (x - m_j)$$

The Covariance

- The variance that we have computed is a *full covariance matrix*
 - And the distance measure requires a matrix inversion

$$C_j = \frac{1}{\sum_k N_k} \sum_k \sum_{i \in segment_k(j)} (x_k(i) - m_j) (x_k(i) - m_j)^T$$

$$d(x, m_j) = (x - m_j)^T C_j^{-1} (x - m_j)$$

$$\Sigma = \begin{pmatrix} \sigma_1^2 & \sigma_{12} & \cdots & \sigma_{1N} \\ \sigma_{21} & \sigma_2^2 & \cdots & \sigma_{2N} \\ \vdots & \vdots & \ddots & \vdots \\ \sigma_{N1} & \sigma_{N2} & \cdots & \sigma_N^2 \end{pmatrix}$$

- In practice we assume that all off-diagonal terms in the matrix are 0**
- This reduces our distance metric to:

$$d(x, m_j) = \sum_l \frac{(x_l - m_{j,l})^2}{\sigma_{j,l}^2}$$

- Where the individual variance terms σ^2 are

$$\sigma_{j,l}^2 = \frac{1}{\sum_k N_k} \sum_k \sum_{i \in segment_k(j)} (x_{k,l}(i) - m_{j,l})^2$$

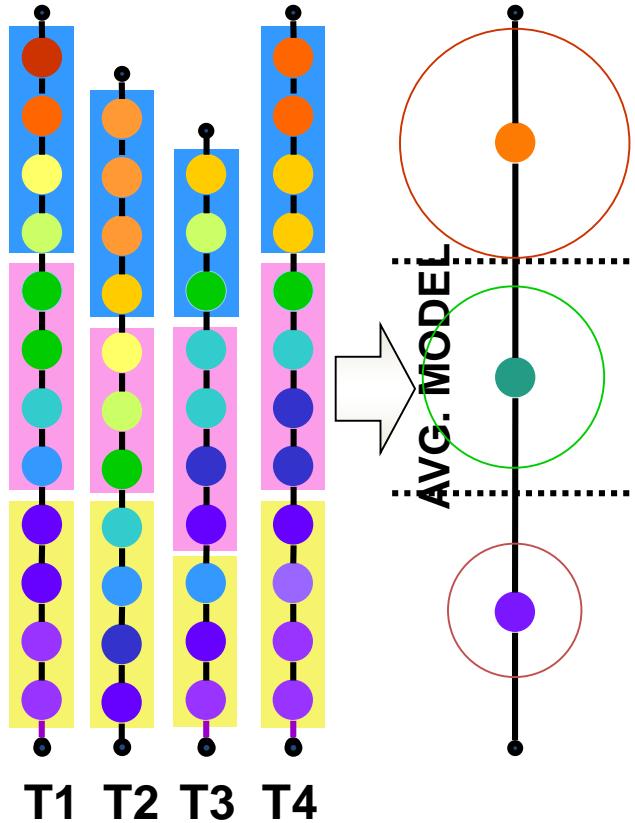
- If we use a negative log Gaussian instead, the modified score (with the diagonal covariance) is

$$d(x, m_j) = 0.5 \sum_l \log(2\pi\sigma_{j,l}^2) + 0.5 \sum_l \frac{(x_l - m_{j,l})^2}{\sigma_{j,l}^2}$$

$$\Sigma = \begin{pmatrix} \sigma_1^2 & 0 & \cdots & 0 \\ 0 & \sigma_2^2 & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & \sigma_N^2 \end{pmatrix}$$

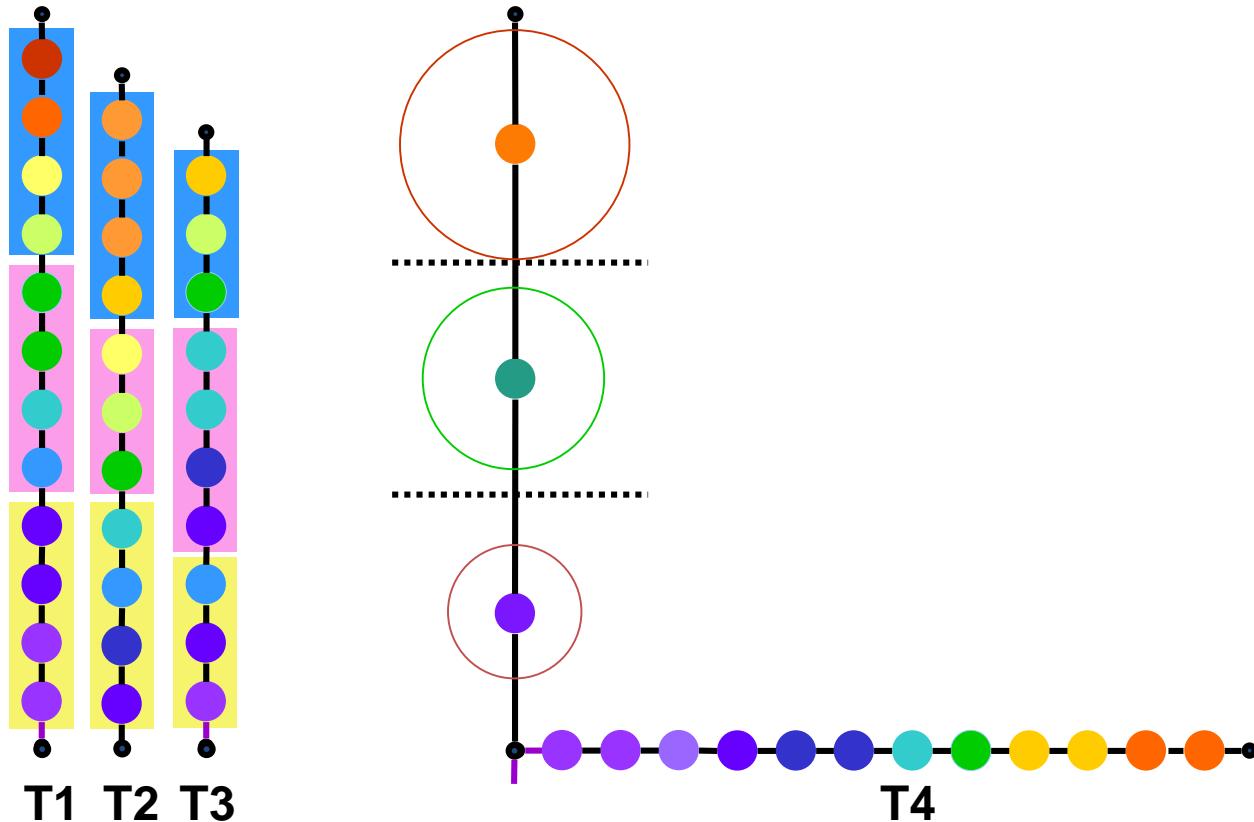
Segmental K-means

MODELS



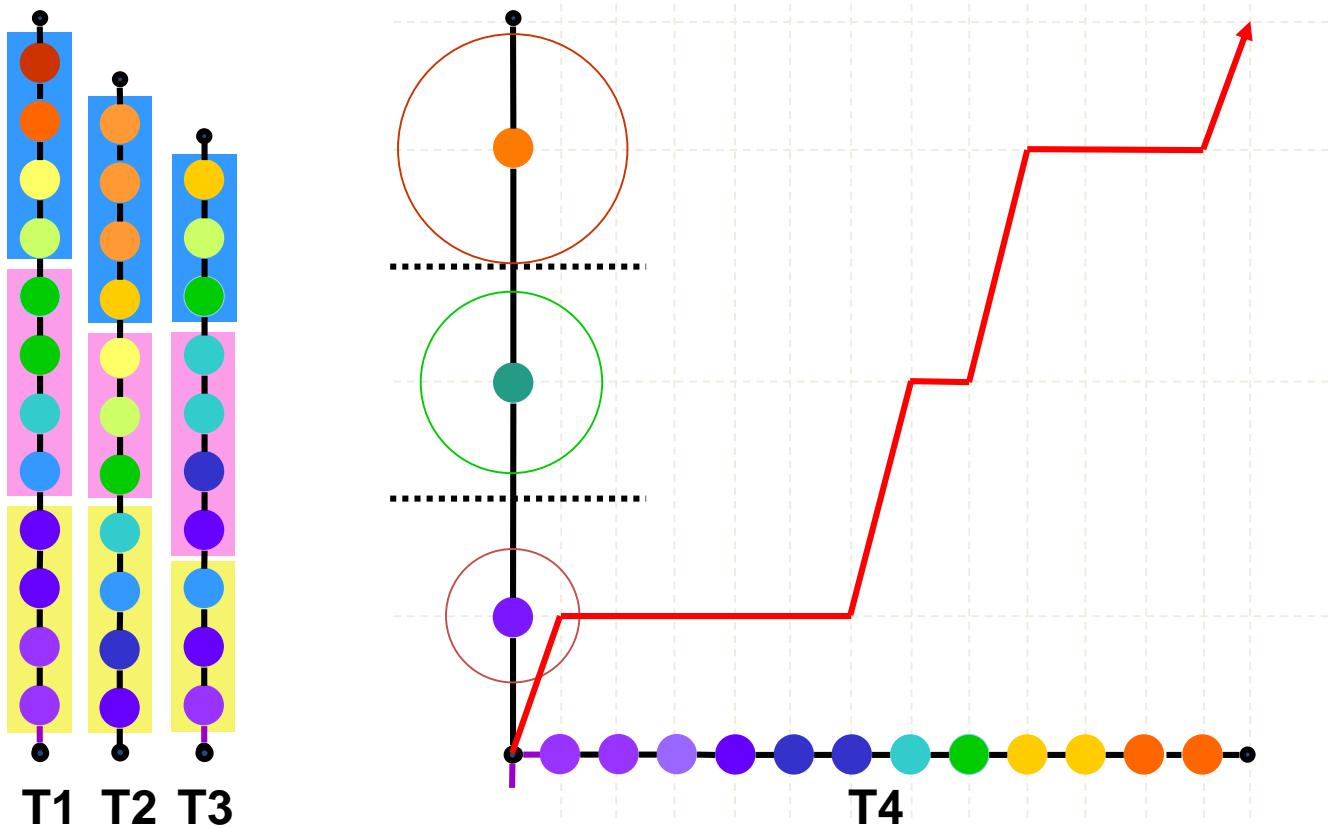
Initialize by uniform segmentation

Alignment for training a model from multiple vector sequences



Initialize by uniform segmentation

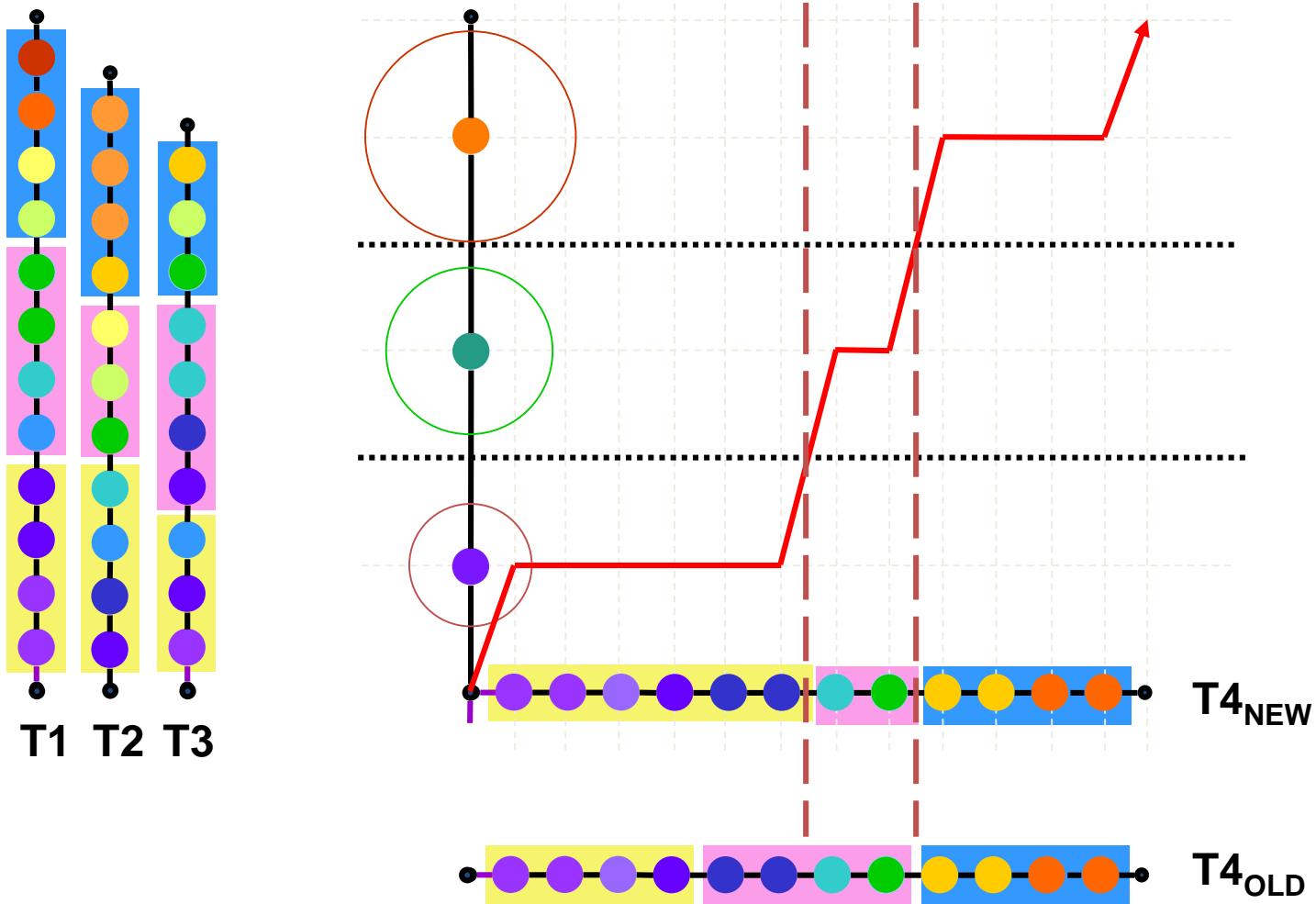
Alignment for training a model from multiple vector sequences



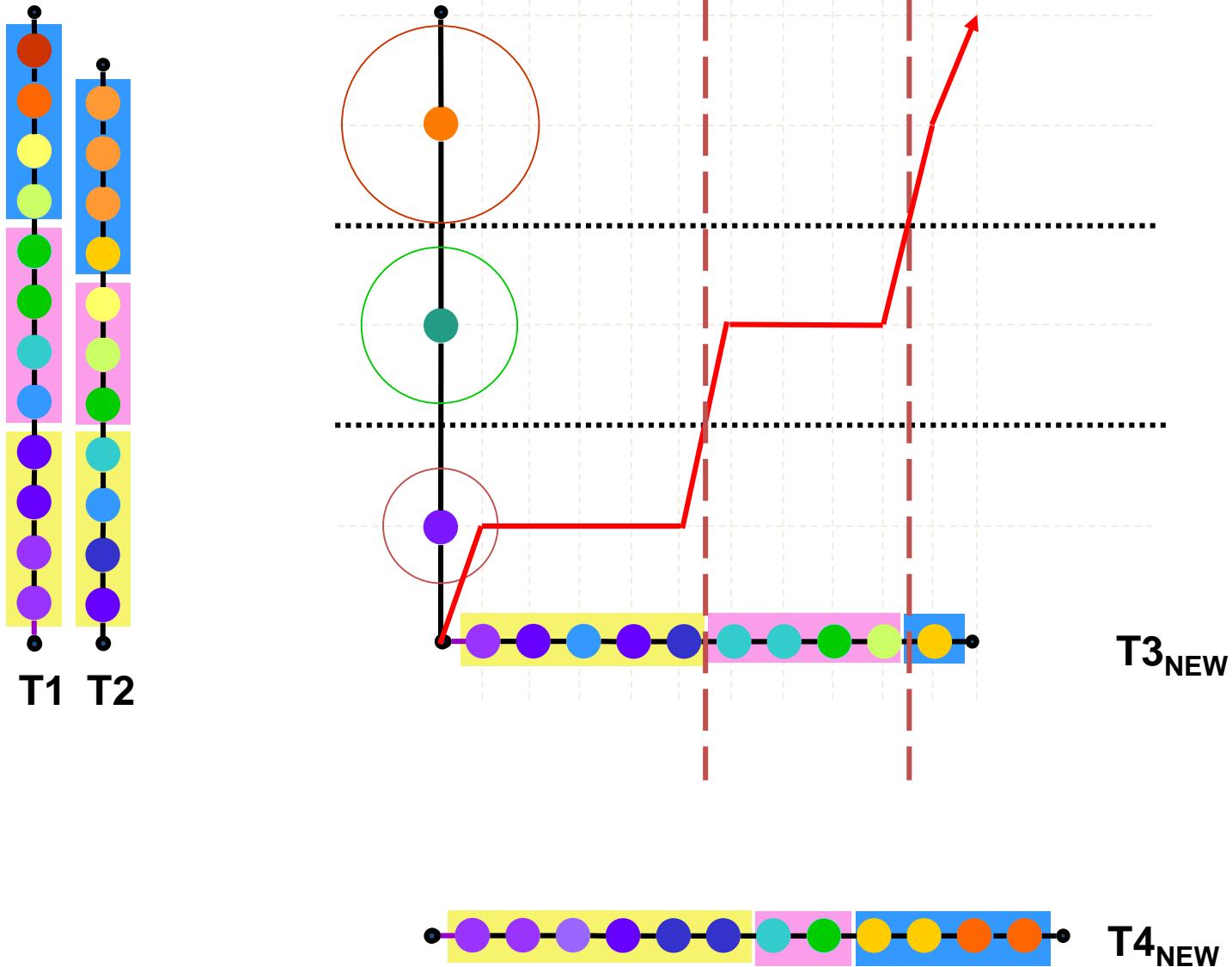
Initialize by uniform segmentation

Align each template to the averaged model to get new segmentations

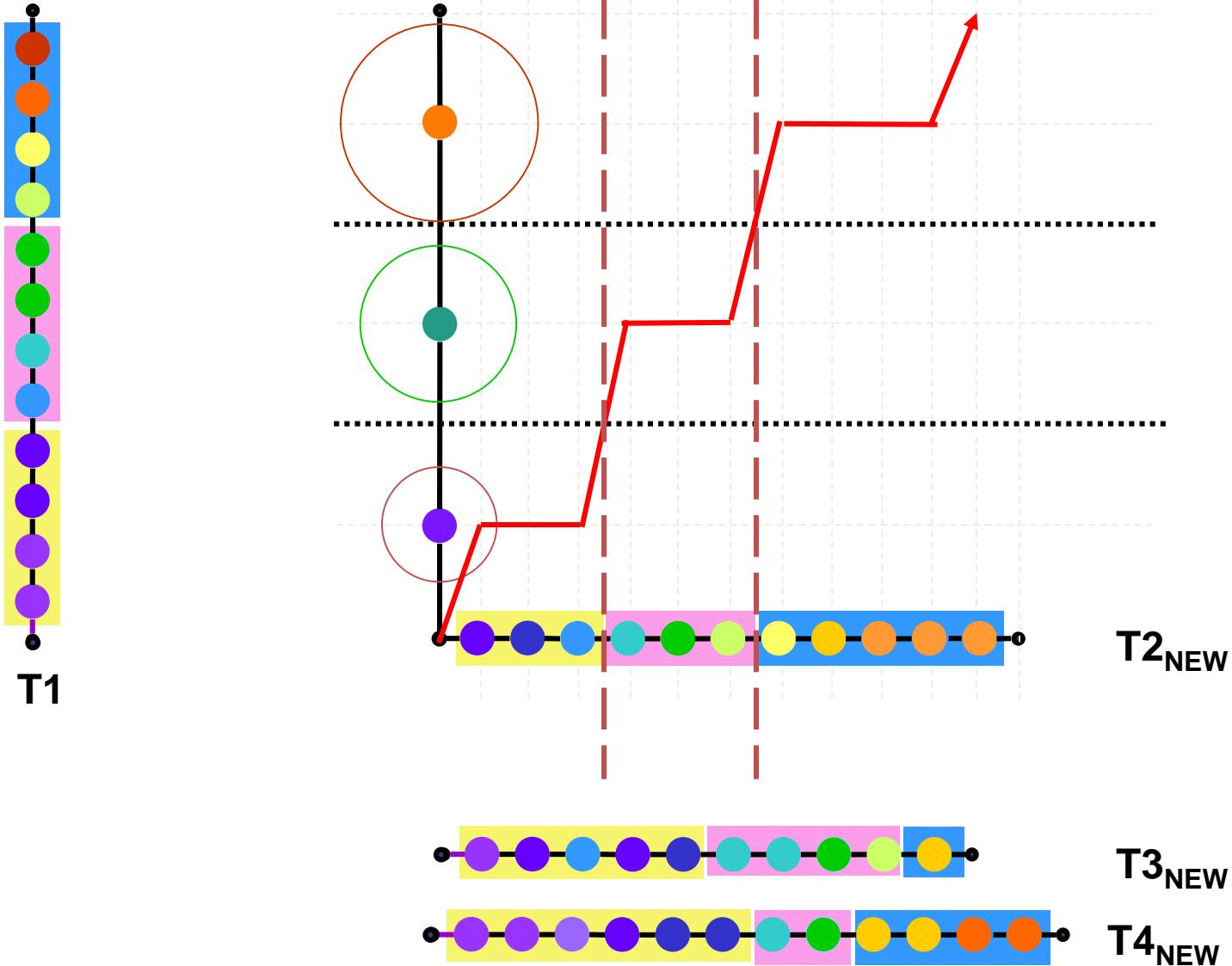
Alignment for training a model from multiple vector sequences



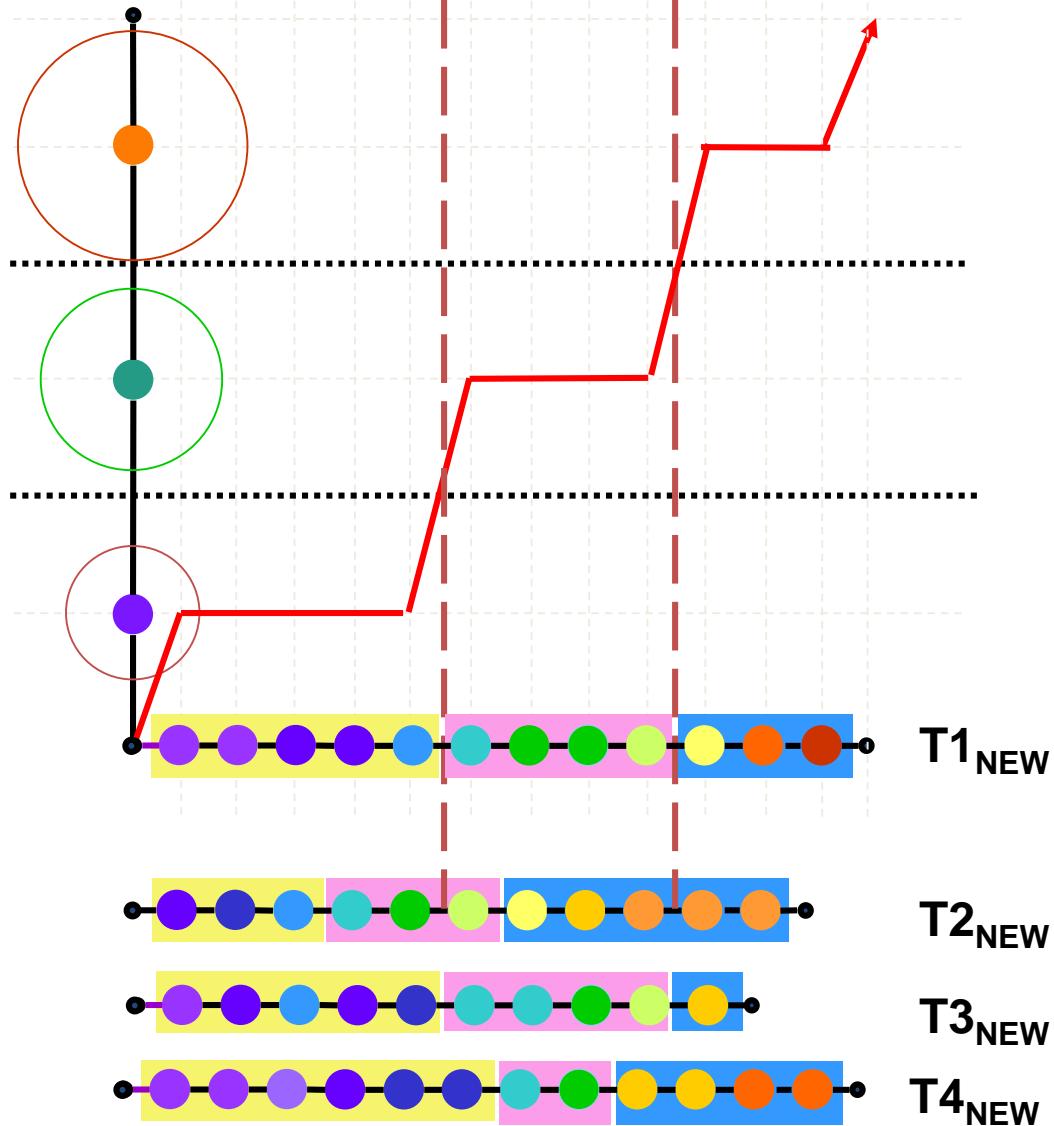
Alignment for training a model from multiple vector sequences



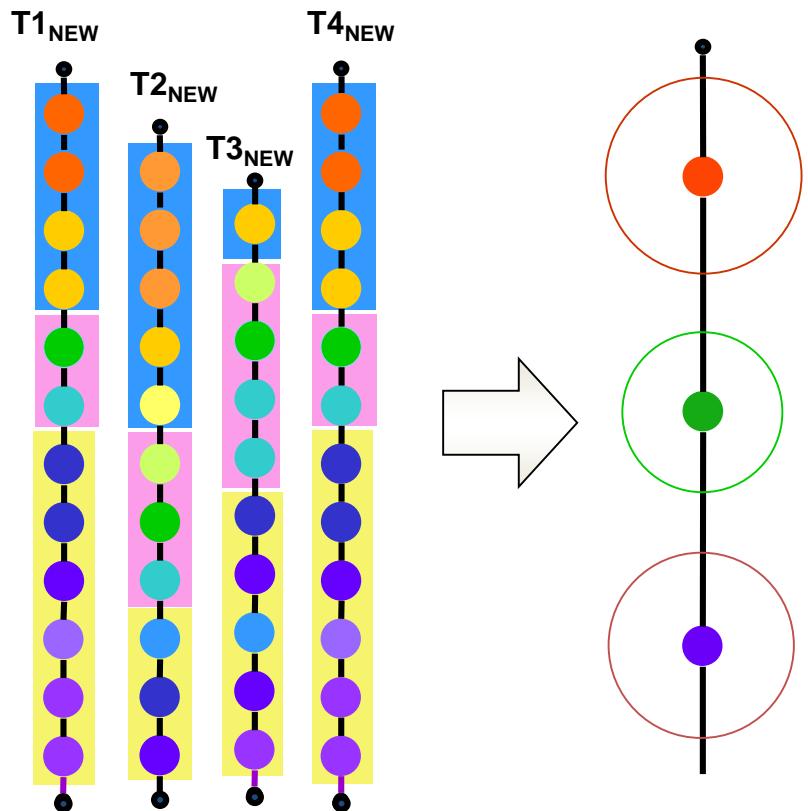
Alignment for training a model from multiple vector sequences



Alignment for training a model from multiple vector sequences



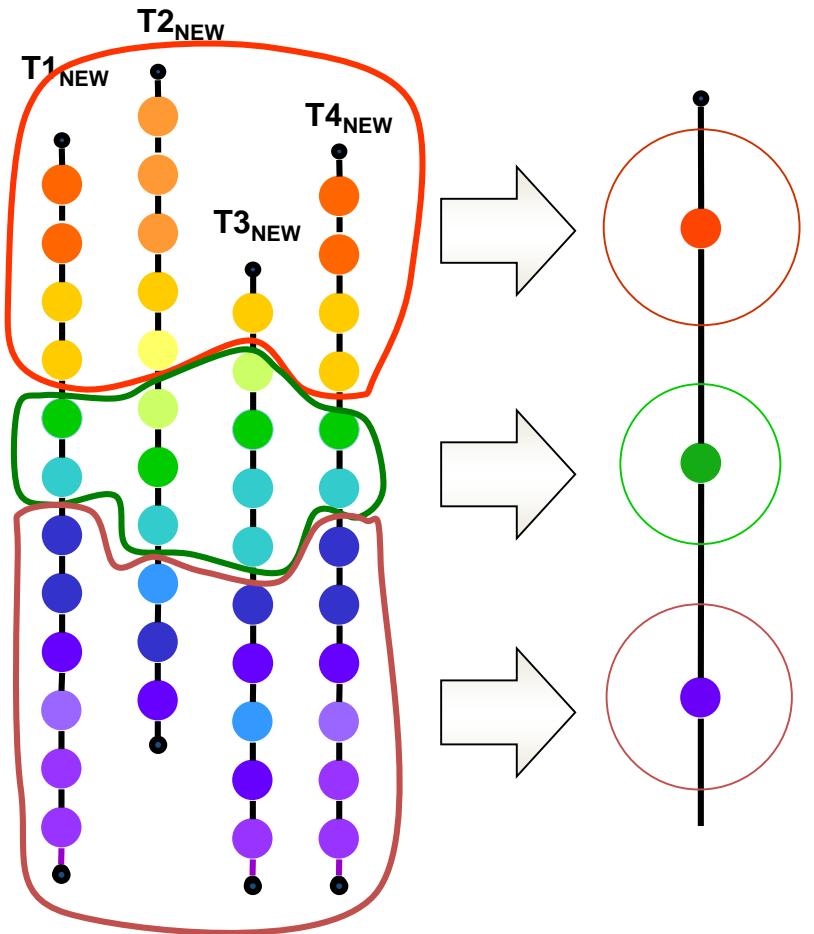
Alignment for training a model from multiple vector sequences



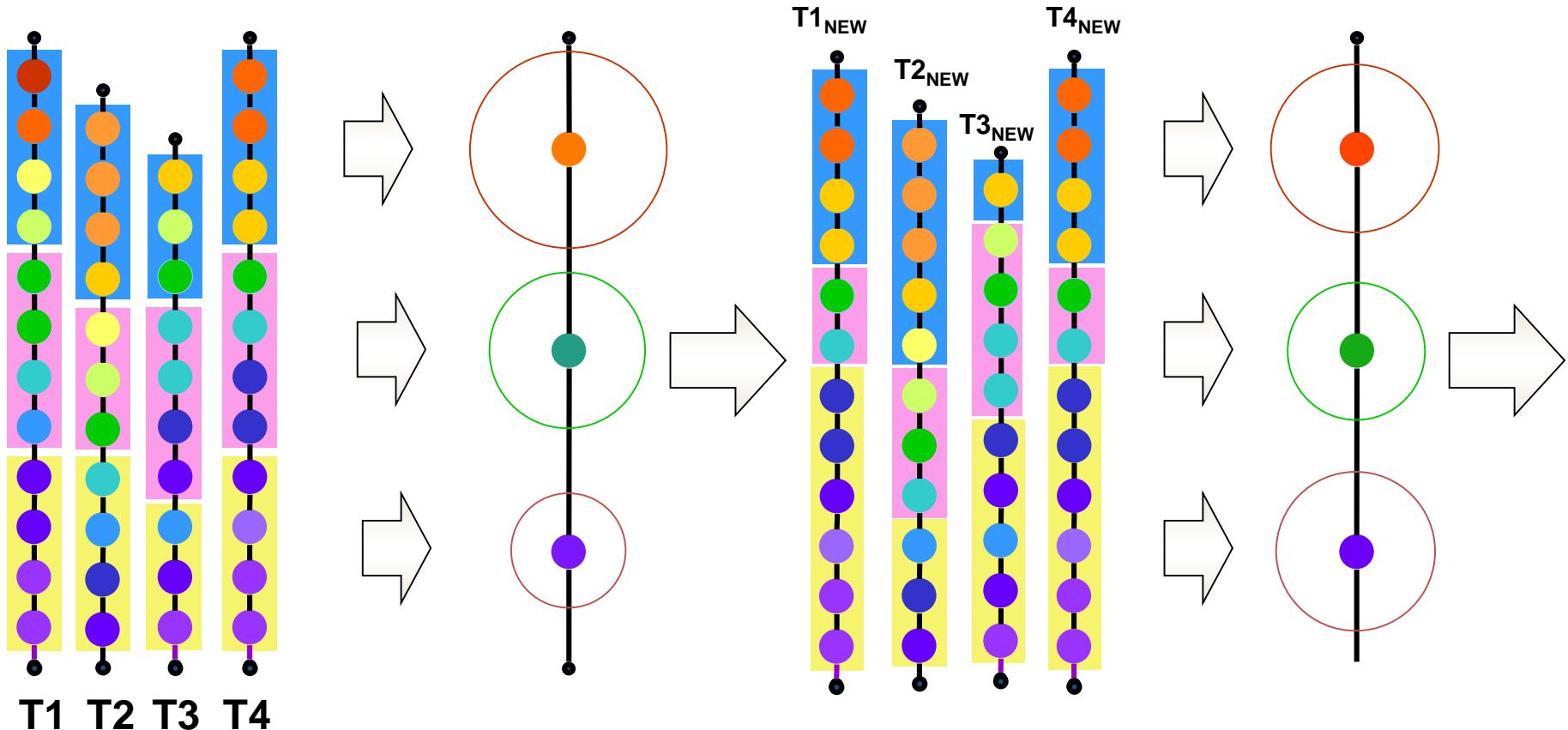
Initialize by uniform segmentation

Align each template to the averaged model to get new segmentations
Recompute the average model from new segmentations

Alignment for training a model from multiple vector sequences



Alignment for training a model from multiple vector sequences



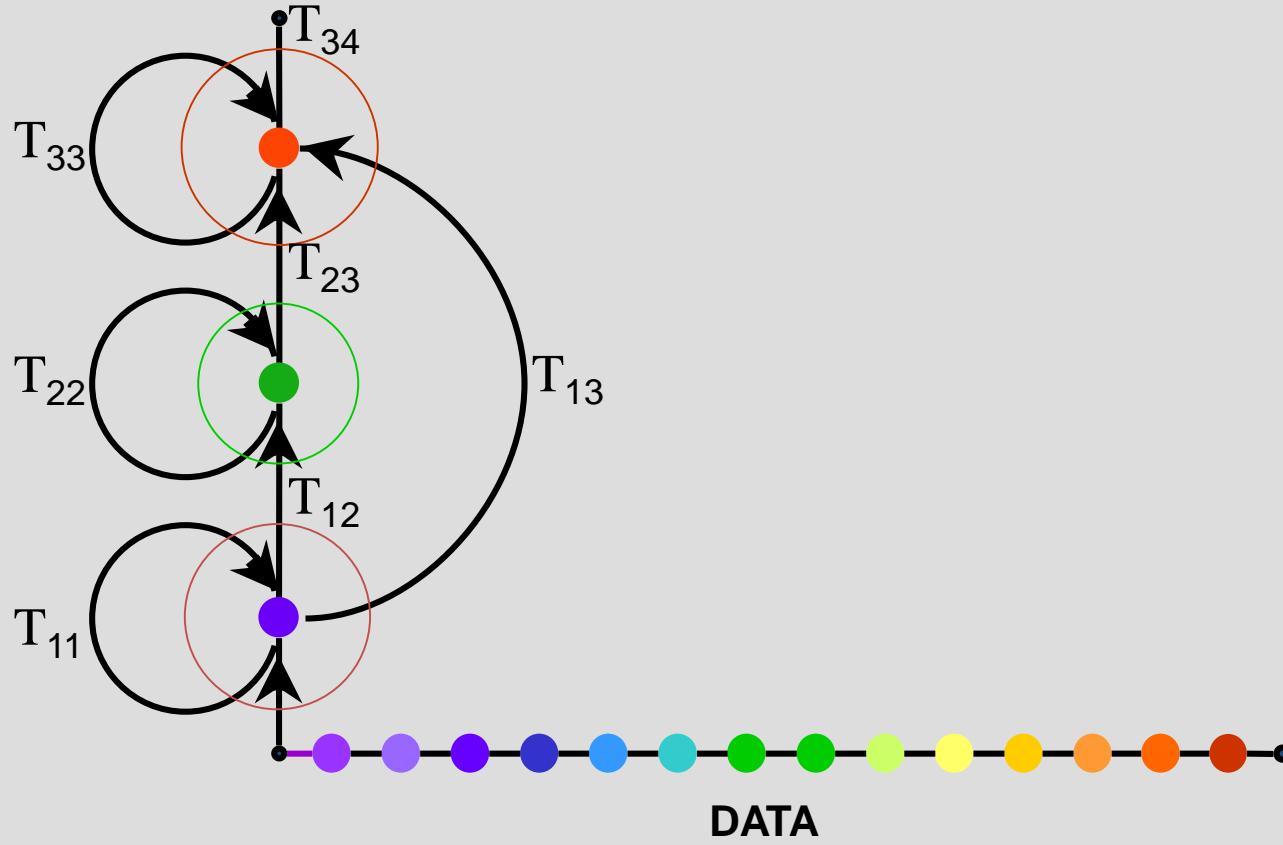
The procedure can be continued until convergence

Convergence is achieved when the total best-alignment error for all training sequences does not change significantly with further refinement of the model

Improving the Templates

- Generalization by averaging the templates
- Generalization by reducing template length
- Accounting for variation within templates represented by the reduced model
- Accounting for varying segment lengths

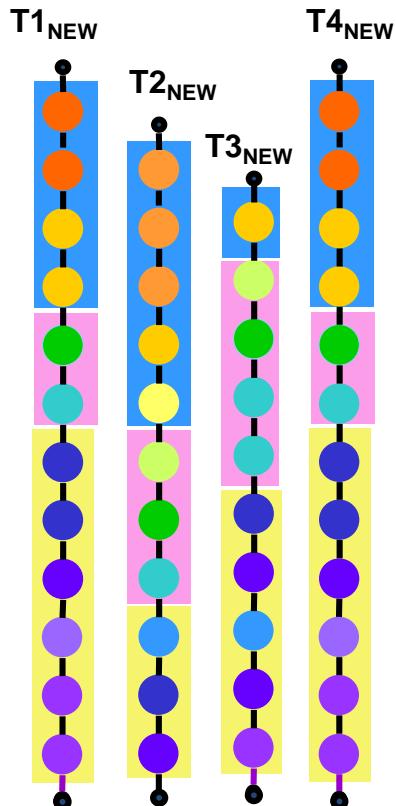
Transition structures in models



This structure also allows the inclusion of arcs that permit the central state to be skipped (deleted)

Other transitions such as returning to the first state from the last state can be permitted by inclusion of appropriate arcs

Modified segmental K-means AKA Viterbi training

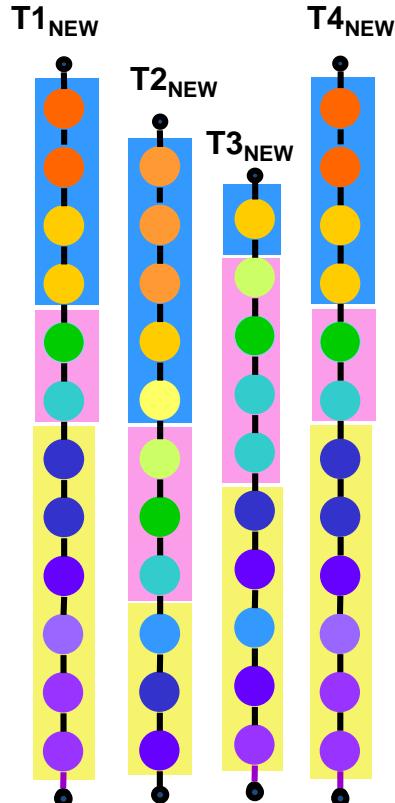


- Transition scores can be easily computed by a simple extension of the segmental K-means algorithm
- Probabilities can be counted by simple counting

$$P_{ij} = \frac{\sum_k N_{k,i,j}}{\sum_k N_{k,i}} \quad T_{ij} = -\log(P_{ij})$$

- $N_{k,i}$ is the number of vectors in the i^{th} segment (state) of the k^{th} training sequence
- $N_{k,i,j}$ is the number of vectors in the i^{th} segment (state) of the k^{th} training sequence that were followed by vectors from the j^{th} segment (state)
 - E.g., No. of vectors in the 1st (yellow) state = 20
No of vectors from the 1st state that were followed by vectors from the 1st state = 16
 $P_{11} = 16/20 = 0.8; \quad T_{11} = -\log(0.8)$

Modified segmental K-means AKA Viterbi training



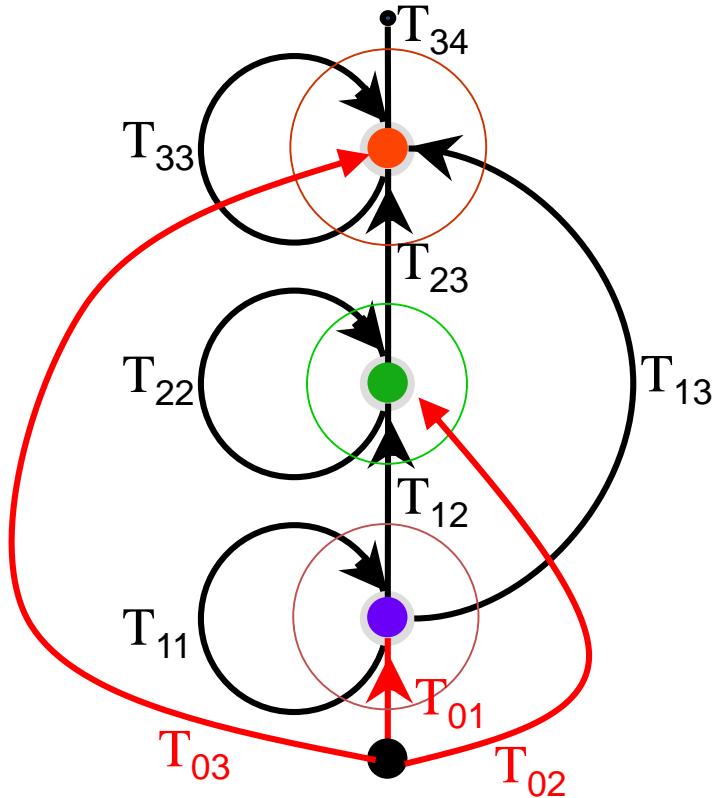
$$\begin{aligned}N &= 4 \\N_{01} &= 4 \\N_{02} &= 0 \\N_{03} &= 0\end{aligned}$$

- A special score is the penalty associated with *starting* at a particular state
- In our examples we always begin at the first state
- Enforcing this is equivalent to setting $T_{01} = 0$, $T_{0j} = \infty$ for $j \neq 1$
- It is sometimes useful to permit entry directly into later states
 - i.e. permit deletion of initial states
- The score for direct entry into any state can be computed as

$$P_j = \frac{N_{0j}}{N} \quad T_{0j} = -\log(P_j)$$

- N is the total number of training sequences
- N_{0j} is the number of training sequences for which the first data vector was in the j^{th} state

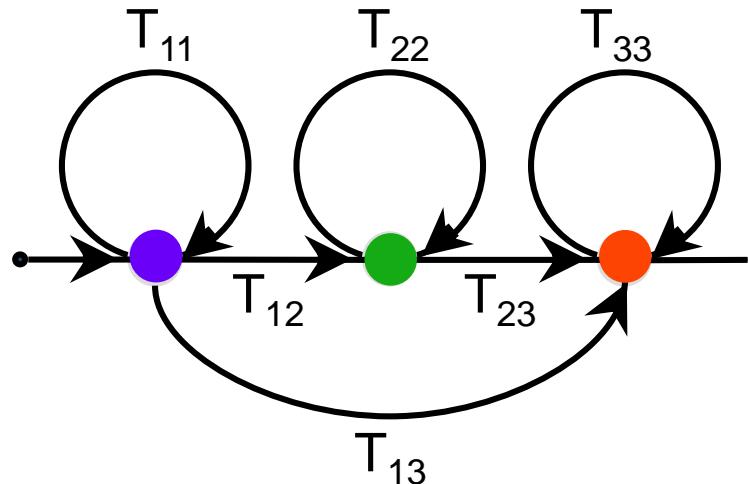
What is the Initial State Probability?



Transition from an initial “dummy” state

The resulting model structure is
also known as an HMM!

