

Algorithms for speech and language processing

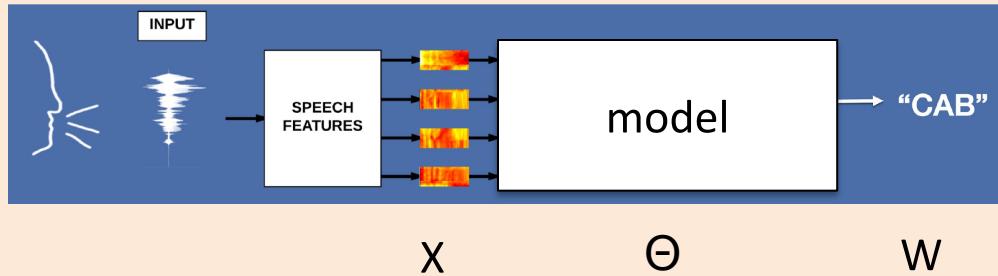
ASR #2: Language Modeling

Course

I. Generative Models

II. Decoding

III. Learning



Notations

θ : model parameters

$X=x_1 \dots x_t$: observed speech frames

$W=w_1..w_n$: word transcriptions

Decoding: given observed X and param. θ ,
find the most likely W .

$$W^* = \arg \max_w p(W|X, \theta)$$

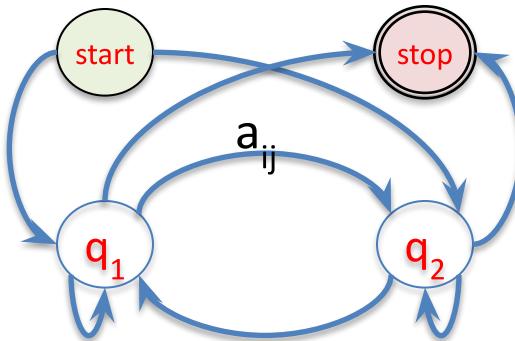
Learning: given an observed X and
model, find the most likely θ

$$\theta^* = \arg \max_{\theta} p(X|\theta)$$

Part I: Generative Models

- Useful finite state models (recap)
- Model decomposition
 - acoustic model
 - pronunciation model
 - language model
- Putting it all together
- Alternative models (RNNs)

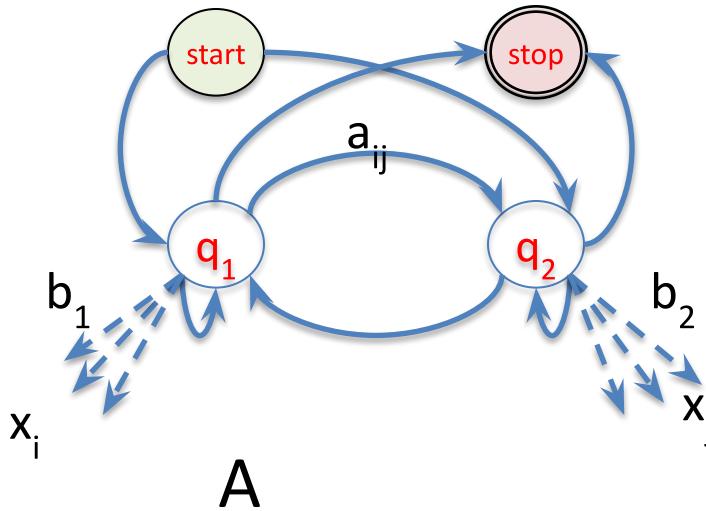
markov chains



A

	q_0 (start)	q_1	q_2	q_F (stop)
q_0 (start)	0	a_{01}	a_{02}	0
q_1	0	a_{11}	a_{12}	a_{1F}
q_2	0	a_{21}	a_{22}	a_{2F}
q_F (stop)	0	0	0	0

hidden markov models (HMMs)



A

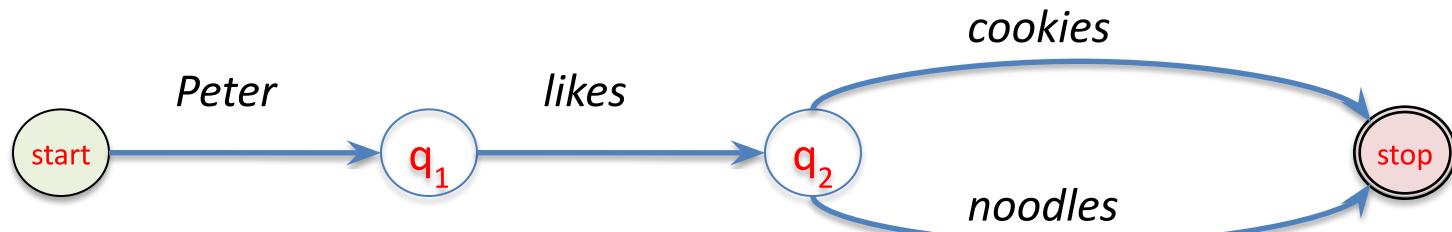
	q_0 (start)	q_1	q_2	q_F (stop)
q_0 (start)	0	a_{01}	a_{02}	0
q_1	0	a_{11}	a_{12}	a_{1F}
q_2	0	a_{21}	a_{22}	a_{2F}
q_F (stop)	0	0	0	0

B

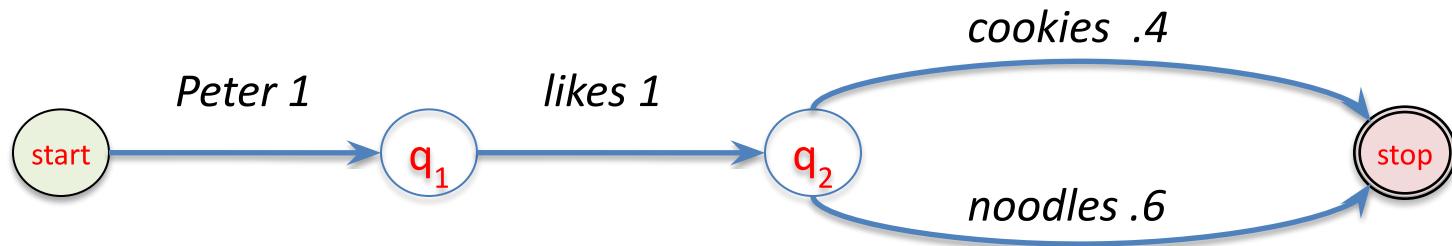
	q_1	q_2
x_1	b_{11}	b_{21}
x_2	b_{12}	a_{22}
x_3	b_{13}	a_{23}

(weighted) finite state transducers

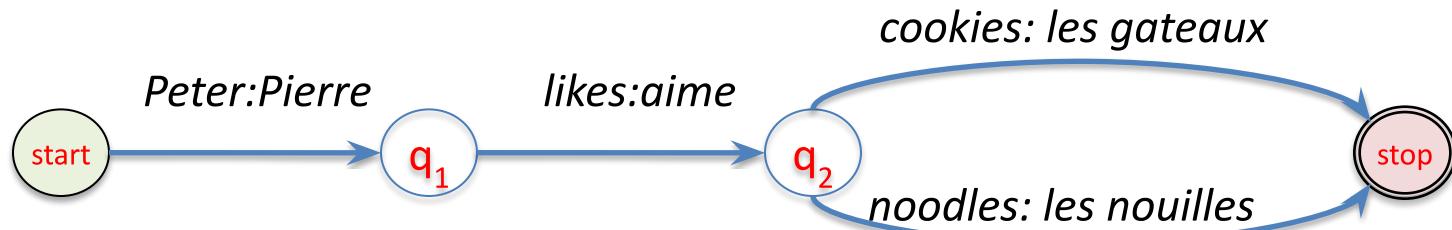
Finite State
Acceptor (FSA)



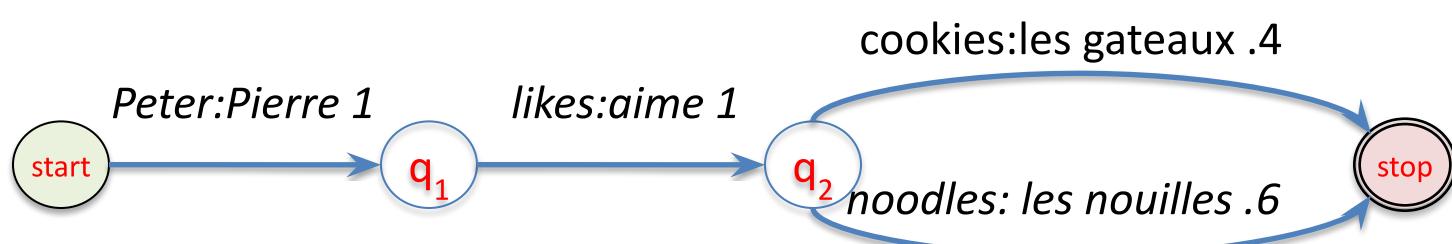
Weighted Finite
State Acceptor
(WFSA)



Finite State
Transducer (FST)

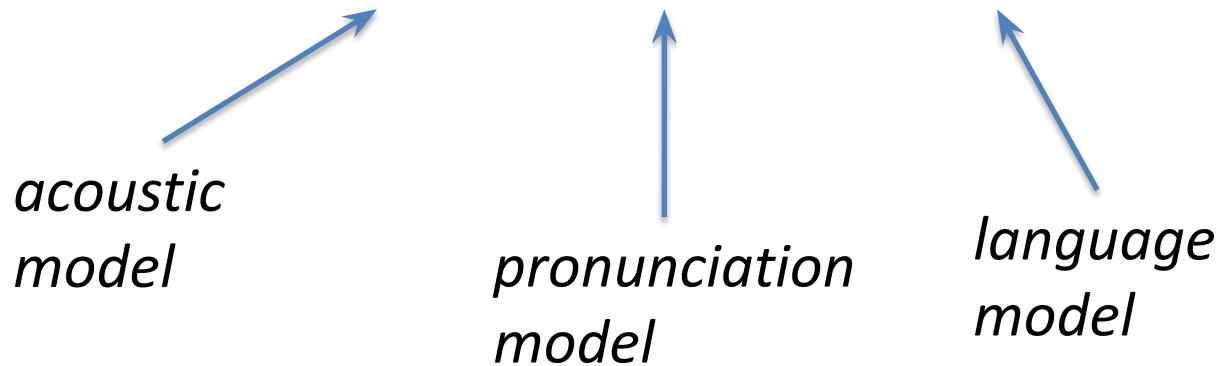


Weighted Finite
State Transducer
(WFST)



Model decomposition

- $W^* = \arg \max_w P_\theta(W|X)$
- $P_\theta(W|X) \propto P_\theta(X|W)P_\theta(W)$
- $= P_\theta(X|Q)P_\theta(Q|W)P_\theta(W)$