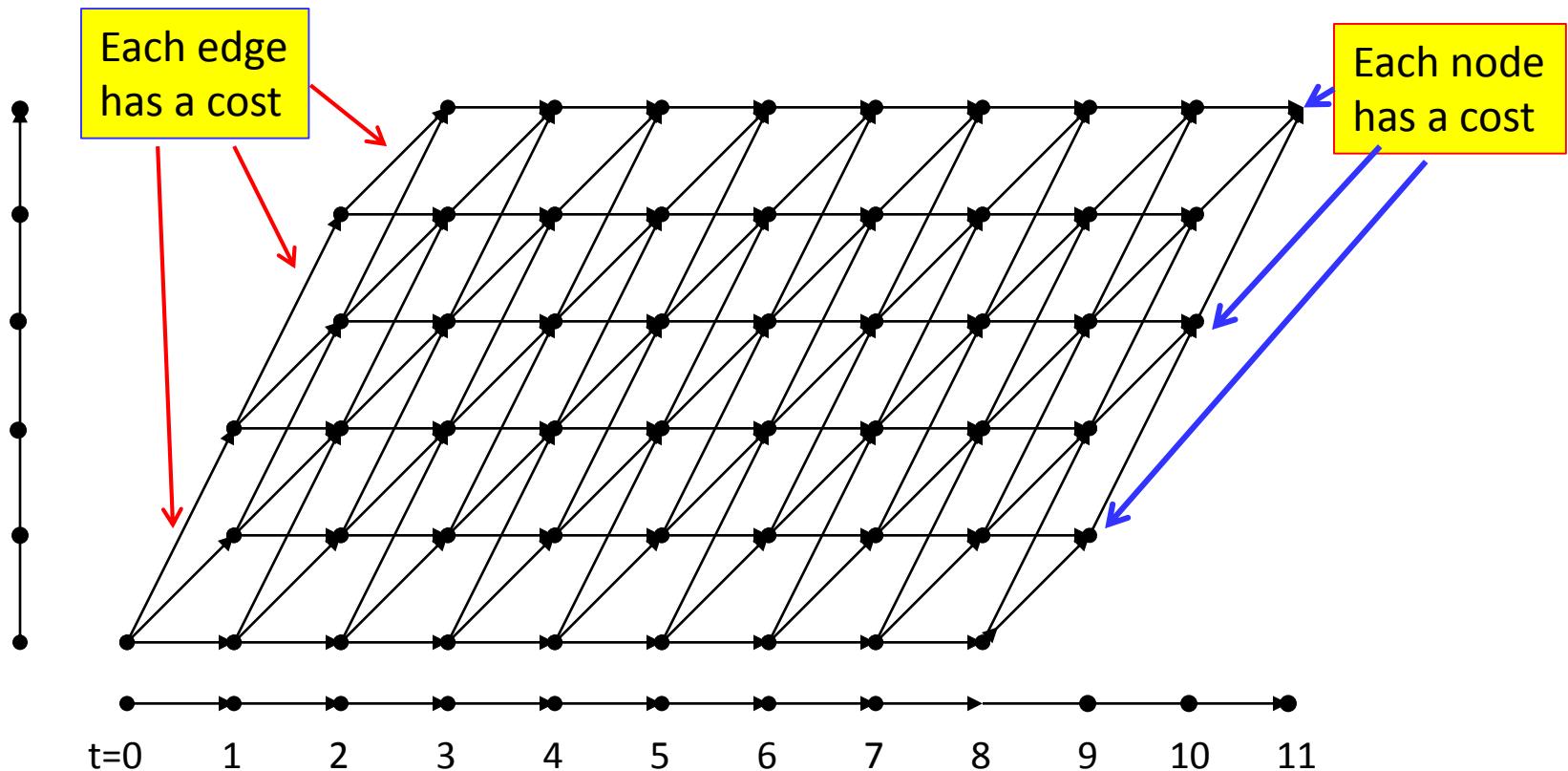
DTW and Hidden Markov Models (HMMs)



- This structure is a generic representation of a statistical model for processes that generate time series
- The “segments” in the time series are referred to as states
 - The process passes through these states to generate time series
- The entire structure may be viewed as *one* generalization of the DTW models we have discussed thus far
- In this example -- strict left-to-right topology
 - Commonly used for speech recognition

DTW with *costs*

- Node and edge *costs* defined for trellis
- Find minimum cost path through trellis..

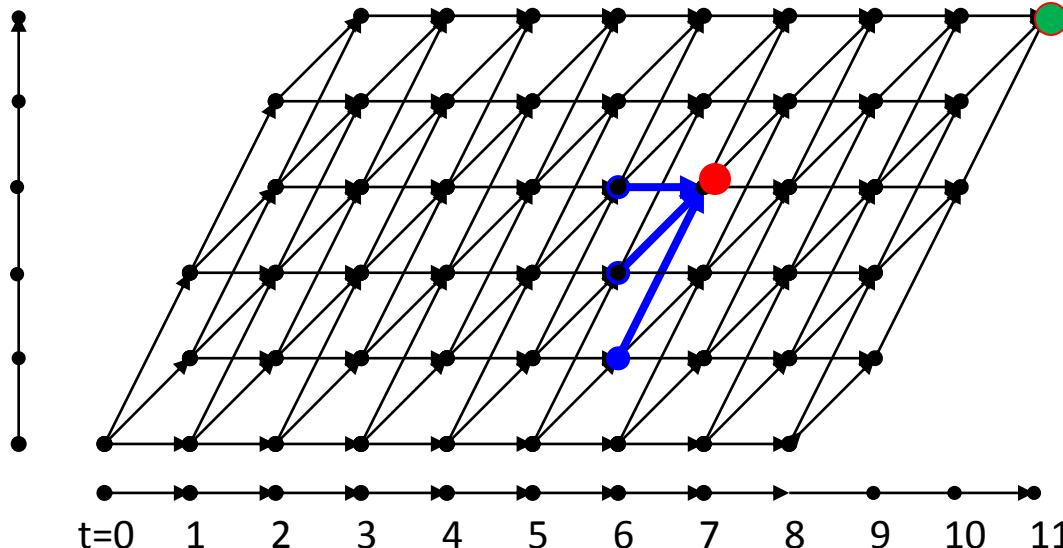


DTW: Finding minimum-cost path

- $P_{i,j}$ = best path cost from origin to node $[i,j]$
- $C_{i,j}$ = local node cost of aligning template frame i to input frame j
- $T_{i,j,k,l}$ = Edge cost from node (i,j) to node (k,l)

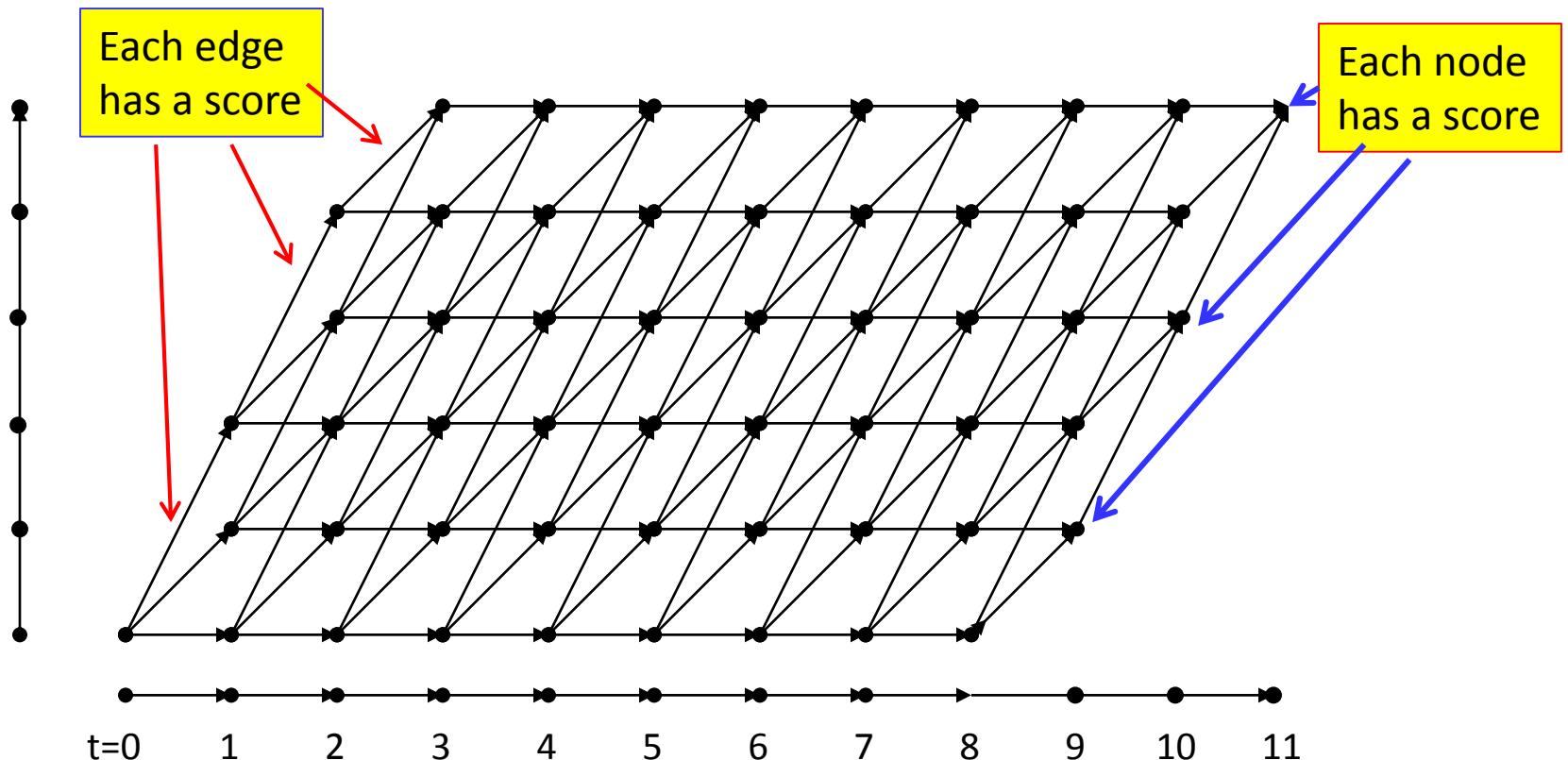
$$P_{i,j} = \min (P_{i,j-1} + T_{i,j-1,i,j}, P_{i-1,j-1} + T_{i-1,j-1,i,j}, P_{i-2,j-1} + T_{i-2,j-1,i,j},) + C_{i,j}$$

– MINIMIZE TOTAL PATH COST



DTW with **scores**

- Node and edge *scores* defined for trellis
- Find **maximum score** path through trellis..

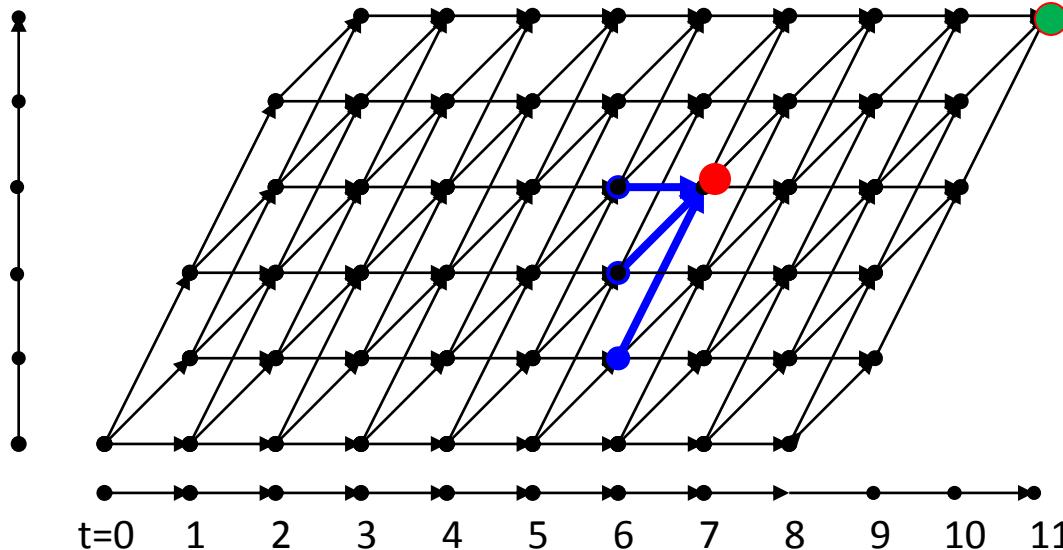


DTW: finding *maximum-score* path

- $P_{i,j}$ = best path score from origin to node $[i,j]$
- $C_{i,j}$ = local node score of aligning template frame i to input frame j
- $T_{i,j,k,l}$ = Edge score from node (i,j) to node (k,l)

$$P_{i,j} = \max (P_{i,j-1} + T_{i,j-1,i,j}, P_{i-1,j-1} + T_{i-1,j-1,i,j}, P_{i-2,j-1} + T_{i-2,j-1,i,j},) + C_{i,j}$$

– MAXIMIZE TOTAL PATH SCORE



Probabilities for Scores

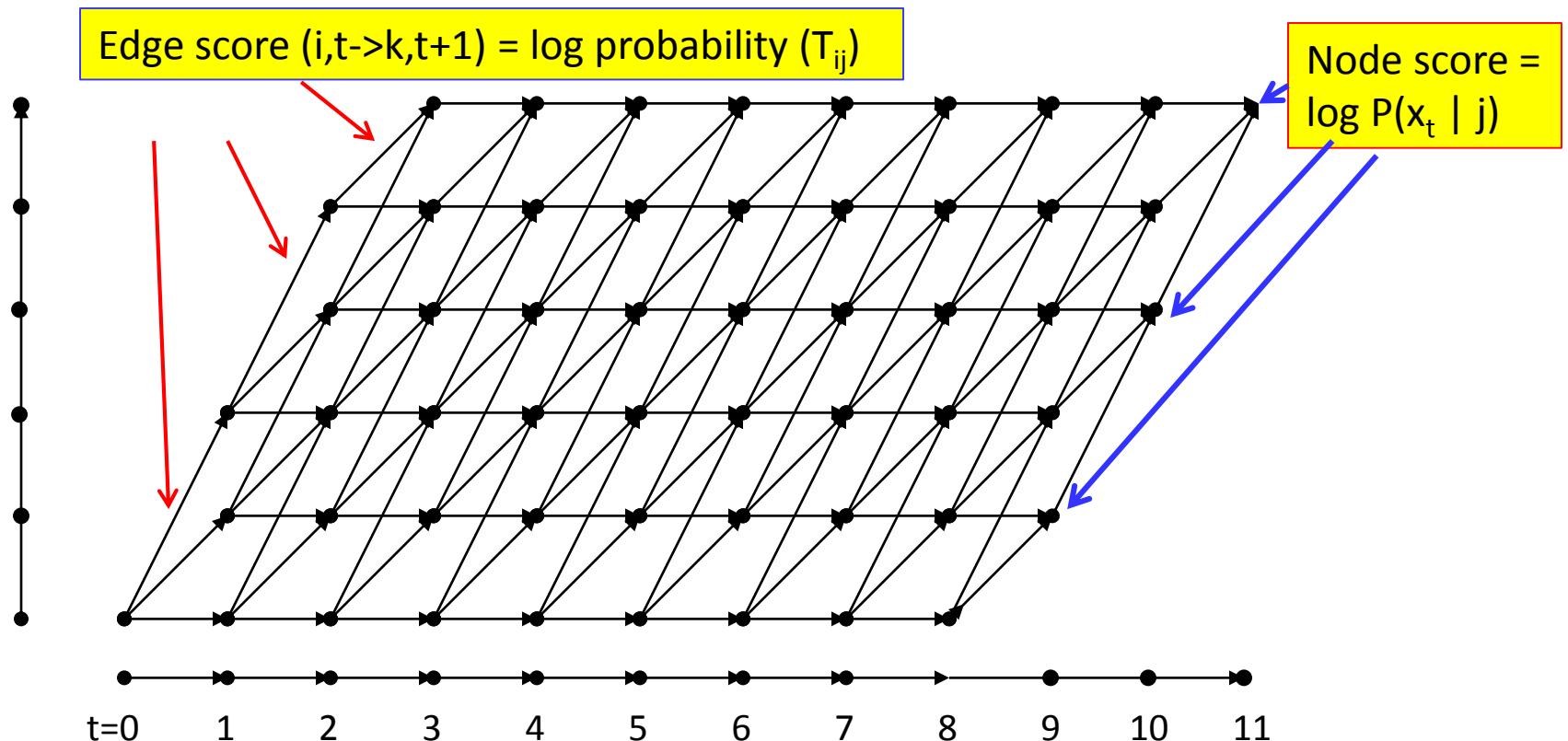
- HMM – inference equivalent to DTW modified to use a *probabilistic* function, for the local node or edge “scores” in the trellis
 - Edges have transition *probabilities*
 - Nodes have *output* or *observation probabilities*
 - They provide the probability of the observed input
 - The output probability may be a Gaussian
 - The goal is to find the template with highest probability of matching the input
- Probability values as “scores” are also called *likelihoods*

Log Likelihoods

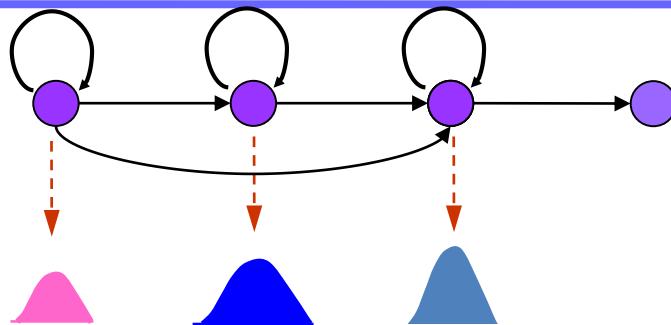
- The problem with probabilities: The probabilities of unrelated variables multiply
 - $P(X, Y) = P(X)P(Y)$
- Probabilities multiply along the path
 - Scores combines multiplicatively along a path
 - score of a path = $\text{Product_over_nodes}(\text{score of node}) \times \text{Product_over_edges}(\text{score of edge})$
- Use *log* probabilities as scores
 - Scores add as in DTW
 - Max instead of Min
- May use *negative* log probabilities
 - *Cost* adds as in DTW

HMM = DTW with **scores**

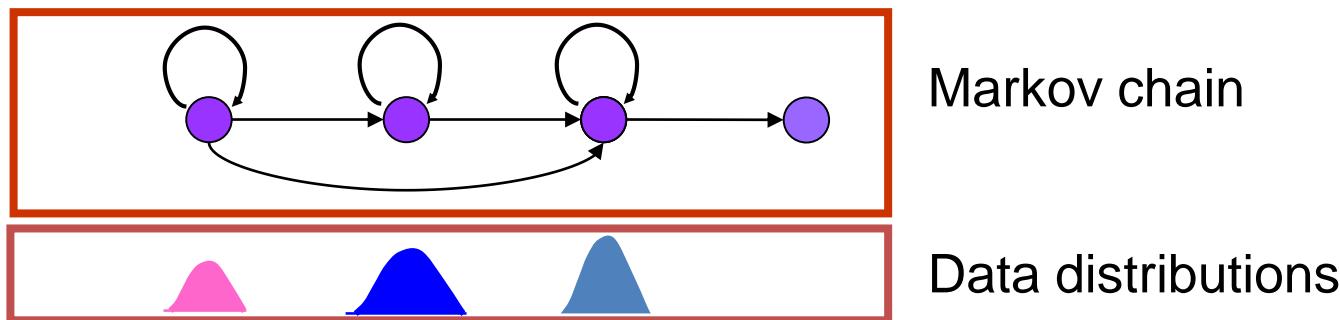
- Node and edge *scores* are log probabilities
- Find **maximum score** path through trellis..



Hidden Markov Models



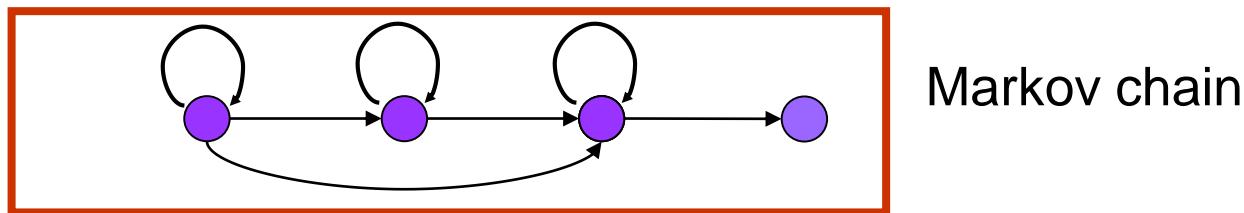
- A Hidden Markov Model consists of two components
 - A state/transition backbone that specifies how many states there are, and how they can follow one another
 - A set of probability distributions, one for each state



- This can be factored into two separate probabilistic entities
 - A probabilistic Markov chain with states and transitions
 - A set of data probability distributions, associated with the states

Basic Structural Questions for HMM

- What is the structure like?
 - What transitions are allowed
 - How many states?



- What are the probability distributions associated with states?

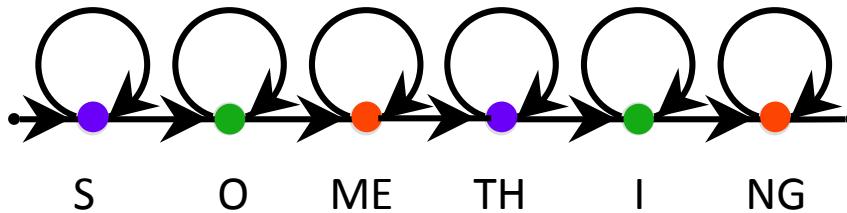


Determining the Number of States

- The correct number of states for any word?
 - We do not know, really
 - Ideally there should be at least one state for each “basic sound” within the word
 - Otherwise widely differing sounds may be collapsed into one state
 - For efficiency, the number of states should be the minimum needed to achieve the desired level of recognition accuracy
 - These two are conflicting requirements, usually solved by making some educated guesses

Determining the Number of States

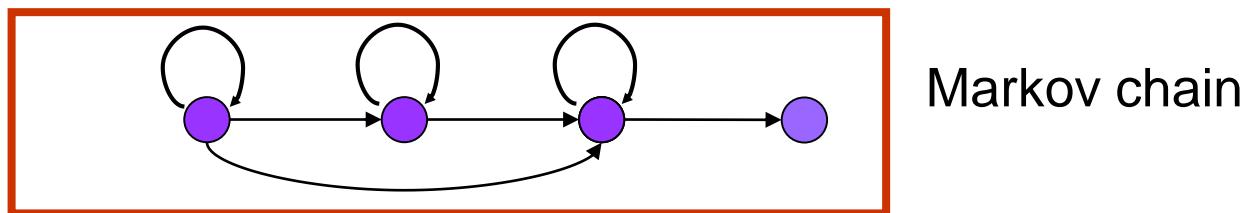
- For small vocabularies, it is possible to examine each word in detail and arrive at reasonable numbers:



- For larger vocabularies, we may be forced to rely on some *ad hoc* principles
 - *E.g.* proportional to the number of letters in the word
 - Works better for some languages than others
 - Spanish, Japanese (Katakana/Hiragana), Indian languages..

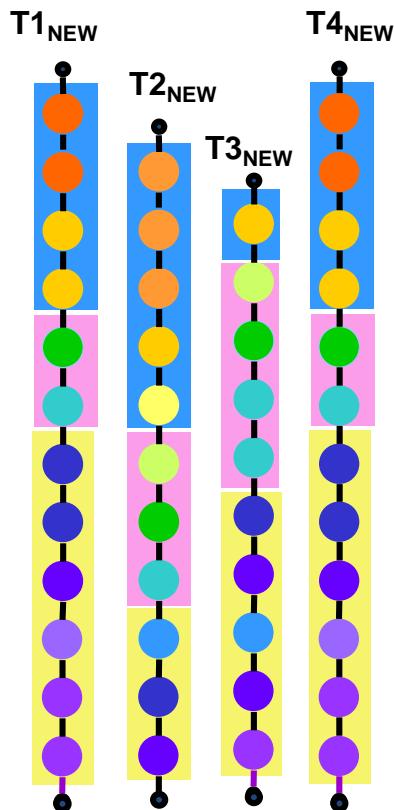
What about the transition structure

- Speech is a left to right process
 - With a deterministic structure to any word
 - Prespecified set of sounds in prespecified order
 - Although some sounds may occasionally be skipped
- This suggests the following kind of structure



- The *Bakis* topology

The Transition *Probabilities?*



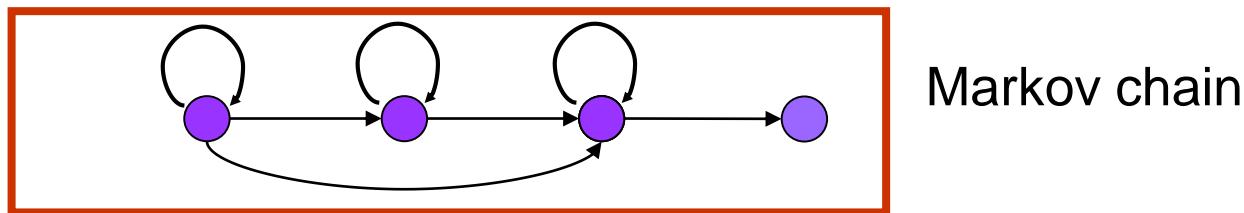
- Transition probabilities are just the probability of making transitions
- As before, can be obtained by simple counting

$$P_{ij} = \frac{\sum_k N_{k,i,j}}{\sum_k N_{k,i}}$$

- $N_{k,i}$ is the number of vectors in the i^{th} segment (state) of the k^{th} training sequence
- $N_{k,i,j}$ is the number of vectors in the i^{th} segment (state) of the k^{th} training sequence that were followed by vectors from the j^{th} segment (state)

Basic Structural Questions for HMM

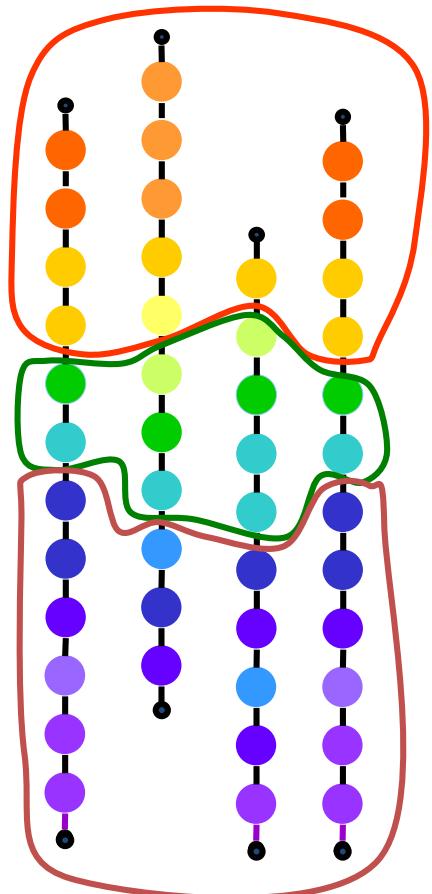
- What is the structure like?
 - What transitions are allowed
 - How many states?



- What are the probability distributions associated with states?



State Output Distributions



- The state output distribution for any state is the distribution of all vectors associated with that state
- It is the distribution of *all* vectors validly associated with the state
 - From all potential instances of the word
 - Not just the training (template) recordings
- We have implicitly assumed this to be Gaussian so far

The State Output Distribution

- The state output distribution is a probability distribution associated with each HMM state
 - The negative log of the probability of any vector as given by this distribution would be the node cost in DTW
- The state output probability distribution could be any distribution at all
- We have assumed so far that state output distributions are Gaussian

$$P(x | j) = \frac{1}{\sqrt{\prod_l 2\pi\sigma_{j,l}^2}} \exp\left(-0.5 \sum_l \frac{(x_l - m_{j,l})^2}{\sigma_{j,l}^2}\right)$$

$$d_j(v) = -\log(P(x | j)) = 0.5 \sum_l \log(2\pi\sigma_{j,l}^2) + 0.5 \sum_l \frac{(x_l - m_{j,l})^2}{\sigma_{j,l}^2}$$

Node cost for
DTW (note change
in notation)

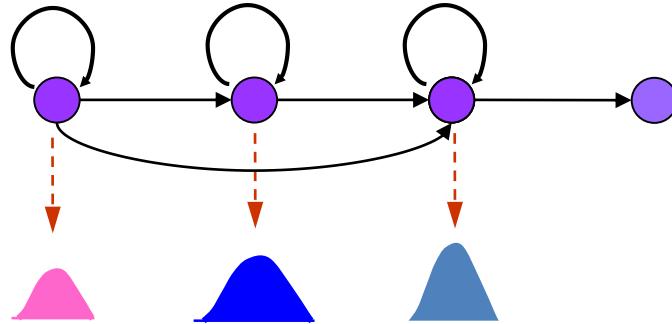
- More generically, we can assume it to be a *mixture of Gaussians*

$$P(x | j) = \sum_k \frac{w_k}{\sqrt{\prod_l 2\pi\sigma_{j,k,l}^2}} \exp\left(-0.5 \sum_l \frac{(x_l - m_{j,k,l})^2}{\sigma_{j,k,l}^2}\right)$$

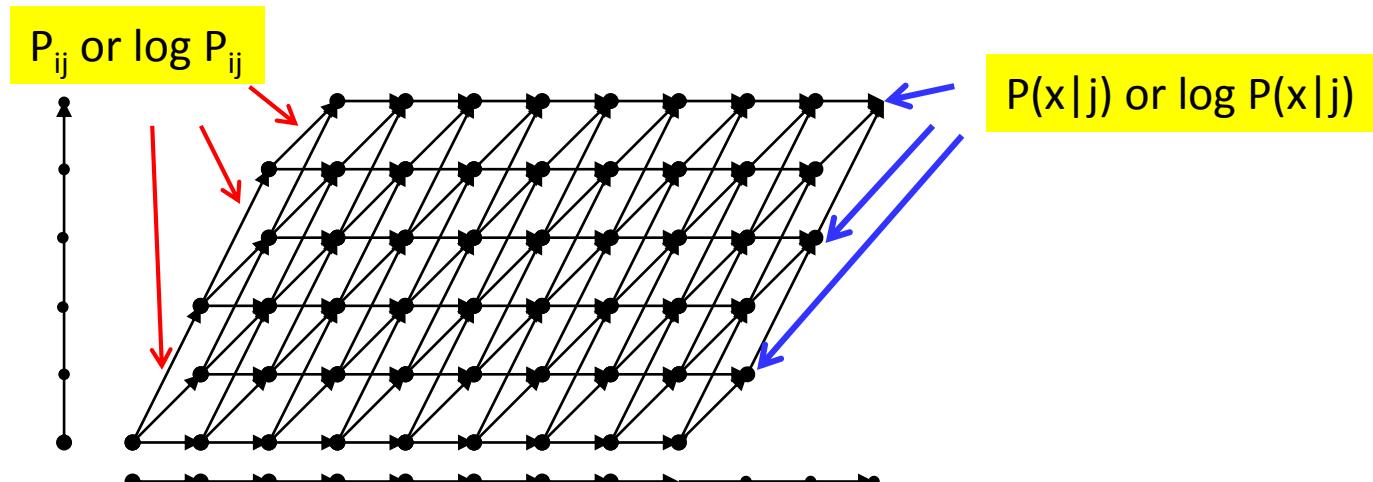
$$d_j(v) = -\log(P(x | j))$$

- More on this later

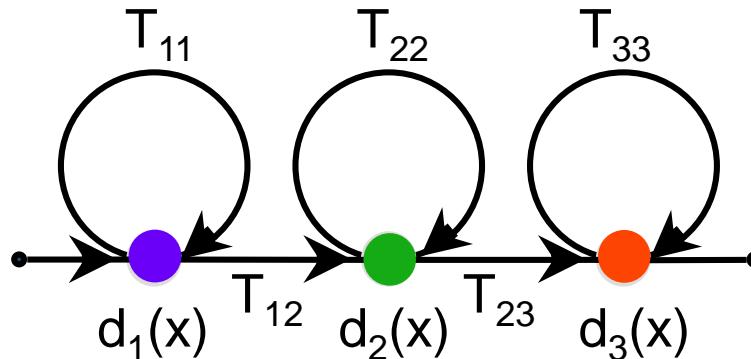
The complete package



- We have seen how to compute all probability terms for this structure
- Sufficient information to compute all terms in a trellis

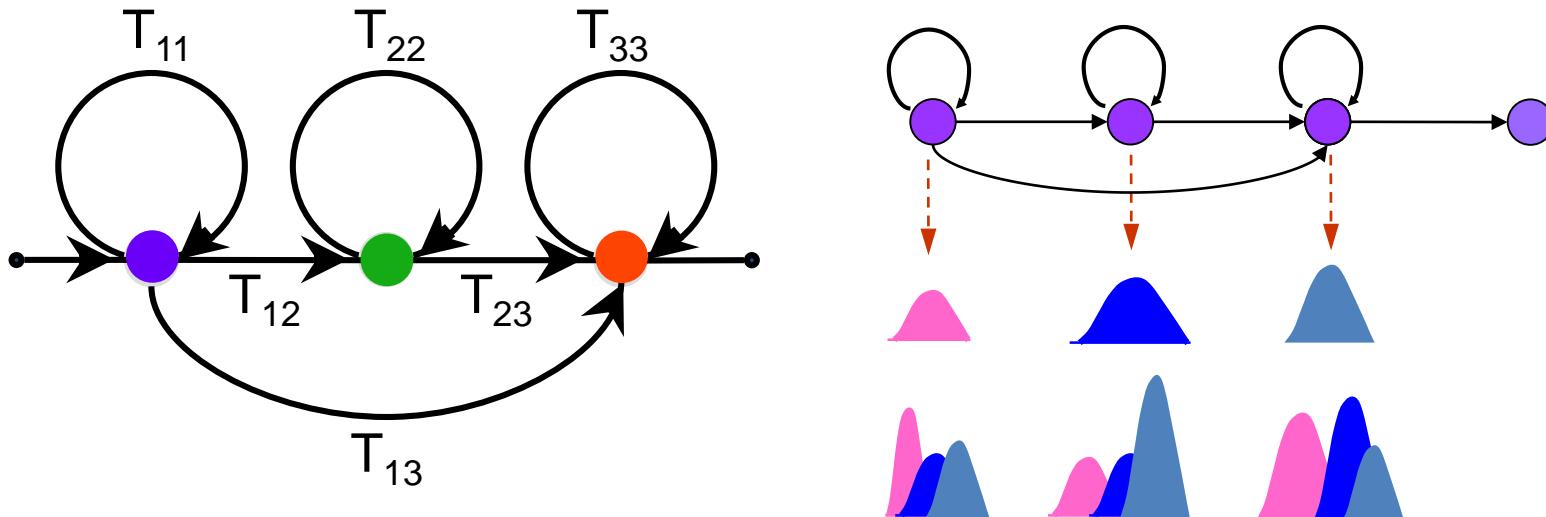


Generalized Templates



- A set of “states”
 - A distance function associated with each state
- A set of transitions
 - Transition-specific penalties

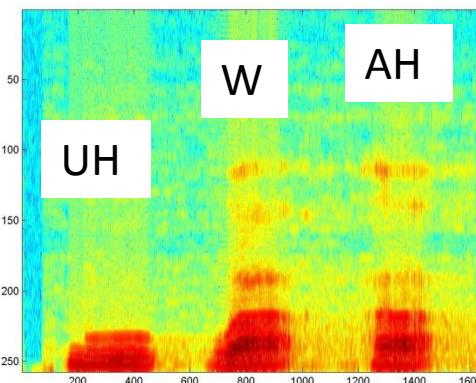
HMMs



- Identical to generalized templates in principle
- “Distance” functions at states replaced by “probability distribution function” for state
- Transition “penalties” replaced by transition probabilities
- Maximize probability of observation
 - Instead of minimizing cost
- The entire structure may be viewed as *one* generalization of the DTW models we have discussed thus far

The HMM Process

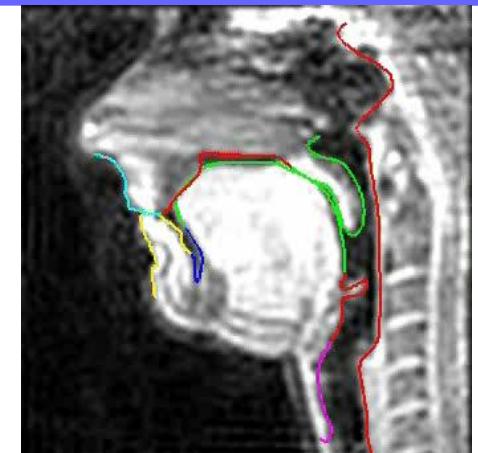
- The HMM models the process underlying the observations as going through a number of states
 - E.g., to produce the sound “Wa”, it first goes through a state where it produces the sound “UH”, then goes into a state where it transitions from “UH” to “AH”, and finally to a state where it produced “AH”



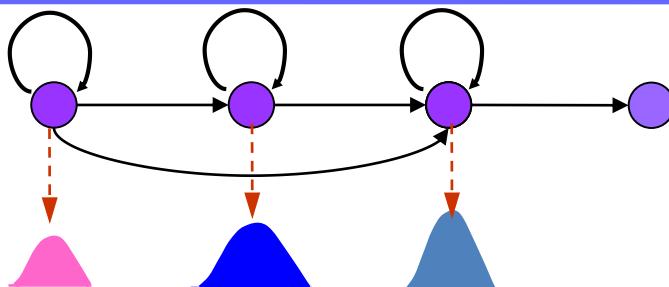
- The true underlying process is the vocal tract here
 - Which roughly goes from the configuration for “UH” to the configuration for “AH”

HMMs are abstractions

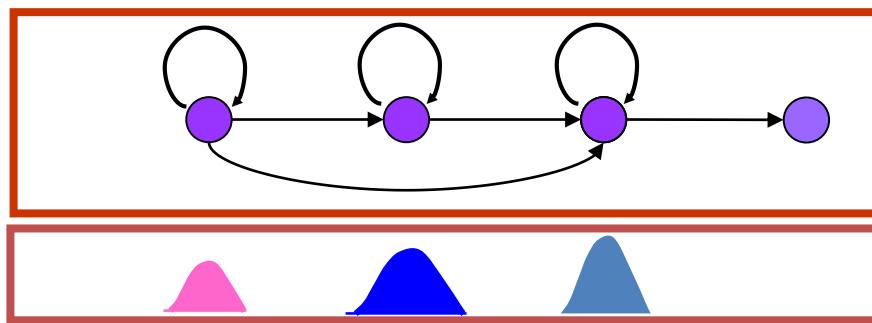
- The states are not directly observed
 - Here states of the process are analogous to configurations of the vocal tract that produces the signal
 - We only hear the speech; we do not see the vocal tract
 - i.e. the states are *hidden*
- The interpretation of states is not always obvious
 - The vocal tract actually goes through a *continuum* of configurations
 - The model represents all of these using only a fixed number of states
- The model *abstracts* the process that generates the data
 - The system goes through a finite number of states
 - When in any state it can either remain at that state, or go to another with some probability
 - When at any states it generates observations according to a distribution associated with that state



Hidden Markov Models



- A Hidden Markov Model consists of two components
 - A state/transition backbone that specifies how many states there are, and how they can follow one another
 - A set of probability distributions, one for each state, which specifies the distribution of all vectors in that state

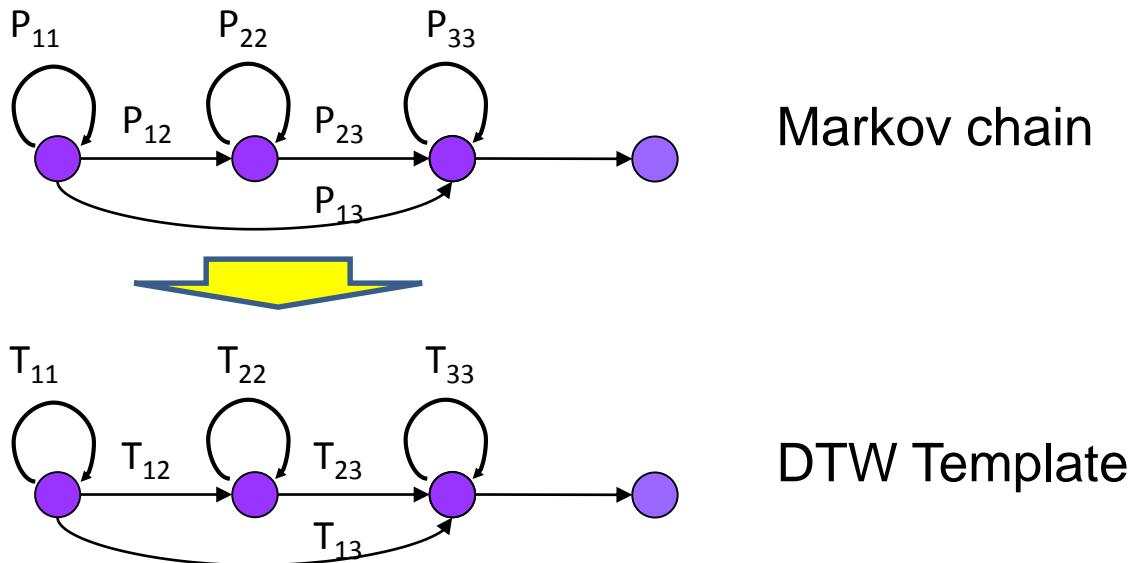


- This can be factored into two separate probabilistic entities
 - A probabilistic Markov chain with states and transitions
 - A set of data probability distributions, associated with the states

Equivalence to DTW templates

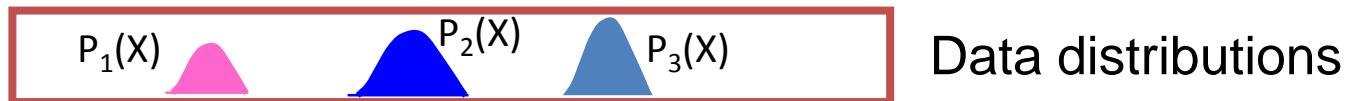
- HMM – inference equivalent to DTW modified to use a *probabilistic* function, for the local node or edge “costs” in the trellis
 - Edges have transition probabilities
 - Nodes have *output* or *observation probabilities*
 - They provide the probability of the observed input
 - The output probability may be a Gaussian or a GMM
 - Goal is to find the template with highest probability of matching the input
- Probability values associated with transitions and edges are called *likelihoods*

Likelihoods and Cost: Transition



- Transitions in the HMM have associated probabilities
 - P_{11} , P_{12} etc
- They can be converted to “scores” through a logarithm
 - $T_{11} = \log(P_{11})$
- Or to “costs” through a negative logarithm
 - $T_{11} = -\log(P_{11})$

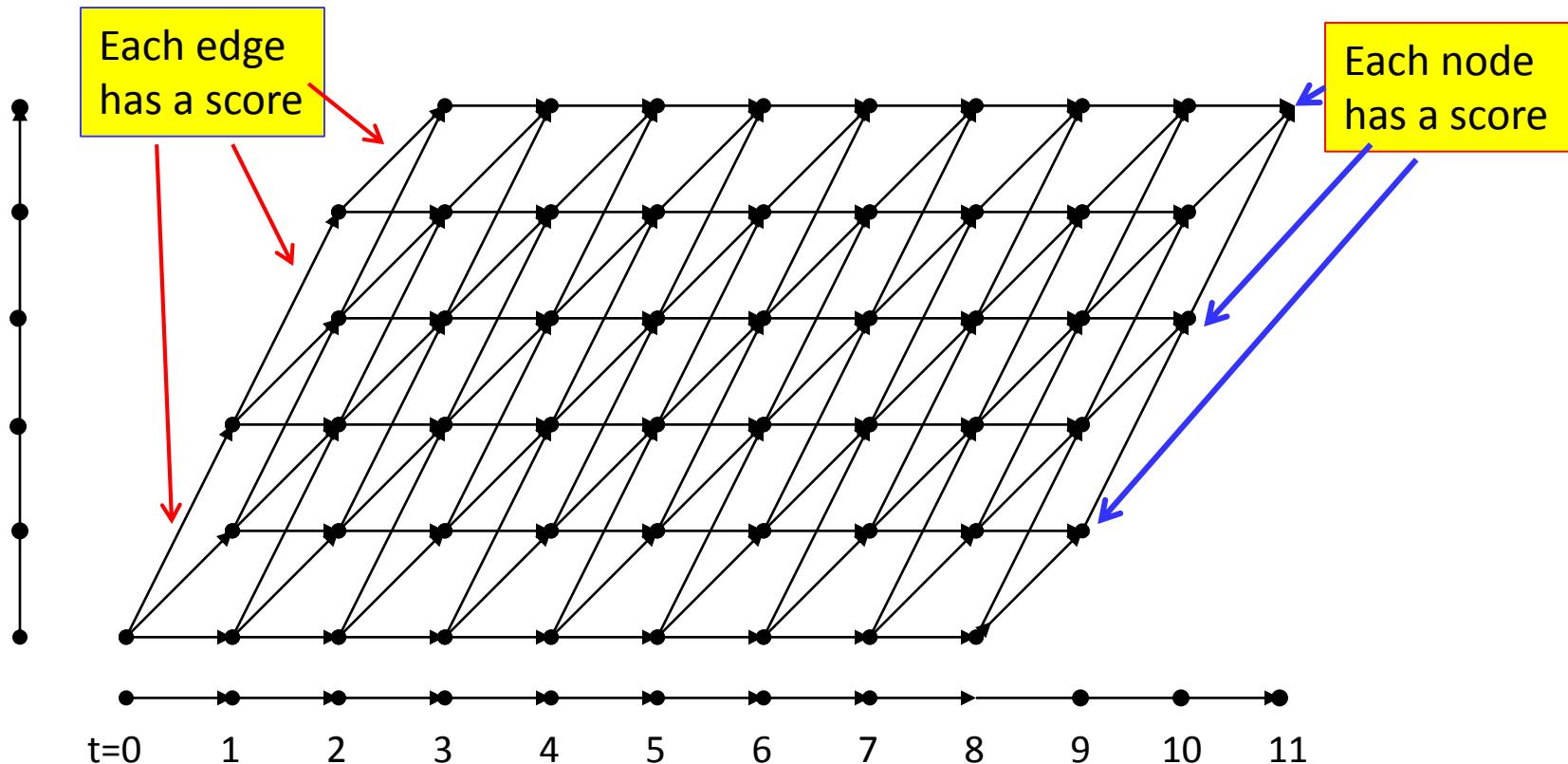
Likelihoods and Cost: Nodes



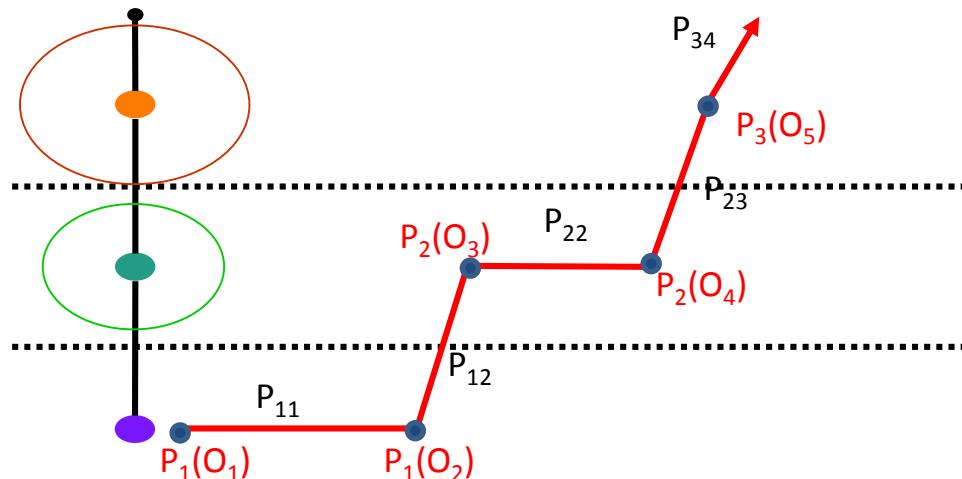
- States in the HMM have probability distributions associated with them
 - E.g Gaussians
 - Whose means and variances have been obtained from the feature frames associated with the nodes/segments
- Nodes in the *trellis* have a probabilities associated with them
 - $P_i(O)$
 - i is the “state” / template node
 - O is the observation associated with any node in the trellis
- Node probabilities may be converted to:
 - Scores: $N_i(O) = \log(P_i(O))$
 - Or Costs: $N_i(O) = -\log(P_i(O))$

Computation is still done with a Trellis

- Node and edge *scores* defined for trellis



Log Likelihoods

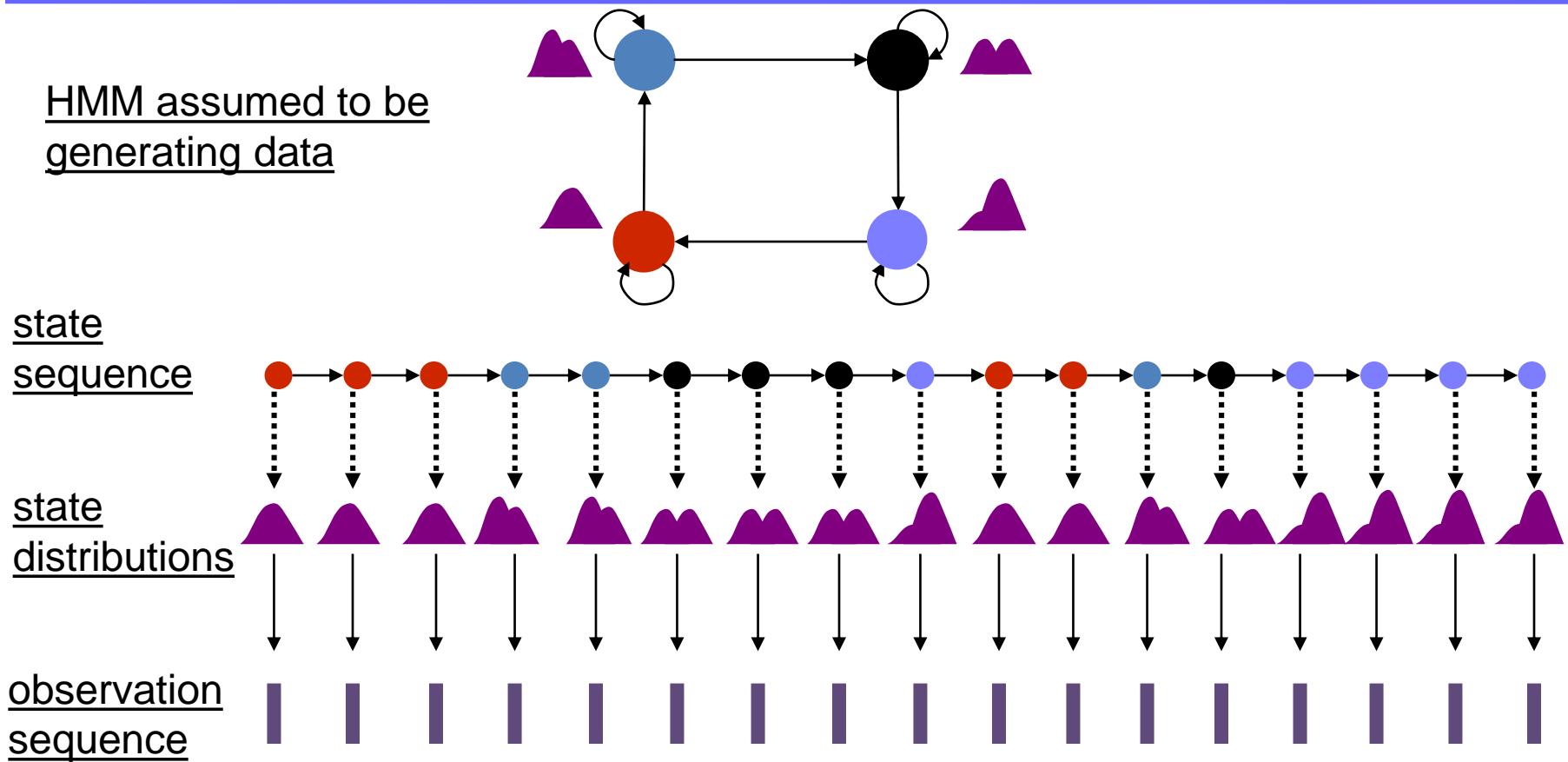


- Use probabilities or likelihoods instead of cost
 - Scores combines multiplicatively along a path
 - Path Score = $P_1(O_1) \cdot P_{11} \cdot P_1(O_2) \cdot P_{12} \cdot P_2(O_3) \cdot P_{22} \cdot P_2(O_4) \cdot P_{23} \cdot P_3(O_5) \cdot P_{23}$
- Alternately use log probabilities as scores: $N_i(O) = \log(P_i(O))$, $T_{11} = \log(P_{11})$
 - Scores add as in DTW
 - Path Score = $N_1(O_1) + T_{11} + N_1(O_2) + T_{12} + N_2(O_3) + T_{22} + N_2(O_4) + T_{23} + N_3(O_5) + T_{23}$
 - Replace all “Min” operations in DTW by “Max”
- Alternately use negative log probabilities as cost: $N_i(O) = -\log(P_i(O))$, $T_{11} = -\log(P_{11})$
 - Cost adds as in DTW
 - Computation remains identical to DTW (with edge costs factored in)

HMM as a statistical model

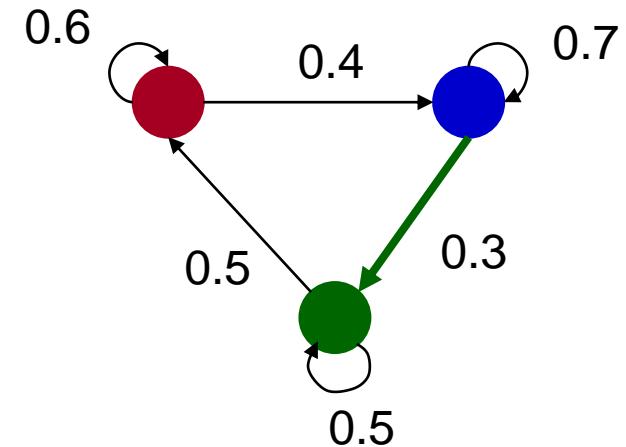
- An HMM is a statistical model for a time-varying process
- The process is always in one of a countable number of states
- When the process visits any state, it generates an observation by a random draw from a distribution associated with that state
- The process constantly moves from state to state. The probability that the process will move to any state is determined solely by the current state
 - i.e. the dynamics of the process are Markovian
- The entire model represents a probability distribution over the sequence of observations
 - It has a specific probability of generating any particular sequence
 - The probabilities of all possible observation sequences sums to 1

How an HMM models a process

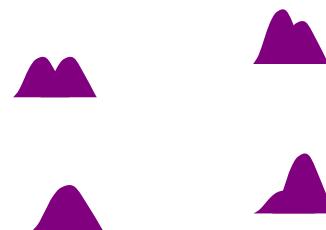


HMM Parameters

- The *topology* of the HMM
 - No. of states and allowed transitions
 - E.g. here we have 3 states and cannot go from the blue state to the red
- The transition probabilities
 - Often represented as a matrix as here
 - T_{ij} is the probability that when in state i , the process will move to j
- The probability of being at a particular state at the first instant
- The *state output distributions*



$$T = \begin{pmatrix} .6 & .4 & 0 \\ 0 & .7 & .3 \\ .5 & 0 & .5 \end{pmatrix}$$



HMM state output distributions

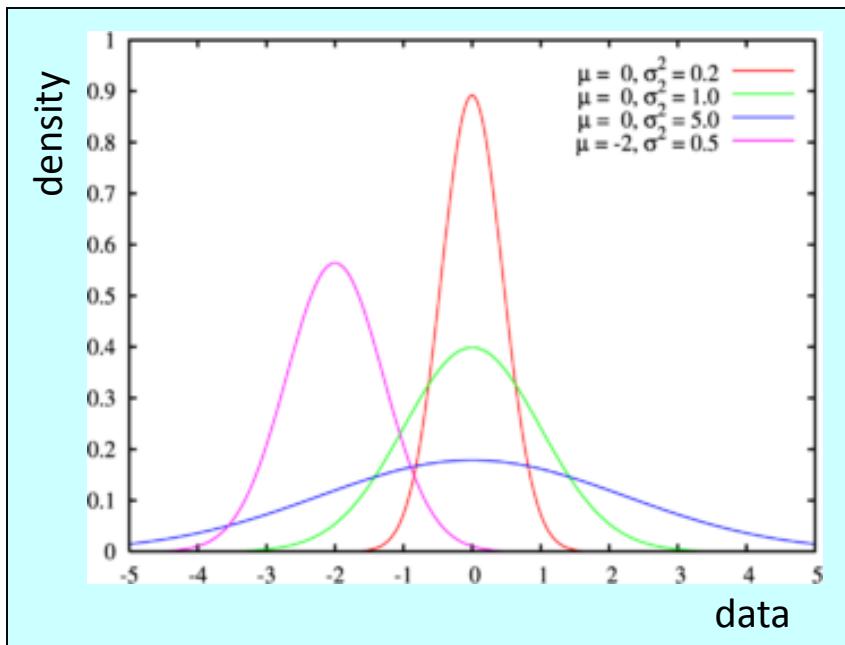
- The state output distribution represents the distribution of data produced from any state
- We have previously considered Gaussian state output distributions

$$P(v) = \text{Gaussian}(v; m, C) = \frac{1}{(2\pi)^{d/2} |C|^{1/2}} \exp\left[-\frac{1}{2}(v-m)^T C^{-1}(v-m)\right]$$

- In reality, the distribution of vectors for any state need not be Gaussian
 - In the most general case it can be arbitrarily complex
 - The Gaussian is only a coarse representation of this distribution
- If we model the output distributions of states better, we can expect the model to be a better representation of the data

Node Score: The Gaussian Distribution

- What does a Gaussian distribution look like?
- For a single (scalar) variable, it is a bell-shaped curve representing the density of data around the mean
- Example:



Four different scalar Gaussian distributions, with different means and variances

The mean is represented by μ , and variance by σ^2

μ and σ are the *parameters* of the Gaussian distribution

(Taken from Wikipedia)

The Scalar Gaussian Function

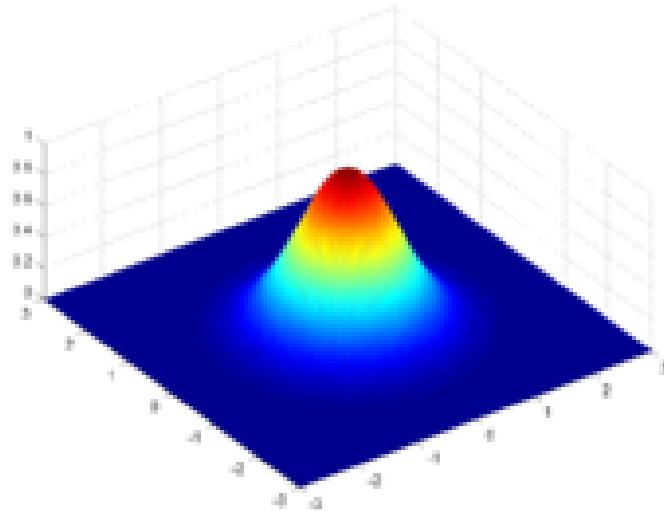
- The Gaussian density function (the bell curve) is:

$$p(x) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{1}{2\sigma^2}(x-\mu)^2}$$

- $p(x)$ is the density function of the variable x , with mean μ and variance σ^2
- The attraction of the Gaussian function (regardless of how appropriate it is!) comes from how easily the mean and variance can be estimated from *sample data* $x_1, x_2, x_3 \dots x_N$
 - $\mu = \sum_i x_i / N$
 - $\sigma^2 = \sum_i (x_i - \mu)^2 / N = \sum_i (x_i^2) / N - \mu^2$

The 2-D Gaussian Distribution

- Speech data are not scalar values, but vectors!(e.g. 39d)
- Needs multi-variate (multi-dimensional) Gaussians
- Figure: A Gaussian for 2-D data
 - Shown as a 3-D plot



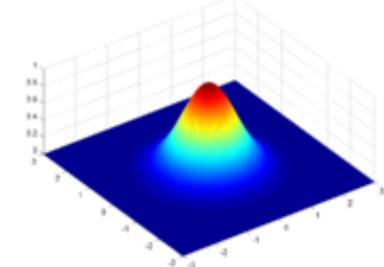
- Distributions for higher dimensions are tough to visualize!

The Multidimensional Gaussian Distribution

- Instead of variance, the multidimensional Gaussian has a *covariance matrix*
- The multi-dimensional Gaussian distribution for a vector variable x with mean μ and covariance Σ is given by:

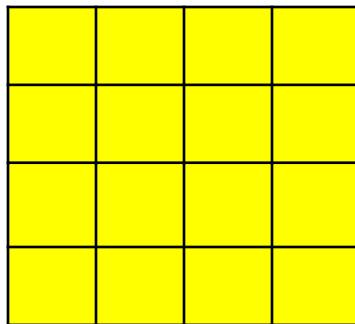
$$f(x) = \frac{1}{\sqrt{(2\pi)^D |C|}} \exp(-0.5(x - \mu)^T C^{-1}(x - \mu))$$

- where D is the vector dimensionality
- The complexity in a full multi-dimensional Gaussian distribution comes from the covariance matrix, which accounts for *dependencies* between the dimensions



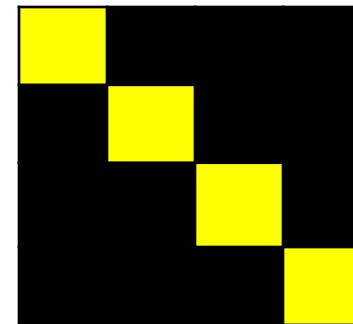
The Diagonal Covariance Matrix

Full covariance:
all elements are
non-zero



$$-0.5(x-\mu)^T C^{-1}(x-\mu)$$

Diagonal covariance:
off-diagonal elements
are zero

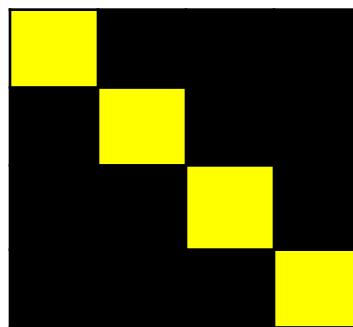


$$-\sum_i (x_i - \mu_i)^2 / 2\sigma_i^2$$

- In speech recognition, we frequently assume that the feature vector dimensions are all *independent* of each other
- *Result:* The covariance matrix is reduced to a diagonal form
 - The exponential term becomes, simply:
 $(\sum_i (x_i - \mu_i)^2 / \sigma_i^2)/2$, i going over all vector dimensions
 - The determinant of the diagonal Σ matrix is easy to compute
- Further, each σ_i^2 (the i -th diagonal element in the covariance matrix) is easily estimated from x_i and μ_i like a scalar

Multivariate Gaussian with Diagonal Covariance

Diagonal covariance:
off-diagonal elements
are zero



$$-\sum_i (x_i - \mu_i)^2 / 2\sigma_i^2$$

$$f(x) = \frac{1}{\sqrt{(2\pi)^D \prod_d \sigma_d^2}} \exp\left(-0.5 \sum_d \frac{(x_d - \mu_d)^2}{\sigma_d^2}\right)$$

- Much simplified

Gaussian Mixtures

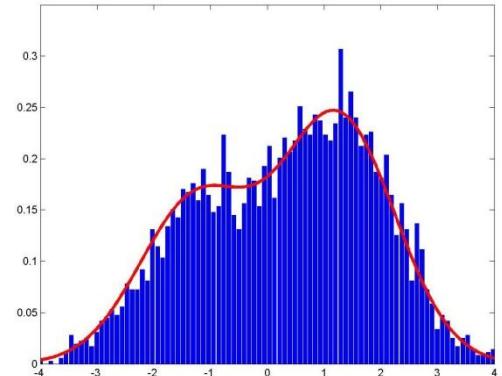
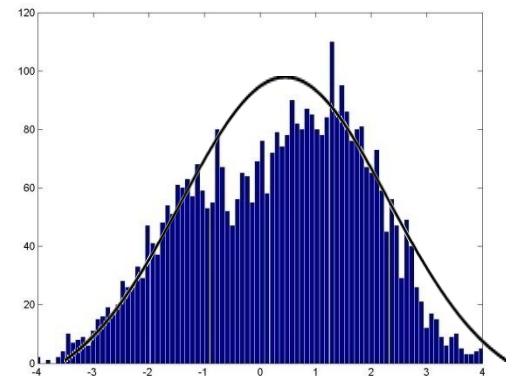
- A Gaussian Mixture is literally a mixture of Gaussians. It is a weighted combination of several Gaussian distributions

$$P(v) = \sum_{i=0}^{K-1} w_i \text{Gaussian}(v; \mu_i, C_i)$$

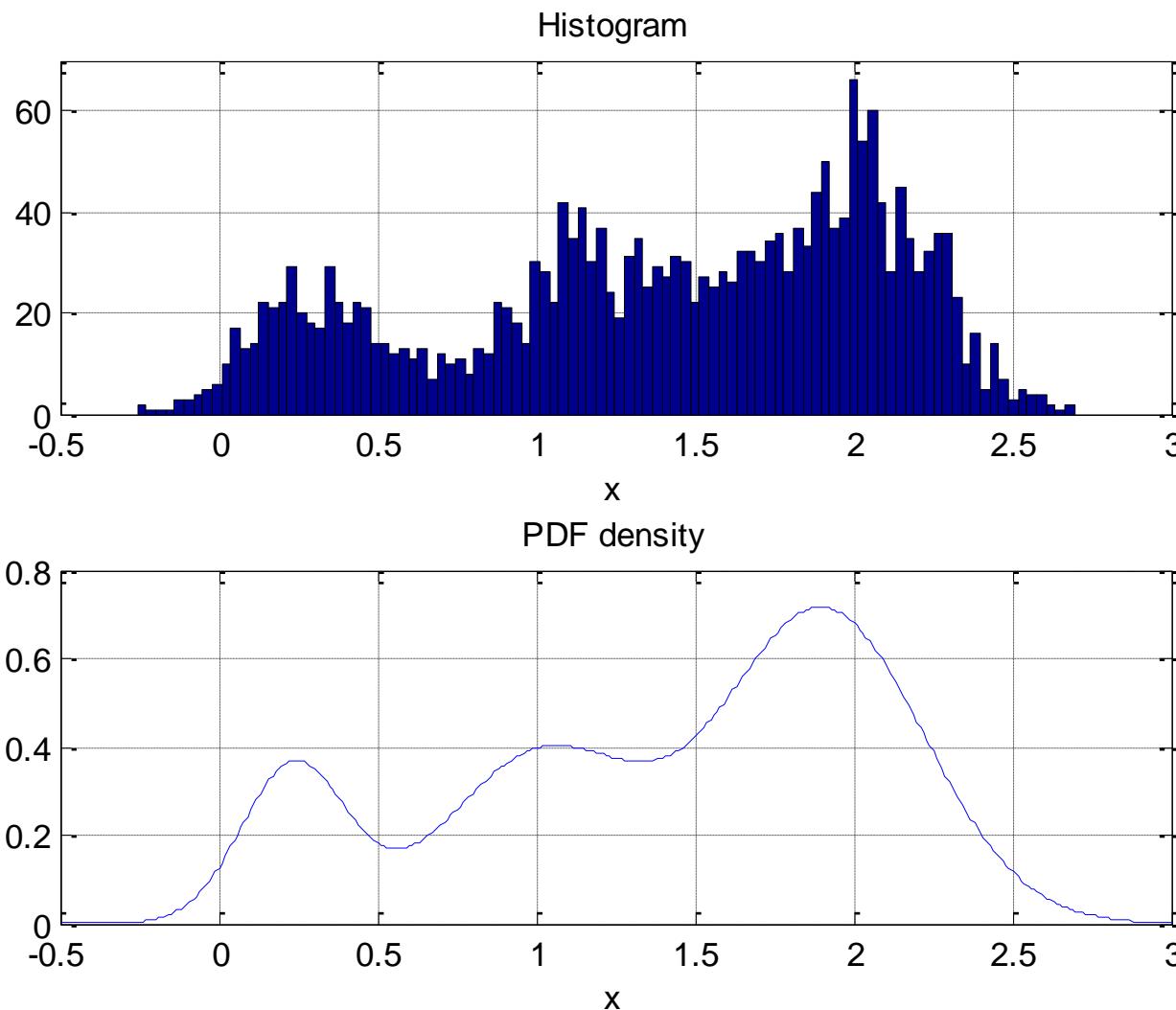
- v is any data vector. $P(v)$ is the probability given to that vector by the Gaussian mixture
- K is the number of Gaussians being mixed
- w_i is the mixture weight of the i^{th} Gaussian. μ_i is its mean and C_i is its covariance
- The Gaussian mixture distribution is also a distribution
 - It is positive everywhere.
 - The total volume under a Gaussian mixture is 1.0.
 - Constraint: the mixture weights w_i must all be positive and sum to 1

Gaussian Mixtures

- A Gaussian mixture can represent data distributions far better than a simple Gaussian
- The two panels show the histogram of an unknown random variable
- The first panel shows how it is modeled by a simple Gaussian
- The second panel models the histogram by a mixture of two Gaussians
- Caveat: It is hard to know the optimal number of Gaussians in a mixture distribution for any random variable



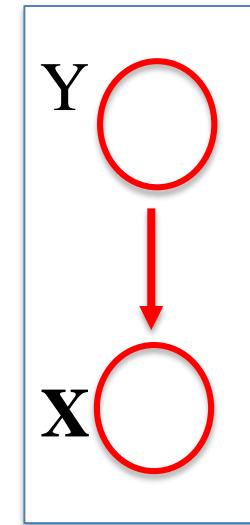
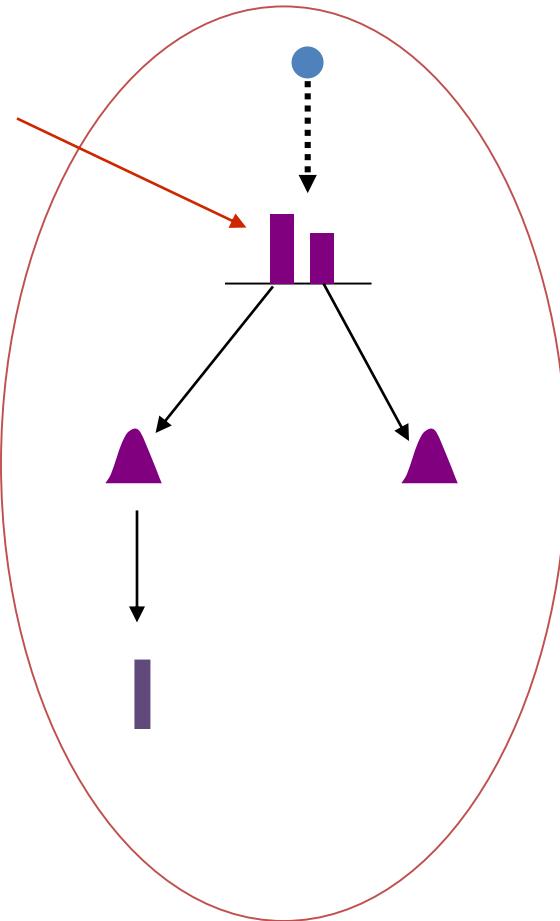
Gaussian Mixtures



Generating an observation from a Gaussian mixture state distribution

First draw the identity of the Gaussian from the a priori probability distribution of Gaussians (mixture weights)

Then draw a vector from the selected Gaussian



HMMs with Gaussian mixture state distributions

- The parameters of an HMM with Gaussian mixture state distributions are:
 - π the set of initial state probabilities for all states
 - T the matrix of transition probabilities
 - A Gaussian mixture distribution for every state in the HMM. The Gaussian mixture for the i^{th} state is characterized by
 - K_i , the number of Gaussians in the mixture for the i^{th} state
 - The set of mixture weights $w_{i,j}$ $0 < j < K_i$
 - The set of Gaussian means $\mu_{i,j}$ $0 < j < K_i$
 - The set of Covariance matrices $C_{i,j}$ $0 < j < K_i$

Three Basic HMM Problems

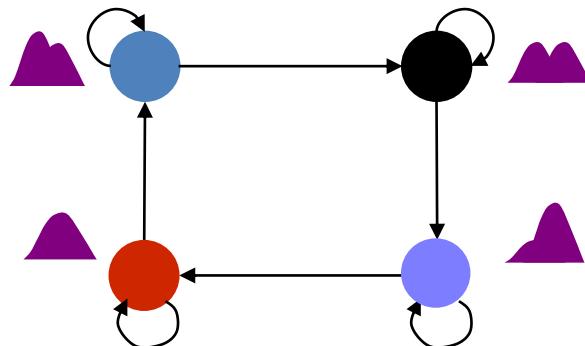
- Given an HMM:
 - What is the probability that it will generate a specific observation sequence
 - Given a observation sequence, how do we determine which observation was generated from which state
 - The state segmentation problem
 - How do we *learn* the parameters of the HMM from observation sequences

Computing the Probability of an Observation Sequence

- Two aspects to producing the observation:
 - Progressing through a sequence of states
 - Producing observations from these states

Progressing through states

HMM assumed to be generating data



state

sequence



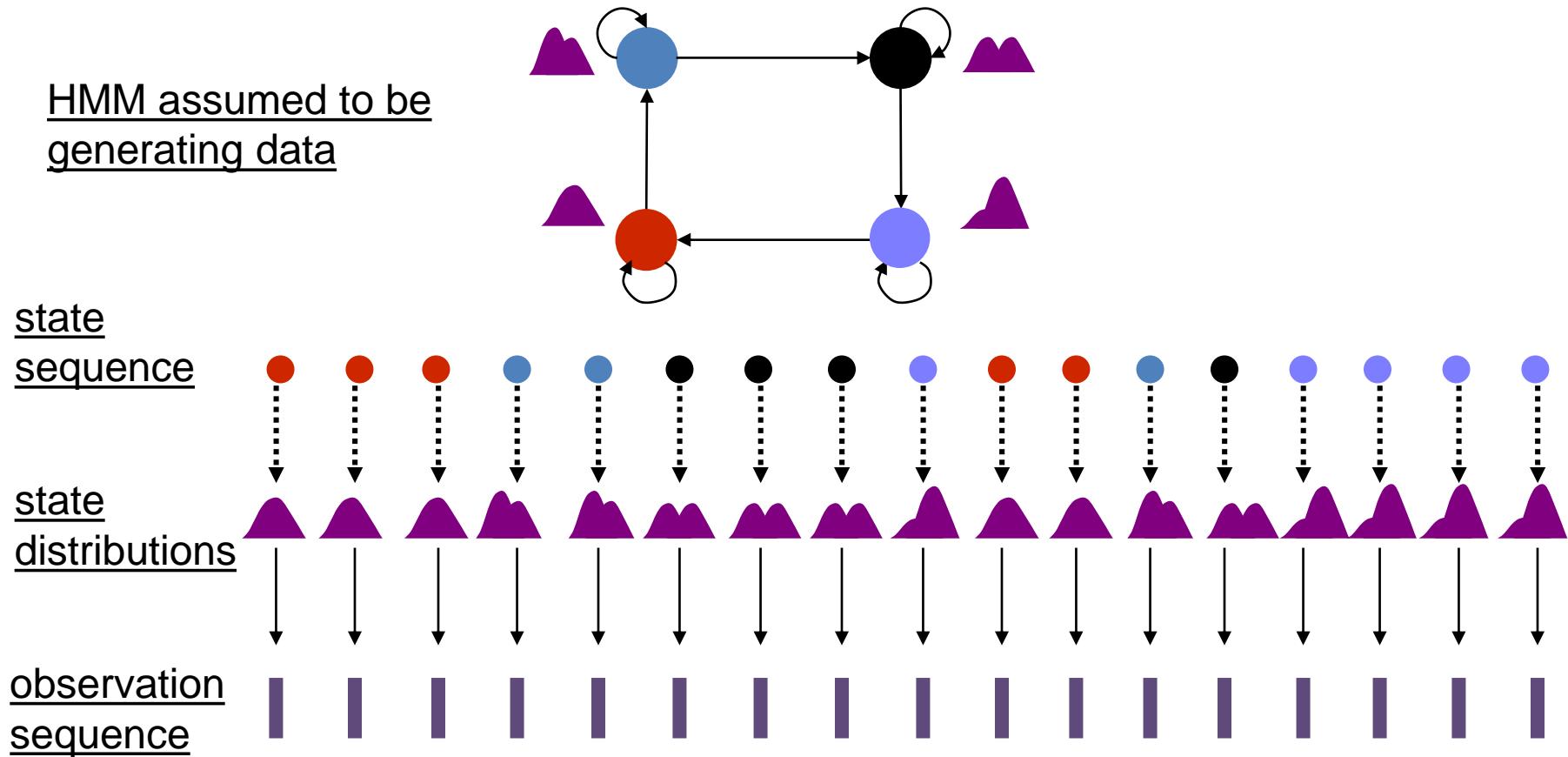
- The process begins at some state (red) here
- From that state, it makes an allowed transition
 - To arrive at the same or any other state
- From that state it makes another allowed transition
 - And so on

Probability that the HMM will follow a particular state sequence

$$P(s_1, s_2, s_3, \dots) = P(s_1) P(s_2|s_1) P(s_3|s_2) \dots$$

- $P(s_1)$ is the probability that the process will initially be in state s_1
- $P(s_i / s_j)$ is the transition probability of moving to state s_i at the next time instant when the system is currently in s_j
 - Also denoted by P_{ij} earlier
 - Related to edge scores in DTW as $T_{ij} = -\log(P(s_i / s_j))$

Generating Observations from States



- At each time it generates an observation from the state it is in at that time

Probability that the HMM will generate a particular observation sequence given a state sequence (state sequence known)

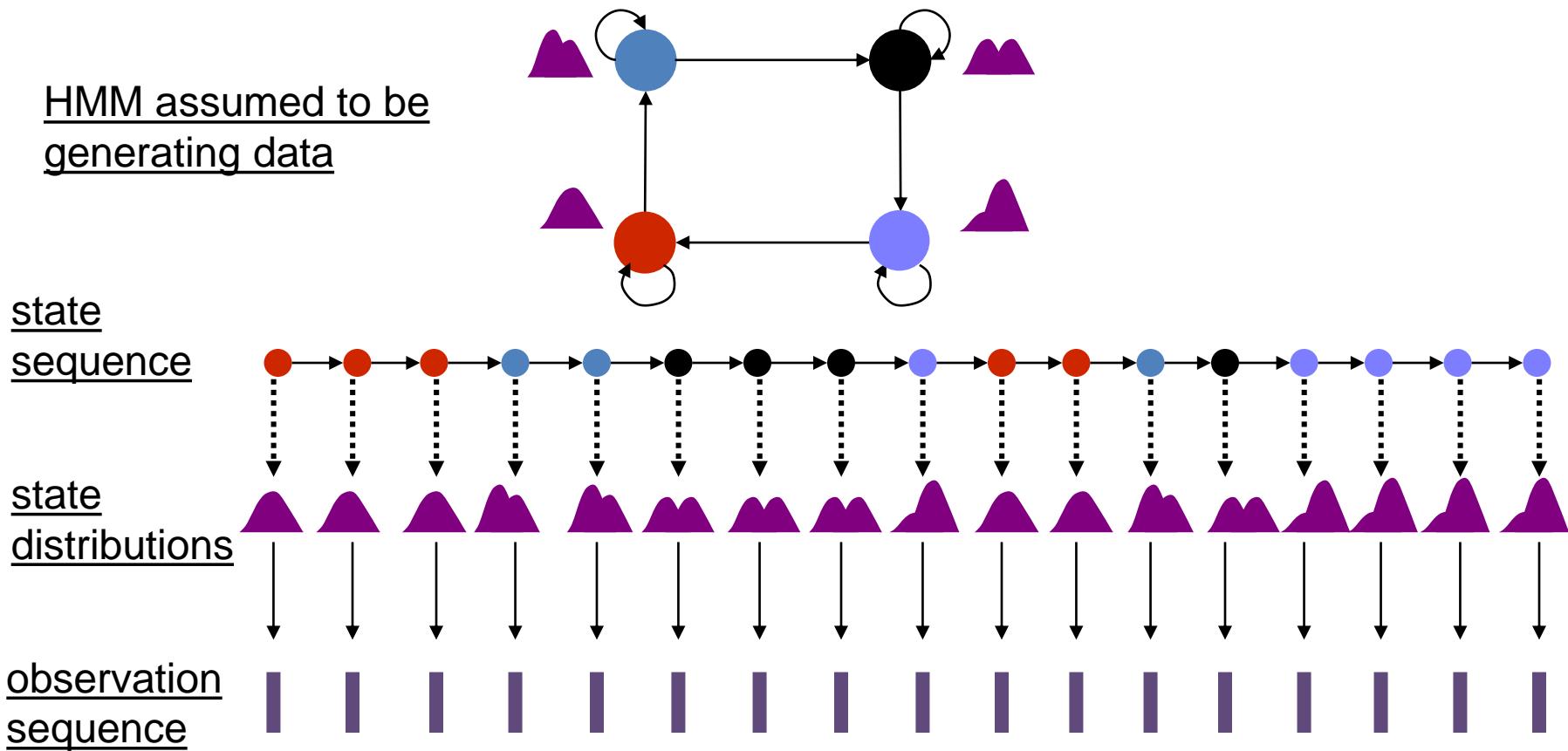
$$P(o_1, o_2, o_3, \dots | s_1, s_2, s_3, \dots) = P(o_1|s_1) P(o_2|s_2) P(o_3|s_3) \dots$$



Computed from the Gaussian or Gaussian mixture for state s_1

- $P(o_i / s_i)$ is the probability of generating observation o_i when the system is in state s_i
 - Related to node scores in DTW trellis as:
 $N_i(O) = -\log(P(o_i / s_i))$

Progressing through States and Producing Observations



- At each time it produces an observation and makes a transition

Probability that the HMM will generate a particular state sequence and, from it, generate a particular observation sequence

$$P(o_1, o_2, o_3, \dots, s_1, s_2, s_3, \dots) =$$

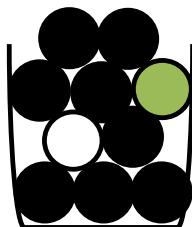
$$P(s_1, s_2, s_3, \dots) P(o_1, o_2, o_3, \dots | s_1, s_2, s_3, \dots) =$$

$$P(s_1) P(s_2 | s_1) P(s_3 | s_2) \dots P(o_1 | s_1) P(o_2 | s_2) P(o_3 | s_3) \dots$$

- Example 1: Marbles in Jars, discrete HMM

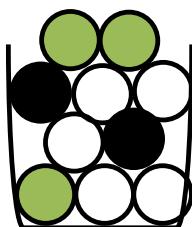
(assume unlimited number of marbles)

State 1



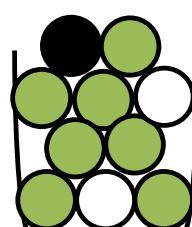
Jar 1

State 2



Jar 2

State 3



Jar 3

$$\begin{aligned} p(b) &= 0.8 \\ p(w) &= 0.1 \\ p(g) &= 0.1 \end{aligned}$$

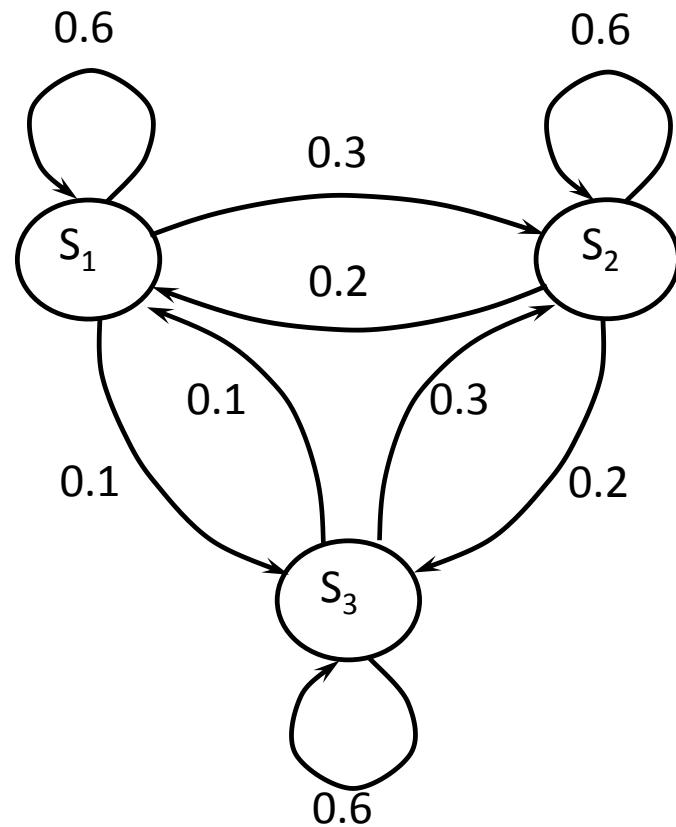
$$\pi_1 = 0.33$$

$$\begin{aligned} p(b) &= 0.2 \\ p(w) &= 0.5 \\ p(g) &= 0.3 \end{aligned}$$

$$\pi_2 = 0.33$$

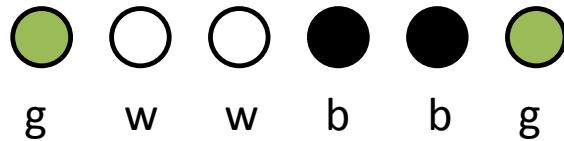
$$\begin{aligned} p(b) &= 0.1 \\ p(w) &= 0.2 \\ p(g) &= 0.7 \end{aligned}$$

$$\pi_3 = 0.33$$



- Example 1: Marbles in Jars (lazy person)
(assume unlimited number of marbles)

- With the following observation:



- What is probability of this observation, given state sequence
 $\{S_3 \ S_2 \ S_2 \ S_1 \ S_1 \ S_3\}$ and the model??

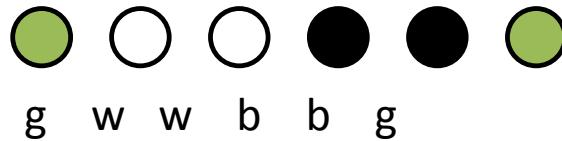
$$= b_3(g) \ b_2(w) \ b_2(w) \ b_1(b) \ b_1(b) \ b_3(g)$$

$$= 0.7 \cdot 0.5 \cdot 0.5 \cdot 0.8 \cdot 0.8 \cdot 0.7$$

$$= 0.0784$$

-
- Example 1: Marbles in Jars (lazy person)
(assume unlimited number of marbles)

- With the *same* observation:



- What is probability of this observation and state sequence
 $\{S_1 S_1 S_3 S_2 S_3 S_1\}$, given the model??

$$= b_1(g) b_1(w) b_3(w) b_2(b) b_3(b) b_1(g)$$

$$= 0.1 \cdot 0.1 \cdot 0.2 \cdot 0.2 \cdot 0.1 \cdot 0.1$$

$$= 4.0 \times 10^{-6}$$

Probability of Generating an Observation Sequence

- If only the observation is known, the precise state sequence followed to produce it is not known
- All possible state sequences must be considered

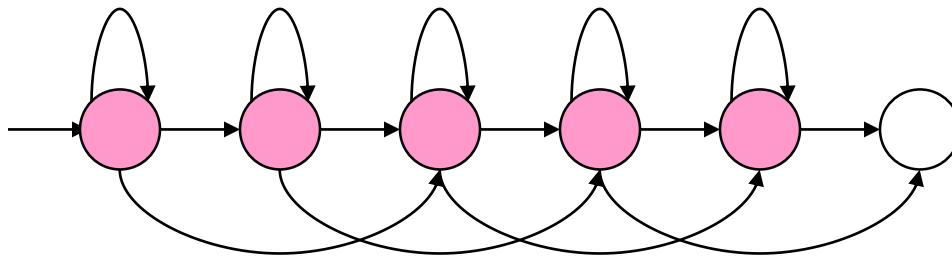
$$P(o_1, o_2, o_3, \dots) = \sum_{\text{all possible state sequences}} P(o_1, o_2, o_3, \dots, s_1, s_2, s_3, \dots) =$$

$$\sum_{\text{all possible state sequences}} P(s_1)P(s_2 | s_1)P(s_3 | s_2)\dots P(o_1 | s_1)P(o_2 | s_2)P(o_3 | s_3)\dots$$

Computing it Efficiently

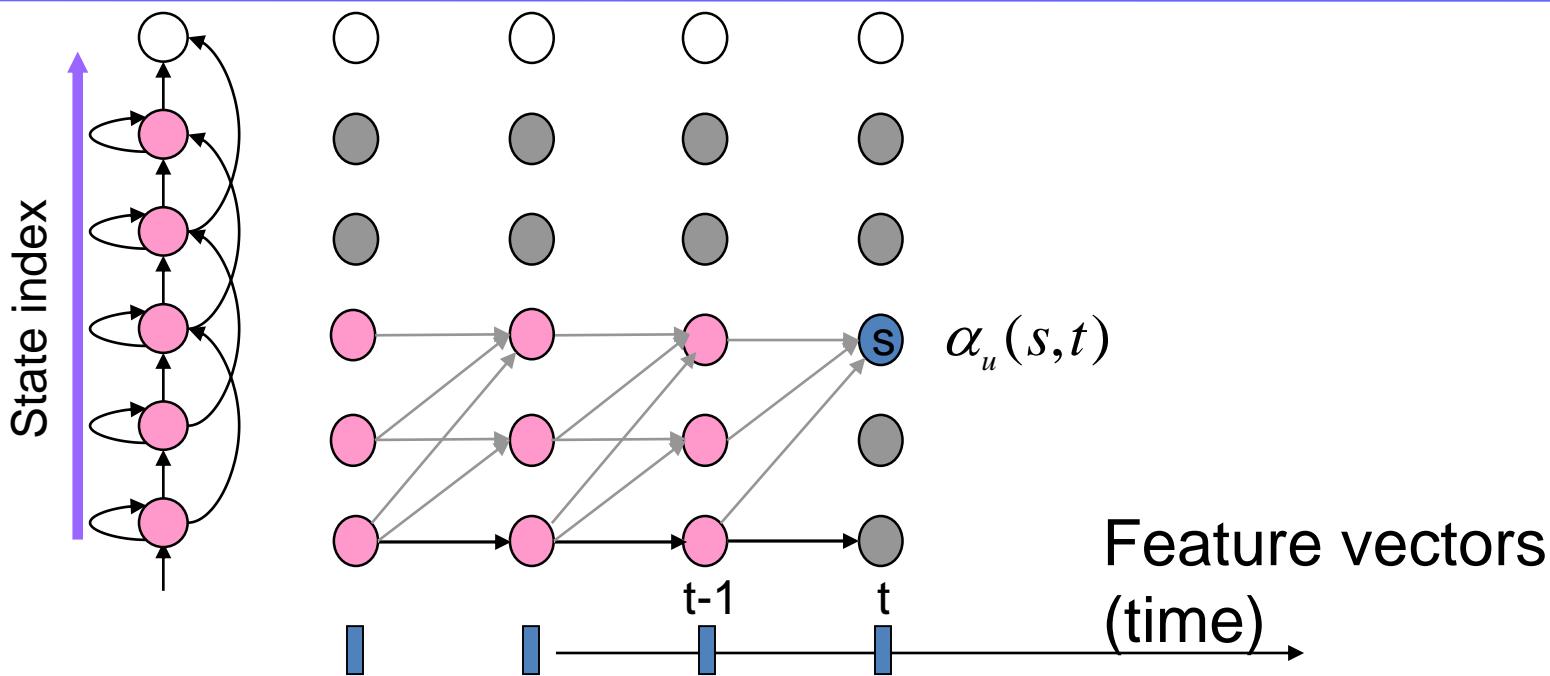
- Explicit summing over all state sequences is not efficient
 - A very large number of possible state sequences
 - For long observation sequences it may be intractable
- Fortunately, we have an efficient algorithm for this:
The forward algorithm
- At each time, for each state compute the total probability of all state sequences that generate observations until that time and end at that state

Illustrative Example



- Consider a generic HMM with 5 states and a “terminating state”. We wish to find the probability of the best state sequence for an observation sequence assuming it was generated by this HMM
 - $P(s_i) = 1$ for state 1 and 0 for others
 - The arrows represent transition for which the probability is not 0.
 - May represent it as $P(s_i | s_j) = a_{ij}$
 - We sometimes also represent the state output probability of s_i as $P(o_t | s_i) = b_i(t)$ for brevity