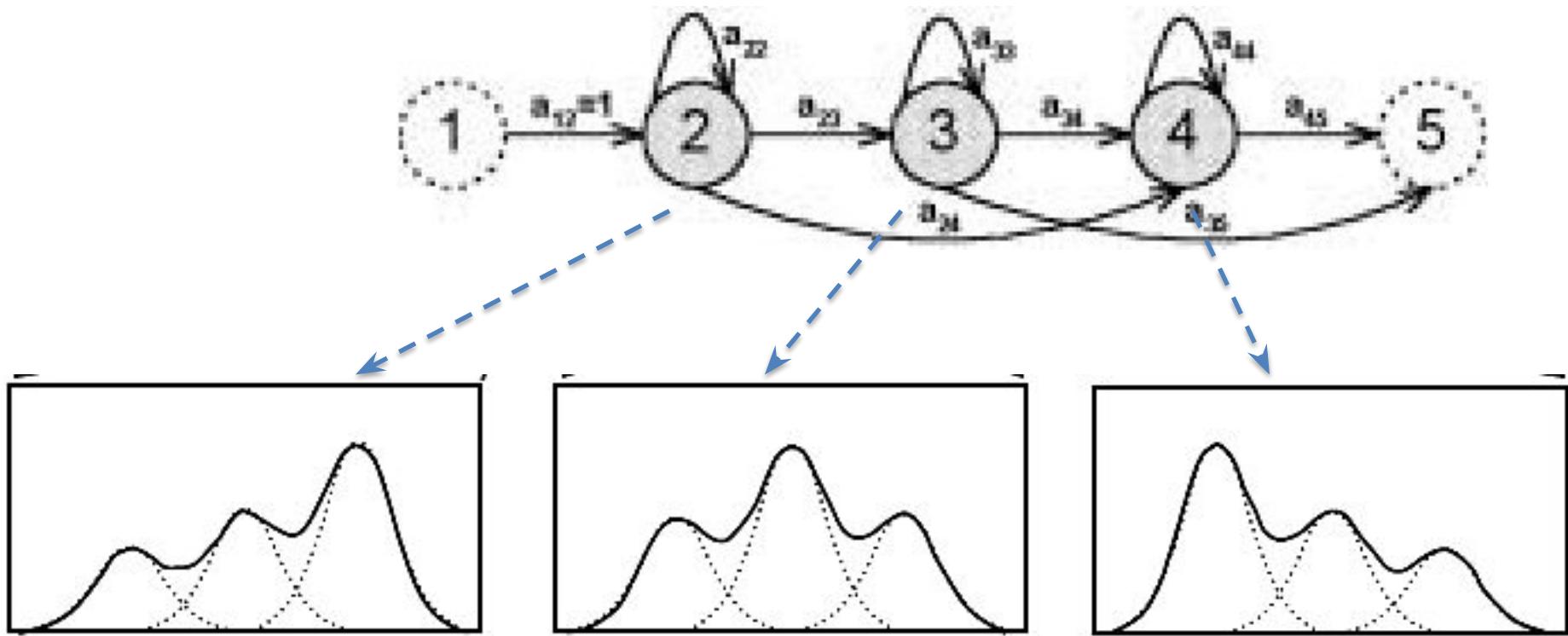


Acoustic models

$$P(X|Q)$$

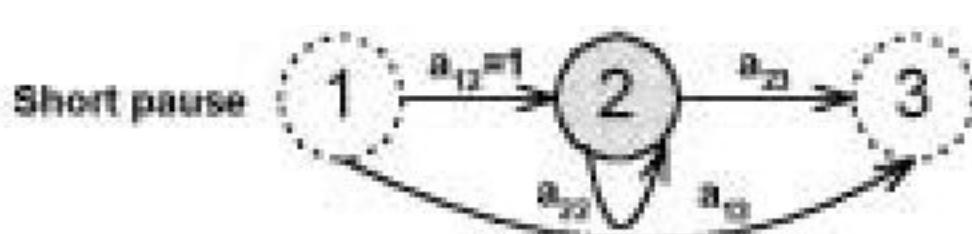
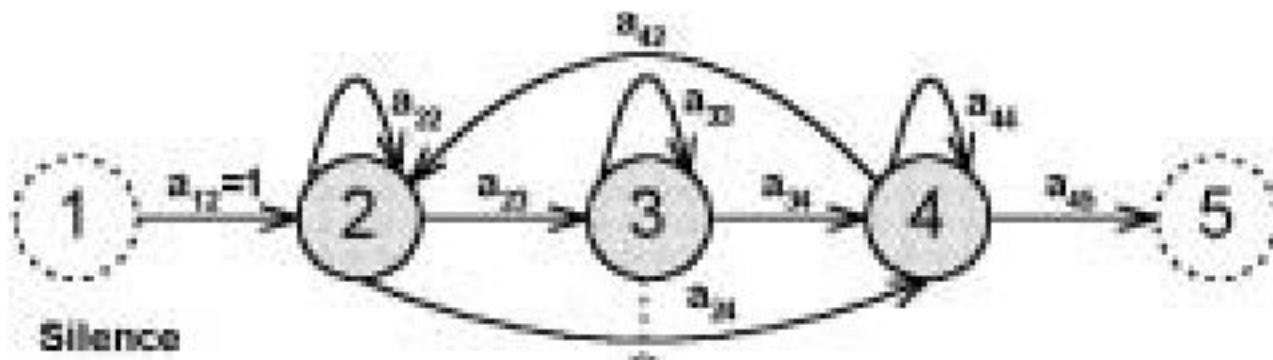
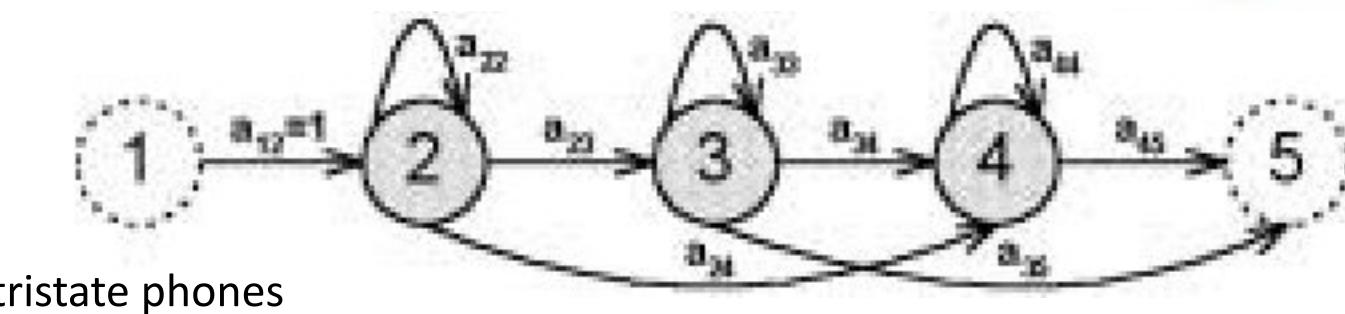


- tri-state phones



Acoustic models

$$P(X|Q)$$



Pronunciation models

$$P(Q|W)$$

Prononciations of the word 'eleven' in the state of Pennsylvania, USA

+ ^ lε > vən
+ ^ lεvən
+ ^ ɻεvən
θ ' lε > vən
+ ^ lεvən
əlεvən
θ ' ɻεvn
əɻεbm
+ ^ ' lεvn
lεvən
əɻε > θ vən
ɻεvn
θ ɻε θ vən
+ ^ ' lεvn
...
(total= 220 variants)