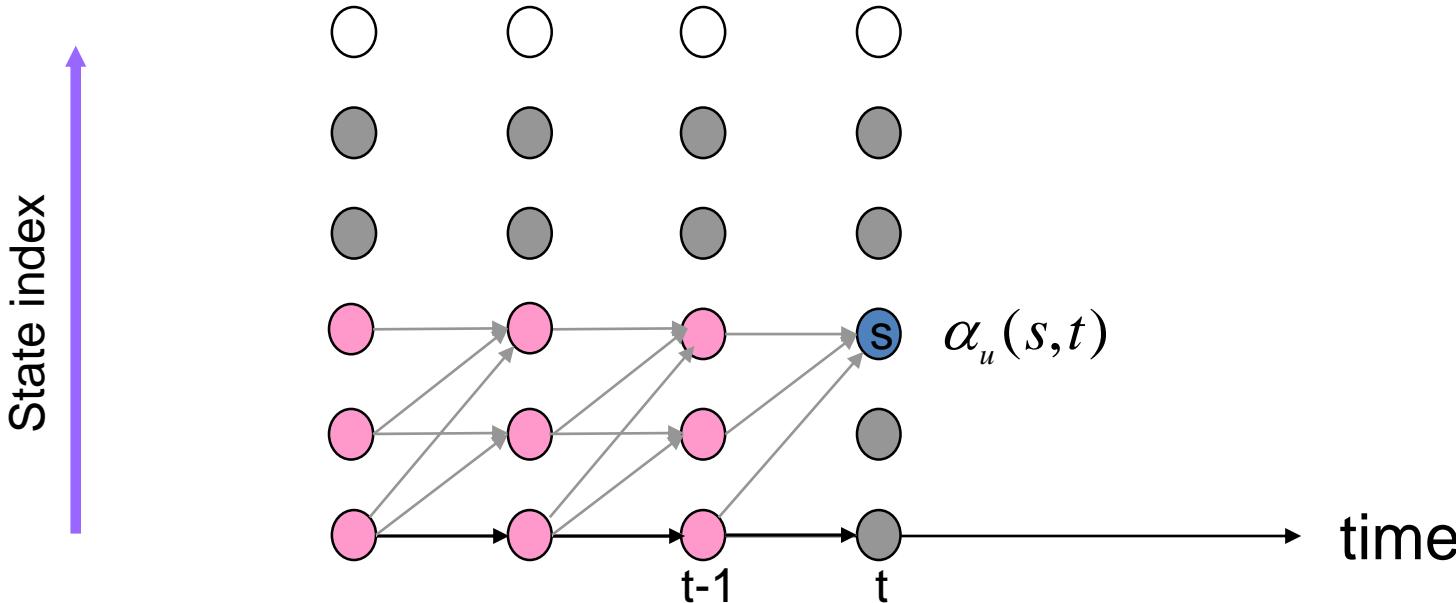
Diversion: The HMM Trellis



- The trellis is a graphical representation of all possible state sequences through the HMM to produce a given observation
 - Analogous to the DTW search graph / trellis
- The Y-axis represents HMM states, X axis represents observations
- Edges in trellis represent valid transitions in the HMM over a single time step
- Every node represents the event of a particular observation being generated from a particular state

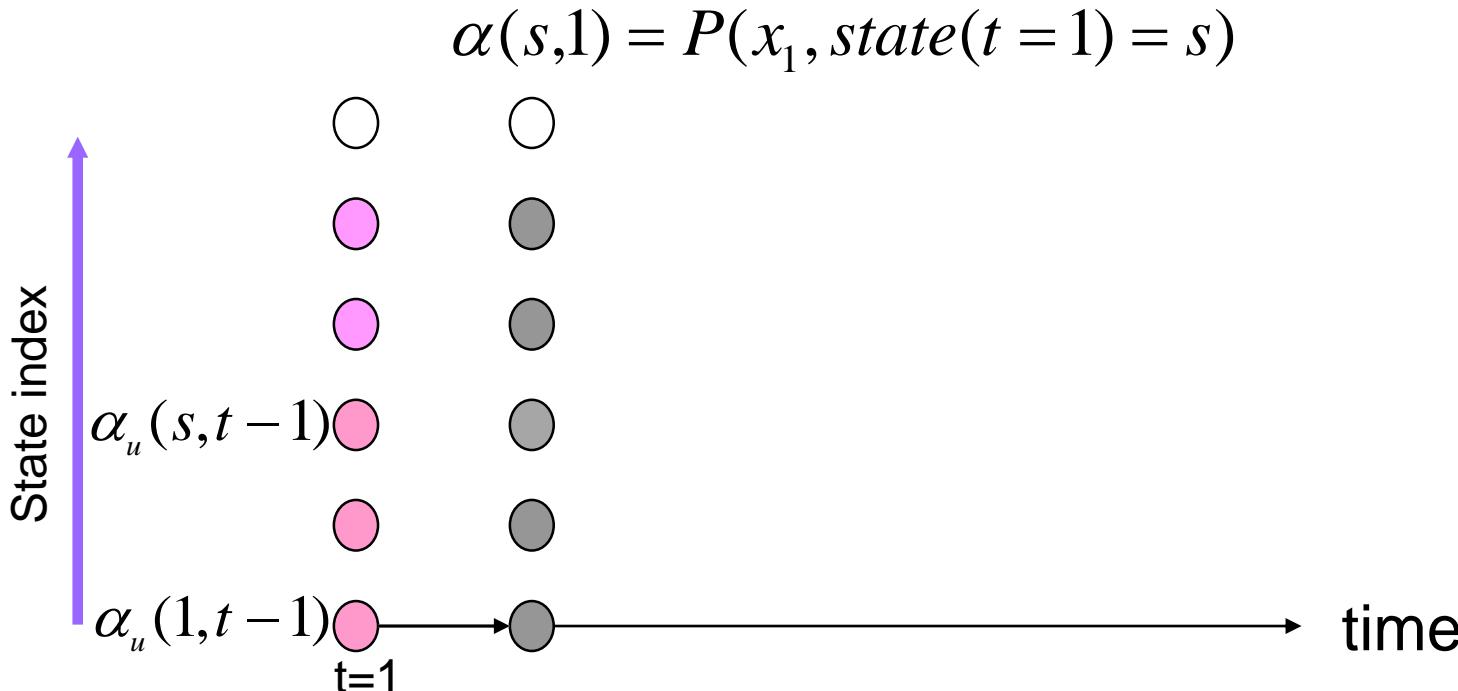
The Forward Algorithm

$$\alpha_u(s, t) = P(x_{u,1}, x_{u,2}, \dots, x_{u,t}, \text{state}(t) = s | \lambda)$$



- $\alpha_u(s, t)$ is the total probability of ALL state sequences that end at state s at time t , and all observations until x_t

The Forward Algorithm at t=1

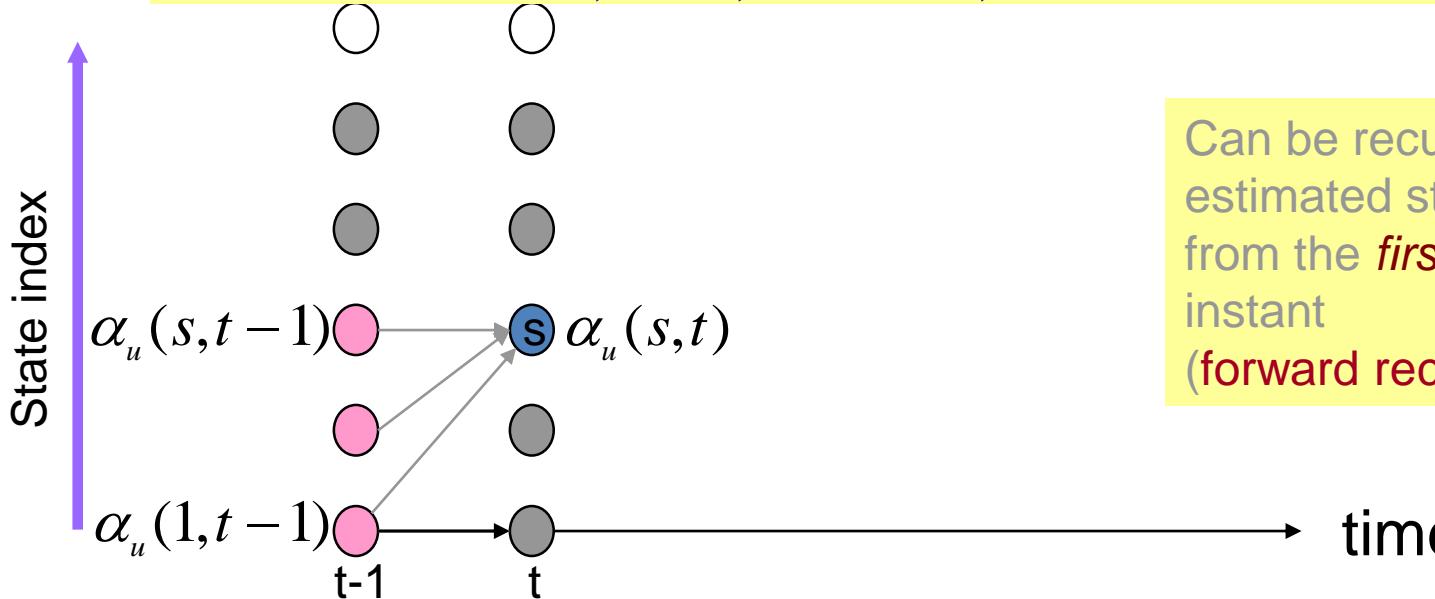


$$\alpha(s,1) = \pi(s)P(x_1 | s)$$

- $\alpha(s,1)$ is simply the probability of being in state s at $t=1$ *and* generating observation x_1 from s

The Forward Algorithm

$$\alpha_u(s, t) = P(x_{u,1}, x_{u,2}, \dots, x_{u,t}, \text{state}(t) = s | \lambda)$$



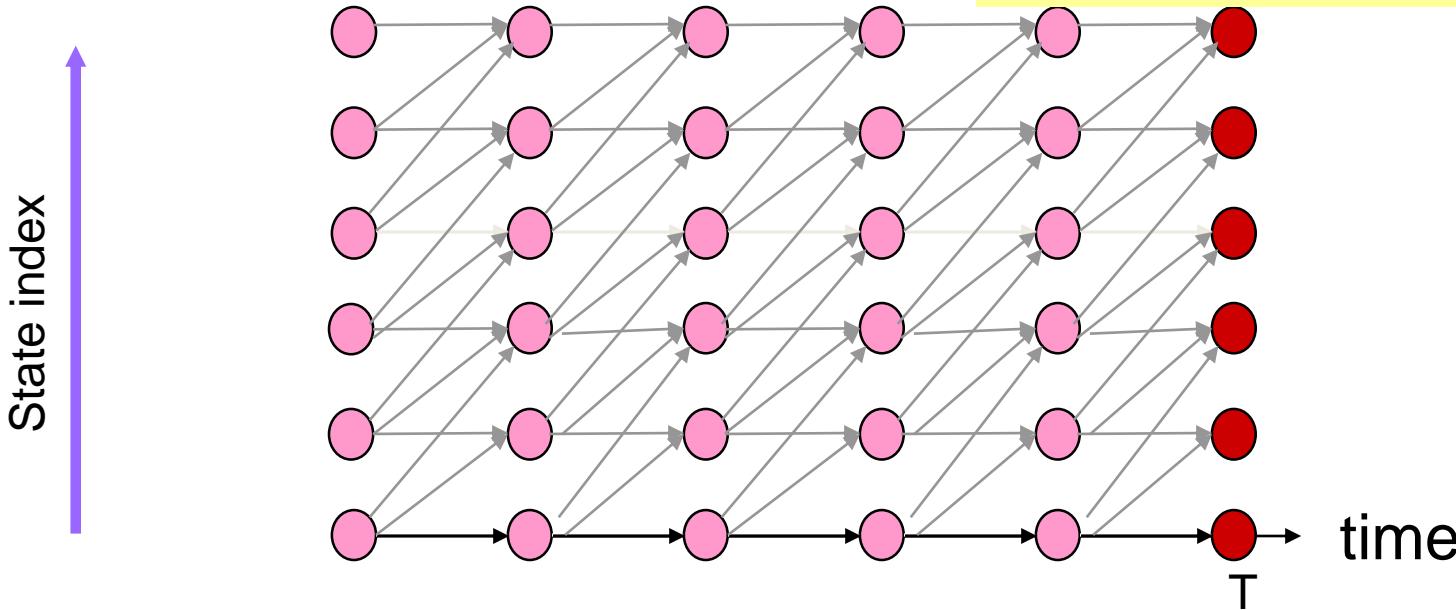
Can be recursively estimated starting from the *first* time instant
(forward recursion)

$$\alpha_u(s, t) = \sum_{s'} \alpha_u(s', t - 1) P(s | s') P(x_{u,t} | s)$$

- $\alpha_u(s, t)$ can be recursively computed in terms of $\alpha_u(s', t')$, the forward probabilities at time $t-1$

The Forward Algorithm

$$Totalprob = \sum_s \alpha_u(s, T)$$



- In the final observation the alpha at each state gives the probability of all state sequences ending at that state
- The total probability of the observation is the sum of the alpha values at all states

The forward algorithm

1. Initialize all alpha terms at t=1:

$$\alpha(s,1) = \pi(s)P(x_1 | s)$$

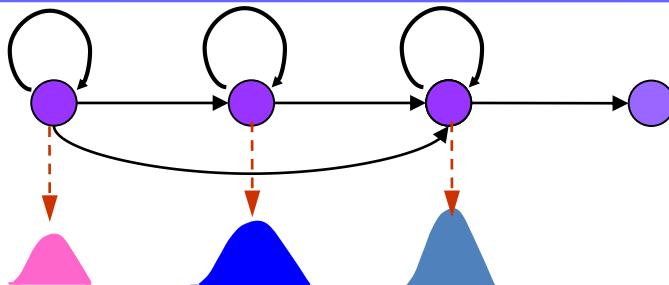
2. Recursive estimate alphas for all subsequent time steps

$$\alpha(s,t) = \sum_{s'} \alpha(s',t-1)P(s | s')P(x_t | s)$$

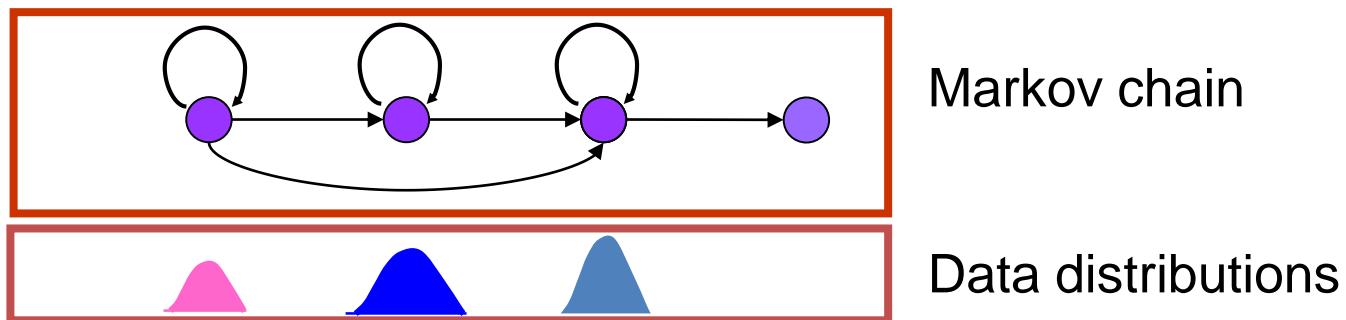
3. Compute overall probability

$$Totalprob = \sum_s \alpha(s,T)$$

RECAP: Hidden Markov Models



- A Hidden Markov Model consists of two components
 - A state/transition backbone that specifies how many states there are, and how they can follow one another
 - A set of probability distributions, one for each state, which specifies the distribution of all vectors in that state



- This can be factored into two separate probabilistic entities
 - A probabilistic Markov chain with states and transitions
 - A set of data probability distributions, associated with the states

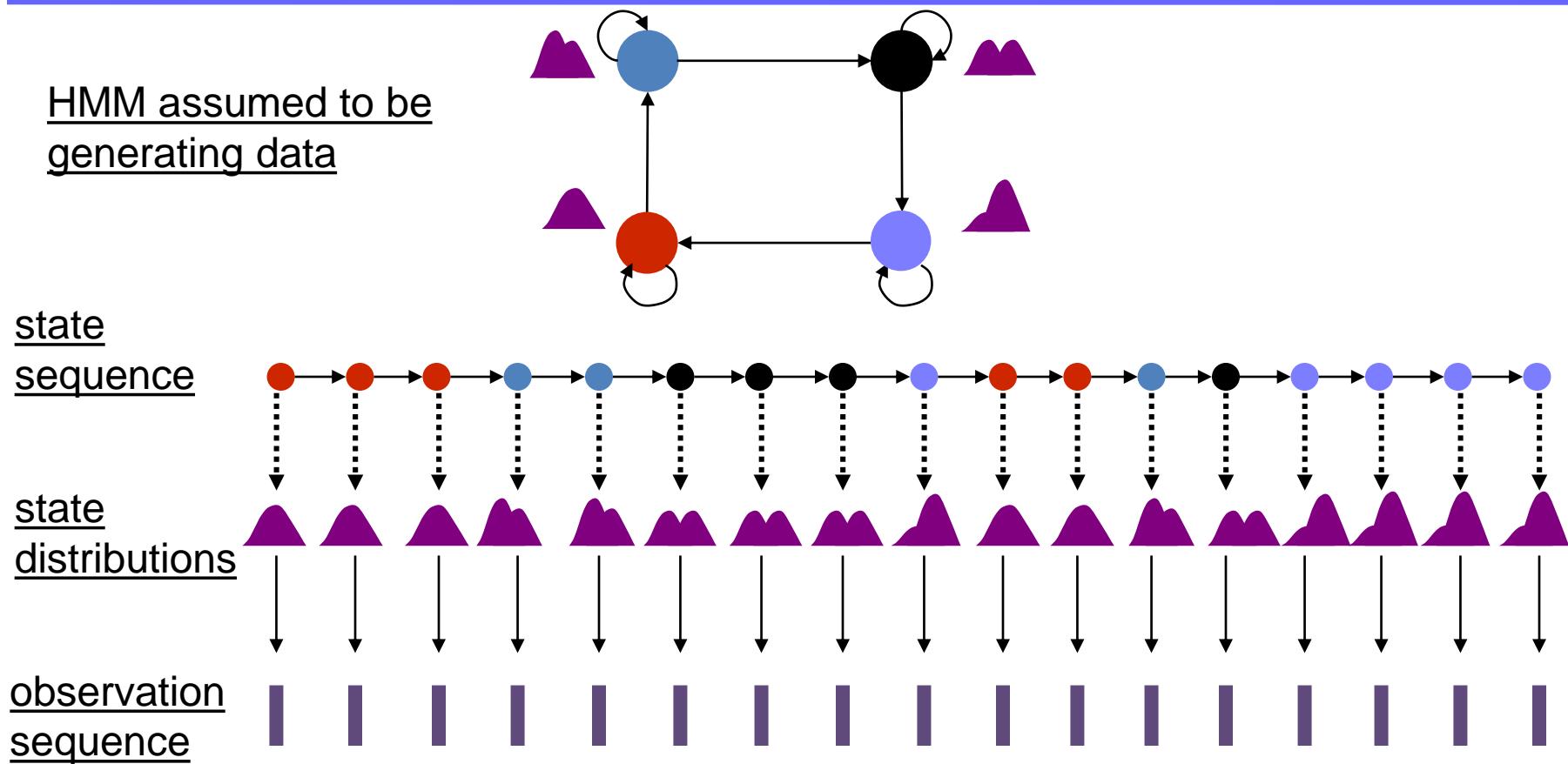
RECAP: Three Basic HMM Problems

- Given an HMM:
 - *What is the probability that it will generate a specific observation sequence (already done)*
 - Given a observation sequence, how do we determine which observation was generated from which state
 - The state segmentation problem
- How do we *learn* the parameters of the HMM from observation sequences

Problem 2: The state segmentation problem

- Given only a sequence of observations, how do we determine which sequence of states was followed in producing it?

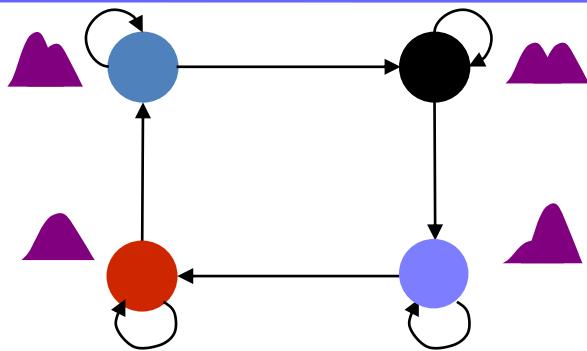
The HMM as a generator



- The process goes through a series of states and produces observations from them

States are Hidden

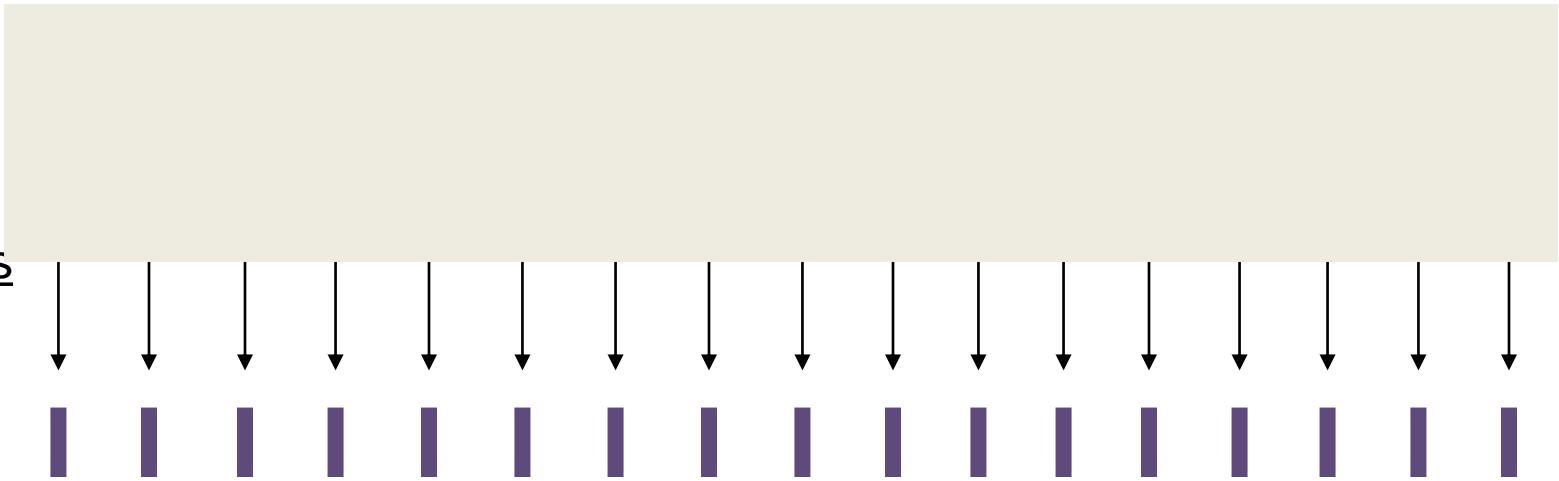
HMM assumed to be generating data



state sequence

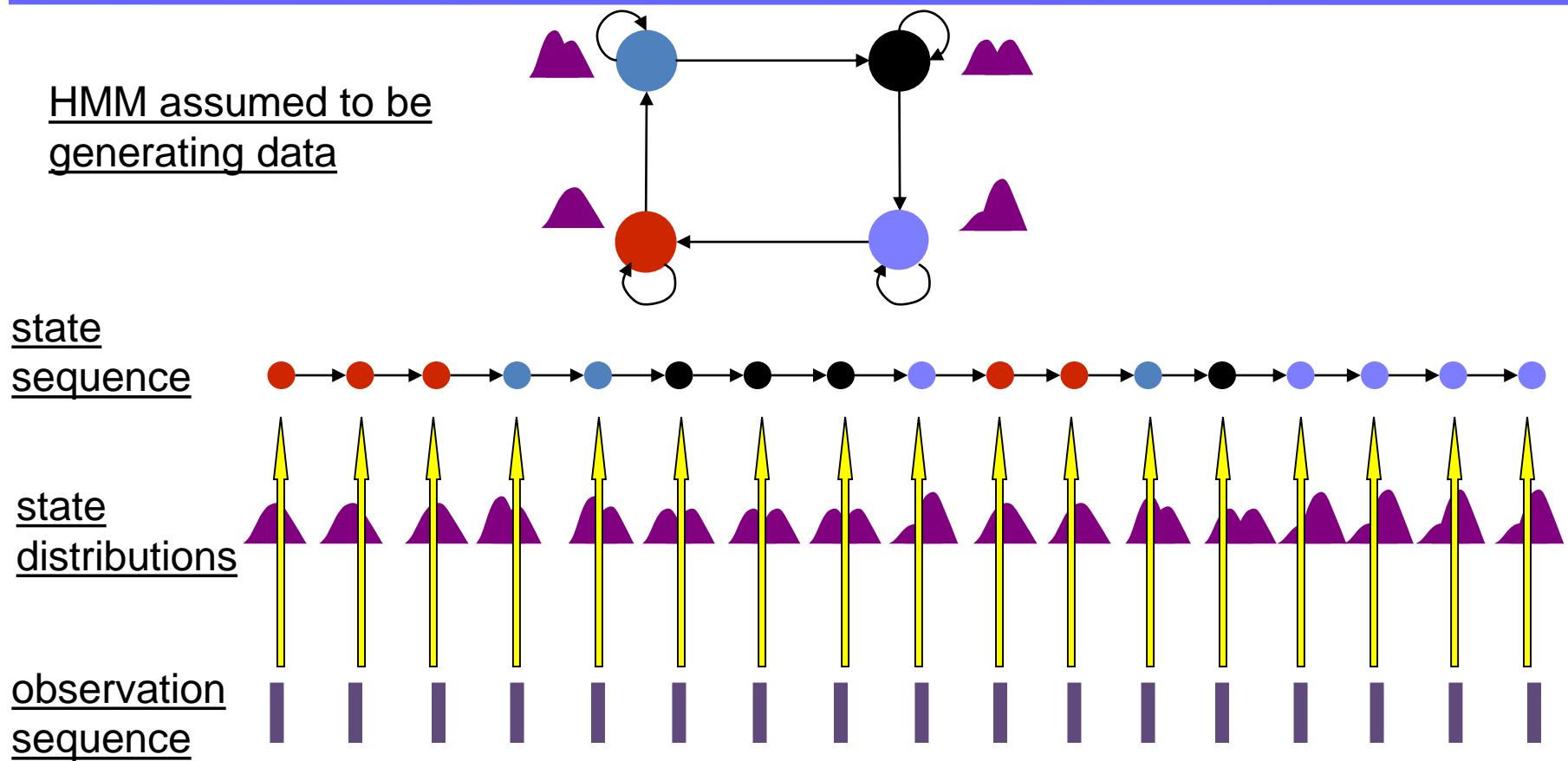
state distributions

observation sequence



- The observations do not reveal the underlying state

The state segmentation problem



- State segmentation: Estimate state sequence given observations

Estimating the State Sequence

- Any number of state sequences could have been traversed in producing the observation
 - In the worst case *every* state sequence may have produced it
- Solution: Identify the most *probable* state sequence
 - The state sequence for which the probability of progressing through that sequence and generating the observation sequence is maximum
 - i.e $P(o_1, o_2, o_3, \dots, s_1, s_2, s_3, \dots)$ is maximum

Estimating the state sequence

- Once again, exhaustive evaluation is impossibly expensive
- But once again a simple dynamic-programming solution is available

$$P(o_1, o_2, o_3, \dots, s_1, s_2, s_3, \dots) =$$

$$P(s_1)P(s_2 | s_1)P(s_3 | s_2) \dots P(o_1 | s_1)P(o_2 | s_2)P(o_3 | s_3) \dots$$

- Needed:

$$\arg \max_{s_1, s_2, s_3, \dots} P(o_1 | s_1)P(s_1)P(o_2 | s_2)P(s_2 | s_1)P(o_3 | s_3)P(s_3 | s_2)$$

Estimating the state sequence

- Once again, exhaustive evaluation is impossibly expensive
- But once again a simple dynamic-programming solution is available

$$P(o_1, o_2, o_3, \dots, s_1, s_2, s_3, \dots) =$$

$$P(s_1)P(s_2 | s_1)P(s_3 | s_2) \dots P(o_1 | s_1)P(o_2 | s_2)P(o_3 | s_3) \dots$$

- Needed:

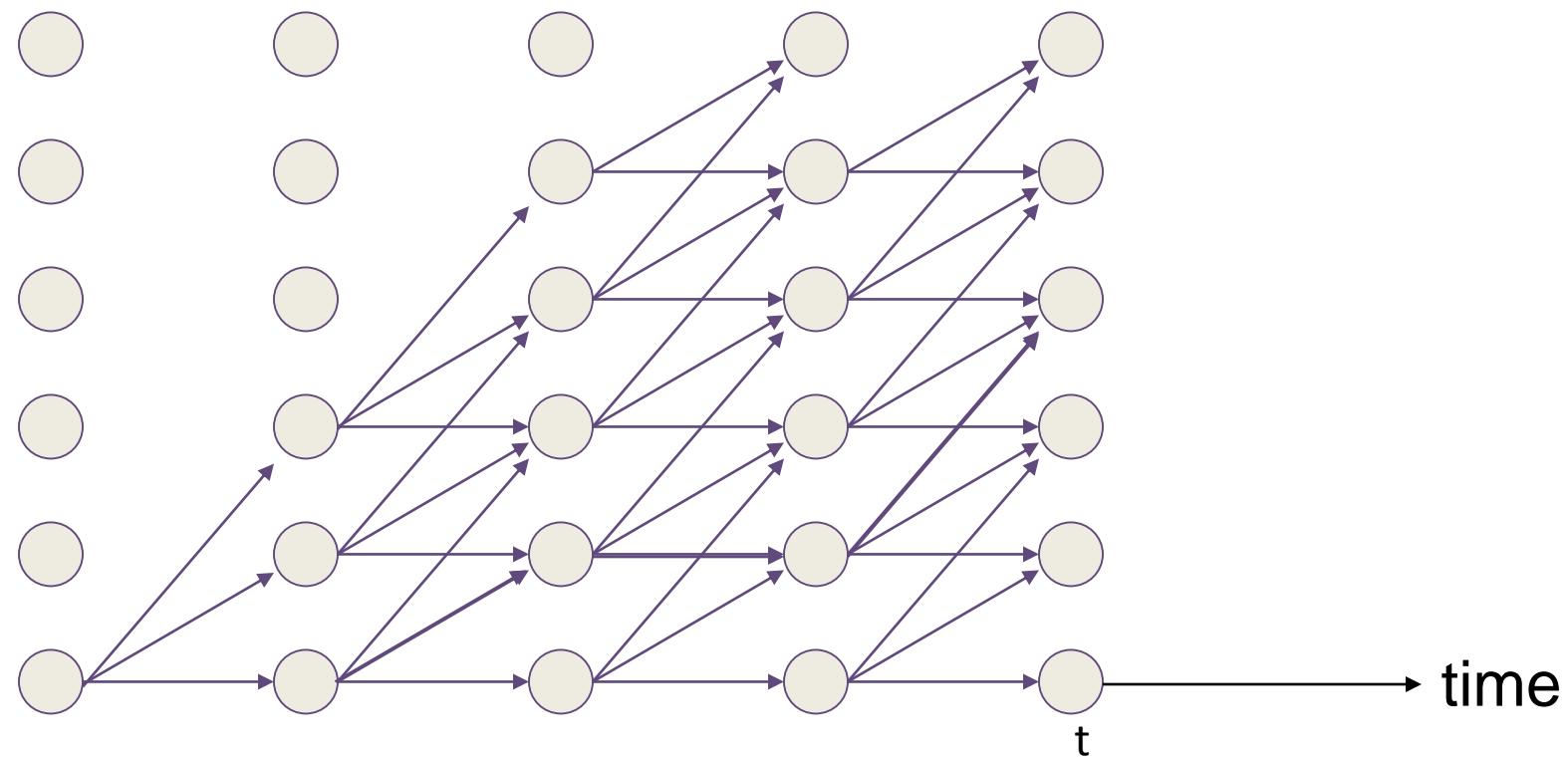
$$\arg \max_{s_1, s_2, s_3, \dots} P(o_1 | s_1)P(s_1)P(o_2 | s_2)P(s_2 | s_1)P(o_3 | s_3)P(s_3 | s_2)$$

The state sequence

- The probability of a state sequence $?, ?, ?, ?, ?, s_x, s_y$ ending at time t is simply
 - $P(?, ?, ?, ?, s_x, s_y) = P(?, ?, ?, ?, s_x) P(o_t|s_y)P(s_y|s_x)$
- The *best* state sequence that ends with s_x, s_y at t will have a probability equal to the probability of the best state sequence ending at $t-1$ at s_x times $P(o_t|s_y)P(s_y|s_x)$
 - Since the last term is independent of the state sequence leading to s_x at $t-1$

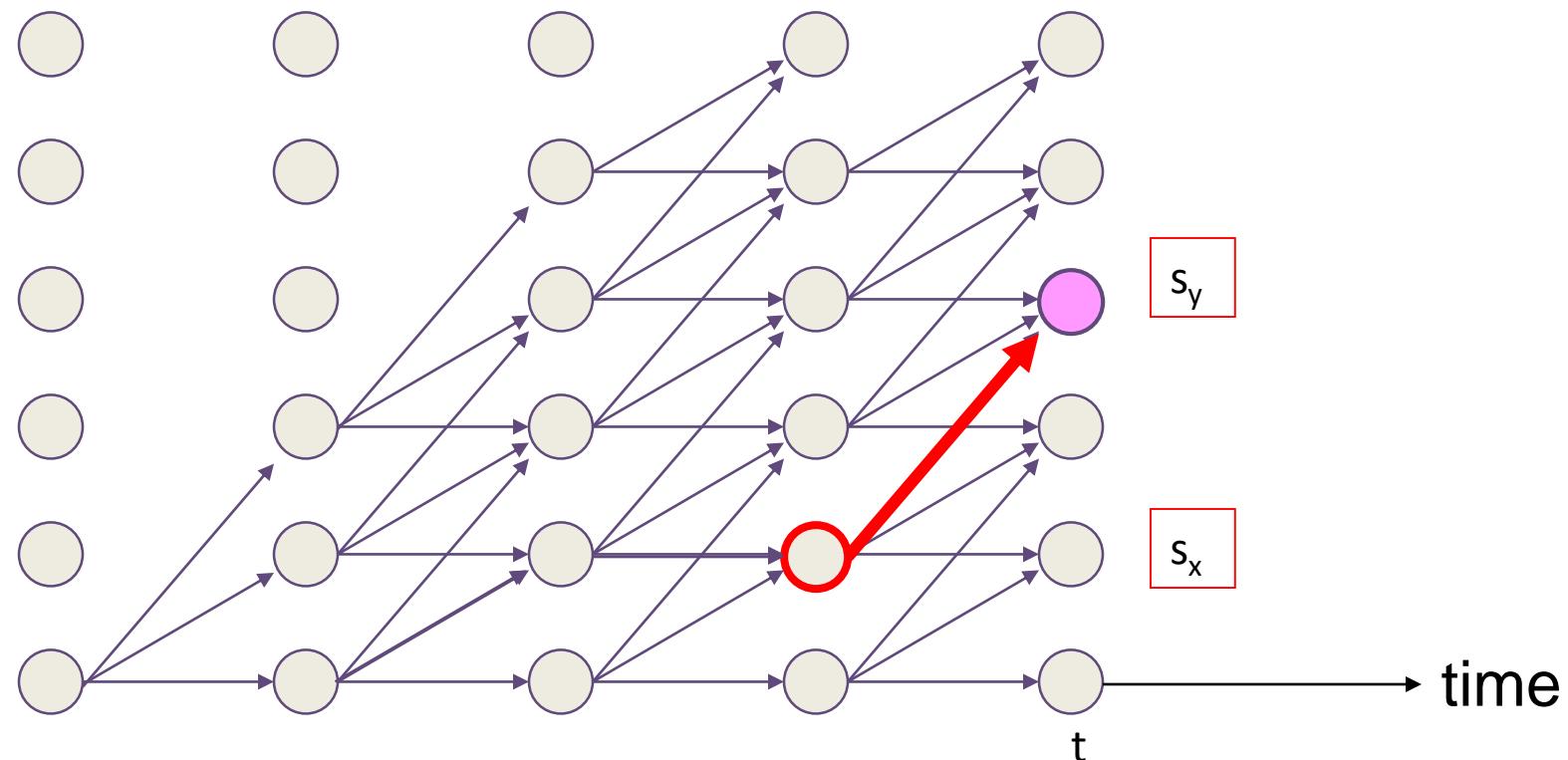
Trellis

- The graph below shows the set of all possible state sequences through this HMM in five time intants



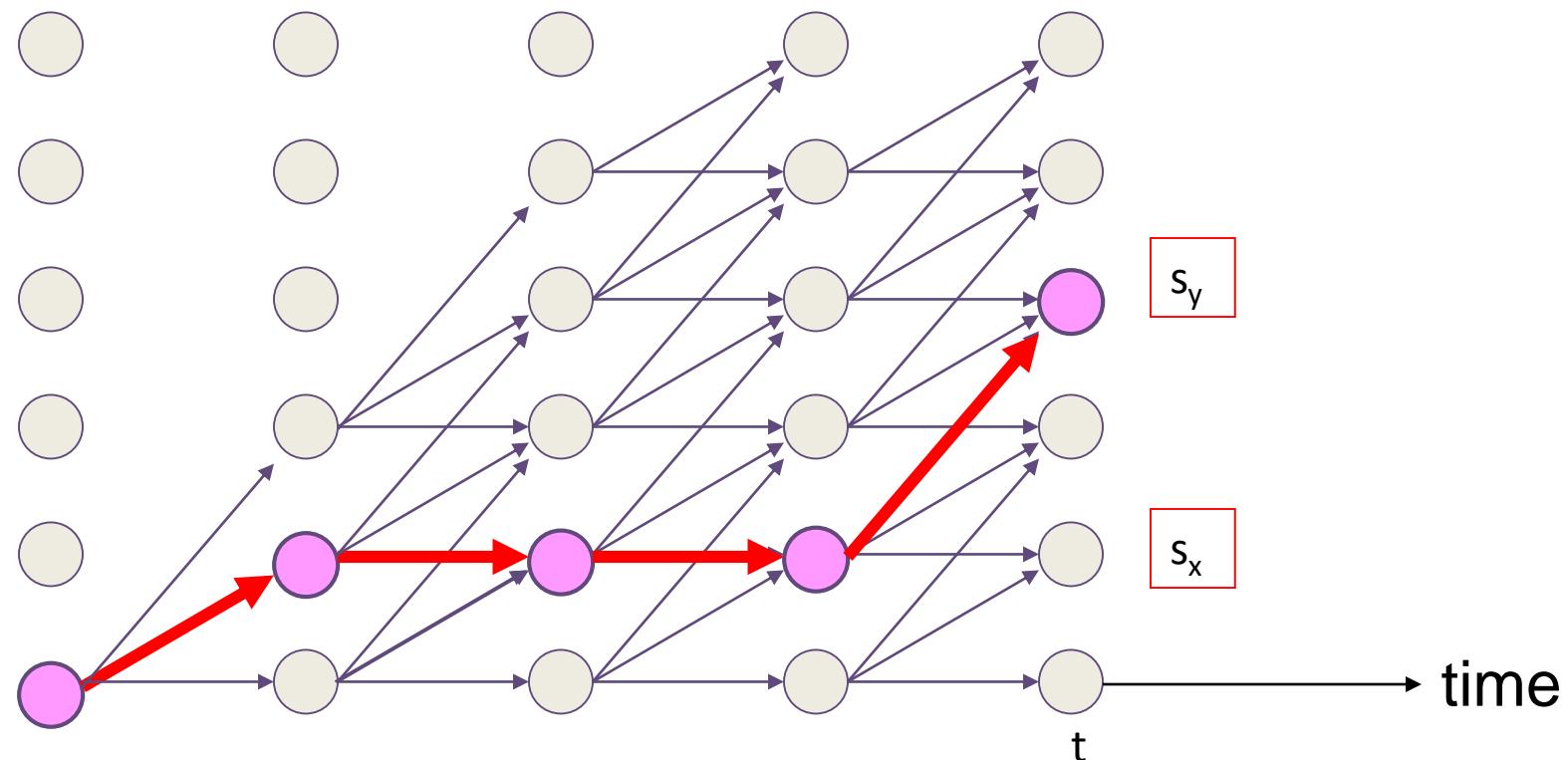
The cost of extending a state sequence

- The cost of extending a state sequence ending at s_x is only dependent on the transition from s_x to s_y , and the observation probability at s_y



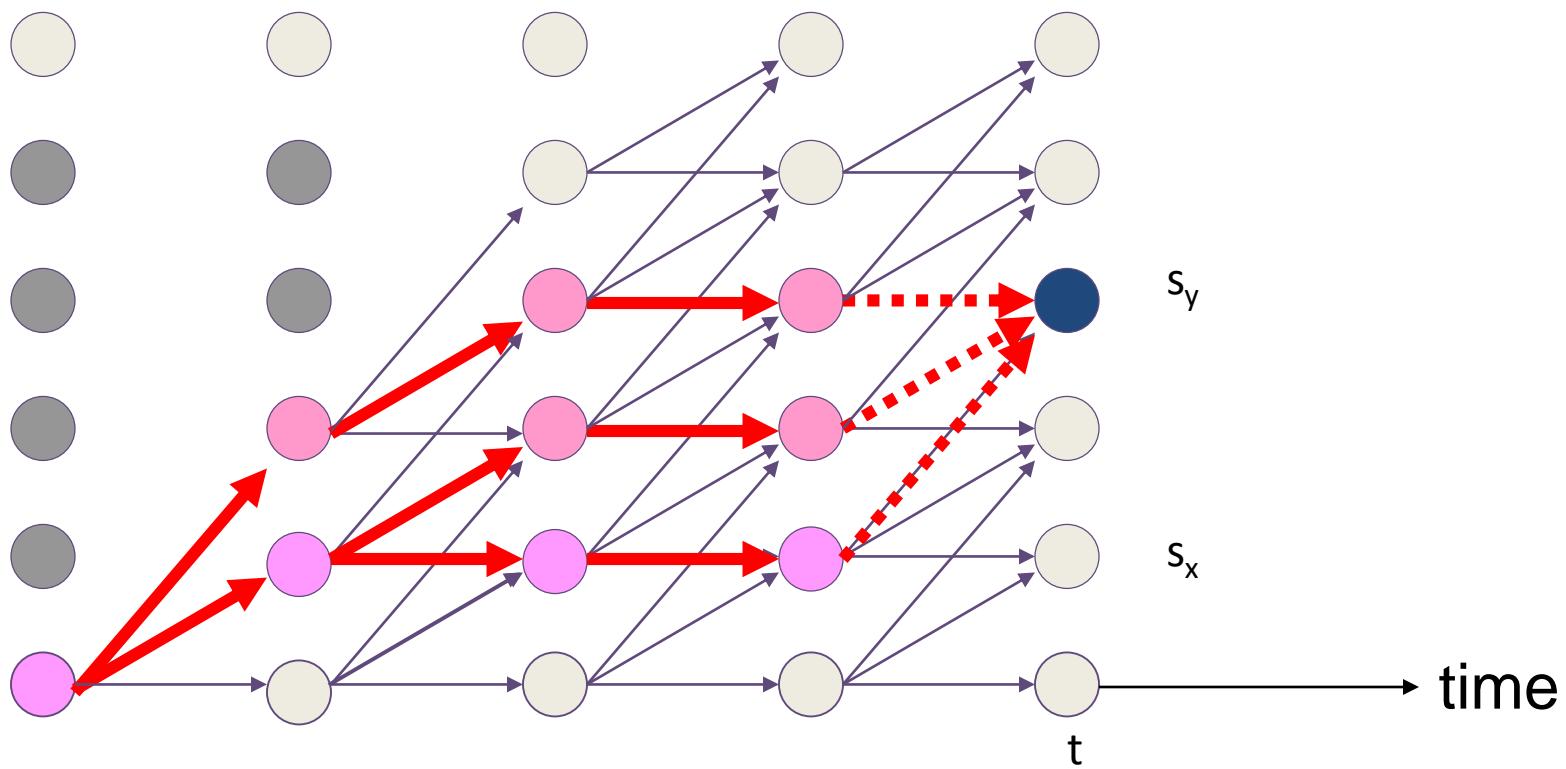
The cost of extending a state sequence

- The best path to s_y through s_x is simply an extension of the best path to s_x