<https://www.youtube.com/watch?v=J3IYLphzAnw>

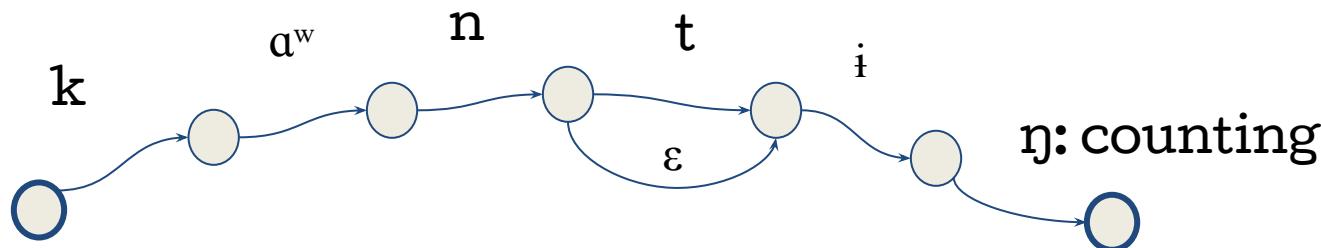
Phonological/Phonetic transformations

- | | |
|------------------------|---------------|
| What are you doing? | 'wʌtʃə'duɪn |
| I can inquire. | 'aɪkɳ'kwaiə |
| Did you eat yet? | 'dʒitjɛ? |
| I don't believe him. | ai'doðba'livm |
| We ought to have come. | wi'ɔf̩y'kʌm |

Pronunciation models

$P(Q|W)$

orthography	SAMPA	IPA
counting	kA_wntlN	kə ^w ntiŋ
counting(1)	kA_wnlN	kə ^w nɪŋ
amortization	@mOrt@zeS@n	əmɔrtəzeʃən
amortization(1)	{mOrt@zeS@n	æmɔrtəzeʃən
amortization(2)	@mOrtA_jzeS@n	əmɔrtə <i>ʒ</i> eʃən
..etc		



Problem: what about OOVs (out-of-vocabulary)?

2. Using G2P to extend the dictionary

- eg: g2p-seq2seq (CMU), phonetisaurus
 - align graphemes and phonemes
 - build an ngram model
 - make a WFST from it
 - iterate

“speaking” = s p ea k ing
[spi:kɪŋ] [s] [p] [i:] [k] [ɪŋ]

ବାସ୍ ⇒ c a: h

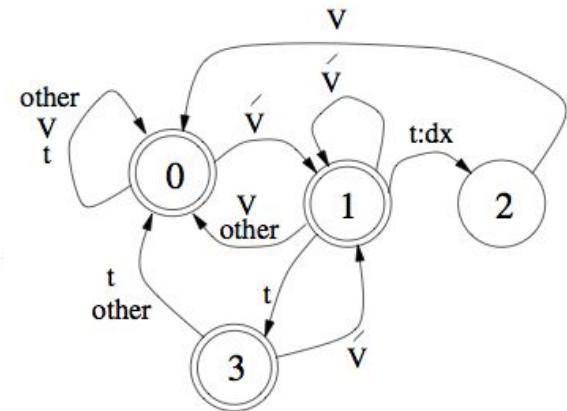
ପ୍ରଯୁକ୍ତ୍ତ ⇒ b r a: j u t

Problem: what about pronunciations that are not listed in the dictionary?

3. phonological rules (hand crafted)

- ex: english flapping

$/t/ \rightarrow [dx] / [+vowel] - \begin{bmatrix} +vowel \\ -stress \end{bmatrix}$



Input symbols: ax t eh n y uw ey t ax d

State sequence: 0 → 0 → 0 → 1 → 0 → 0 → 0 → 1 → 2 → 0 → 0

Output symbols: ax t eh n y uw ey dx ax d

Name	Rule	Example	Prob
Syllabic Rules*			
Syllabic n	[ax ix] n → en	button	.35
Syllabic m	[ax ix] m → em	bottom	.32
Syllabic l	[ax ix] l → el	bottle	.72
Syllabic r	[ax ix] r → axr	butter	.77
Flapping	[tcl dcl] [t d] → dx /V ____ [ax ix axr]	button	.87
Flapping-r	[tcl dcl] [t d] → dx /V r ____ [ax ix axr]	barter	.92
H-voicing	hh → hv / [+voice] ____ [+voice]	ahead	.92
L-deletion	l → Ø/ ____ y [ax ix axr]	million	n/a
Gliding	iy → y / ____ [ax ix axr]	colonial	n/a
Nasal-deletion	[n m ng] → Ø/ ____ [-voice -consonant]	rant	n/a
Function words			
h-deletion	h → Ø/ # ____	he, him	n/a
w-deletion	w → Ø/ # ____	will, would	n/a
dh-deletion	dh → Ø/ # ____	this, those	n/a
Dental-deletion	[tcl dcl] [t d] → Ø/ [+vowel] ____ [th dh]	breadth	n/a
Final dental-deletion	([tcl dcl]) [t d] → Ø/ [+cons +continuant] ____ # soft (as)		n/a
Slur	ax → Ø/ [+consonant] ____ [r l n] [+vowel]	camera	n/a
Stressed slur	[+vowel +stress] r → er	warts	n/a
Pre-stress contraction	ax → Ø/ [+cons] ____ [+cons] [+vowel +stress]	senility	n/a
Ruh-reduction	r ax → er / [-word bdry] ____ [-word bdry]	separable	n/a
Transitional stops			
t-introduction	Ø → tcl / [+dental +nasal] ____ [+fricative]	prin[t]ce	n/a
t-deletion	[tcl] → Ø/ [+dental +nasal] ____ [+fricative]	prints	n/a

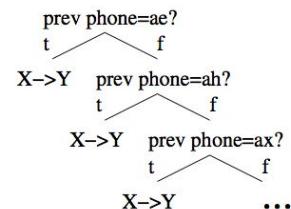
4. phonological rules (data driven)

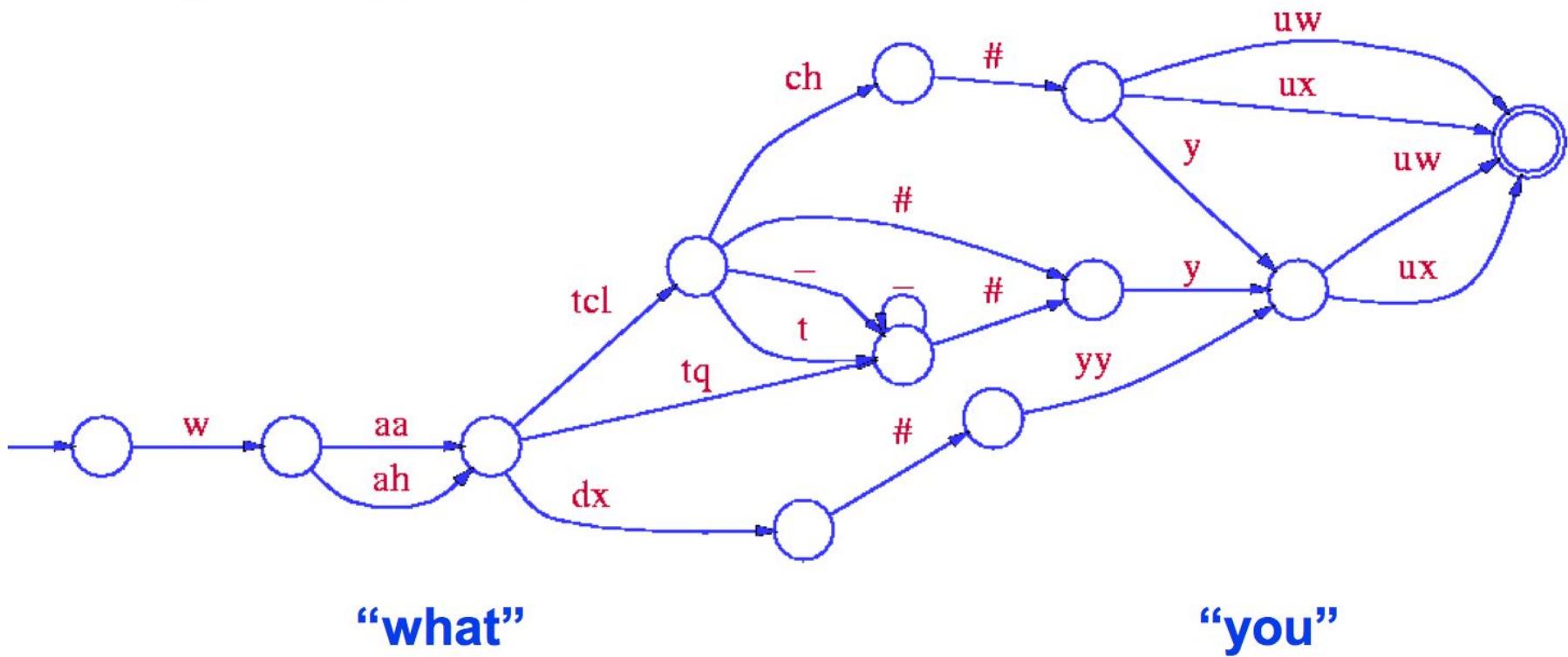
- using a phone recognizer to learn a pronunciation dictionary

a	baseball	game	word sequence
ax	b ey s b ao l	g ey m	canonical phone sequence
ax	b eh ey s b el	g eh m	aligned surface phones
ax	b eh ey s b el	g eh m	proposed new pronunciations

- inducing pronunciation rules
 - trigram models
 - decision trees

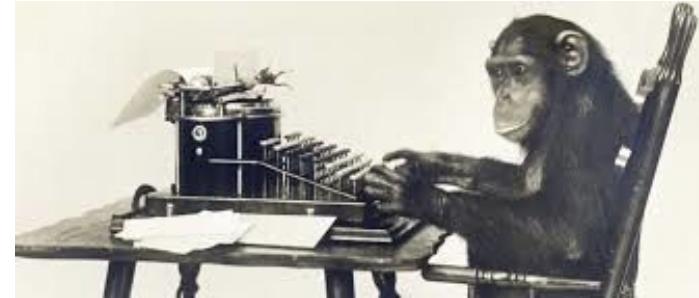
$$P(S^i) = \prod_{j=1}^n (P(s_j^i | c_{j-1}, c_j, c_{j+1})) .$$





Language models

P(W)



$$p(w_1..w_n) = p(w_1)p(w_2|w_1) \dots p(w_n|w_1..w_{n-1})$$

- N-gram models
 - truncate $p(w_i|w_1..w_{i-N+1}..w_{i-1})$ to $p(w_i|w_{i-N+1}..w_{i-1})$
- 1-gram: $p(w_i)$
- 2-gram: $p(w_i|w_{i-1})$
- 3-gram: $p(w_i|w_{i-2}w_{i-1})$
- etc.

N-1

1 - To him swallowed confess hear both. Which. Of save on trail
for are ay device and rote life have

Gram - Hill he late speaks; or! a more to leg less first you enter

2 - Why dost stand forth thy canopy, forsooth; he is this palpable
hit the King Henry. Live king. Follow.

Gram - What means, sir. I confess she? then all sorts, he is trim,
captain.

3 - Fly, and will rid me these news of price. Therefore the sadness
of parting, as they say, 'tis done.

Gram - This shall forbid it should be branded, if renown made it
empty.

4 - King Henry. What! I will go seek the traitor Gloucester. Exeunt
some of the watch. A great banquet serv'd in;

Gram - It cannot be but so.