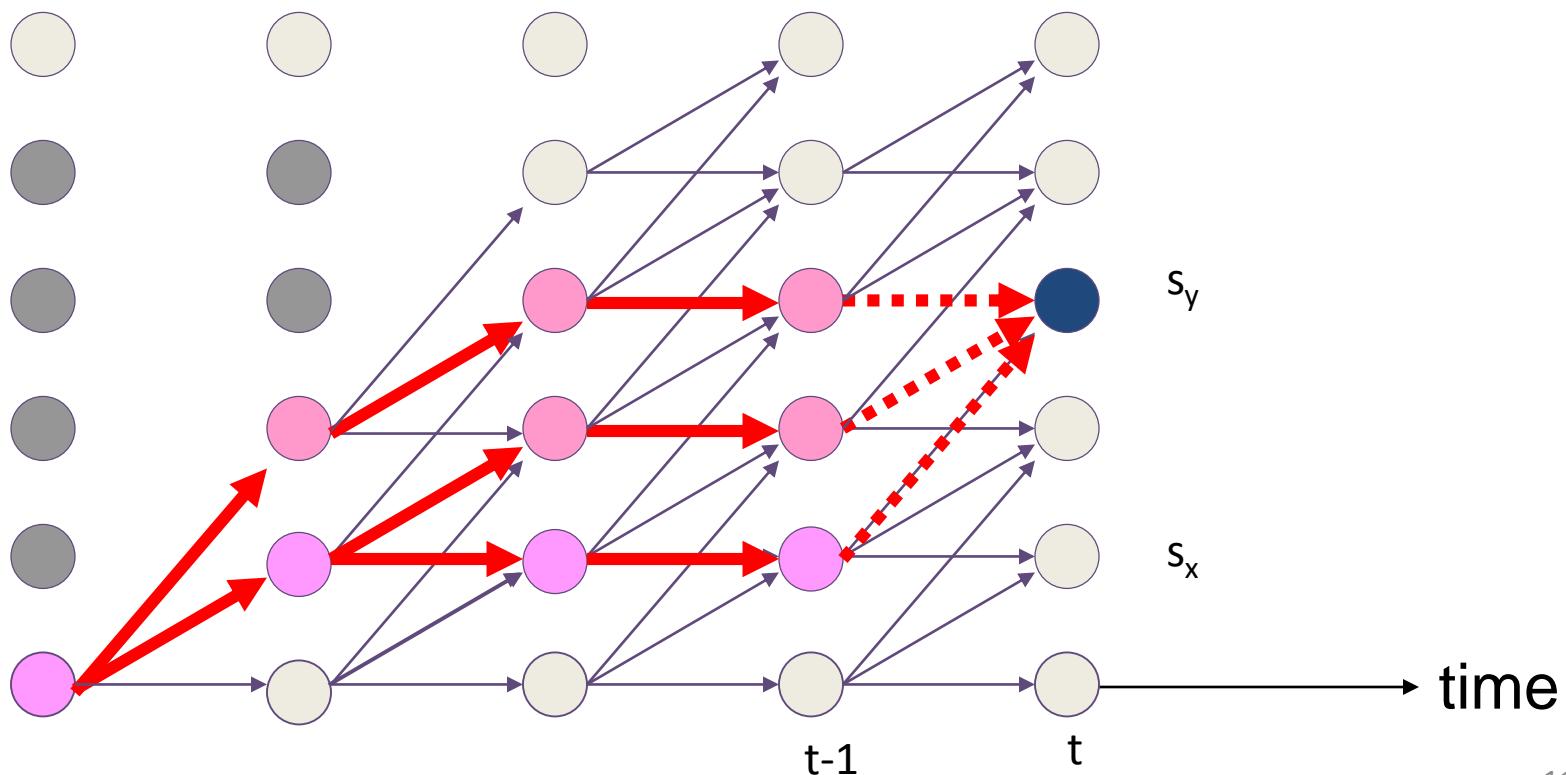
The Recursion

- The overall best path to s_y is an extension of the best path to one of the states at the previous time



The Recursion

- Bestpath prob(s_y, t) =
Best_? (Bestpath prob($s_?, t$) $P(s_y | s_?) P(o_t | s_y)$)



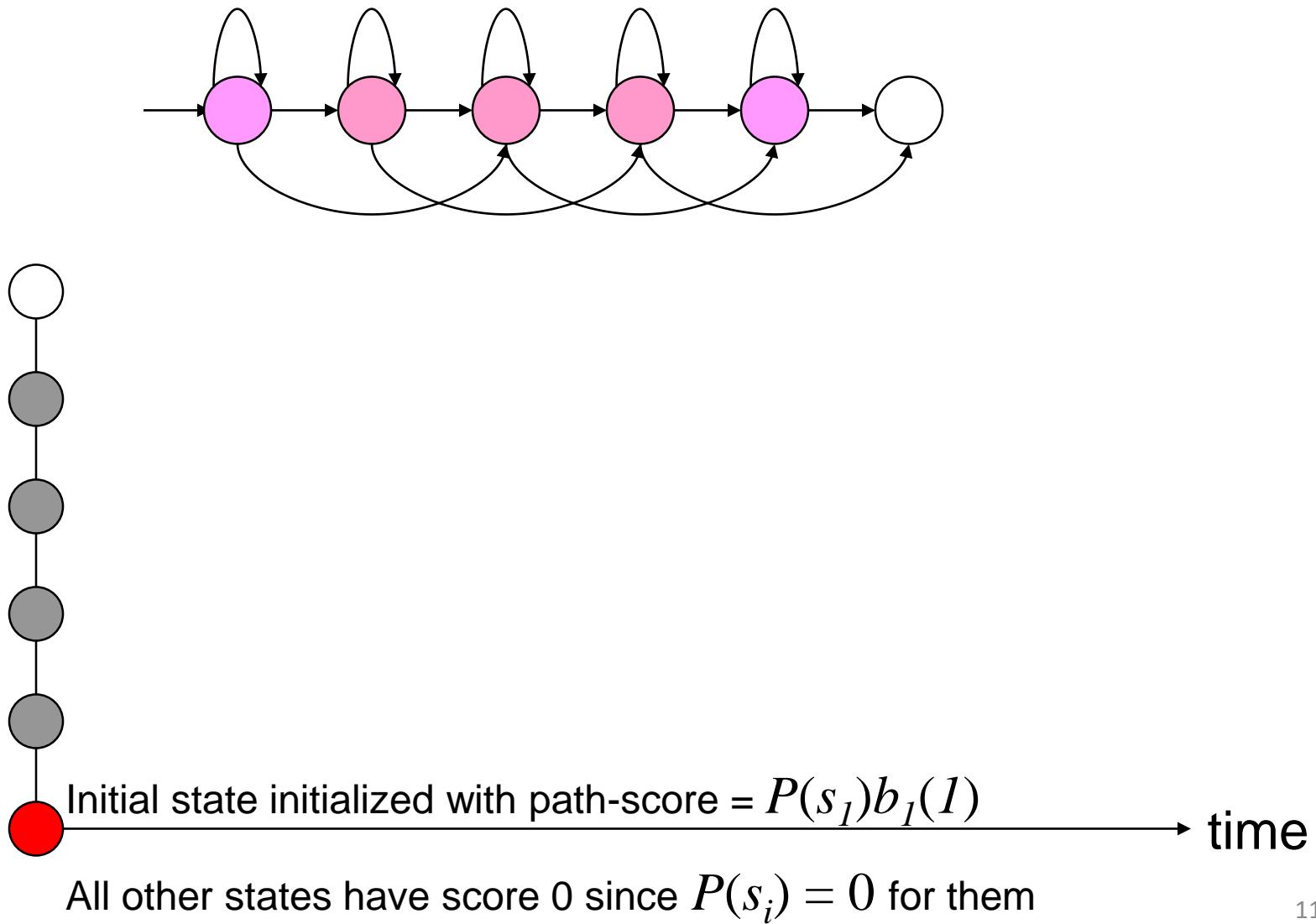
Finding the best state sequence

- This gives us a simple recursive formulation to find the overall best state sequence:
 1. The best state sequence $X_{1,i}$ of length 1 ending at state s_i is simply s_i .
 - The probability $C(X_{1,i})$ of $X_{1,i}$ is $P(o_1 | s_i) P(s_i)$
 2. The best state sequence of length $t+1$, terminating at state j , is simply given by
 - $(\text{argmax}_{X_{t,i}} C(X_{t,i}) P(o_{t+1} | s_j) P(s_j | s_i))$
 - The best overall state sequence for an utterance of length T is given by
$$\text{argmax}_{X_{t,i}} C(X_{T,i})$$
 - The state sequence of length T with the highest overall probability

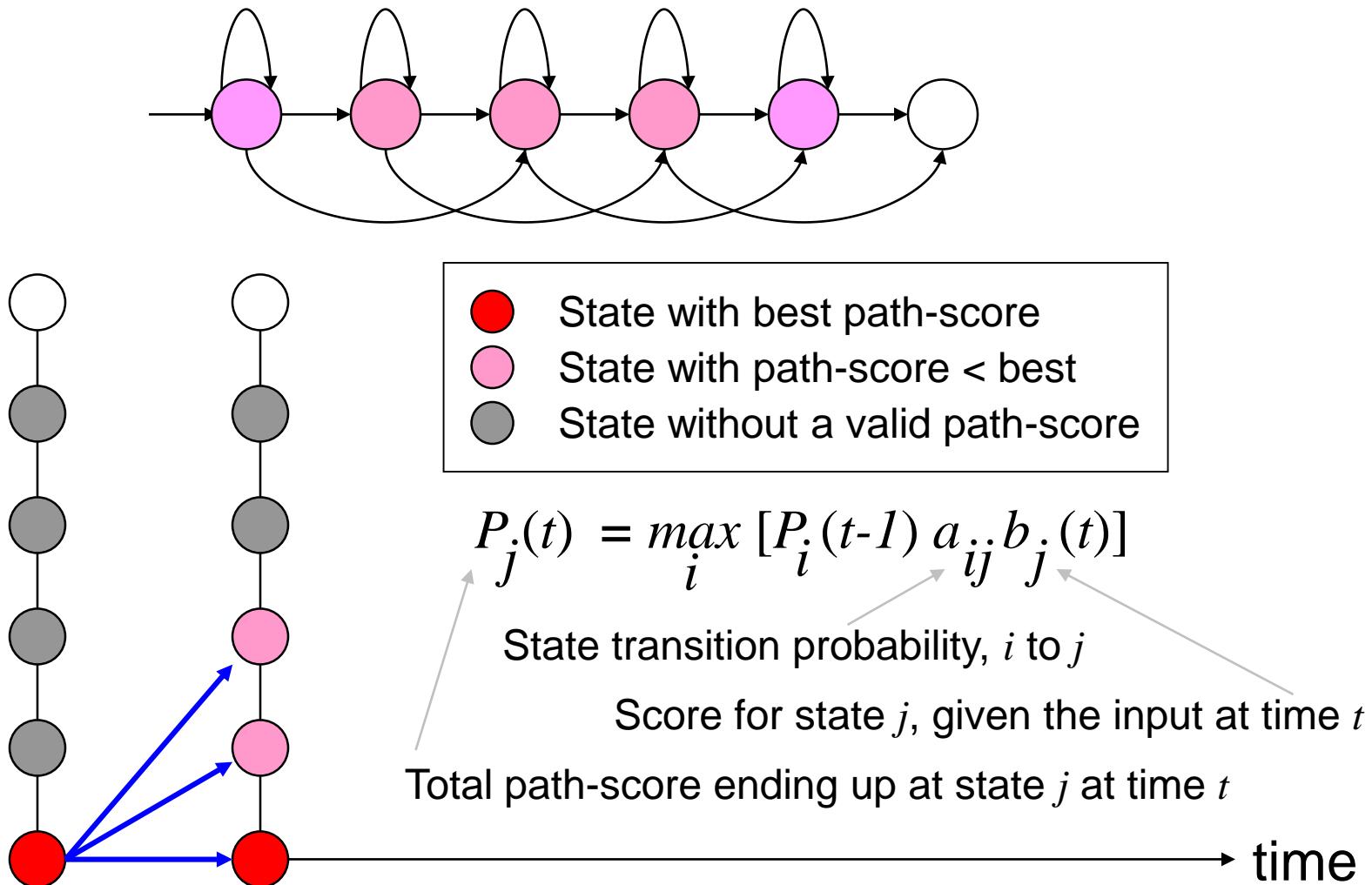
Finding the best state sequence

- The simple algorithm just presented is called the VITERBI algorithm in the literature
 - After A.J.Viterbi, who invented this dynamic programming algorithm for a completely different purpose: decoding error correction codes!
- The Viterbi algorithm can also be viewed as a breadth-first graph search algorithm

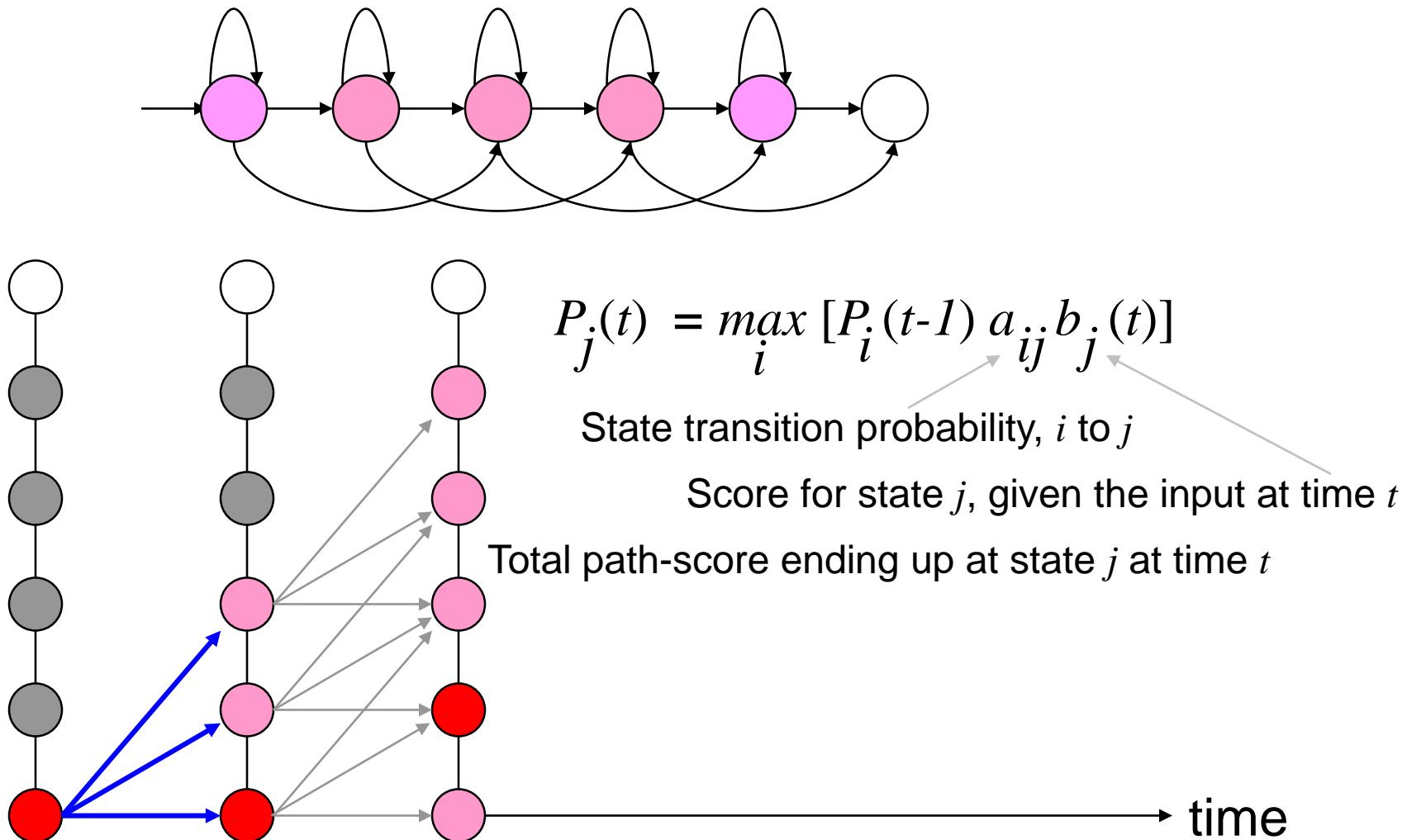
Viterbi Search (contd.)



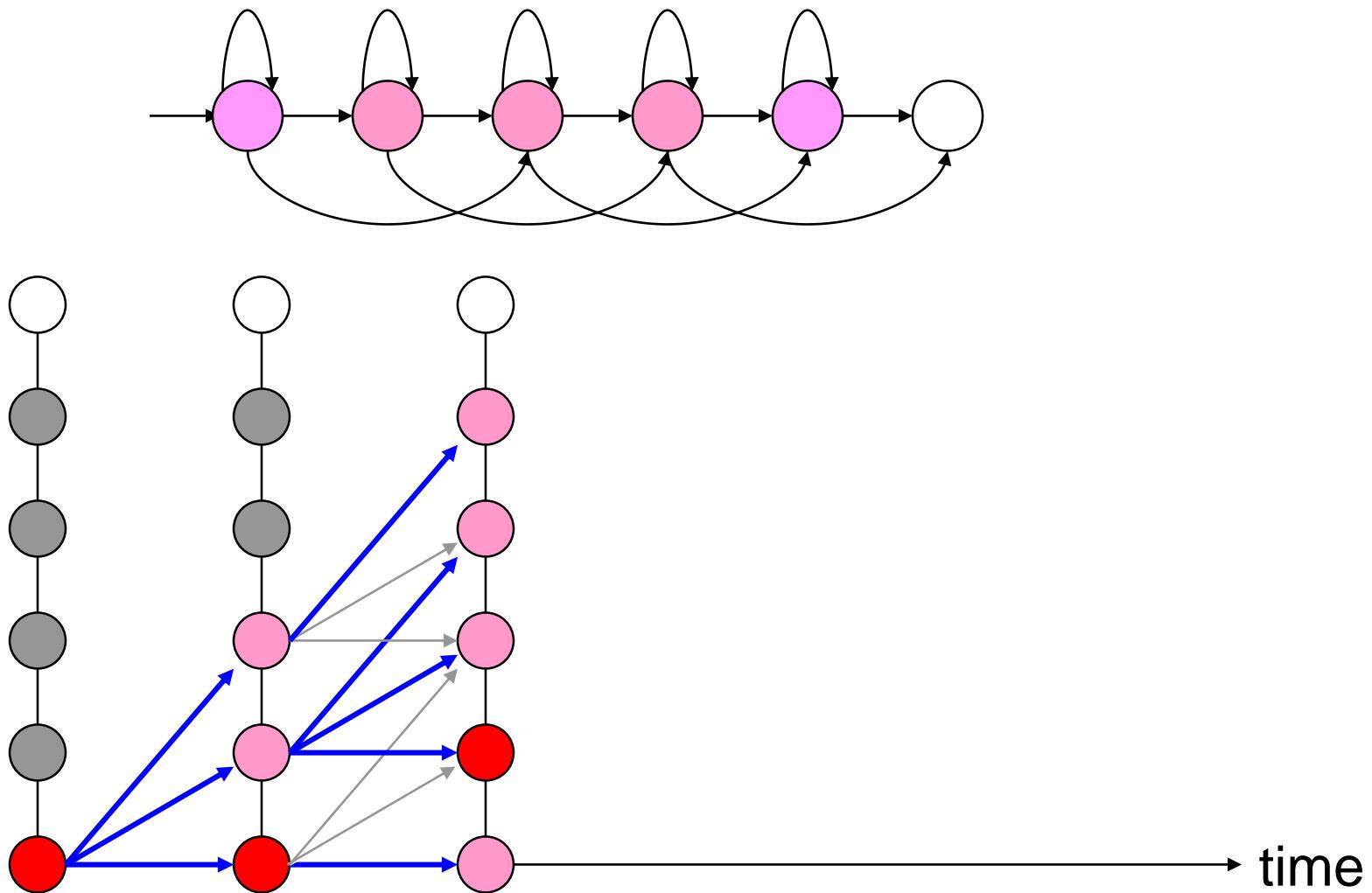
Viterbi Search (contd.)



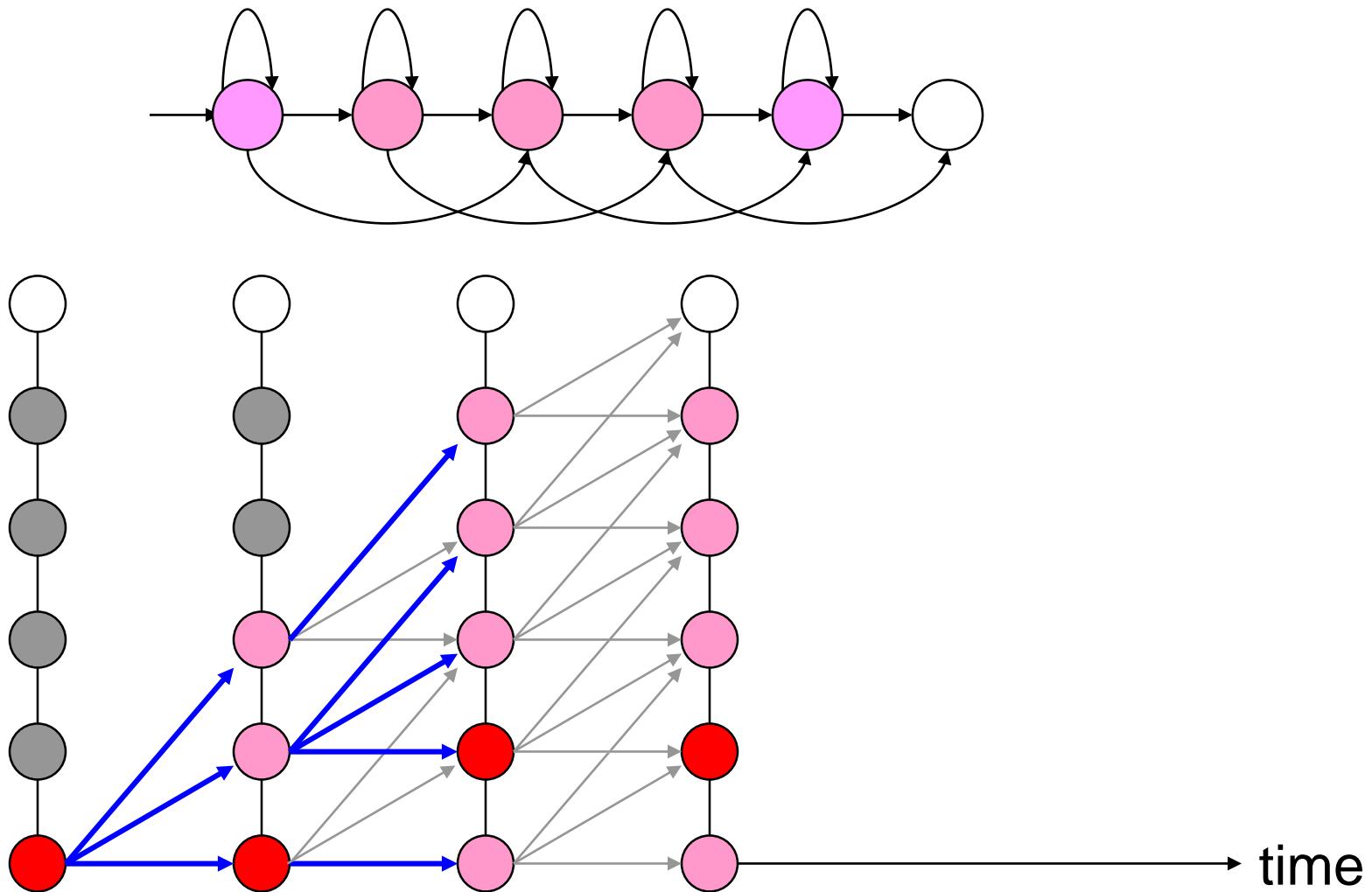
Viterbi Search (contd.)



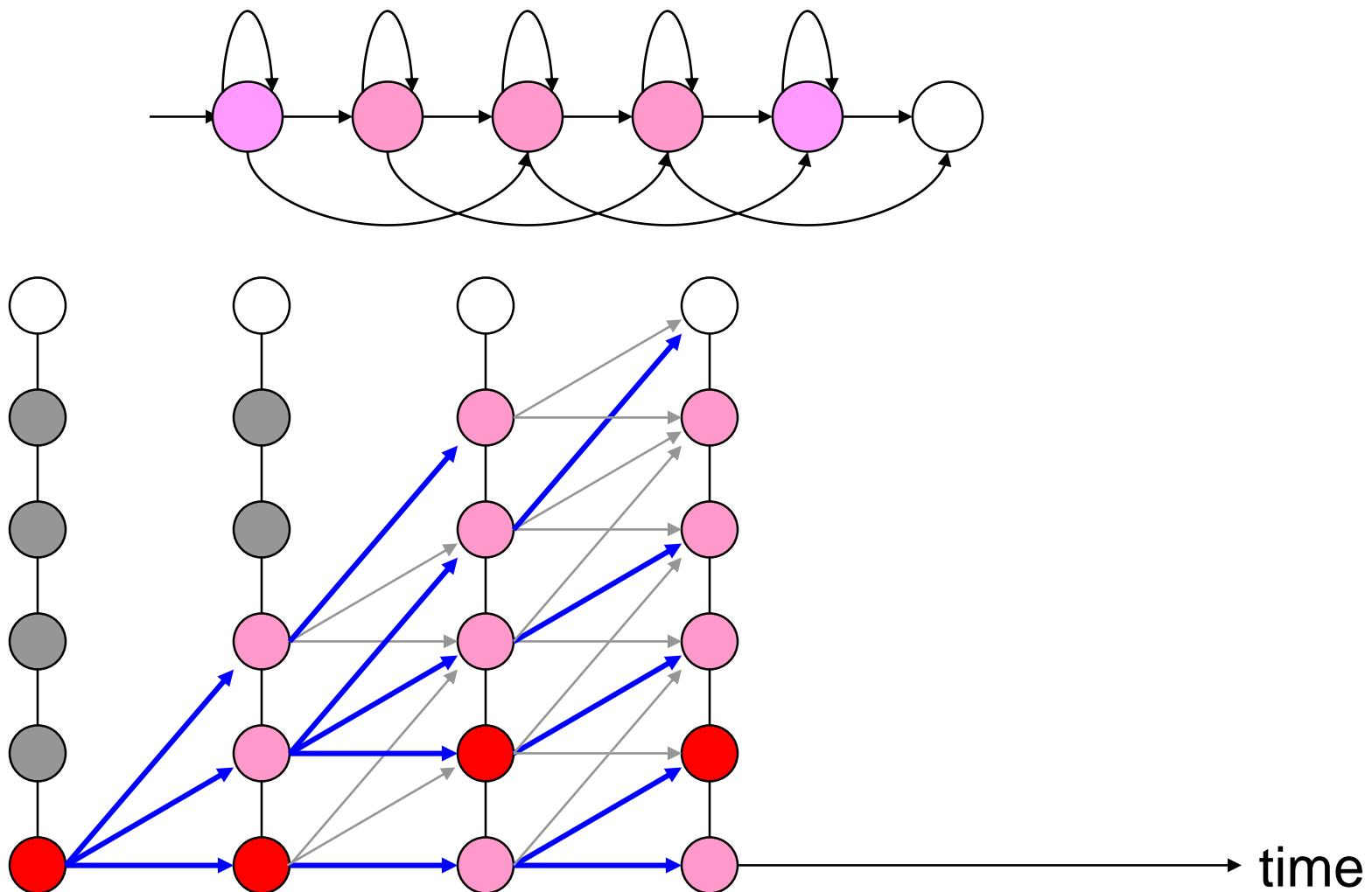
Viterbi Search (contd.)



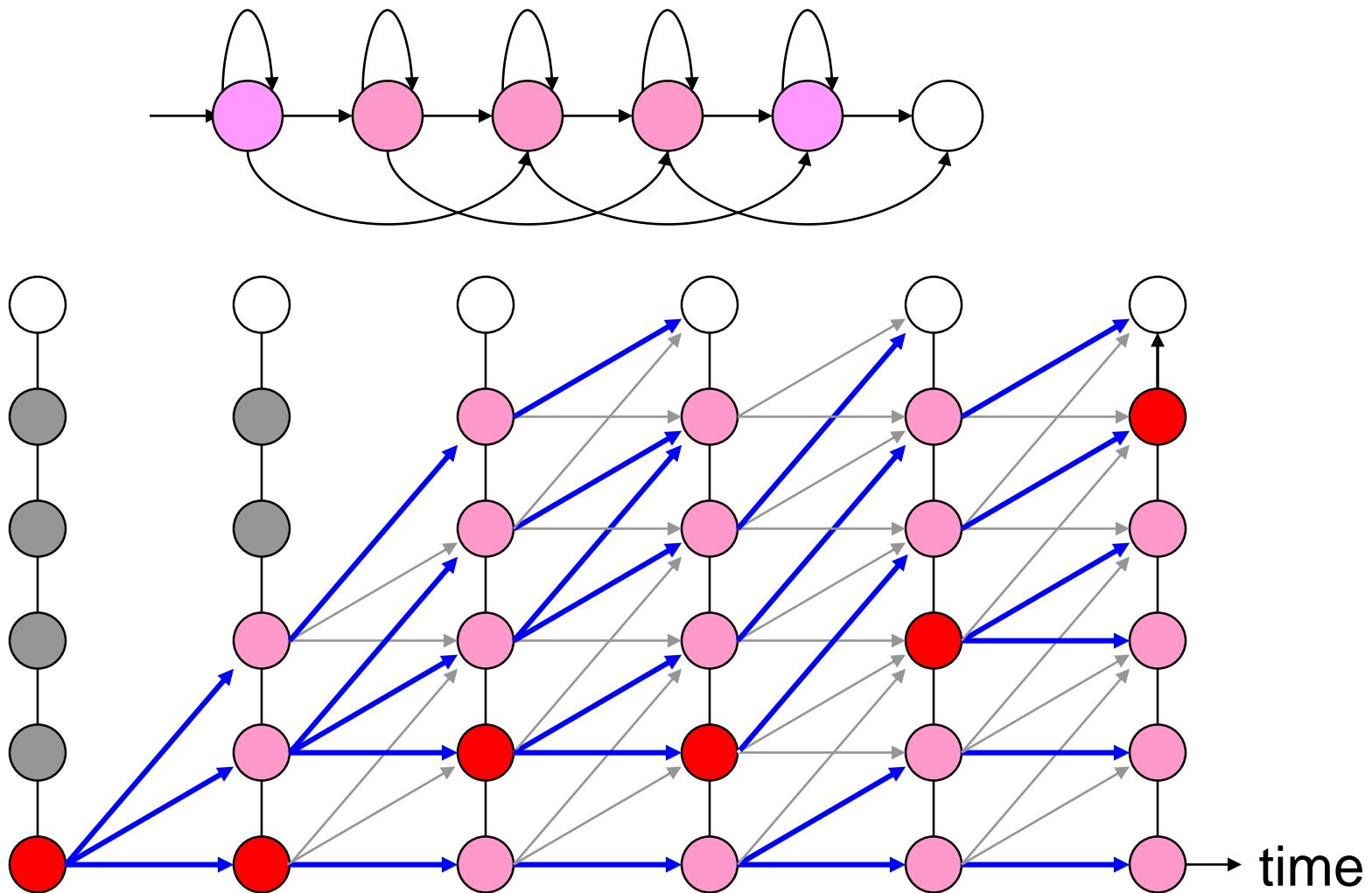
Viterbi Search (contd.)



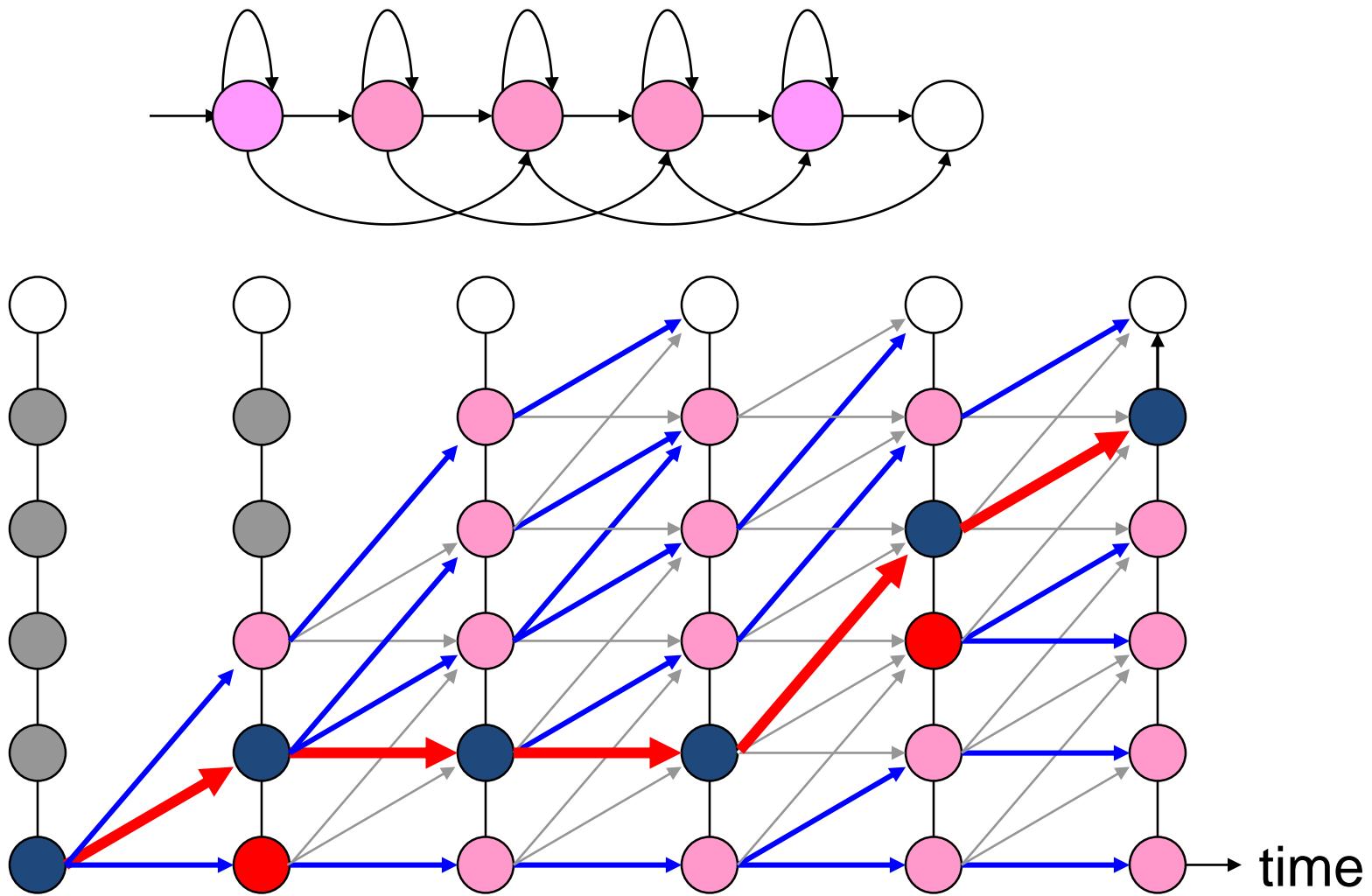
Viterbi Search (contd.)



Viterbi Search (contd.)

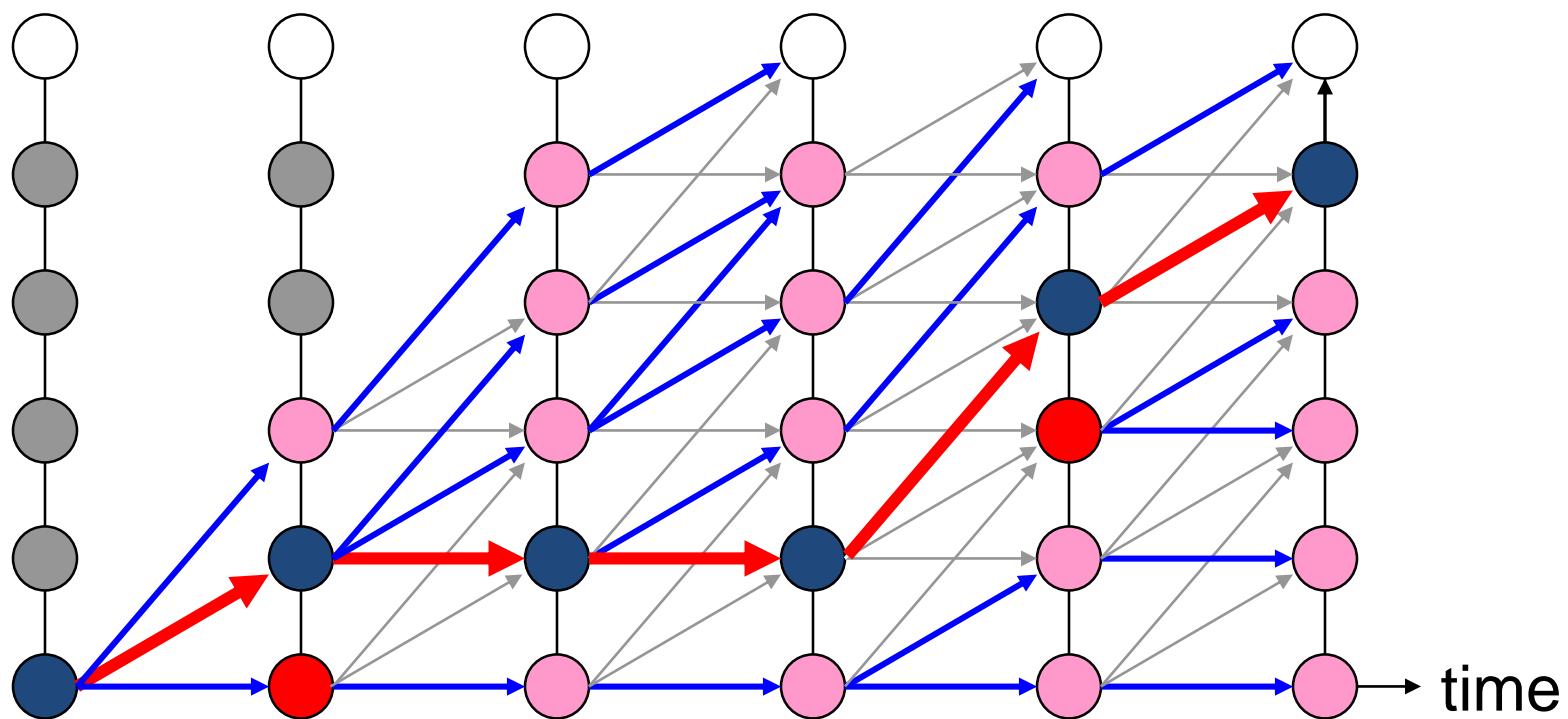


Viterbi Search (contd.)



Viterbi Search (contd.)

THE BEST STATE SEQUENCE IS THE ESTIMATE OF THE STATE SEQUENCE FOLLOWED IN GENERATING THE OBSERVATION



Viterbi and DTW

- The Viterbi algorithm is identical to the string-matching procedure used for DTW that we saw earlier
- It computes an estimate of the state sequence followed in producing the observation
- *It also gives us the probability of the best state sequence*

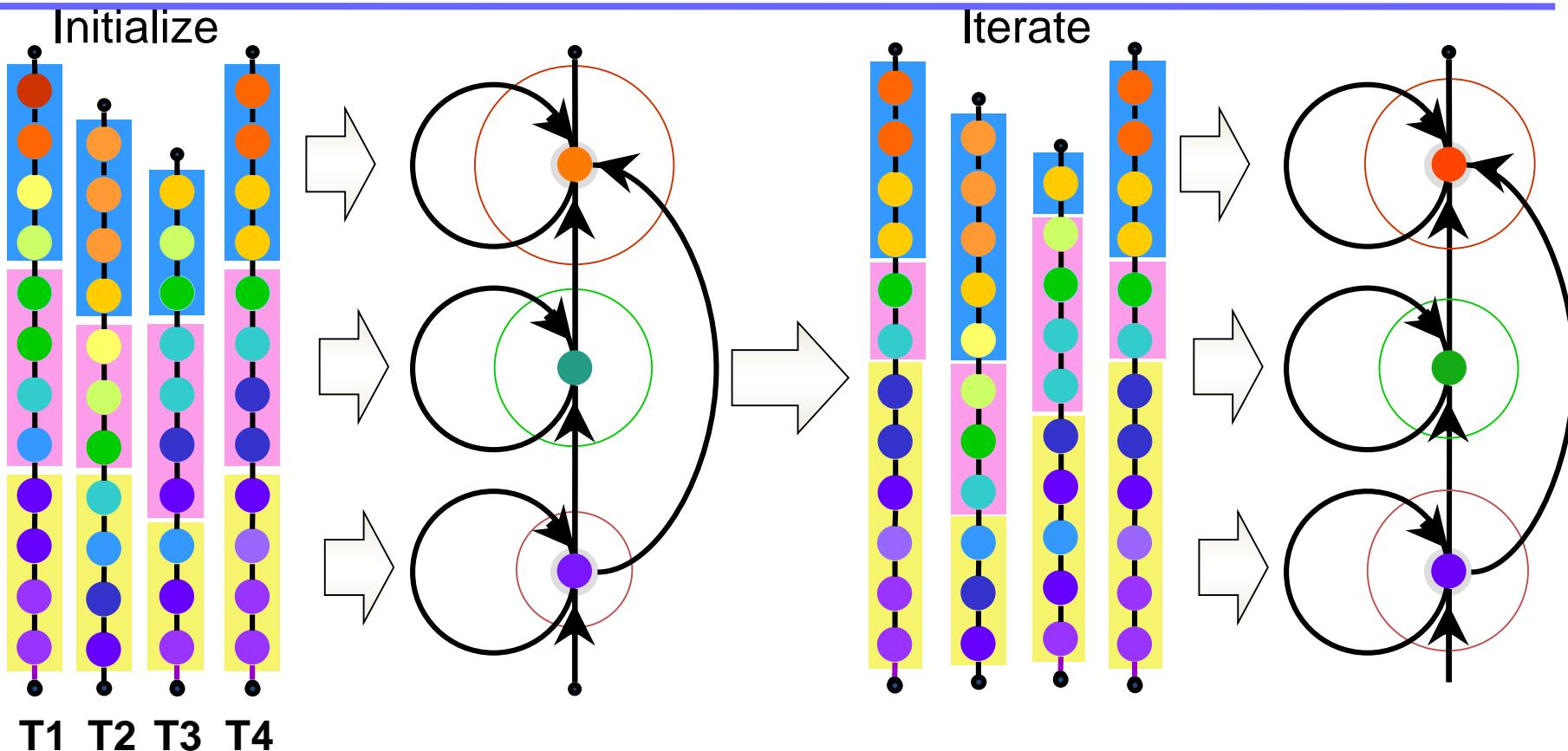
Problem3: Training HMM parameters

- We can compute the probability of an observation, and the best state sequence given an observation, using the HMM's parameters
- But where do the HMM parameters come from?
- They must be learned from a collection of observation sequences
- We have already seen one technique for training HMMs: The segmental K-means procedure

Modified segmental K-means AKA Viterbi training

- The entire segmental K-means algorithm:
 1. Initialize all parameters
 - State means and covariances
 - Transition probabilities
 - Initial state probabilities
 2. Segment all training sequences
 3. Reestimate parameters from segmented training sequences
 4. If not converged, return to 2

Segmental K-means

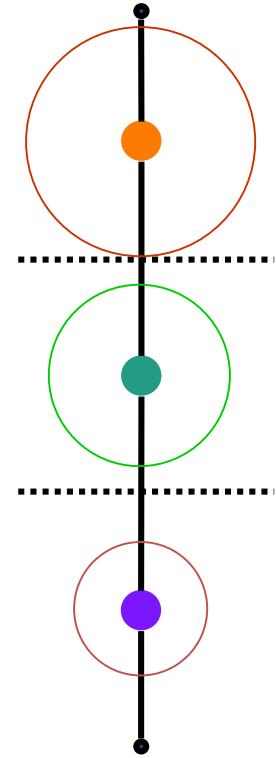
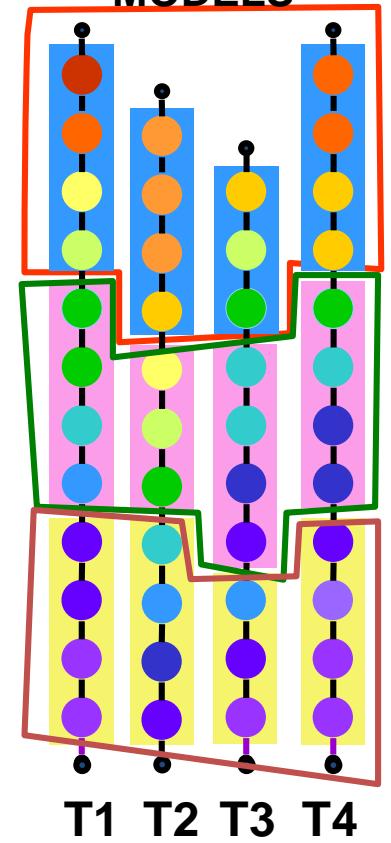


The procedure can be continued until convergence

Convergence is achieved when the total best-alignment error for all training sequences does not change significantly with further refinement of the model

Training by segmentation: Hard Assignment

MODELS



$$m_j = \frac{1}{\sum_{i \in \text{segment}(j)} 1} \sum_{i \in \text{segment}(j)} x_i$$

$$C_j = \frac{1}{\sum_{i \in \text{segment}(j)} 1} \sum_{i \in \text{segment}(j)} (x_i - m_j)(x_i - m_j)^T$$

Each vector belongs uniquely to a segment

$$P(x | j) = \frac{1}{\sqrt{2\pi |C|}} \exp(-0.5(x - m_j)^T C^{-1} (x - m_j))$$

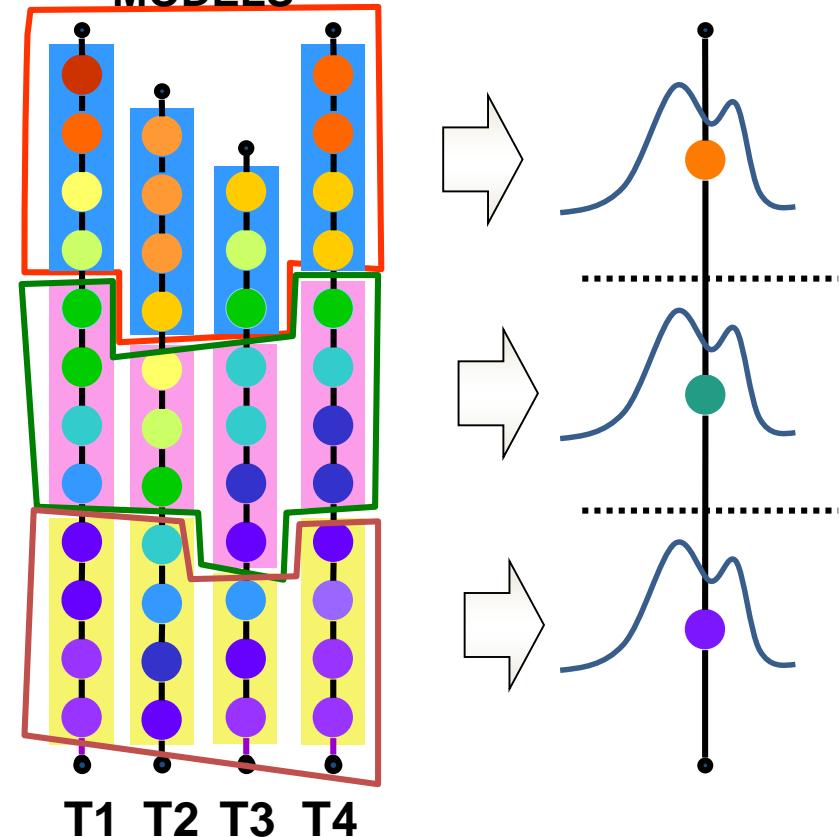
Assumes state output distribution is Gaussian

$$d_j(v) = -\log(P(x | j))$$

$$P(x | j) = \frac{1}{\sqrt{2\pi |C|}} \exp\left(-\sum_d \frac{(x_d - m_{j,d})^2}{2\sigma_d^2}\right)$$

Training by segmentation: Hard Assignment with Gaussian Mixtures

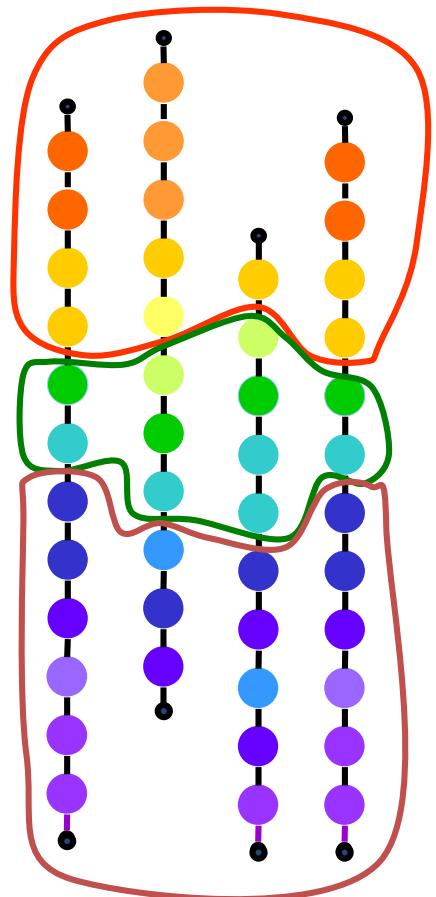
MODELS



Dealing with Gaussian Mixtures
at states?