- estimating Ngram models from data
 - counts

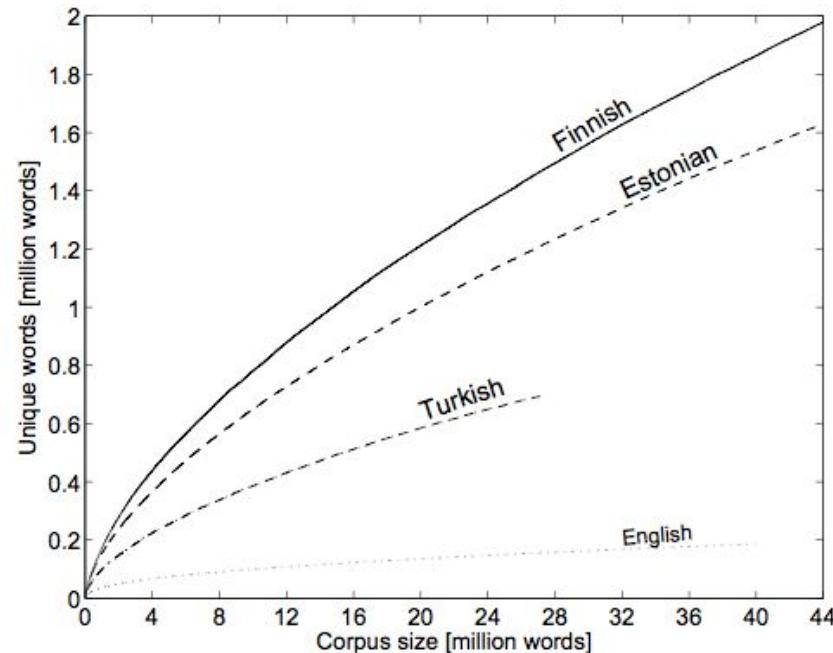
$$P(w_i) = \frac{c_i}{N}$$

- estimating Ngram models from data

- counts

$$P(w_i) = \frac{c_i}{N}$$

- unknown words (OOVs)
 - the lexicon is infinite!
 - a new corpus will always have new words
 - a fast and dirty solution: recoding with <UNK>



	i	want	to	eat	chinese	food	lunch	spend
i	5	827	0	9	0	0	0	2
want	2	0	608	1	6	6	5	1
to	2	0	4	686	2	0	6	211
eat	0	0	2	0	16	2	42	0
chinese	1	0	0	0	0	82	1	0
food	15	0	15	0	1	4	0	0
lunch	2	0	0	0	0	1	0	0
spend	1	0	1	0	0	0	0	0

- unobserved ngrams
 - most ngrams have been seen 0 time
 - problem grows exponentially with n
 - true 0s versus sampling
 - solutions:
 - smoothing
 - backoff, interpolation,
 - clustering

see Jurafsky & Martin 2017

$$\hat{P}(w_n | w_{n-2} w_{n-1}) = \lambda_1 P(w_n | w_{n-2} w_{n-1}) + \lambda_2 P(w_n | w_{n-1}) + \lambda_3 P(w_n)$$

- evaluating language models
 - extrinsic: improved WER
 - intrinsic: perplexity

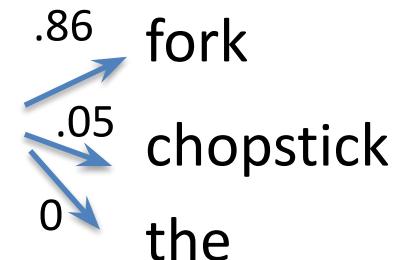
$$PP(W) = P(W)^{-1/N}$$

$$PP(W) = \prod_1^N P(w_i | w_1..w_{i-1})^{-1/N}$$

$$PP(W) = 2^{H(W)}$$

with $H(W) = -1/N \log_2 P(W)$

He is eating a steak with a knife and a ...



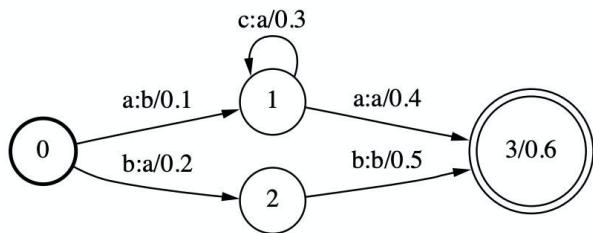
Putting all together

- acoustic model (an wFST): A
- pronunciation model (an wFST): P
- language model (an wFST): G

the combined model (a large wFST!): $A \square P \square G$

wFST composition

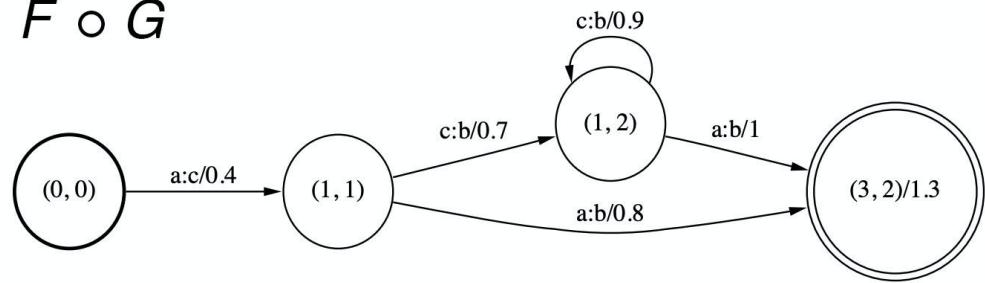
F



a a → b a

a c a → b a a

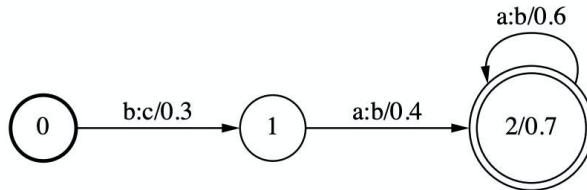
$F \circ G$



a a → c b

a c a → c b b

G



b a → c b

b a a → c b b

from:

<https://jonathan-hui.medium.com/speech-recognition-weighted-finite-state-transducers-wfst-a4ece08a89b7>

Putting all together

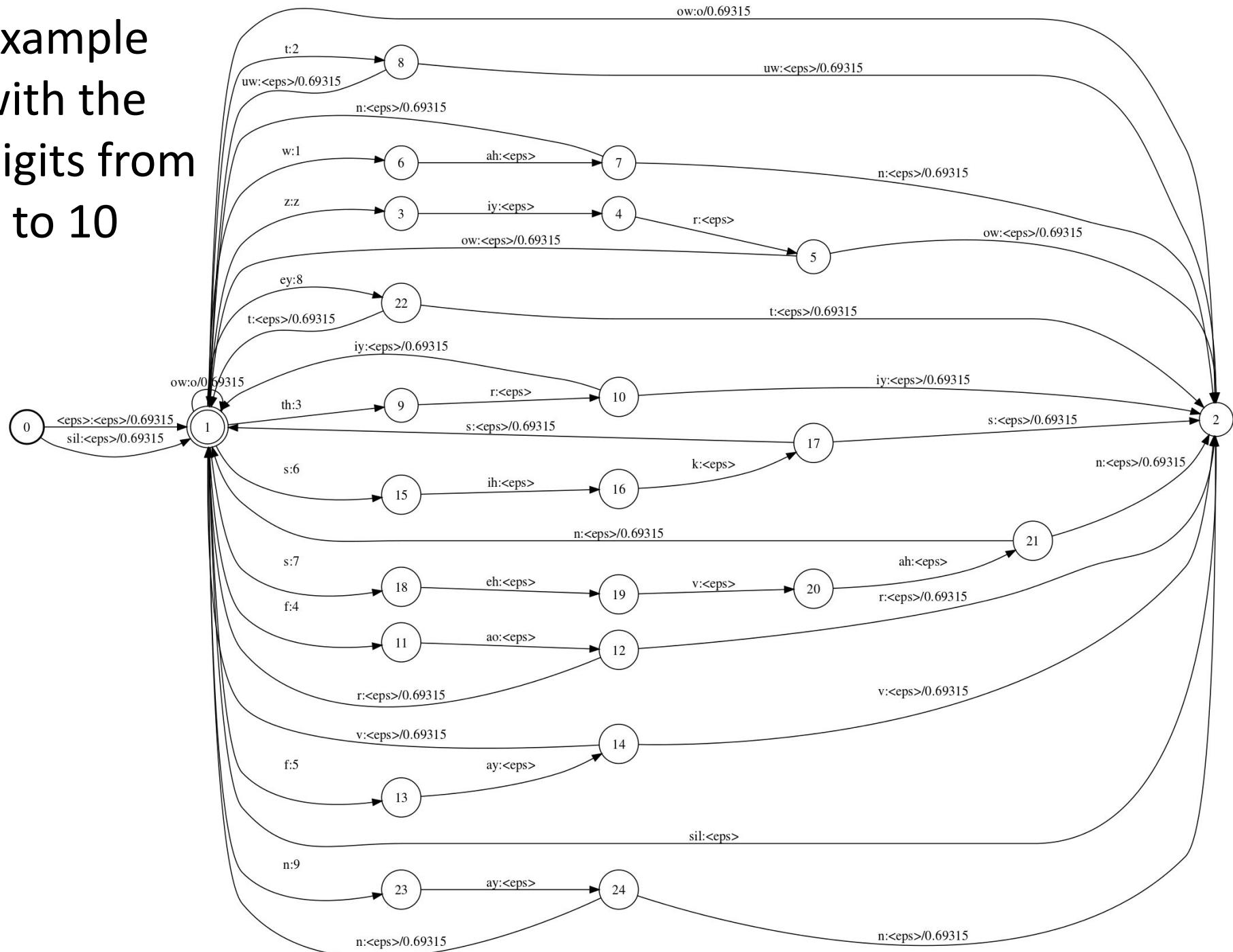
- acoustic model (an wFST): A
- pronunciation model (an wFST): P
- language model (an wFST): G

the combined model (a large wFST!): $A \square P \square G$

→ applying the wFST to acoustic input: a large decoding graph.

→ task: finding the optimal path

Example with the digits from 0 to 10



in practice

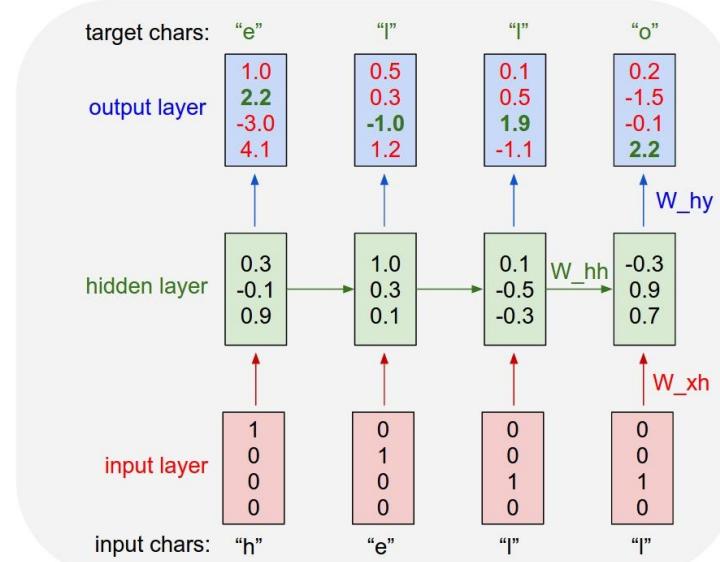
- FST tools
 - openFST (<http://www.openfst.org>)
 - composing & simplifying FSTs
- decomposing processing in two stages
 - a first decoding lattice, followed by lattice rescoring with a large LM

to know more

- Mohri, M. (1997). Finite-state transducers in language and speech processing. *Computational linguistics*, 23(2):269–311.
- Mohri, M., Pereira, F., and Riley, M. (2002). Weighted finite-state transducers in speech recognition. *Computer Speech & Language*, 16(1):69–88.