The distribution of vectors in
any state is assumed to be a
Gaussian mixture

Training by segmentation: Hard Assignment with Gaussian Mixtures



Assume the distribution of each collection of vectors is a Gaussian mixture

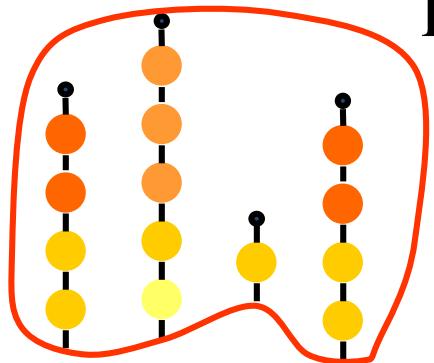
$$P(x | j) = \sum_k \frac{w_k}{\sqrt{\prod_l 2\pi\sigma_{j,k,l}^2}} \exp\left(-0.5 \sum_l \frac{(x_l - m_{j,k,l})^2}{\sigma_{j,k,l}^2}\right)$$

$$d_j(v) = -\log(P(x | j))$$

Above equation assumes Gaussian covariance matrices are diagonal

Training a GMM by hard counting

How does one learn the parameters of the GMM?



$$P(x | j) = \sum_k \frac{w_k}{\sqrt{\prod_l 2\pi\sigma_{j,k,l}^2}} \exp\left(-0.5 \sum_l \frac{(x_l - m_{j,k,l})^2}{\sigma_{j,k,l}^2}\right)$$

$$d_j(v) = -\log(P(x | j))$$

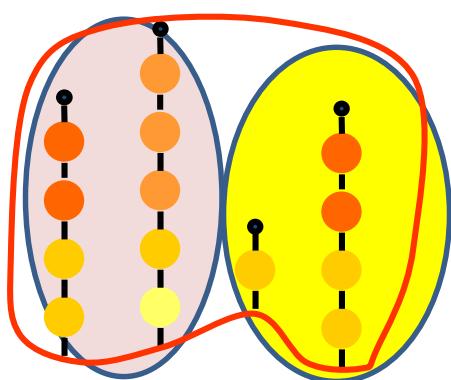
Gaussian Mixtures

- A Gaussian Mixture is literally a mixture of Gaussians. It is a weighted combination of several Gaussian distributions

$$P(x) = \sum_{i=0}^{K-1} w_i \text{Gaussian}(x; m_i, C_i)$$

- x is any data vector. $P(x)$ is the probability given to that vector by the Gaussian mixture
- K is the number of Gaussians being mixed
- w_i is the mixture weight of the i^{th} Gaussian. m_i is its mean and C_i is its covariance

Gaussian Mixtures: A “hard” perspective



$$P(x) = \sum_{i=0}^{K-1} w_i \text{Gaussian}(x; m_i, C_i)$$

Fraction of vectors
in i^{th} group

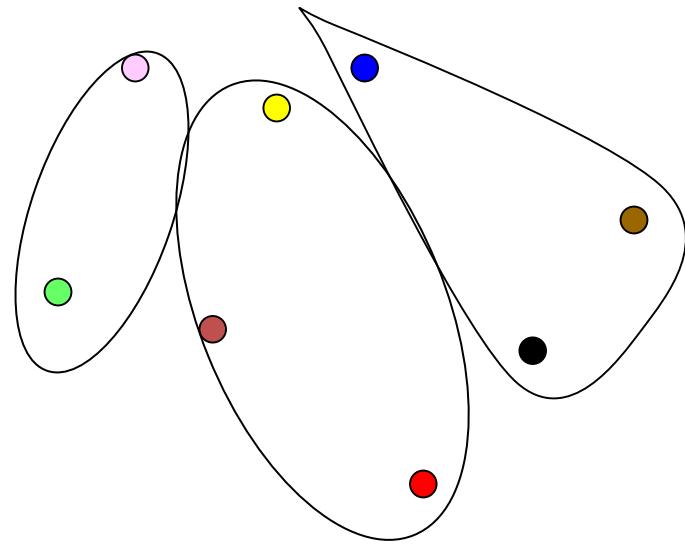
Mean of
 i^{th} group

Covariance
of i^{th} group

- Data from each Gaussian can be clearly grouped into clusters
- The parameters of the GMM are the parameters of individual clusters
- We can use a clustering algorithm to find the clusters
 - K-means

The K-means algorithm

- The K-means algorithm is an iterative algorithm for clustering similar data from a data set
 - Similarity defined in terms of distance between clusters and data
 - E.g. distance from cluster mean
 - Negative log probability of the vector given by the distribution of the cluster
- The algorithm tries to find the most consistent clusters of data
 - Consistency in terms of specified distance measure



K-Means training Gaussian Mixtures

- The K-means algorithm can be used to estimate Gaussian mixture distributions for a data set
- Each of the K Gaussians is assumed to represent a separate cluster of the data
- The j^{th} cluster is characterized by
 - Its covariance C_j
 - Its mean vector m_j
 - A mixture weight w_j that specifies what portion of the total data belongs to that cluster
- Define the distance between a vector and the j^{th} cluster as

$$d(v, j) = 0.5 \log((2\pi)^D |C_j|) + 0.5(v - m_j)^T C_j^{-1} (v - m_j) - \log(w_j)$$

$-\log P(v, j) = -\log P(j) - \log P(v|j)$, $P()$ is a Gaussian

D = dimension of vectors

K-Means: Estimating parameters for a cluster

- The parameters for a cluster are its mixture weight, mean vector and covariance matrix. These are computed as follows:

$$m_j = \frac{1}{N_j} \sum_{v:j(v)=j} v$$

$j(v)$ is the cluster that vector v is assigned to

- N_j is the number of vectors that have been tagged as belonging to cluster j
- The summation is over all vectors who have been tagged as belonging to j

$$C_j = \frac{1}{N_j} \sum_{v:j(v)=j} (v - m_j)(v - m_j)^T$$

$$w_j = \frac{N_j}{N}$$

- N is the total number of training vectors for all clusters

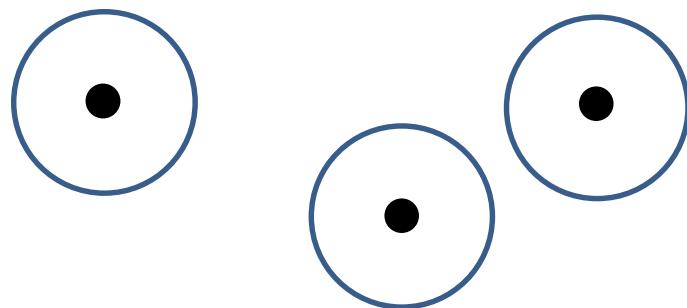
The K-means algorithm

- Initialize all clusters somehow (the number of clusters is assumed)
- For each training vector, find the closest cluster
- Reassign training vectors to their closest clusters
- Iterate the above two steps until the total distance of all training vectors from their clusters converges
 - Convergence can be proved for most distance measures

K-means

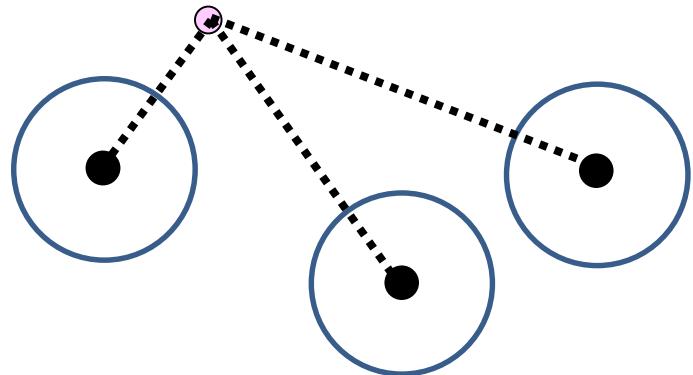
1. Initialize cluster parameters

How? We'll return to this shortly



K-means

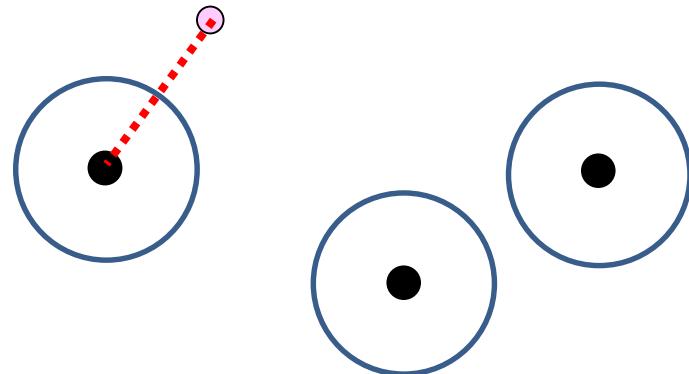
1. Initialize cluster parameters
2. For each data point x , find the distance from each cluster
 - $d_{cluster} = d(x, cluster)$



$$d(v, j) = 0.5 \log((2\pi)^D |C_j|) + 0.5(v - m_j)^T C_j^{-1} (v - m_j) - \log(w_i)$$

K-means

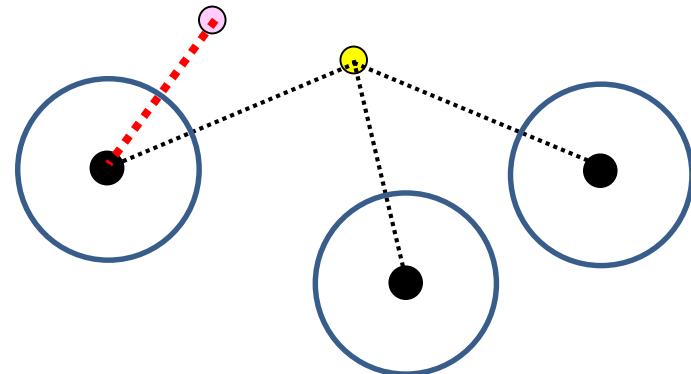
1. Initialize cluster parameters
2. For each data point x , find the distance from each cluster
 - $d_{cluster} = d(x, cluster)$
3. Put data point in the cluster of the closest centroid
 - Cluster for which $d_{cluster}$ is minimum



$$d(v, j) = 0.5 \log((2\pi)^D |C_j|) + 0.5(v - m_j)^T C_j^{-1} (v - m_j) - \log(w_i)$$

K-means

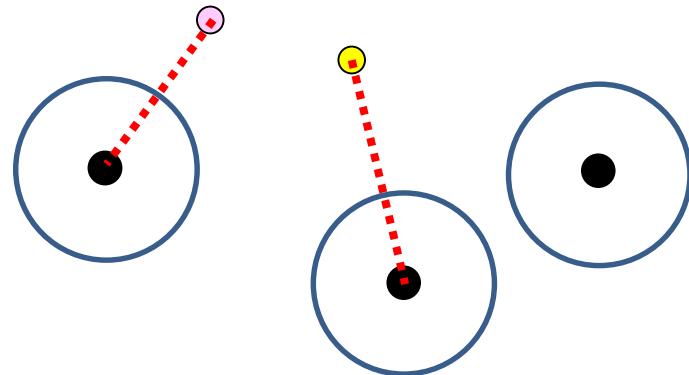
1. Initialize cluster parameters
2. For each data point x , find the distance from each cluster
 - $d_{cluster} = d(x, cluster)$
3. Put data point in the cluster of the closest centroid
 - Cluster for which $d_{cluster}$ is minimum



$$d(v, j) = 0.5 \log((2\pi)^D |C_j|) + 0.5(v - m_j)^T C_j^{-1} (v - m_j) - \log(w_i)$$

K-means

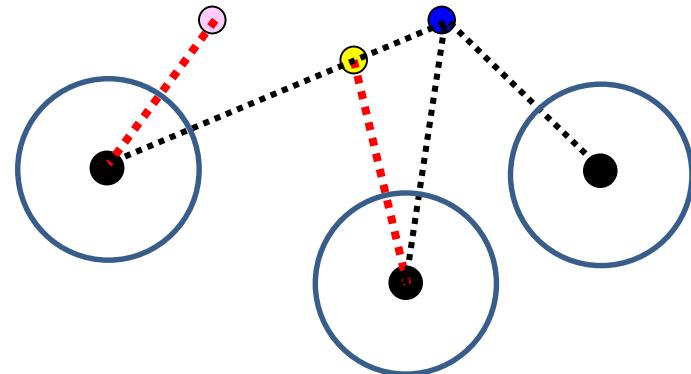
1. Initialize cluster parameters
2. For each data point x , find the distance from each cluster
 - $d_{cluster} = d(x, cluster)$
3. Put data point in the cluster of the closest centroid
 - Cluster for which $d_{cluster}$ is minimum



$$d(v, j) = 0.5 \log((2\pi)^D |C_j|) + 0.5(v - m_j)^T C_j^{-1} (v - m_j) - \log(w_i)$$

K-means

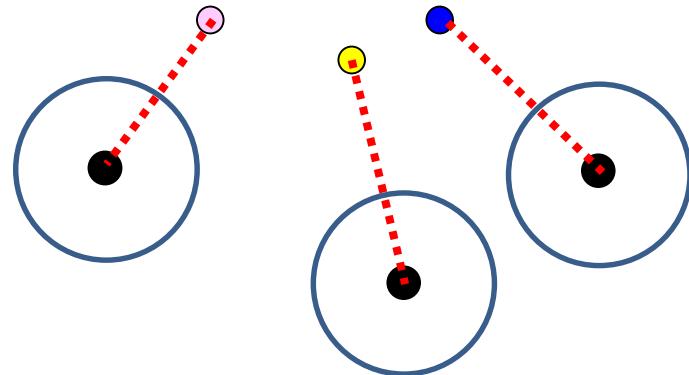
1. Initialize cluster parameters
2. For each data point x , find the distance from each cluster
 - $d_{cluster} = d(x, cluster)$
3. Put data point in the cluster of the closest centroid
 - Cluster for which $d_{cluster}$ is minimum



$$d(v, j) = 0.5 \log((2\pi)^D |C_j|) + 0.5(v - m_j)^T C_j^{-1} (v - m_j) - \log(w_i)$$

K-means

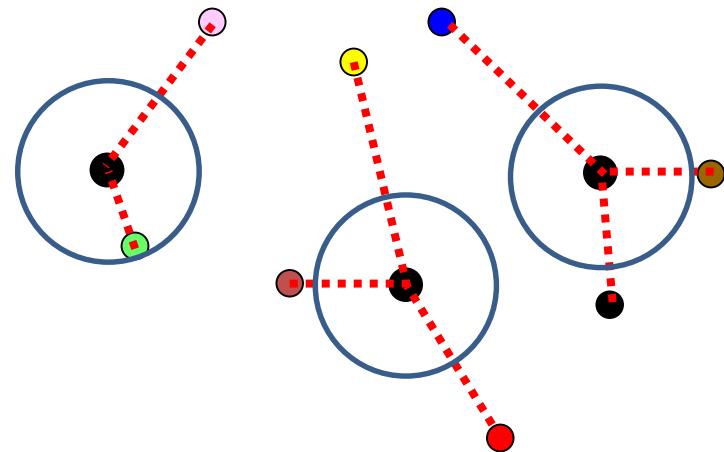
1. Initialize cluster parameters
2. For each data point x , find the distance from each cluster
 - $d_{cluster} = d(x, cluster)$
3. Put data point in the cluster of the closest centroid
 - Cluster for which $d_{cluster}$ is minimum



$$d(v, j) = 0.5 \log((2\pi)^D |C_j|) + 0.5(v - m_j)^T C_j^{-1} (v - m_j) - \log(w_i)$$

K-means

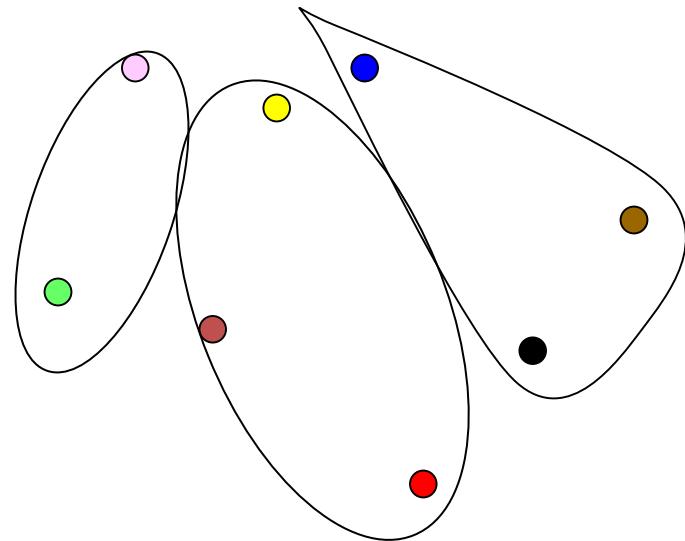
1. Initialize cluster parameters
2. For each data point x , find the distance from each cluster
 - $d_{cluster} = d(x, cluster)$
3. Put data point in the cluster of the closest centroid
 - Cluster for which $d_{cluster}$ is minimum



$$d(v, j) = 0.5 \log((2\pi)^D |C_j|) + 0.5(v - m_j)^T C_j^{-1} (v - m_j) - \log(w_i)$$

K-means

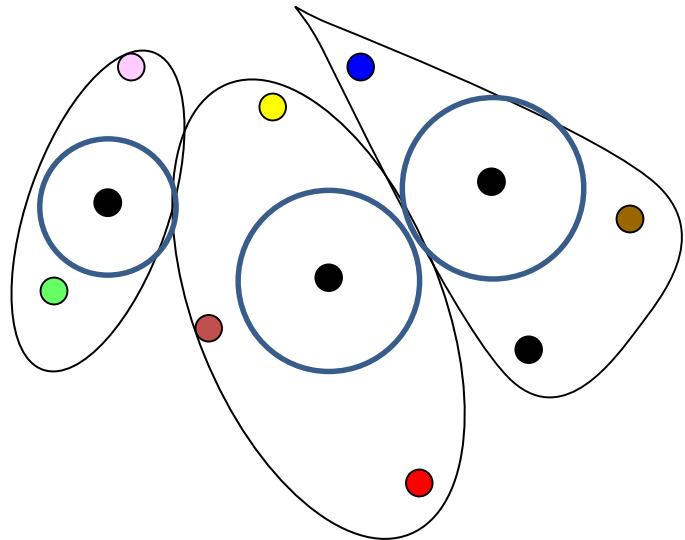
1. Initialize cluster parameters
2. For each data point x , find the distance from each cluster
 - $d_{cluster} = d(x, cluster)$
3. Put data point in the cluster of the closest centroid
 - Cluster for which $d_{cluster}$ is minimum



$$d(v, j) = 0.5 \log((2\pi)^D |C_j|) + 0.5(v - m_j)^T C_j^{-1} (v - m_j) - \log(w_i)$$

K-means

1. Initialize cluster parameters
2. For each data point x , find the distance from each cluster
 - $d_{cluster} = d(x, cluster)$
3. Put data point in the cluster of the closest centroid
 - Cluster for which $d_{cluster}$ is minimum
4. When all data points clustered, recompute cluster parameters
 - Means, variances, weights



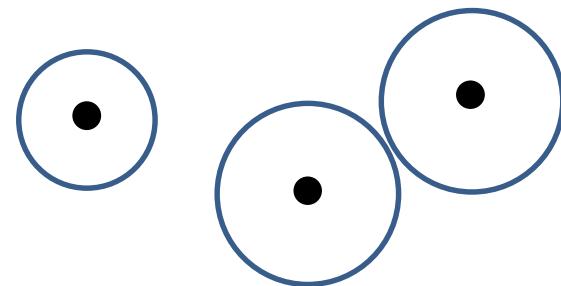
$$m_j = \frac{1}{N_j} \sum_{v:j(v)=j} v$$

$$C_j = \frac{1}{N_j} \sum_{v:j(v)=j} (v - m_j)(v - m_j)^T$$

$$w_j = \frac{N_j}{N}$$

K-means

1. Initialize cluster parameters
2. For each data point x , find the distance from each cluster
 - $d_{cluster} = d(x, cluster)$
3. Put data point in the cluster of the closest centroid
 - Cluster for which $d_{cluster}$ is minimum
4. When all data points clustered, recompute cluster parameters
 - Means, variances, weights
5. If not converged, go back to 2



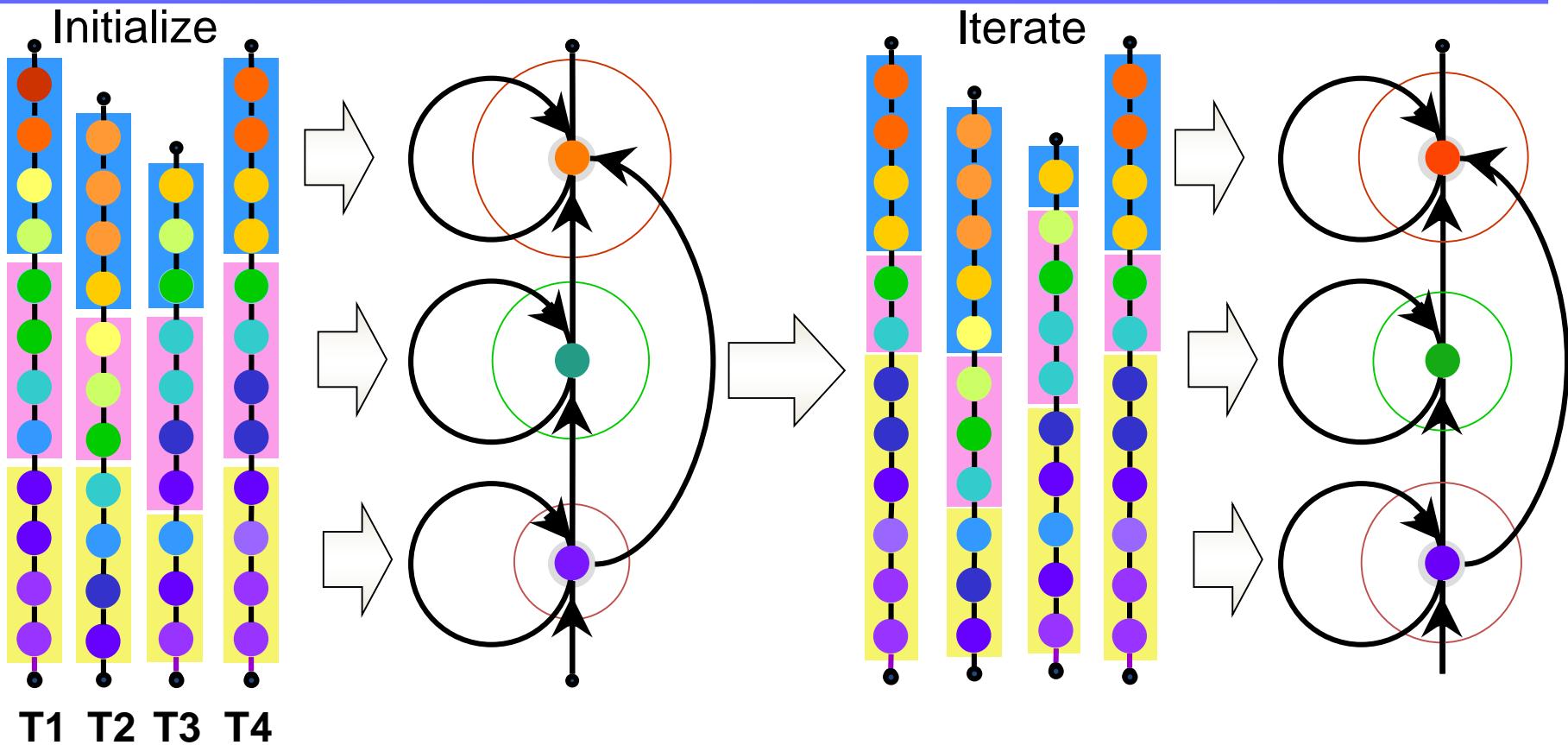
K-Means for Gaussian Mixtures

- Initialization: There are several ways of initializing the K-means algorithm. One common method is:
 - Set $w_i = 1/K$ for all clusters
 - Set $m_i = \text{random}(\{v\})$: i.e. a randomly selected vector from the training data
 - Set $C_j =$ the global covariance of the entire training data
- There are other ways of initializing the K-means algorithm
 - Some of these may be better than the initialization procedure described above, but are more complicated
 - Revisit later

The set of cluster means, variances and mixture weights constitute the parameters of the Gaussian mixture distribution for the data

$$P(v) = \sum_j w_j \text{Gaussian}(v; m_j, C_j)$$

Overall Segmental Kmeans



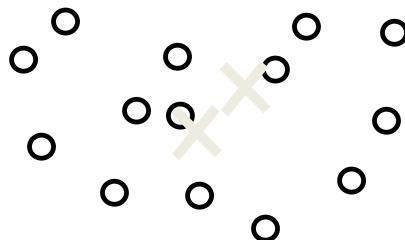
Identical to what we learned before, with one change:
At each iteration, learn a Gaussian *mixture* distribution for each state.

More on Gaussian Mixtures

- More common procedure:
 1. Train HMMs with 1-Gaussian per state HMMs using a first-pass of segmental K-means
 - This is a trivial Gaussian mixture with only one Gaussian
 2. *Split* the Gaussians in the state output distributions to obtain a larger Gaussian mixture at each state
 3. Run segmental K-means to convergence with updated Gaussian Mixtures
 4. Use EM (expectation maximization) algorithm to fine training
 5. If desired number of Gaussians not obtained for each state, return to 2
- What is “splitting”? What is the ideal no. of Gaussians?
 - We get to this shortly

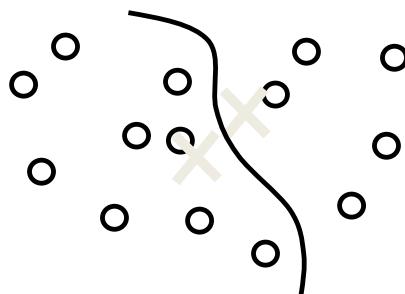
- Binary Split Algorithm

Given a set of points in a codebook:

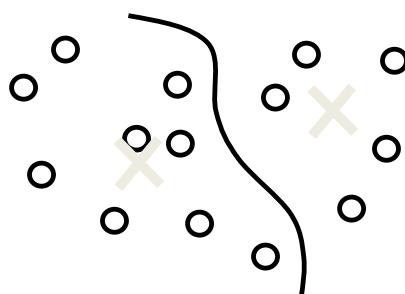


$$\mathbf{y}_n^+ = \mathbf{y}_n(1 + \varepsilon)$$

$$\mathbf{y}_n^- = \mathbf{y}_n(1 - \varepsilon)$$



use K -means to
cluster split codebook



and compute new
centroids

EM algorithm for GMM

Given an initial guess Θ^g , find Θ^{new} as follows

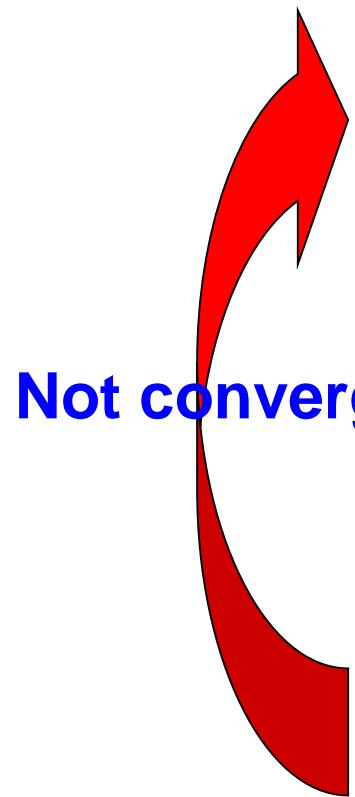
$$p(l | \mathbf{x}_i, \Theta^g) = \frac{\alpha_l^g p_l(\mathbf{x}_i | \theta_l^g)}{\sum_{j=1}^M \alpha_j^g p_j(\mathbf{x}_i | \theta_j^g)}$$

$$\alpha_l^{new} = \frac{1}{N} \sum_{i=1}^N p(l | \mathbf{x}_i, \Theta^g)$$

$$\boldsymbol{\mu}_l^{new} = \frac{\sum_{i=1}^N \mathbf{x}_i p(l | \mathbf{x}_i, \Theta^g)}{\sum_{i=1}^N p(l | \mathbf{x}_i, \Theta^g)}$$

$$\Sigma_l^{new} = \frac{\sum_{i=1}^N p(l | \mathbf{x}_i, \Theta^g) (\mathbf{x}_i - \boldsymbol{\mu}_l^{new})(\mathbf{x}_i - \boldsymbol{\mu}_l^{new})^T}{\sum_{i=1}^N p(l | \mathbf{x}_i, \Theta^g)}$$

$$\Theta^g \leftarrow \Theta^{new}$$

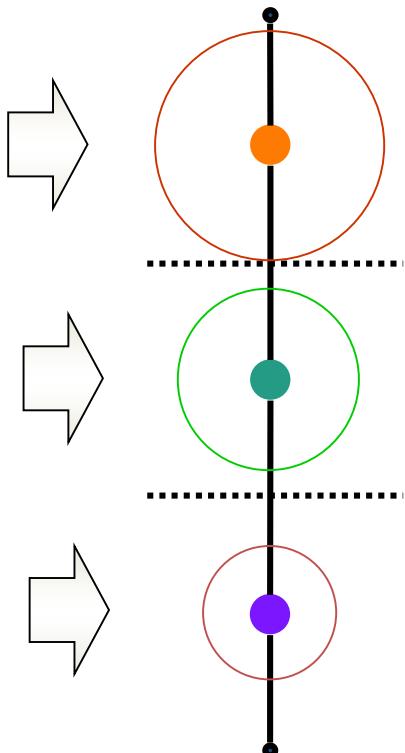
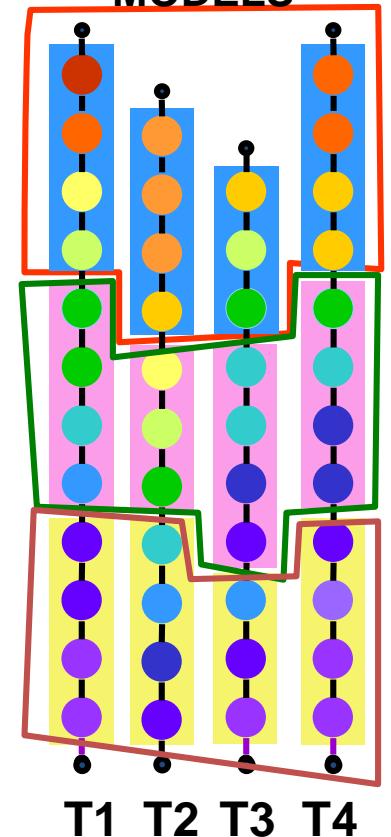


A Better Technique

- The Segmental K-means technique uniquely assigns each observation to one state
- However, this is only an estimate and may be wrong
- A better approach is to take a “soft” decision
 - Assign each observation to *every* state with a probability

Training by segmentation: Soft Assignment (1 Gaussian case)

MODELS



$$\mu_j = \frac{1}{\sum_{i \in \text{All vectors}} f_{i,j}} \sum_{i \in \text{All vectors}} f_{i,j} x_i$$

$$\sum_{j \in \text{all segments}} f_{i,j} = 1$$

Assignment is *fractioned*:
Every segment gets a piece of
every vector

Means and variances are computed
from fractioned vectors

Where do the fractions come from?

The “probability” of a state

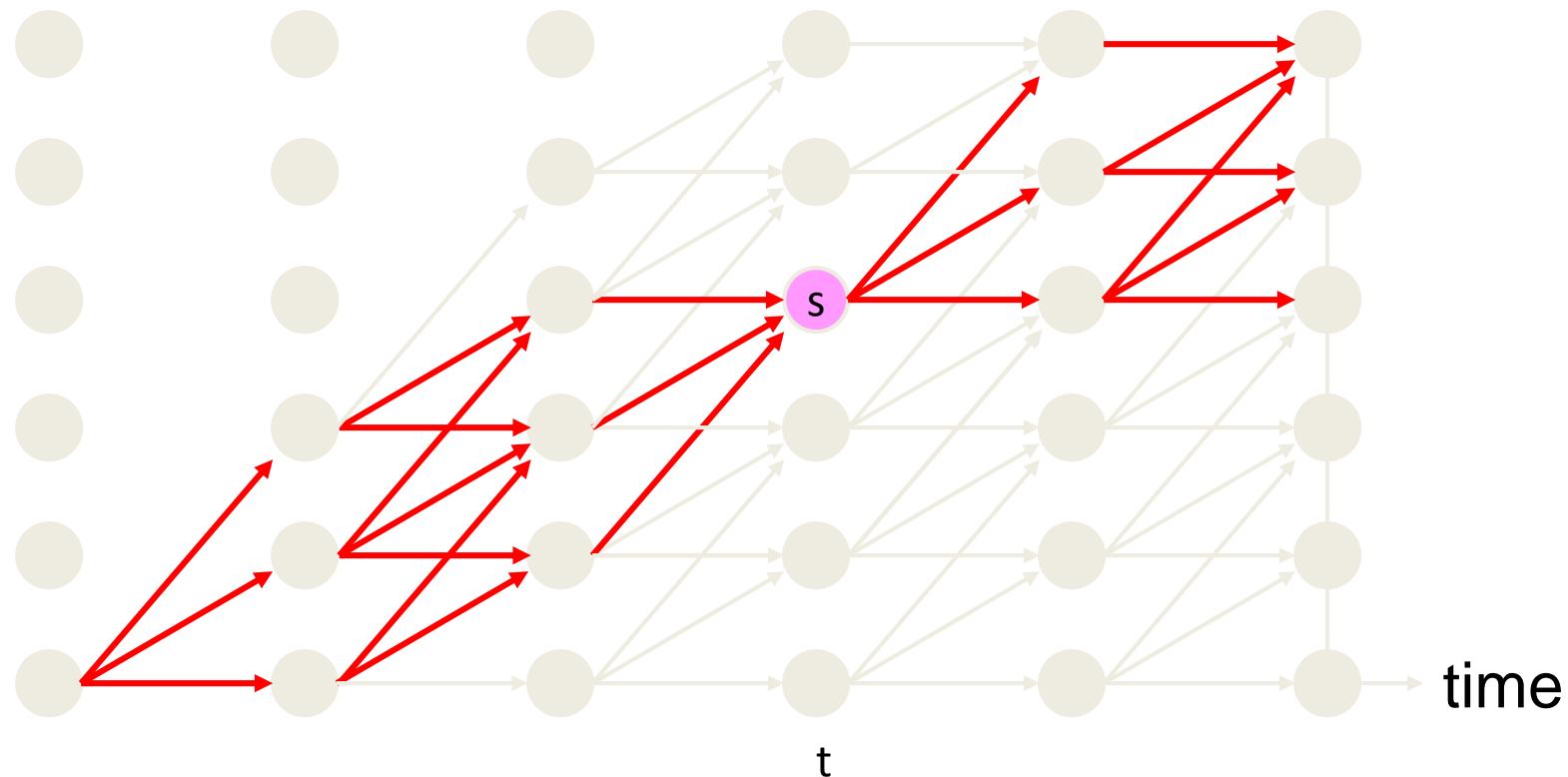
- The probability assigned to any state s , for any observation x_t is the probability that the process was at s when it generated x_t
- We want to compute

$$P(state(t) = s \mid x_1, x_2, \dots, x_T) \propto P(state(t) = s, x_1, x_2, \dots, x_T)$$

- We will compute $P(state(t) = s, x_1, x_2, \dots, x_T)$ first
 - This is the probability that the process visited s at time t while producing the entire observation

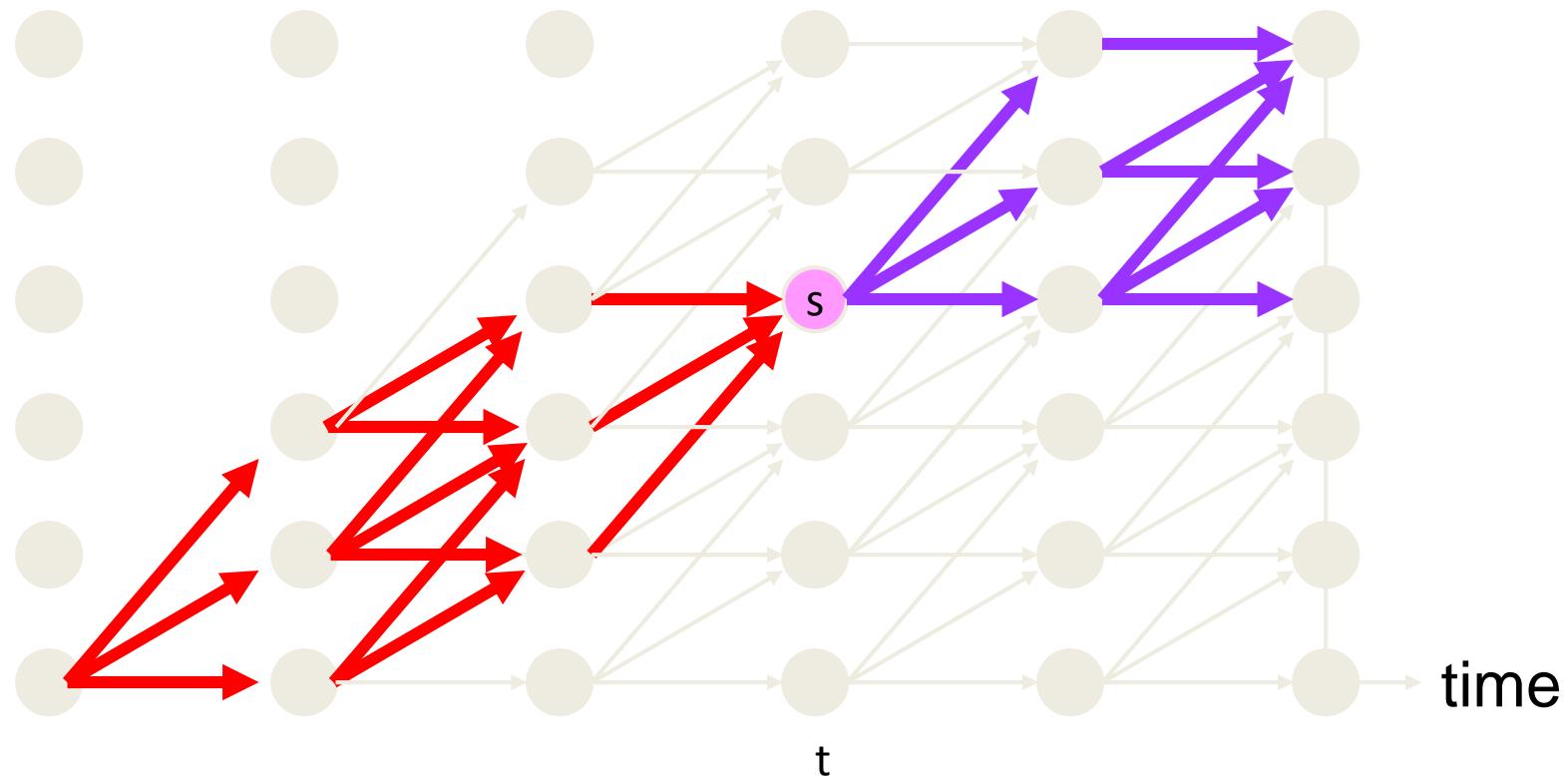
Probability of Assigning an Observation to a State

- The probability that the HMM was in a particular state s when generating the observation sequence is the probability that it followed a state sequence that passed through s at time t



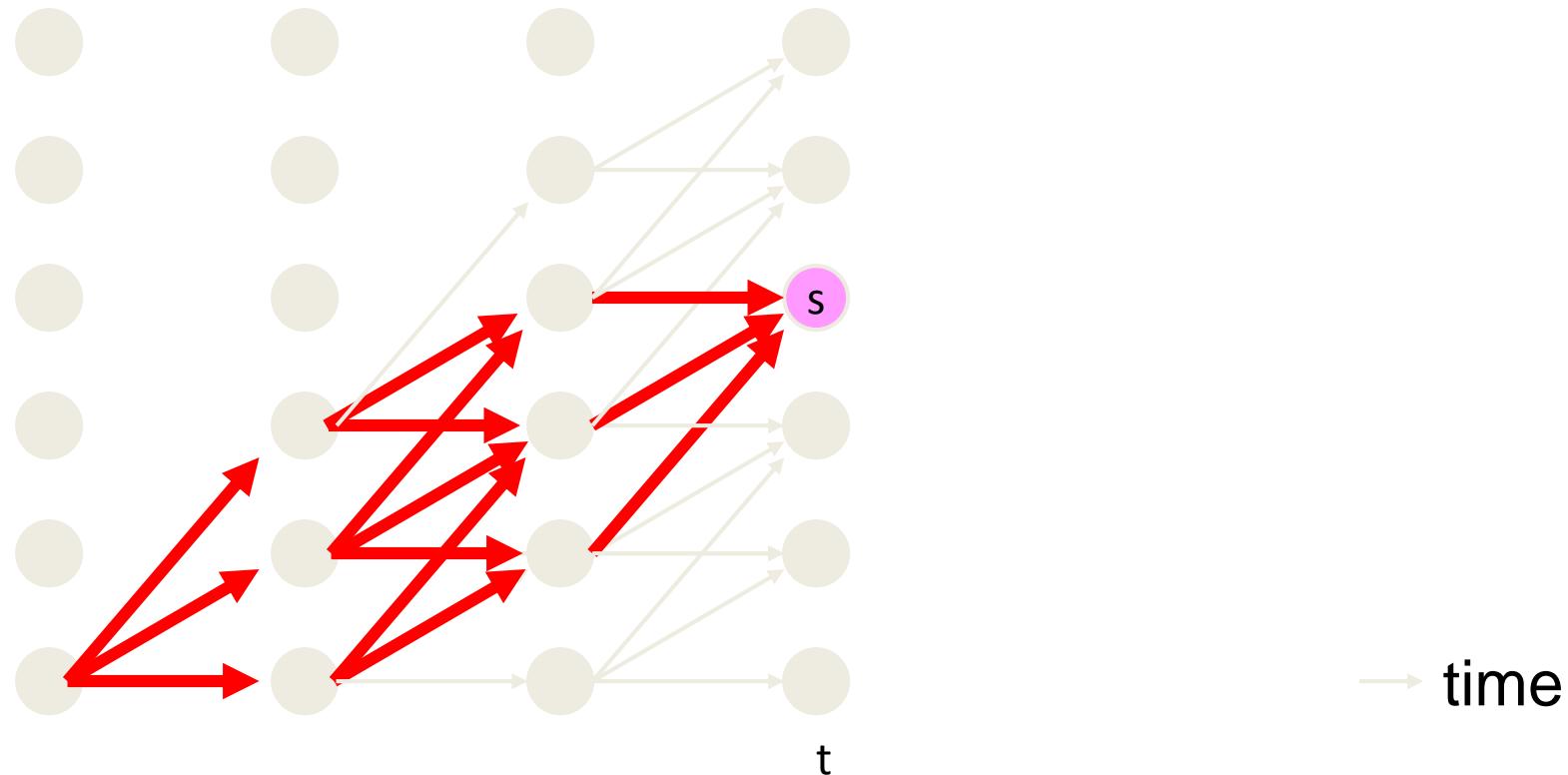
Probability of Assigning an Observation to a State

- This can be decomposed into two multiplicative sections
 - The section of the lattice leading into state s at time t and the section leading out of it



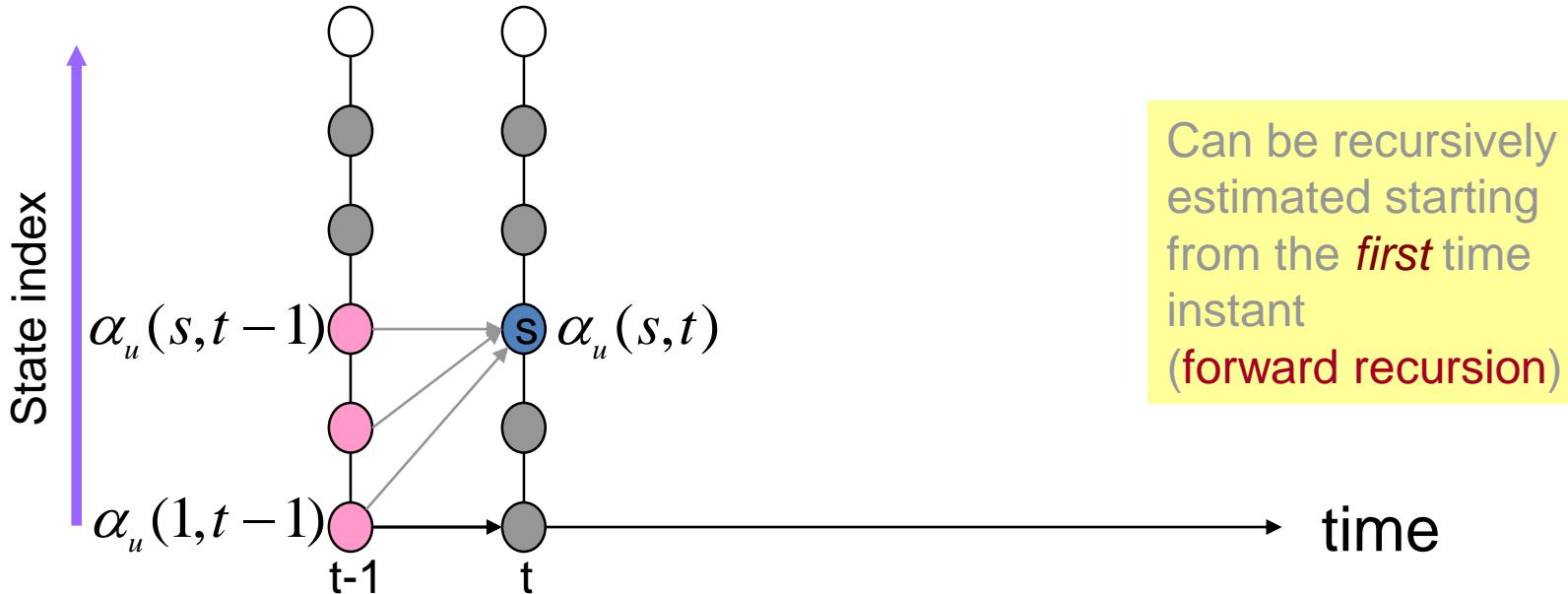
Probability of Assigning an Observation to a State

- The probability of the red section is the total probability of all state sequences ending at state s at time t
 - This is simply $\alpha(s,t)$
 - Can be computed using the forward algorithm



The forward algorithm

$$\alpha_u(s, t) = P(x_{u,1}, x_{u,2}, \dots, x_{u,t}, \text{state}(t) = s | \lambda)$$

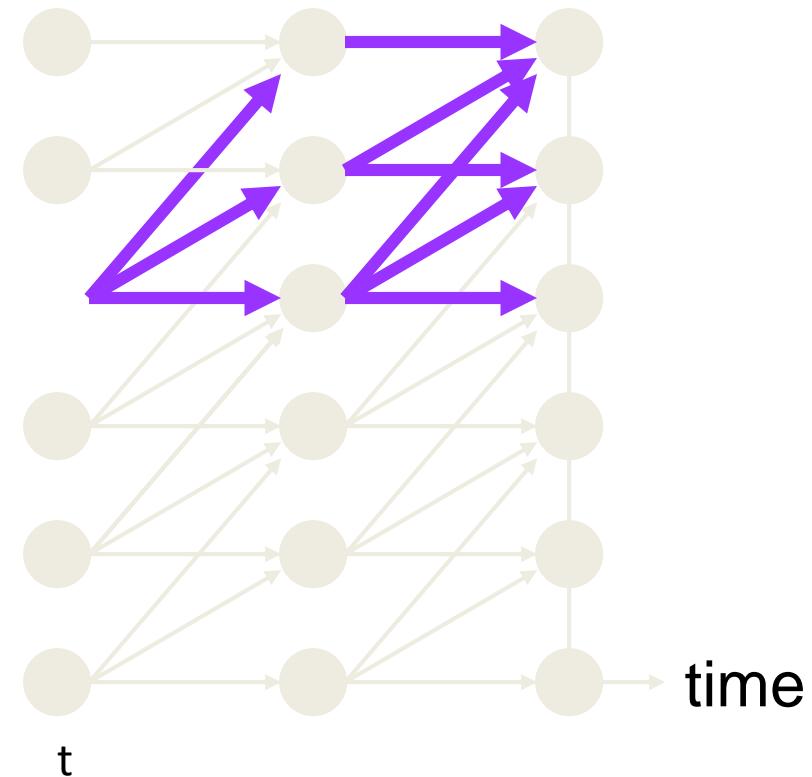


$$\alpha_u(s, t) = \sum_{s'} \alpha_u(s', t-1) P(s|s') P(x_{u,t} | s)$$

λ represents the complete current set of HMM parameters

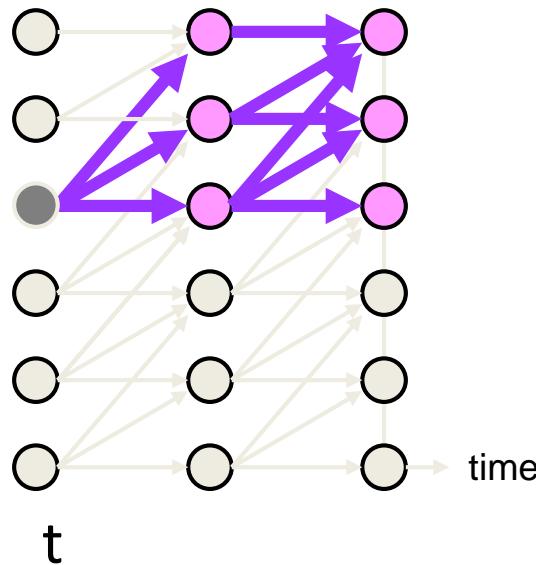
The Future Paths

- The blue portion represents the probability of all state sequences that began at state s at time t
 - Like the red portion it can be computed using a *backward recursion*



The Backward Recursion

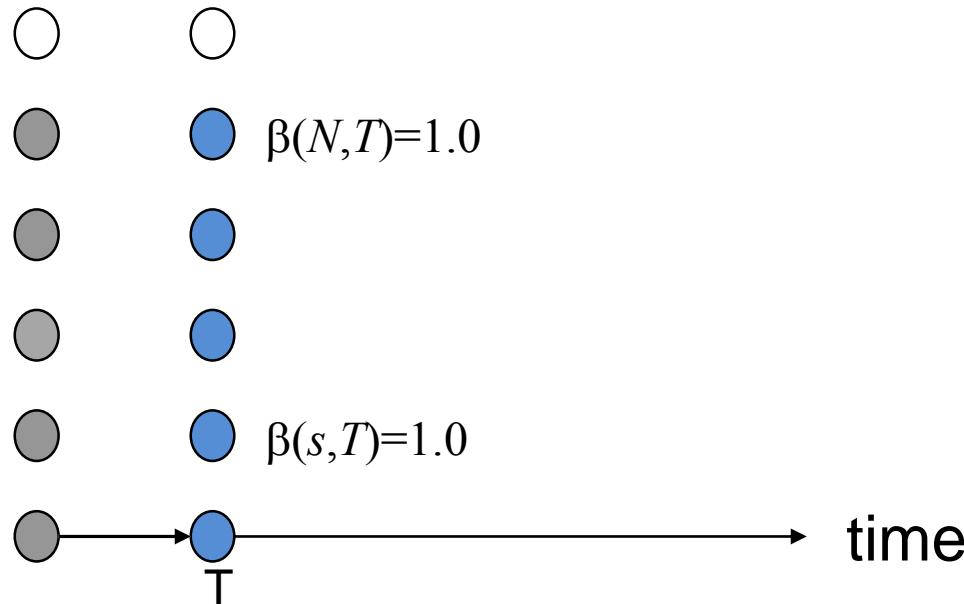
$$\beta_u(s,t) = P(x_{u,t+1}, x_{u,t+2}, \dots, x_{u,T} | state(t) = s, \lambda)$$



- $\beta_u(s,t)$ is the total probability of ALL state sequences that depart from s at time t , and all observations after x_t

The Backward Recursion

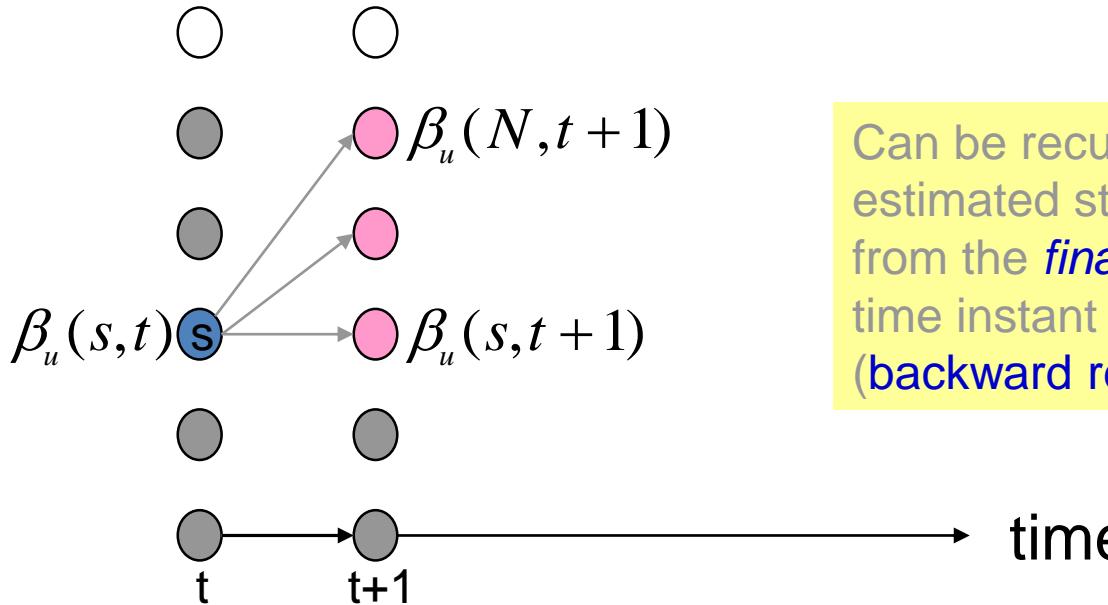
$$\beta_u(s, t) = P(x_{u,t+1}, x_{u,t+2}, \dots, x_{u,T} | state(t) = s, \lambda)$$



- $\beta(s, T) = 1$ at the final time instant for all valid final states
 - Since the future is a deterministic nothing..

The Backward Recursion

$$\beta_u(s, t) = P(x_{u,t+1}, x_{u,t+2}, \dots, x_{u,T} | state(t) = s, \lambda)$$



$$\beta_u(s, t) = \sum_{s'} \beta_u(s', t+1) P(s'|s) P(x_{u,t+1}|s')$$

- Note: Beta for any time t does *not* include the contribution of the observation at that time
 - $\beta(s, t)$ does not factor in $P(x_t|s)$

The backward algorithm

1. Initialize all beta terms at t=T:

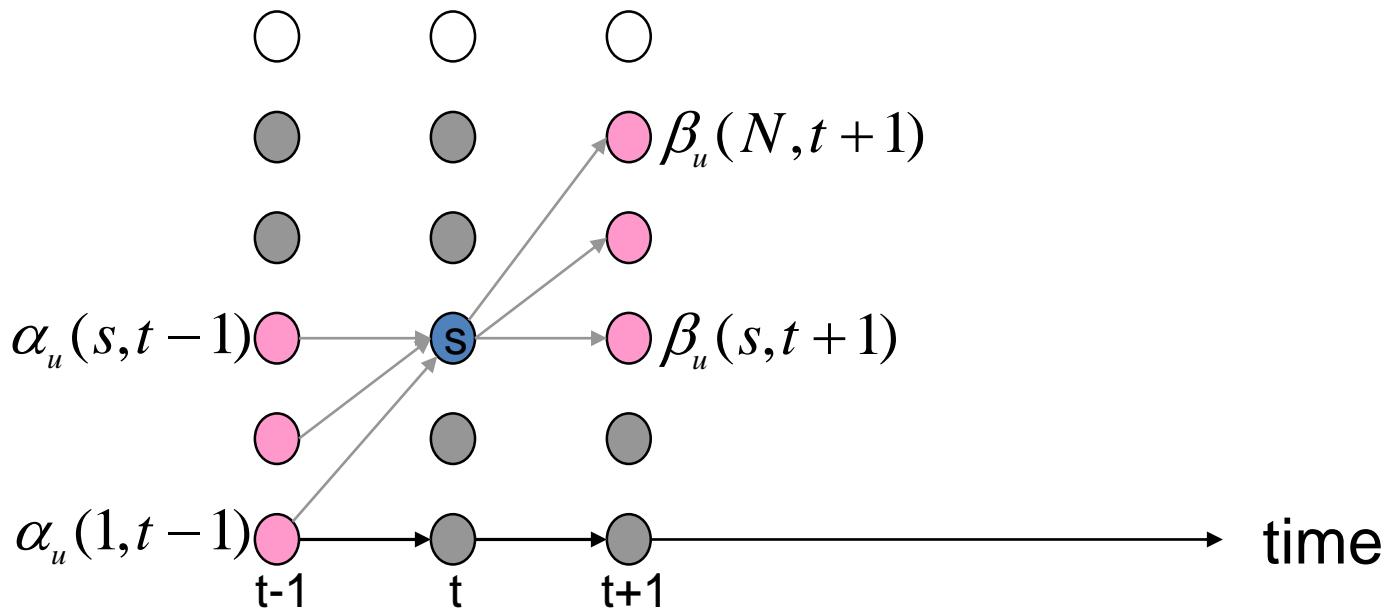
$$\beta(s, T) = 1$$

2. Recursively estimate betas for all prior time steps

$$\beta_u(s, t) = \sum_{s'} \beta_u(s', t+1) P(s'|s) P(x_{u,t+1}|s')$$

The complete probability

$$\alpha_u(s,t)\beta_u(s,t) = P(x_{u,1}, x_{u,2}, \dots, x_{u,T}, state(t) = s | \lambda)$$



$$= P(\mathbf{X}_u, state(t) = s | \lambda)$$

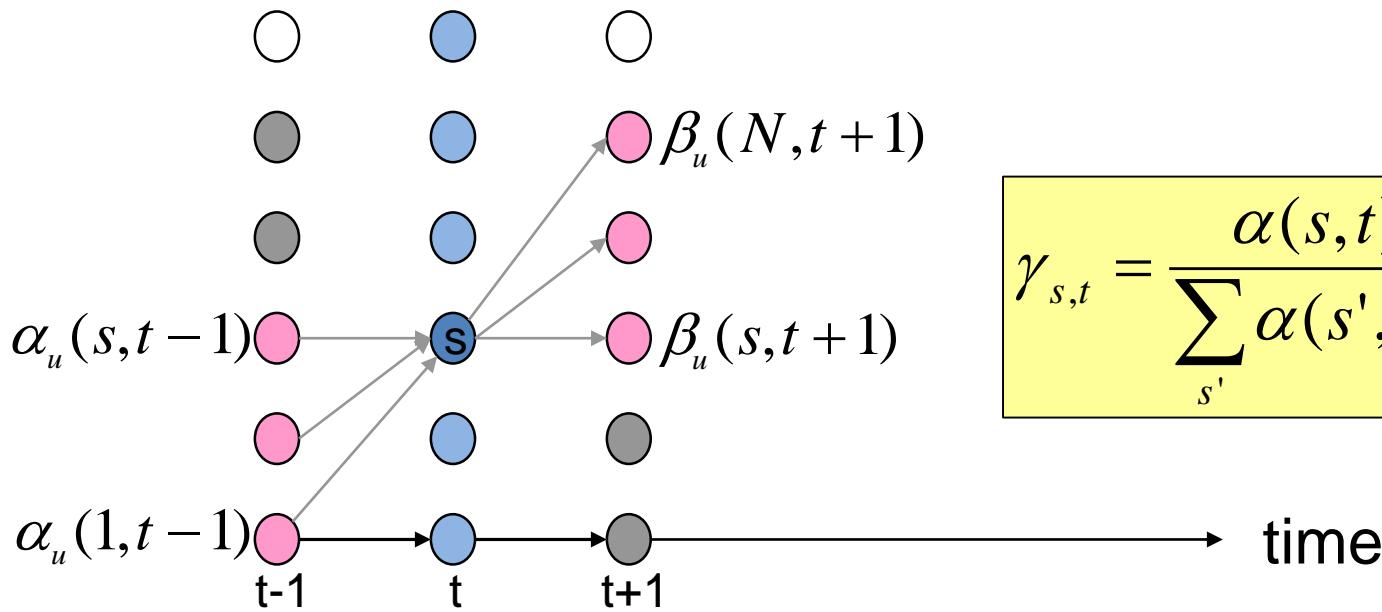
Posterior probability of a state

- The probability that the process was in state s at time t , given that we have observed the data is obtained by simple normalization

$$P(state(t) = s | \mathbf{X}_u, \lambda) = \frac{P(\mathbf{X}_u, state(t) = s | \lambda)}{\sum_{s'} P(\mathbf{X}_u, state(t) = s' | \lambda)} = \frac{\alpha_u(s, t) \beta_u(s, t)}{\sum_{s'} \alpha_u(s', t) \beta_u(s', t)}$$

- This term is often referred to as the gamma term and denoted by $\gamma_{s,t}$

The complete probability



- The gamma at any state at any time is obtained by normalizing the product of alphas and betas to sum to 1.0 over the corresponding column of the trellis

Update Rules

- Once we have the state probabilities (the gammas) the update rules are obtained through a simple modification of the formulae used for segmental K-means
 - This new learning algorithm is known as the Baum-Welch learning procedure
- Case1: State output densities are Gaussians

Update Rules

$$\mu_s = \frac{1}{N_s} \sum_{x \in s} x$$

$$C_s = \frac{1}{N_s} \sum_{x \in s} (x - \mu_s)^T (x - \mu_s)$$

Segmental K-means

$$\mu_s = \frac{\sum_u \sum_t \gamma_{u,s,t} x_{u,t}}{\sum_u \sum_t \gamma_{u,s,t}}$$

$$C_s = \frac{\sum_u \sum_t \gamma_{u,s,t} (x - \mu_s)^T (x - \mu_s)}{\sum_u \sum_t \gamma_{s,u,t}}$$

Baum Welch

- A similar update formula reestimates transition probabilities
- The *initial* state probabilities $P(s)$ also have a similar update rule

Case 2: State output densities are Gaussian Mixtures

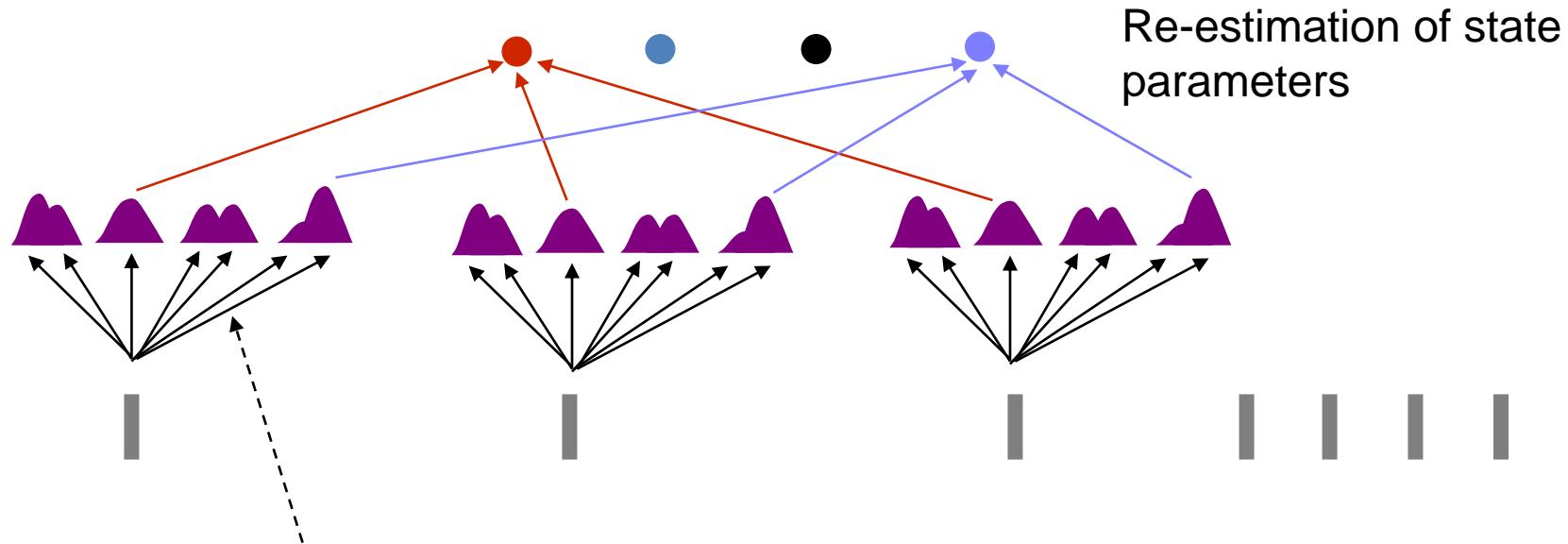
- When state output densities are Gaussian *mixtures*, more parameters must be estimated

$$P(x | s) = \sum_{i=0}^{K-1} w_{s,i} \text{Gaussian}(x; \mu_{s,i}, C_{s,i})$$

- The mixture weights $w_{s,i}$, mean $\mu_{s,i}$ and covariance $C_{s,i}$ of every Gaussian in the distribution of each state must be estimated

Splitting the Gamma

We split the gamma for any state among all the Gaussians at that state



A posteriori probability that the t^{th} vector was generated by the k^{th} Gaussian of state s

$$\gamma_{k,s,u,t} = P(\text{state}(t) = s | \mathbf{X}_u, \lambda) P(k^{\text{th}}. \text{Gaussian} | \text{state}(t) = s, x_{u,t}, \lambda)$$

Splitting the Gamma among Gaussians

A posteriori probability that the t^{th} vector was generated by the k^{th} Gaussian of state s

$$\gamma_{k,s,t} = P(\text{state}(t) = s \mid \mathbf{X}, \lambda) P(k^{\text{th}}.\text{Gaussian} \mid \text{state}(t) = s, x_t, \lambda)$$

$$\gamma_{k,s,t} = \gamma_{s,t} \frac{w_{k,s} \frac{1}{\sqrt{(2\pi)^D |C_{k,s}|}} e^{-\frac{1}{2}(x_t - \mu_{k,s})^T C_{k,s}^{-1} (x_t - \mu_{k,s})}}{\sum_{k'} w_{k',s} \frac{1}{\sqrt{(2\pi)^D |C_{k',s}|}} e^{-\frac{1}{2}(x_t - \mu_{k',s})^T C_{k',s}^{-1} (x_t - \mu_{k',s})}}$$

EM algorithm for GMM

Given an initial guess Θ^g , find Θ^{new} as follows

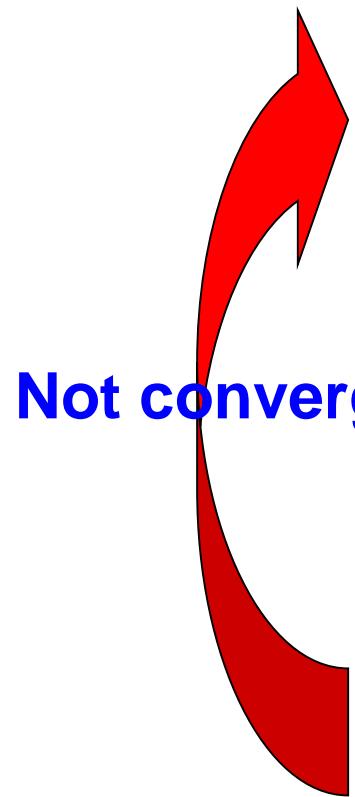
$$p(l | \mathbf{x}_i, \Theta^g) = \frac{\alpha_l^g p_l(\mathbf{x}_i | \theta_l^g)}{\sum_{j=1}^M \alpha_j^g p_j(\mathbf{x}_i | \theta_j^g)}$$

$$\alpha_l^{new} = \frac{1}{N} \sum_{i=1}^N p(l | \mathbf{x}_i, \Theta^g)$$

$$\boldsymbol{\mu}_l^{new} = \frac{\sum_{i=1}^N \mathbf{x}_i p(l | \mathbf{x}_i, \Theta^g)}{\sum_{i=1}^N p(l | \mathbf{x}_i, \Theta^g)}$$

$$\Sigma_l^{new} = \frac{\sum_{i=1}^N p(l | \mathbf{x}_i, \Theta^g) (\mathbf{x}_i - \boldsymbol{\mu}_l^{new})(\mathbf{x}_i - \boldsymbol{\mu}_l^{new})^T}{\sum_{i=1}^N p(l | \mathbf{x}_i, \Theta^g)}$$

$$\Theta^g \leftarrow \Theta^{new}$$



Updating HMM Parameters

$$\tilde{\mu}_{k,s} = \frac{\sum_u \sum_t \gamma_{k,s,u,t} x_{u,t}}{\sum_u \sum_t \gamma_{k,s,u,t}}$$

$$\tilde{\mathbf{C}}_{k,s} = \frac{\sum_u \sum_t \gamma_{k,s,u,t} (x_{u,t} - \tilde{\mu}_{k,s}) (x_{u,t} - \tilde{\mu}_{k,s})^T}{\sum_u \sum_t \gamma_{k,s,u,t}}$$

$$\tilde{w}_{k,s} = \frac{\sum_u \sum_t \gamma_{k,s,u,t}}{\sum_u \sum_t \sum_j \gamma_{j,s,u,t}}$$

- Note: Every observation contributes to the update of parameter values of every Gaussian of every state

Overall Training Procedure: Single Gaussian PDF

- Determine a topology for the HMM
 - Initialize all HMM parameters
 - Initialize all allowed transitions to have the same probability
 - Initialize all state output densities to be Gaussians
 - We'll revisit initialization
1. Over all utterances, compute the “sufficient” statistics $\sum_u \sum_t \gamma_{u,s,t}$ $\sum_u \sum_t \gamma_{u,s,t} x_{u,t}$ $\sum_u \sum_t \gamma_{u,s,t} (x - \mu_s)^T (x - \mu_s)$
 2. Use update formulae to compute new HMM parameters
 3. If the overall probability of the training data has not converged, return to step 1

An Implementational Detail

- Step 1 computes “buffers” over all utterance

$$\sum_u \sum_t \gamma_{u,s,t} = \sum_{u \in U_1} \sum_t \gamma_{u,s,t} + \sum_{u \in U_2} \sum_t \gamma_{u,s,t} + \dots$$

Assuming
1 gaussian/stt
on this slide

$$\sum_u \sum_t \gamma_{u,s,t} x_{u,t} = \sum_{u \in U_1} \sum_t \gamma_{u,s,t} x_{u,t} + \sum_{u \in U_2} \sum_t \gamma_{u,s,t} x_{u,t} + \dots$$

$$\sum_u \sum_t \gamma_{u,s,t} (x - \mu_s)^T (x - \mu_s) = \sum_{u \in U_1} \sum_t \gamma_{u,s,t} (x - \mu_s)^T (x - \mu_s) + \sum_{u \in U_2} \sum_t \gamma_{u,s,t} (x - \mu_s)^T (x - \mu_s) + \dots$$

- This can be split and parallelized
 - U_1, U_2 etc. can be processed on separate machines

Machine 1

$$\sum_{u \in U_1} \sum_t \gamma_{u,s,t} \quad \sum_{u \in U_1} \sum_t \gamma_{u,s,t} x_{u,t}$$

$$\sum_{u \in U_1} \sum_t \gamma_{u,s,t} (x - \mu_s)(x - \mu_s)^T$$

Machine 2

$$\sum_{u \in U_2} \sum_t \gamma_{u,s,t} \quad \sum_{u \in U_2} \sum_t \gamma_{u,s,t} x_{u,t}$$

$$\sum_{u \in U_2} \sum_t \gamma_{u,s,t} (x - \mu_s)(x - \mu_s)^T$$

An Implementational Detail

- Step2 aggregates and adds buffers before updating the models

$$\sum_u \sum_t \gamma_{u,s,t} = \sum_{u \in U_1} \sum_t \gamma_{u,s,t} + \sum_{u \in U_2} \sum_t \gamma_{u,s,t} + ..$$

$$\sum_u \sum_t \gamma_{u,s,t} x_{u,t} = \sum_{u \in U_1} \sum_t \gamma_{u,s,t} x_{u,t} + \sum_{u \in U_2} \sum_t \gamma_{u,s,t} x_{u,t} + ...$$

$$\sum_u \sum_t \gamma_{u,s,t} (x - \mu_s)^T (x - \mu_s) = \sum_{u \in U_1} \sum_t \gamma_{u,s,t} (x - \mu_s)^T (x - \mu_s) + \sum_{u \in U_2} \sum_t \gamma_{u,s,t} (x - \mu_s)^T (x - \mu_s) + ..$$

$$\tilde{\mu}_{k,s} = \frac{\sum_u \sum_t \gamma_{k,s,u,t} x_{u,t}}{\sum_u \sum_t \gamma_{k,s,u,t}}$$

$$\tilde{\mathbf{C}}_{k,s} = \frac{\sum_u \sum_t \gamma_{k,s,u,t} (x_{u,t} - \tilde{\mu}_{k,s}) (x_{u,t} - \tilde{\mu}_{k,s})^T}{\sum_u \sum_t \gamma_{k,s,u,t}}$$

$$\tilde{w}_{k,s} = \frac{\sum_u \sum_t \gamma_{k,s,u,t}}{\sum_u \sum_t \sum_j \gamma_{j,s,u,t}}$$

An Implementational Detail

- Step2 aggregates and adds buffers before updating the models

$$\sum_u \sum_t \gamma_{u,s,t} = \sum_{u \in U_1} \sum_t \gamma_{u,s,t} + \sum_{u \in U_2} \sum_t \gamma_{u,s,t} + \dots$$

$$\sum_u \sum_t \gamma_{u,s,t} x_{u,t} = \sum_{u \in U_1} \sum_t \gamma_{u,s,t} x_{u,t} + \sum_{u \in U_2} \sum_t \gamma_{u,s,t} x_{u,t} + \dots$$

$$\sum_u \sum_t \gamma_{u,s,t} (x - \mu_s)^T (x - \mu_s) = \sum_{u \in U_1} \sum_t \gamma_{u,s,t} (x - \mu_s)^T (x - \mu_s) + \sum_{u \in U_2} \sum_t \gamma_{u,s,t} (x - \mu_s)^T (x - \mu_s) + \dots$$

$$\tilde{\mu}_{k,s} = \frac{\sum_u \sum_t \gamma_{k,s,u,t} x_{u,t}}{\sum_u \sum_t \gamma_{k,s,u,t}}$$

$$\tilde{\mathbf{C}}_{k,s} = \frac{\sum_t \gamma_{k,s,u,t} (x_{u,t} - \tilde{\mu}_{k,s})}{\sum_u \sum_t \gamma_{k,s,u,t}} (x_{u,t} - \tilde{\mu}_{k,s})^T$$

$$\tilde{w}_{k,s} = \frac{\sum_u \sum_t \gamma_{k,s,u,t}}{\sum_u \sum_t \sum_j \gamma_{j,s,u,t}}$$

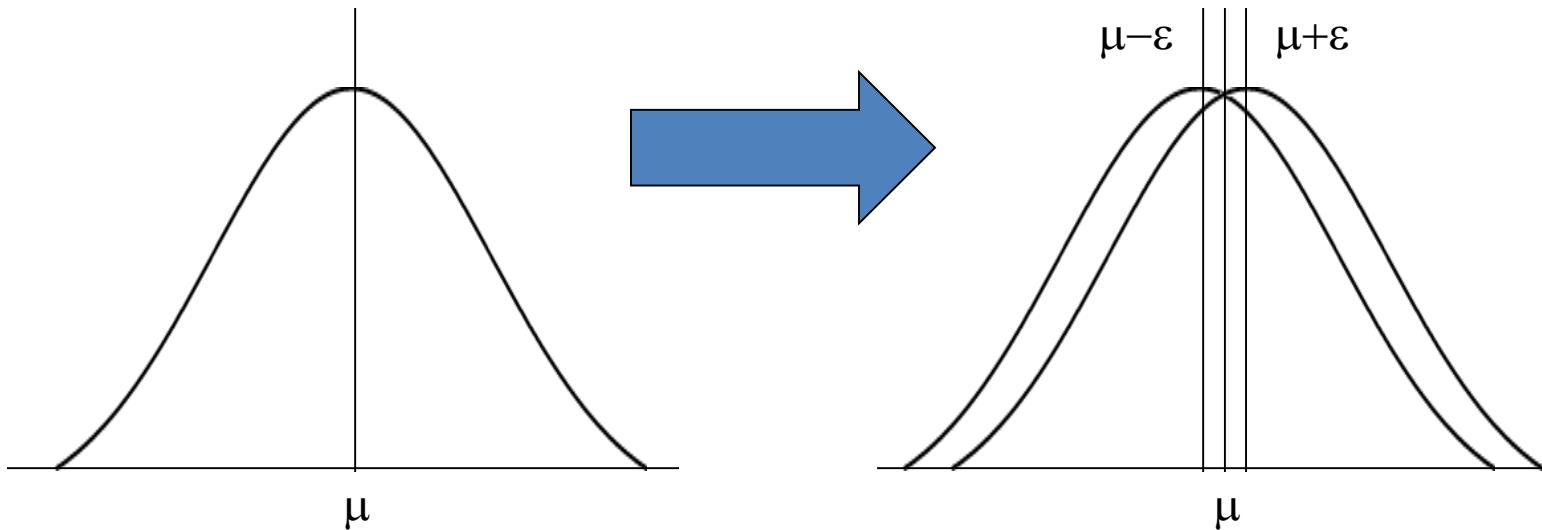
Computed by
machine 1

Computed by
machine 2

Training for HMMs with *Gaussian Mixture* State Output Distributions

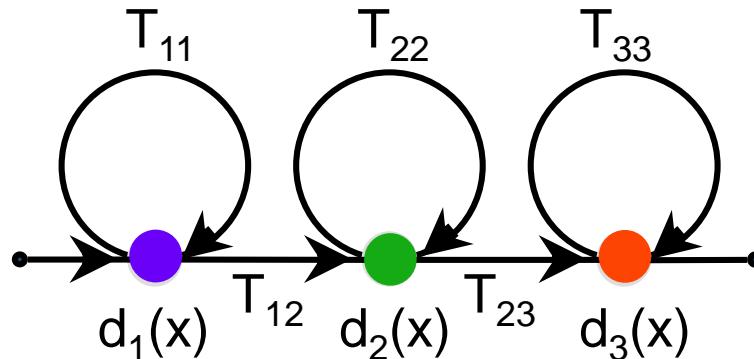
- Gaussian *Mixtures* are obtained by splitting
1. Train an HMM with (single) Gaussian state output distributions
 2. Split the Gaussian with the largest variance
 - Perturb the mean by adding and subtracting a small number
 - This gives us 2 Gaussians. Partition the mixture weight of the Gaussian into two halves, one for each Gaussian
 - A mixture with N Gaussians now becomes a mixture of N+1 Gaussians
 3. Iterate BW to convergence
 4. If the desired number of Gaussians not obtained, return to 2

Splitting a Gaussian

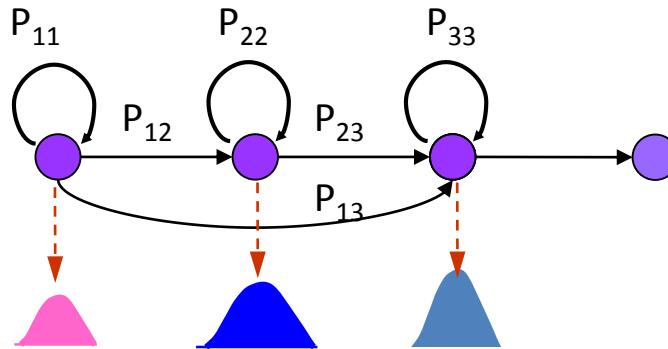


- The mixture weight w for the Gaussian gets shared as $0.5w$ by each of the two split Gaussians

Transition Probabilities

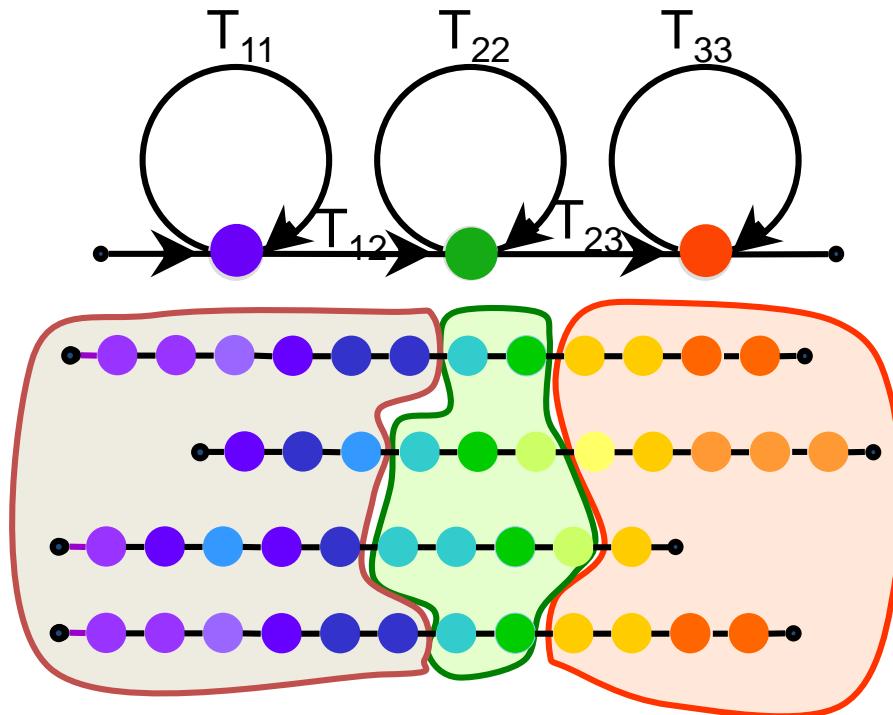


- We have seen how to compute transition penalties for templates



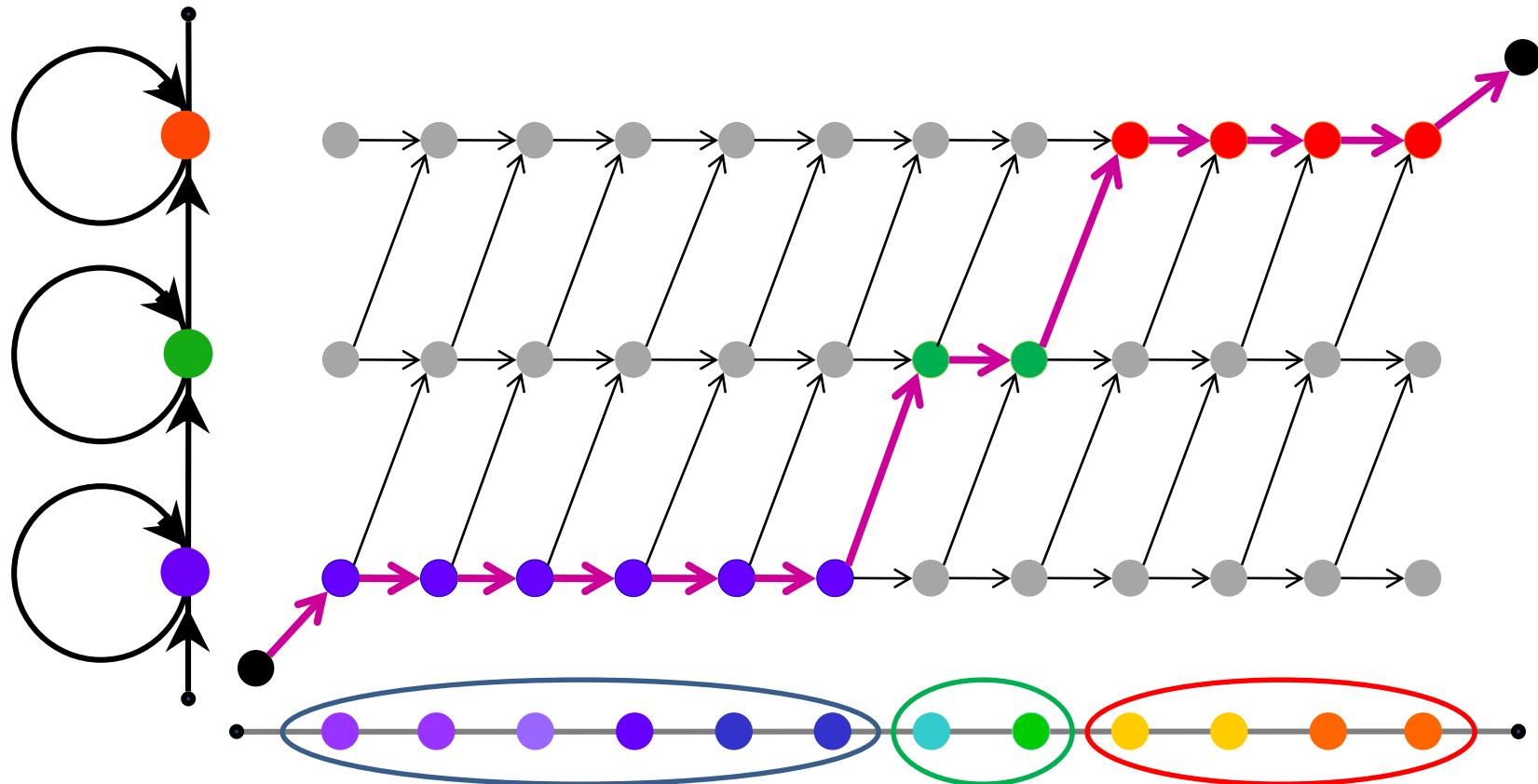
- How about transition probabilities in an HMM?
 - “Hard” estimation – by counting, as for templates
 - “Soft” estimation – need soft counts

Transition penalties by counting



- 20 vectors in state 1
 - 16 are followed by vectors in state 1
 - 4 are followed by vectors in state 2
- $P_{11} = 16/20 = 0.8 \rightarrow T_{11} = -\log(P_{11}) = -\log(0.8)$
- $P_{12} = 4/20 = 0.2 \rightarrow T_{12} = -\log(P_{12}) = -\log(0.2)$

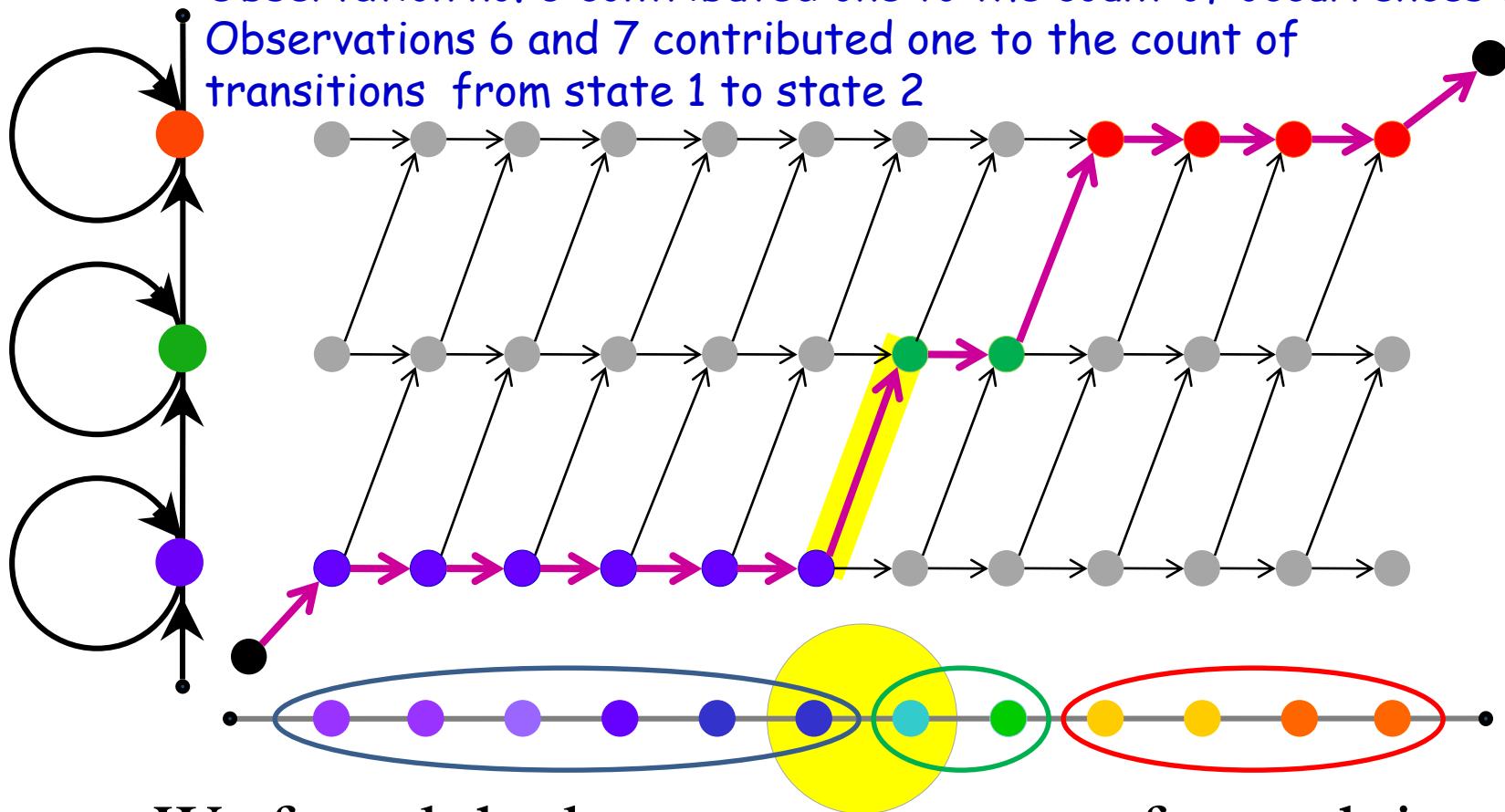
Transitions by counting



- We found the best state sequence for each input
 - And counted transitions

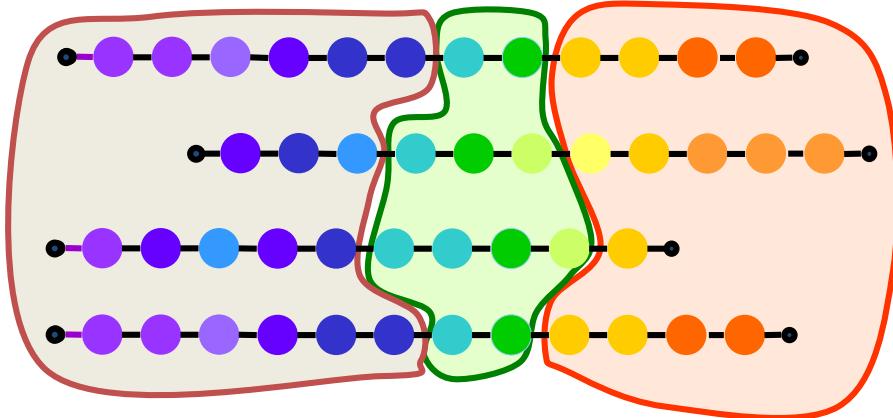
Transitions by counting

Observation no. 6 contributed one to the count of occurrences of state 1
Observations 6 and 7 contributed one to the count of transitions from state 1 to state 2



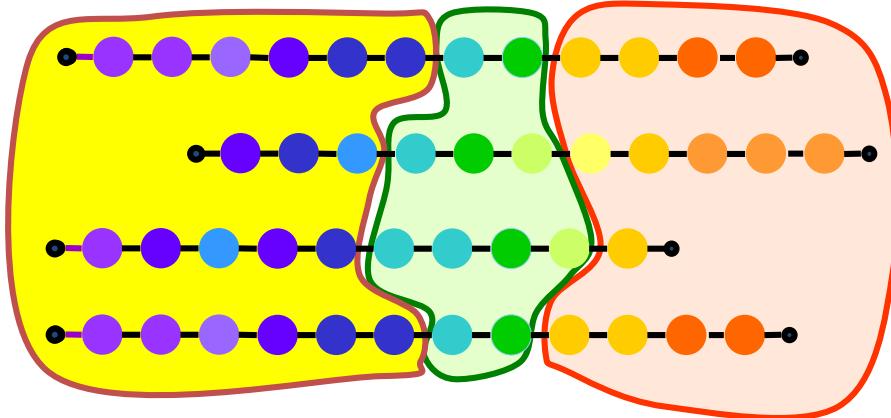
- We found the best state sequence for each input
 - And counted transitions

Probability of transitions



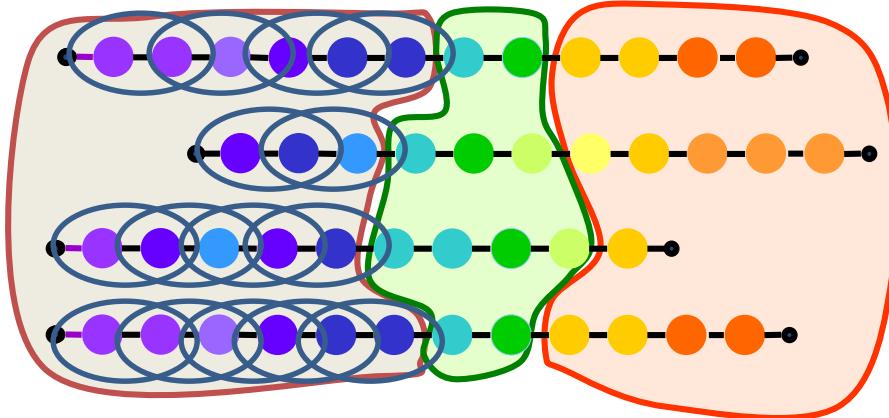
- $P(\text{transition state I} \rightarrow \text{state J}) =$
 - Count transitions(I,J) / count instances(I)

Probability of transitions



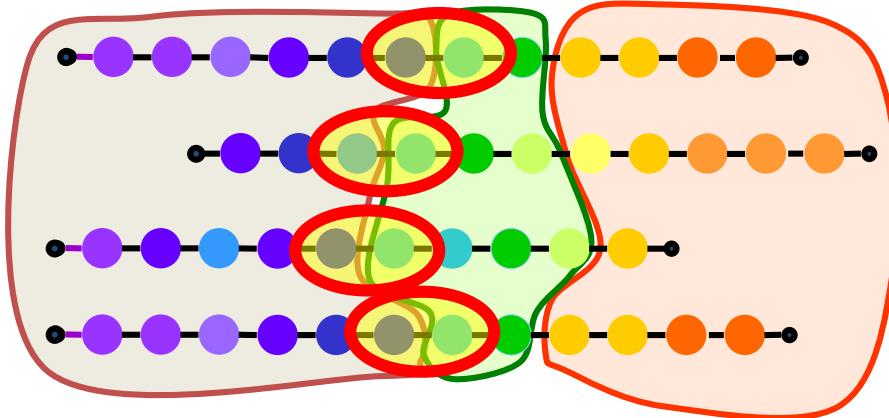
- $P(\text{transition state } I \rightarrow \text{state } J) =$
 - Count transitions(I, J) / count instances(I)
 - Count instances(I) = 20

Probability of transitions



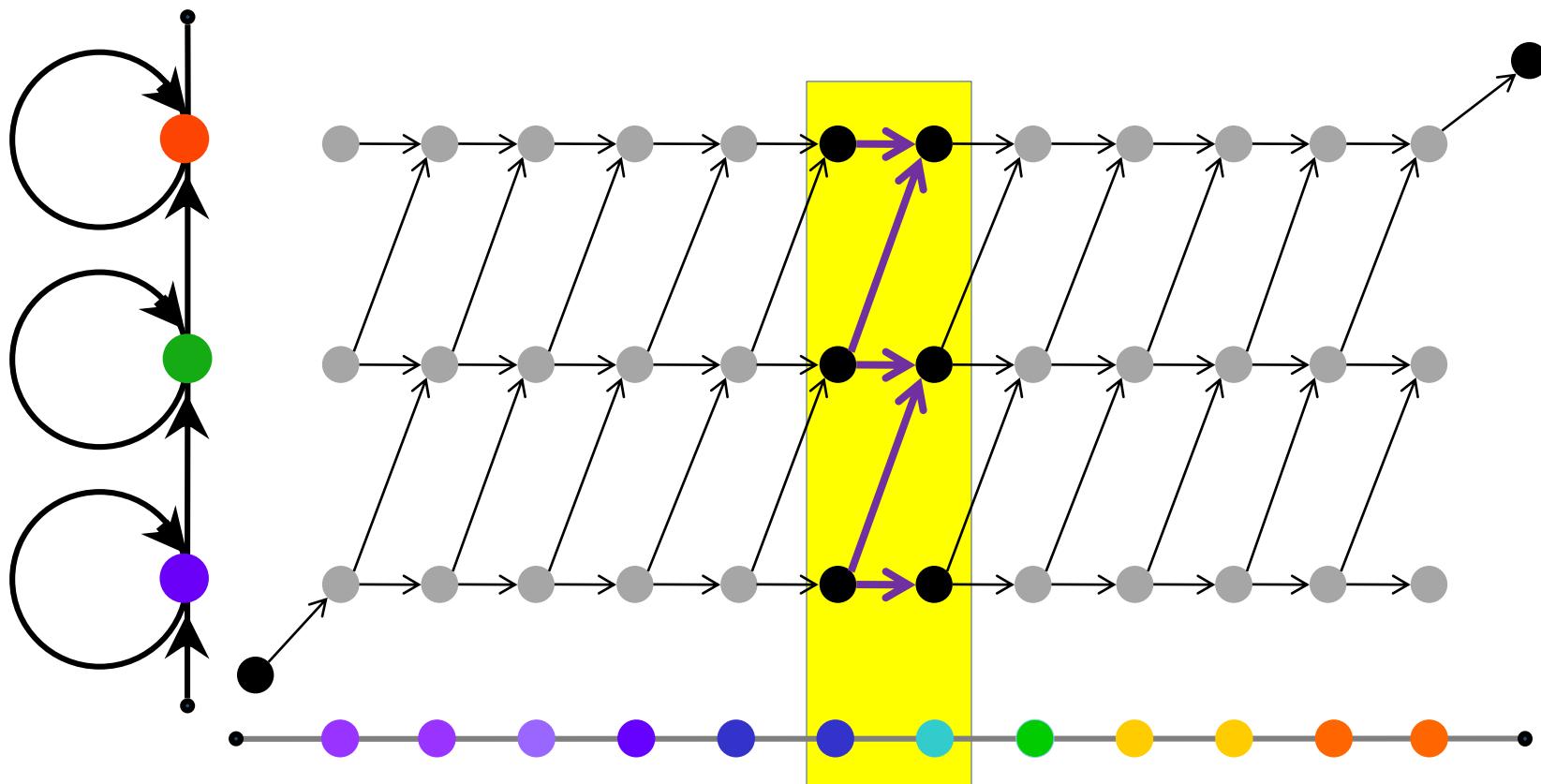
- $P(\text{transition state } I \rightarrow \text{state } J) =$
 - Count transitions(I,J) / count instances(I)
 - Count instances(I) = 20
 - Count transitions ($1,1$) = 16
 - $P(\text{transition state } 1 \rightarrow \text{state } 1) = 0.8$

Probability of transitions



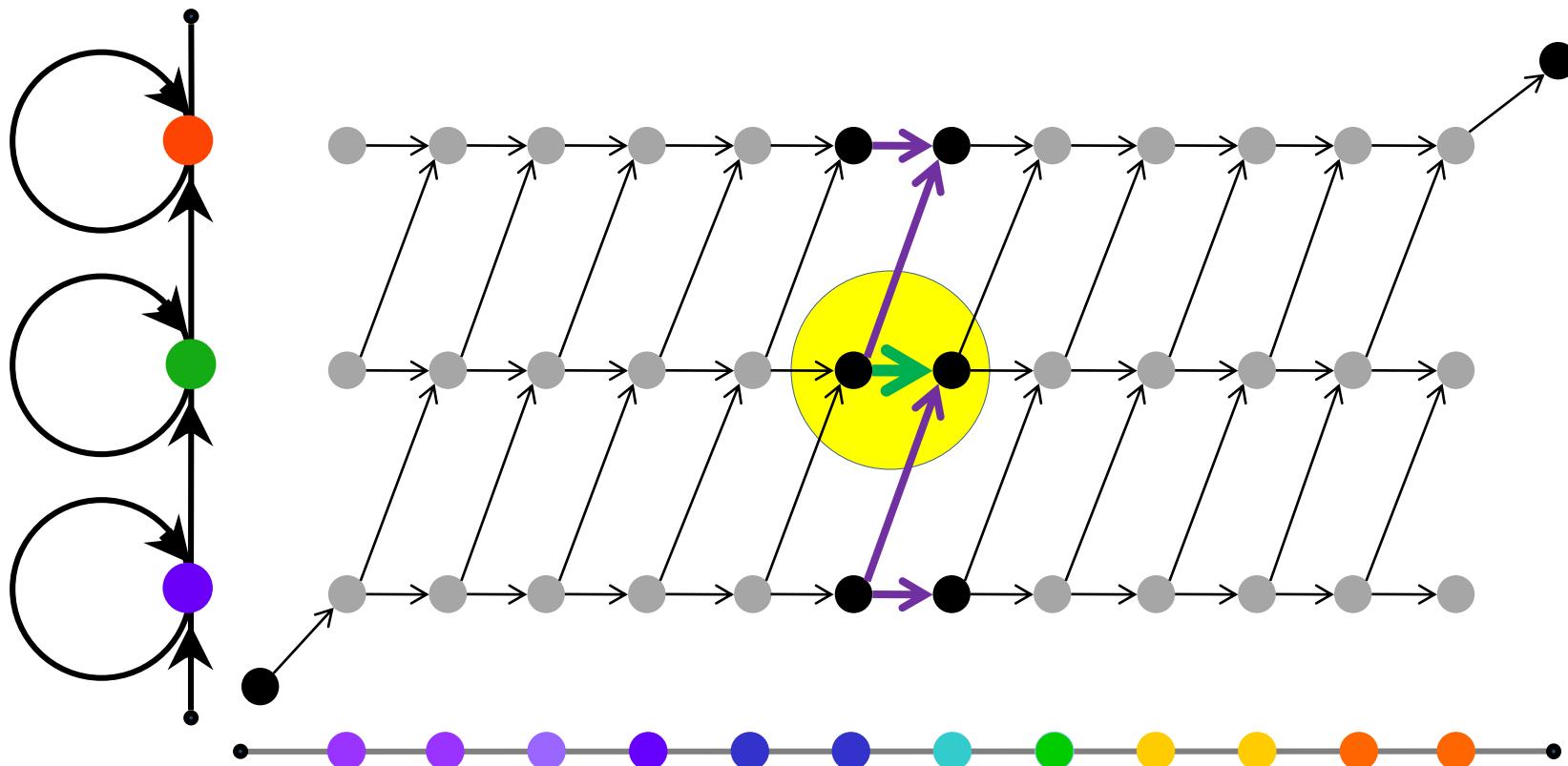
- $P(\text{transition state } I \rightarrow \text{state } J) =$
 - Count transitions(I,J) / count instances(I)
 - Count instances(1) = 20
 - Count transitions (1,2) = 4
 - $P(\text{transition state } 1 \rightarrow \text{state } 2) = 0.2$

Transitions by *soft* counting



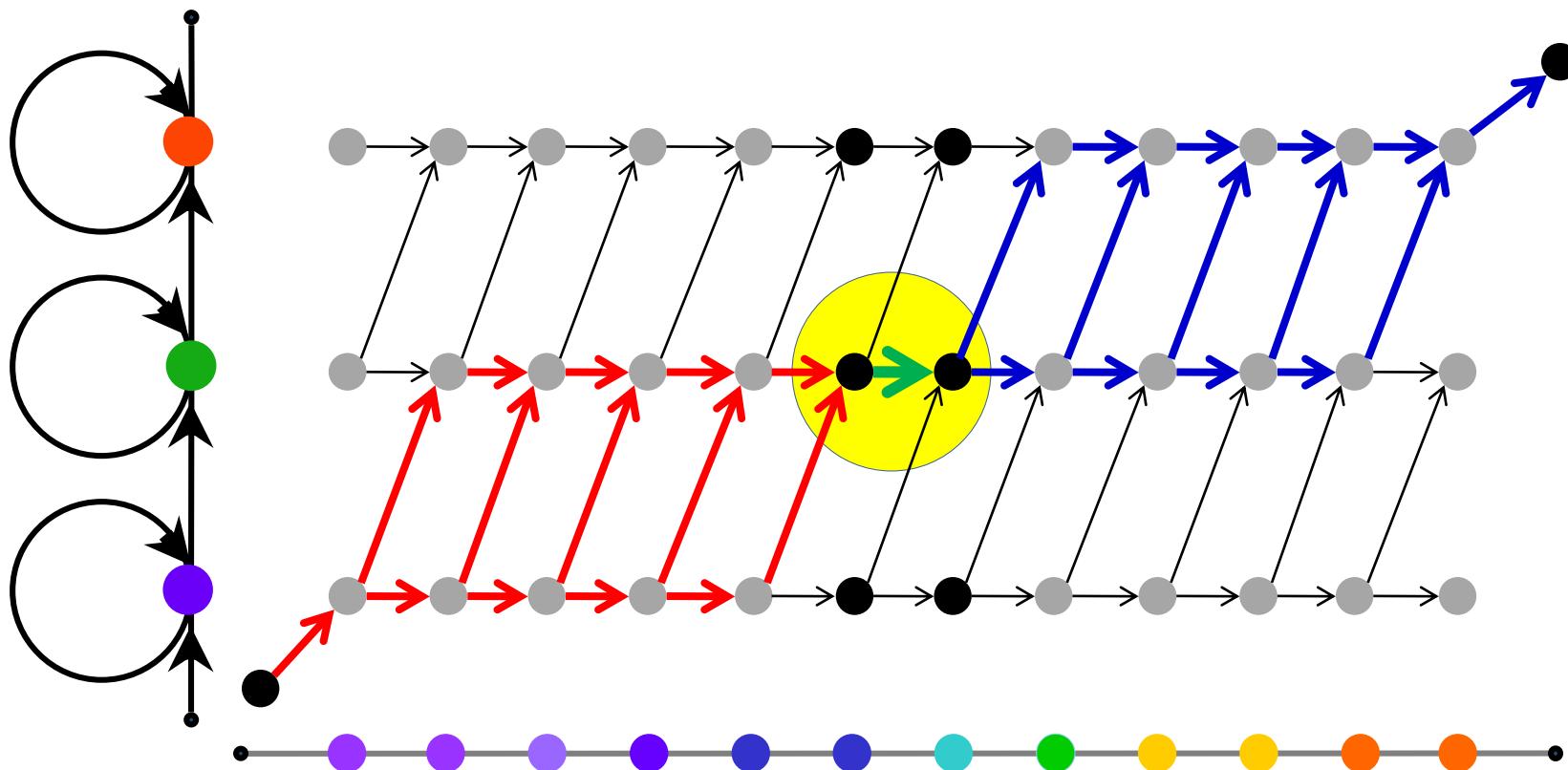
- Each observation pair contributes to every transition
 - E.g. observations 6,7 contribute counts to all of the following:
 - Transition $(1 \rightarrow 1)$, Transition $(1 \rightarrow 2)$, Transition $(2 \rightarrow 2)$, Transition $(2 \rightarrow 3)$, Transition $(3 \rightarrow 3)$

Transitions by *soft* counting



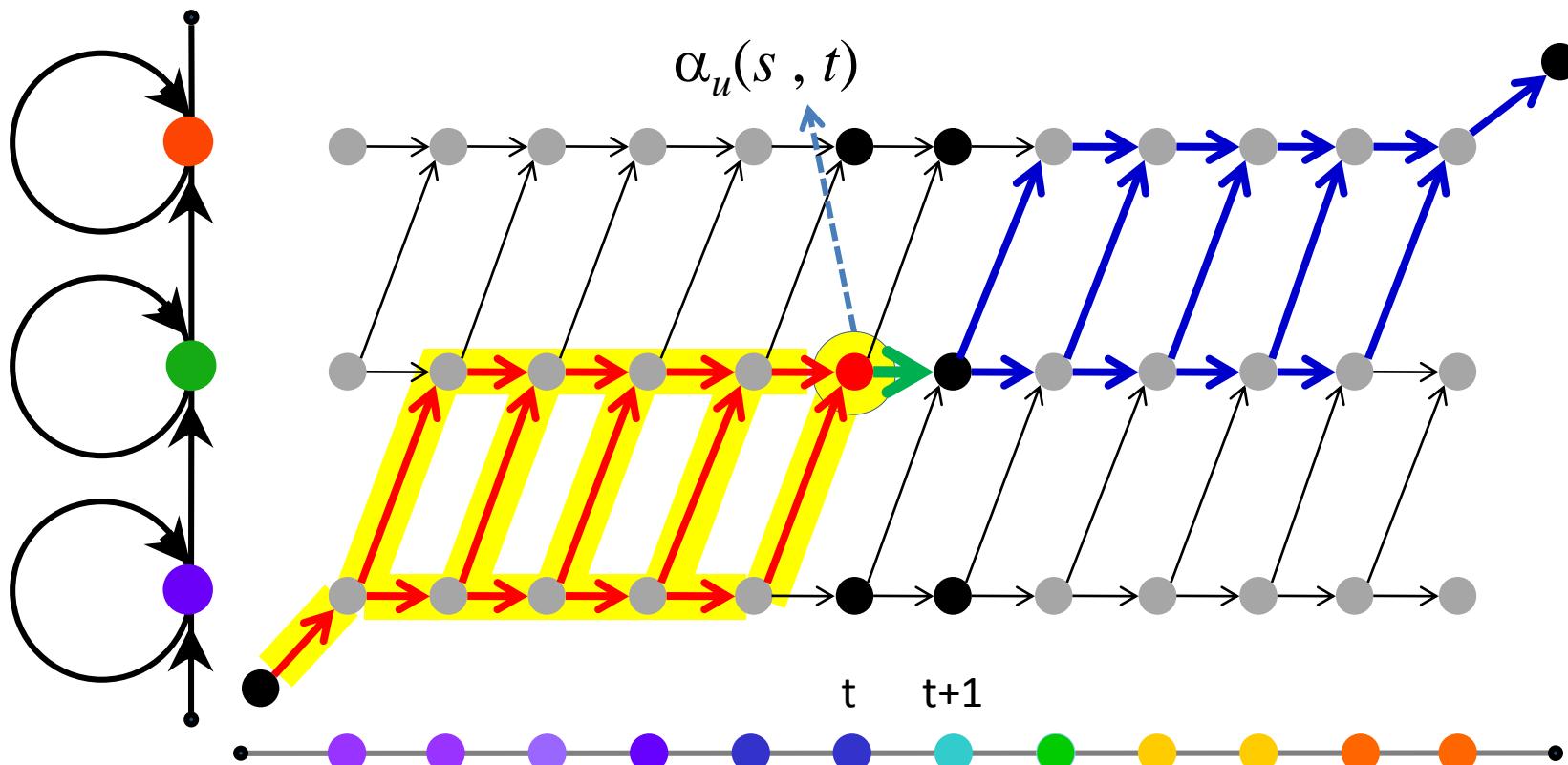
- Contribution of any transition to the count is the *a posteriori* probability of the count
 - This is a fraction
 - The fractions for all possible transitions at any time sum to 1

Transitions by *soft* counting



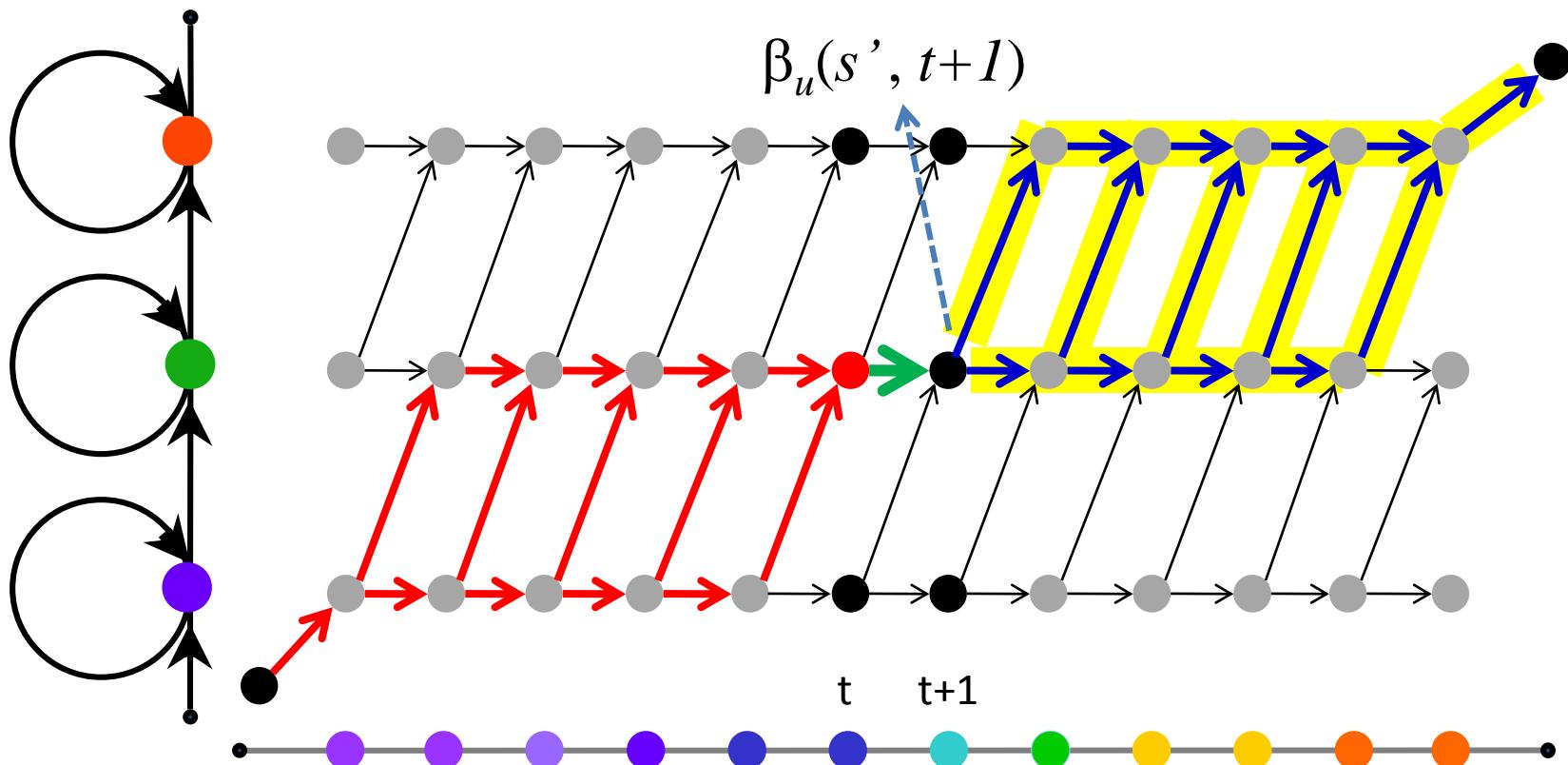
- Probability of a transition is the total probability of all paths that include the transition

Transitions by *soft* counting



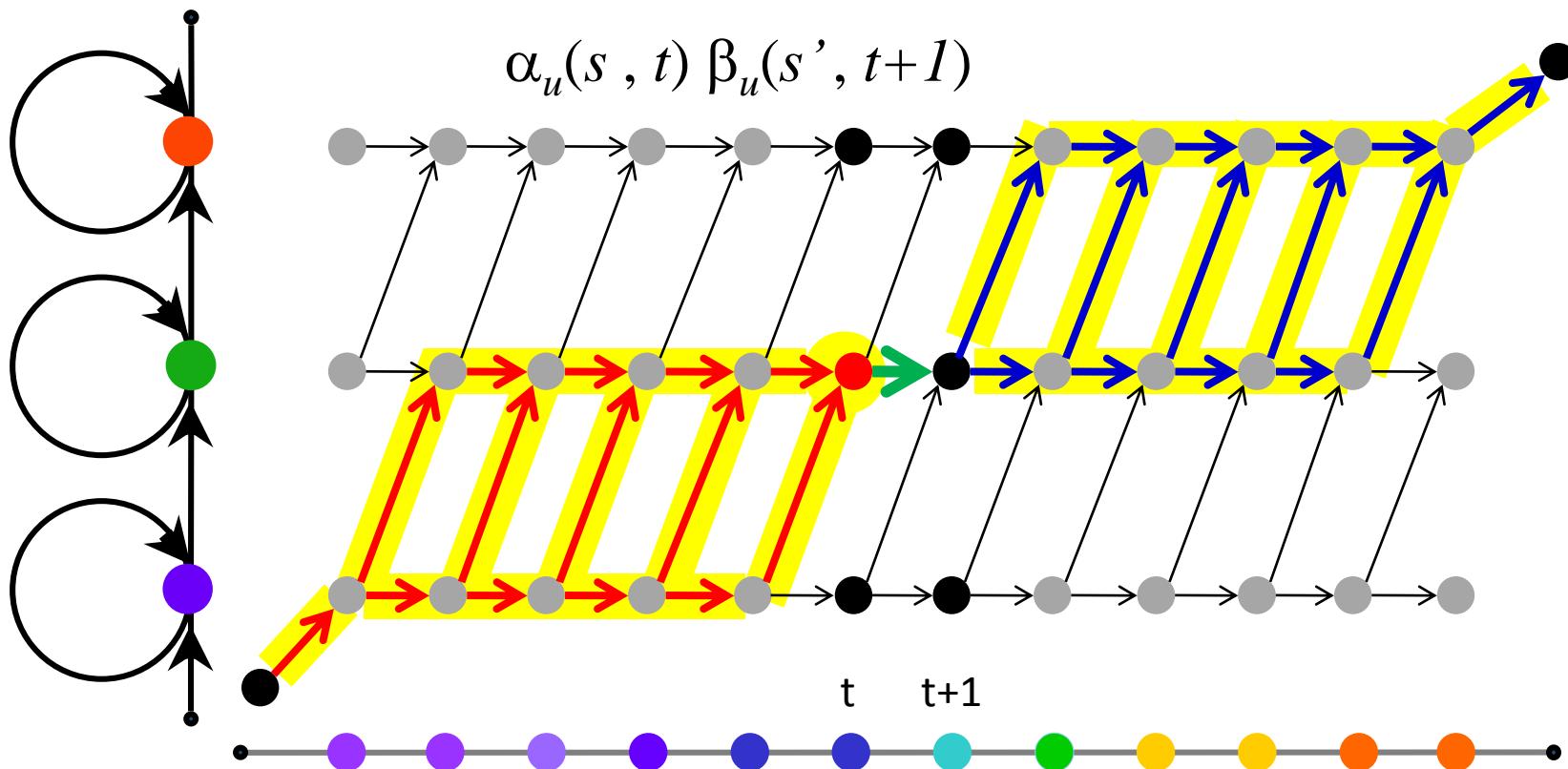
- The forward probability of the source state at t accounts for all incoming paths at time t
 - *including the t -th observation x_t*

Transitions by *soft* counting



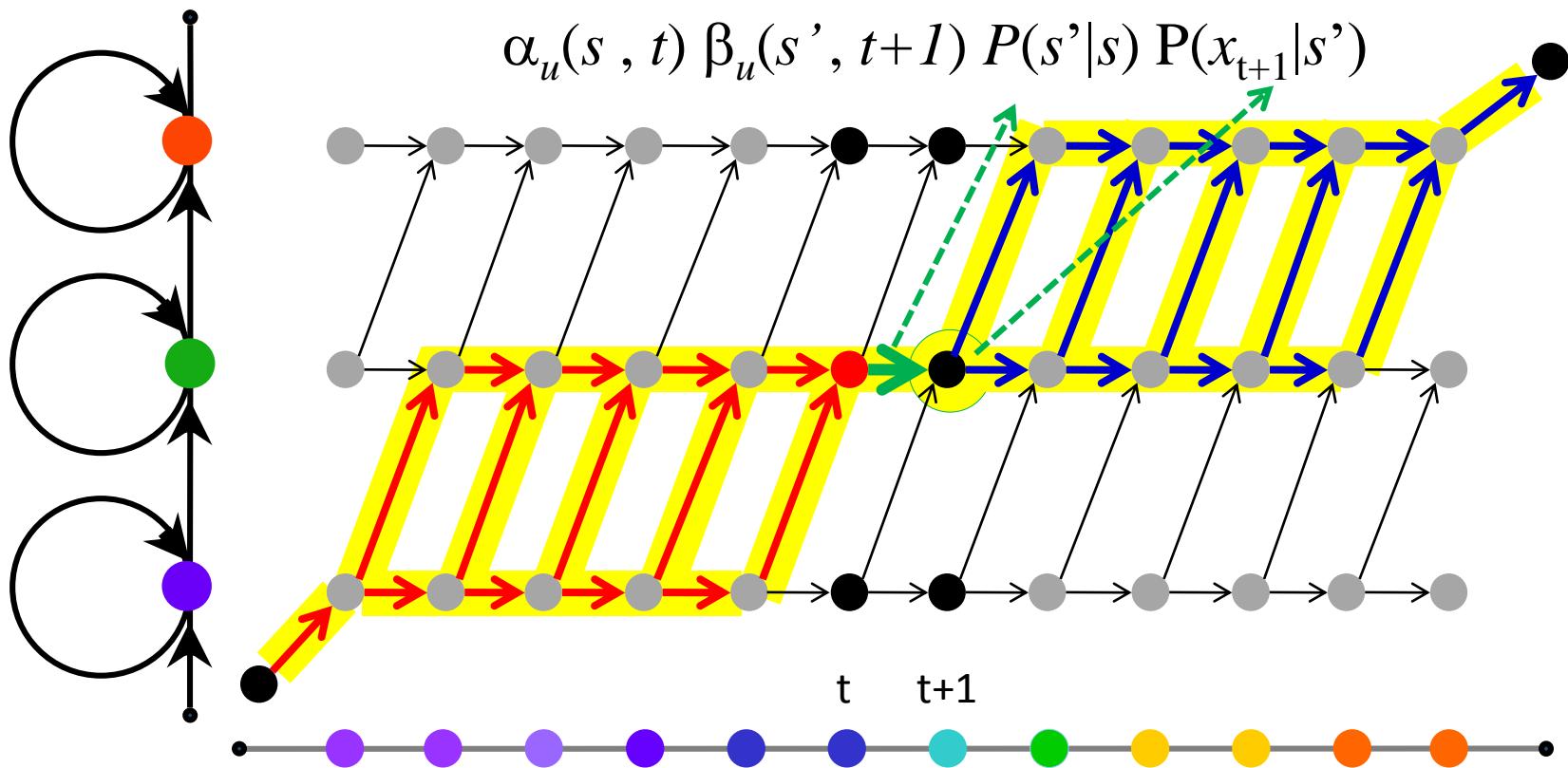
- The *backward* probability of the *destination* state at $t+1$ accounts for all *outgoing* paths from the state at time $t+1$
 - ***NOT including*** the $t+1$ -th observation x_{t+1}

Transitions by *soft* counting



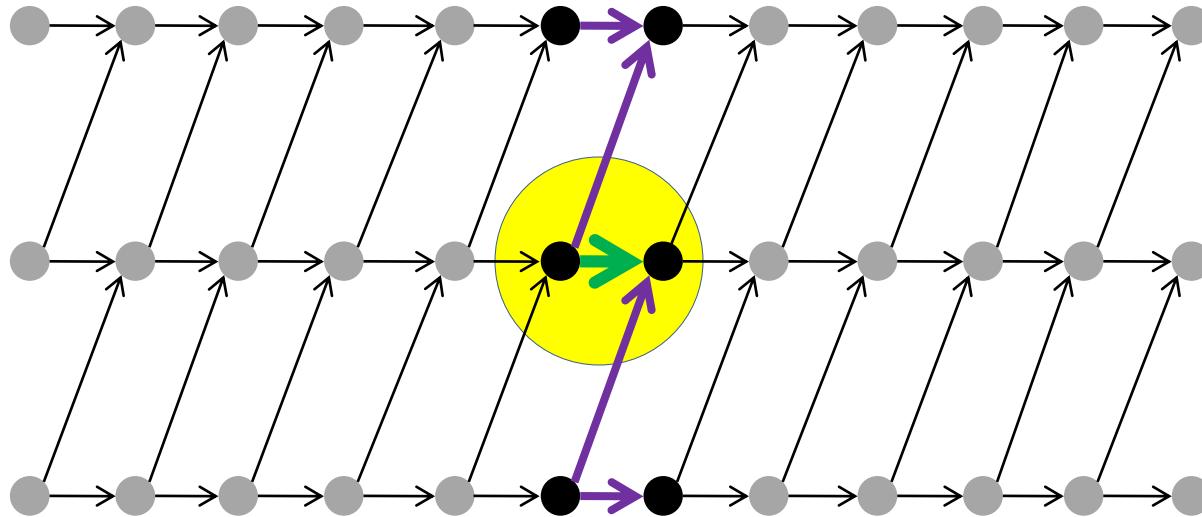
- The product of the forward probability of s at t and s' at $t+1$ accounts for all paths TO state s at t , and all paths FROM s' at $t+1$
 - But not the *transition* from s to s' or the observation at $t+1$

Transitions by *soft* counting



- By factoring in the transition probability and observation probabilities, the total *probability* is obtained

From probability to *a posteriori* probability



- The *a posteriori* probability of a transition is the ratio of its probability to the sum of all transitions at the same time

A posteriori probability of a transition

- Probability of a transition

$$P(state(t) = s, state(t+1) = s', x_1, x_2, \dots, x_N) = \alpha_u(s, t) P(s' | s) P(x_t + 1 | s') \beta_u(s', t+1)$$

- *A posteriori* probability of a transition

$$\gamma_{u,s,t,s',t+1} = \frac{P(state(t) = s, state(t+1) = s', x_1, x_2, \dots, x_N)}{\sum_{S,S'} P(state(t) = S, state(t+1) = S', x_1, x_2, \dots, x_N)}$$

$$\gamma_{u,s,t,s',t+1} = \frac{\alpha_u(s, t) P(s' | s) P(x_t + 1 | s') \beta_u(s', t+1)}{\sum_{S,S'} \alpha_u(S, t) P(S' | S) P(x_t + 1 | S') \beta_u(S', t+1)}$$

Estimate of Transition Probabilities

$$P(s'|s) = \frac{\sum_{u}^u \sum_{t}^t \gamma_{u,s,t,s',t+1}}{\sum_u^u \sum_t^t \gamma_{u,s,t}}$$

- Numerator is total “soft” count of transitions from state s to s'
- Denominator is total “soft” count of instances of state s

Implementation of BW: underflow

- **Arithmetic underflow is a problem**

$$\alpha_u(s, t) = \sum_{s'} \underbrace{\alpha_u(s', t-1)}_{\text{probability term}} \underbrace{P(s|s') P(x_{u,t}|s)}_{\text{probability terms}}$$

- The alpha terms are a recursive product of probability terms
 - As t increases, an increasingly greater number probability terms are factored into the alpha
- All probability terms are less than 1
- With increasing time, alpha values decrease
- Within a few time instants, they underflow to 0
 - Every alpha goes to 0 at some time t . All future alphas remain 0
 - As the dimensionality of the data increases, alphas goes to 0 faster

Underflow: Solution

- One method of avoiding underflow is to scale all alphas at each time instant
 - Scale with respect to the largest alpha to make sure the largest scaled alpha is 1.0
 - Scale with respect to the sum of the alphas to ensure that all alphas sum to 1.0
 - Scaling constants must be appropriately considered when computing the final probabilities of an observation sequence
- An alternate method: Compute alphas and betas in log domain
 - How? (Not obvious)

Implementation of BW: underflow

- Similarly, arithmetic underflow can occur during beta computation

$$\beta_u(s, t) = \sum_{s'} \beta_u(s', t+1) P(s' | s) P(x_{u,t+1} | s')$$

- The beta terms are also a recursive product of probability terms and can underflow
- Underflow can be prevented by
 - Scaling: Divide all beta terms by a constant that prevents underflow
 - By performing beta computation in the log domain (now? Not obvious..)
- QUESTION: HOW DOES SCALING AFFECT THE ESTIMATION OF GAMMA TERMS**
 - For Gaussian parameter updates?
 - For transition probability updates?

Log-Domain Mathematics

there may be some lingering doubt as to what, exactly, logAdd() in the Viterbi/Forward-Backward templates is doing.

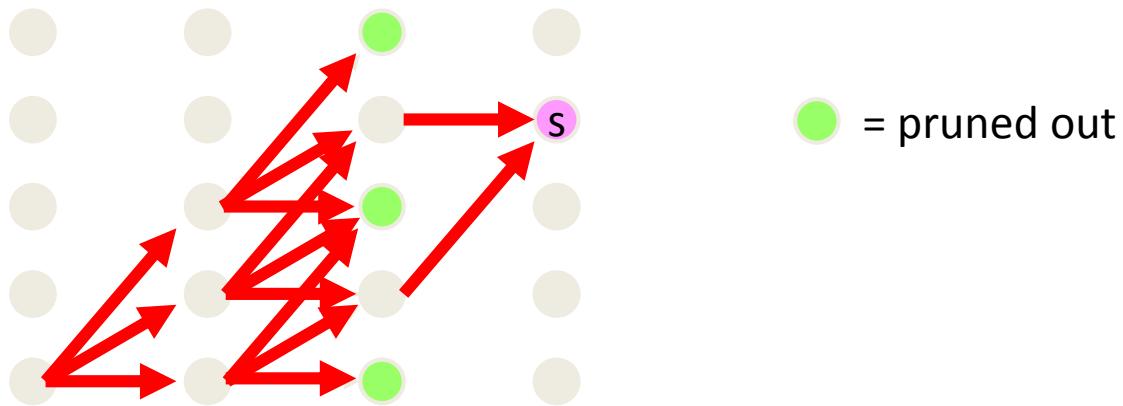
<i>“normal” domain</i>	<i>log domain</i>
x^y	$x \cdot y$
$x \cdot y$	$x+y$
$x+y$	$\text{logAdd}(x,y)$

One method: $\log_{10}(10^x + 10^y)$

But the reason we’re using logs is to avoid underflow... so

$$\begin{aligned}\log(x+y) = & \quad \log\left(x + \left[\frac{y}{x} \cdot x\right]\right) \\ & \quad \log\left(x\left[1 + \frac{y}{x}\right]\right) \\ & \quad \log(x) + \log\left(1 + \frac{y}{x}\right) \\ & \quad \log(x) + \log\left(1 + e^{\log(y) - \log(x)}\right)\end{aligned}$$

Implementation of BW: pruning



- The forward backward computation can get very expensive
- Solution: Prune
- Pruning in the forward backward algorithm raises some additional issues
 - Pruning from forward pass can be employed by backward pass
 - Convergence criteria and tests may be affected
 - More later

Building a recognizer for isolated words

- Now have all necessary components to build an HMM-based recognizer for isolated words
 - Where each word is spoken by itself in isolation
 - E.g. a simple application, where one may either say “Yes” or “No” to a recognizer and it must recognize what was said

Isolated Word Recognition with HMMs

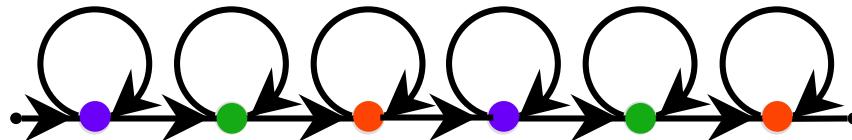
- Assuming all words are equally likely
- Training
 - Collect a set of “training” recordings for each word
 - Compute feature vector sequences for the words
 - Train HMMs for each word
- Recognition:
 - Compute feature vector sequence for test utterance
 - Compute the forward probability of the feature vector sequence from the HMM for each word
 - Alternately compute the best state sequence probability using Viterbi
 - Select the word for which this value is highest

Issues

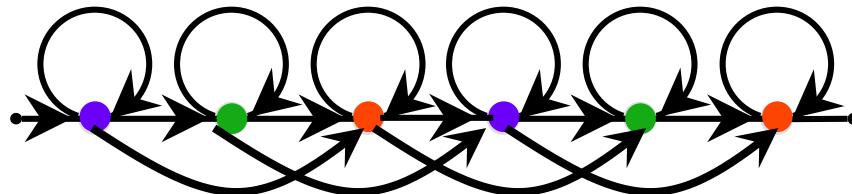
- What is the topology to use for the HMMs
 - How many states
 - What kind of transition structure
 - If state output densities have Gaussian Mixtures:
how many Gaussians?

HMM Topology

- For speech a left-to-right topology works best
 - The “Bakis” topology
 - Note that the initial state probability $P(s)$ is 1 for the 1st state and 0 for others. This need not be *learned*



- States may be skipped

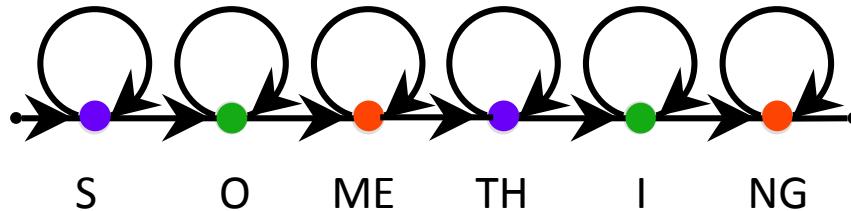


Determining the Number of States

- How do we know the number of states to use for any word?
 - We do not, really
 - Ideally there should be at least one state for each “basic sound” within the word
 - Otherwise widely differing sounds may be collapsed into one state
 - The average feature vector for that state would be a poor representation
- For computational efficiency, the number of states should be small
 - These two are conflicting requirements, usually solved by making some educated guesses

Determining the Number of States

- For small vocabularies, it is possible to examine each word in detail and arrive at reasonable numbers:



- For larger vocabularies, we may be forced to rely on some *ad hoc* principles
 - E.g.* proportional to the number of letters in the word
 - Works better for some languages than others
 - Spanish and Indian languages are good examples where this works as almost every letter in a word produces a sound

How many Gaussians

- No clear answer for this either
- The number of Gaussians is usually a function of the amount of training data available
 - Often set by trial and error
 - A minimum of 4 Gaussians is usually required for reasonable recognition

Implementation of BW: initialization of alphas and betas

- Initialization for alpha: $\alpha_u(s, 1)$ set to 0 for all states except the first state of the model. $\alpha_u(s, 1)$ set to 1 for the first state
 - All observations *must* begin at the first state
- Initialization for beta: $\beta_u(s, T)$ set to 0 for all states except the terminating state. $\beta_u(s, t)$ set to 1 for this state
 - All observations *must* terminate at the final state

Initializing State Output Density Parameters

1. Initially only a single Gaussian per state assumed
 - Mixtures obtained by splitting Gaussians
2. For Bakis-topology HMMs, a good initialization is the “flat” initialization
 - Compute the *global* mean and variance of all feature vectors in all training instances of the word
 - Initialize *all Gaussians* (i.e all state output distributions) with this mean and variance
 - Their means and variances will converge to appropriate values automatically with iteration
 - Gaussian splitting to compute Gaussian mixtures takes care of the rest

Isolated word recognition: Final thoughts

- All relevant topics covered
 - How to compute features from recordings of the words
 - We will not explicitly refer to feature computation in future lectures
 - How to set HMM topologies for the words
 - How to train HMMs for the words
 - Baum-Welch algorithm
 - How to select the most probable HMM for a test instance
 - Computing probabilities using the forward algorithm
 - Computing probabilities using the Viterbi algorithm
 - Which also gives the state segmentation

Questions

- ?