fast forward: RNN language models

- character language models



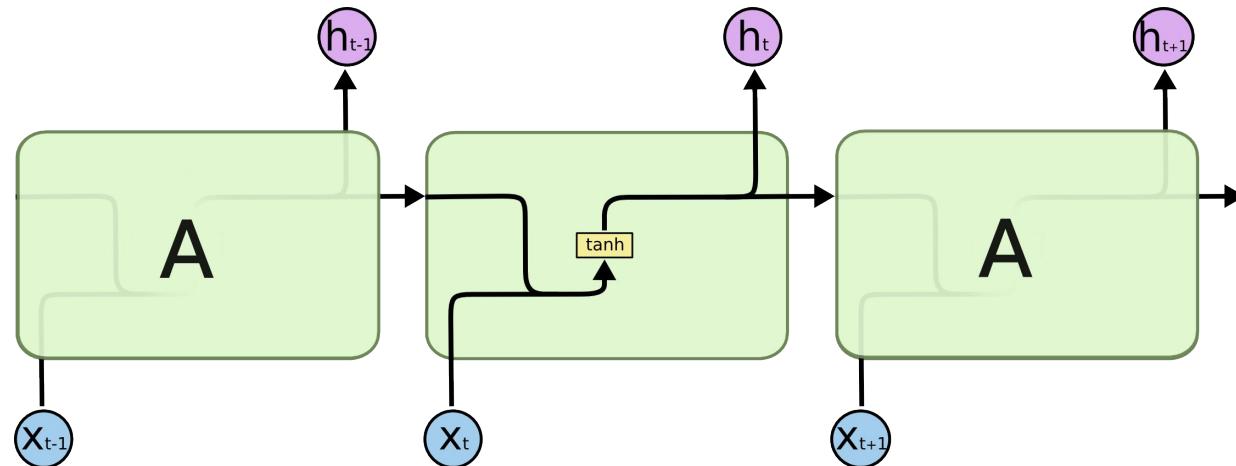
<http://karpathy.github.io/2015/05/21/rnn-effectiveness/>

- RNNs can approximate $P(W)$ better than ngram models by relying on indefinite amount of past information

$$p(w_1..w_n) = p(w_1)p(w_2|w_1) \dots p(w_n|w_1..w_{n-1})$$

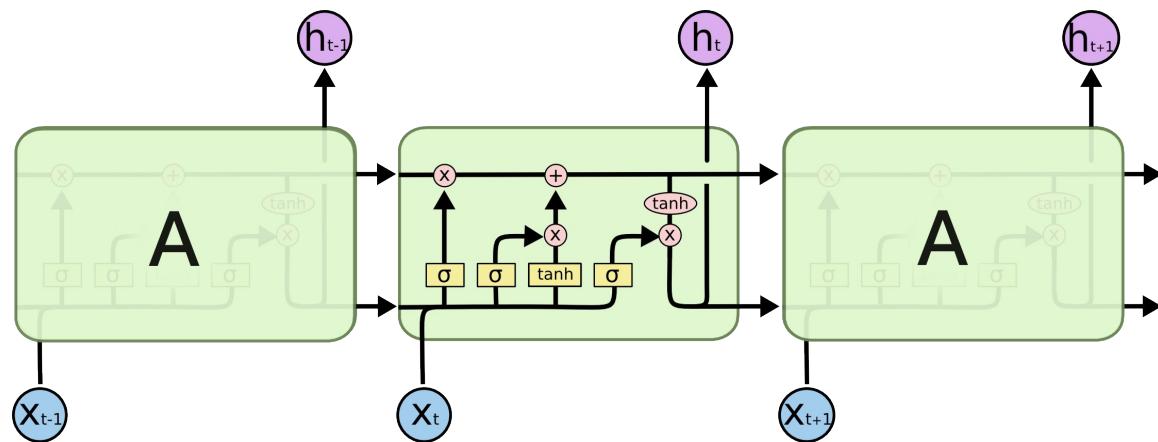
- but, the Markov assumption no longer holds

Vanilla RNN



$$h_t = \tanh(W_C[h_{t-1}, x_t] + b)$$

LSTM



$$f_t = \sigma(W_f \cdot [h_{t-1}, x_t] + b_f)$$

$$i_t = \sigma(W_i \cdot [h_{t-1}, x_t] + b_i)$$

$$\tilde{C}_t = \tanh(W_C \cdot [h_{t-1}, x_t] + b_C)$$

$$C_t = f_t * C_{t-1} + i_t * \tilde{C}_t$$

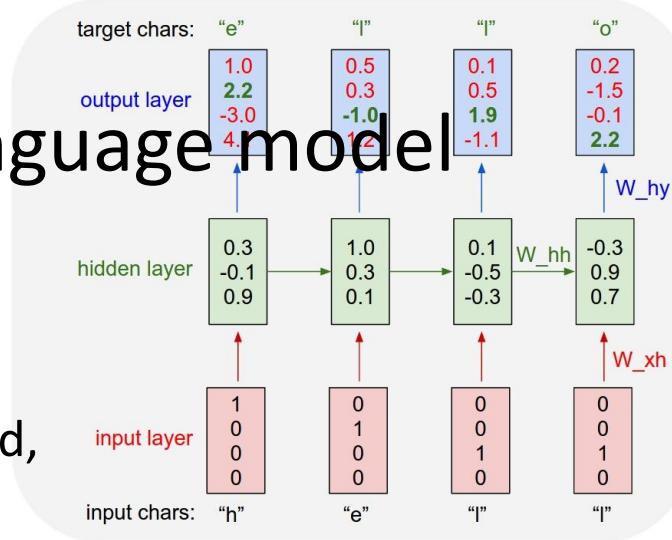
$$o_t = \sigma(W_o \cdot [h_{t-1}, x_t] + b_o)$$

$$h_t = o_t * \tanh(C_t)$$

- sample output of a character language model

PANDARUS:

Alas, I think he shall be come approached and the day
 When little strain would be attain'd into being never fed,
 And who is but a chain and subjects of his death,
 I should not sleep.



Second Senator:

They are away this miseries, produced upon my soul,
 Breaking and strongly should be buried, when I perish
 The earth and thoughts of many states.

DUKE VINCENTIO:

Well, your wit is in the care of side and that.

Second Lord: They would be ruled after this chamber, and
 my fair nues begun out of the fact, to be conveyed,
 Whose noble souls I'll have the heart of the wars.

- word RNNs
 - in practice word RNNs work better than char RNNs, but grow very large
 - Typically, the input layer is as large as the nb of words; a linear embedding layer is used to reduce dimensionality
- wordpiece RNNs
 - use Byte Pair Encoding (BPE) to define a lexicon of useful word parts

Part II. decoding

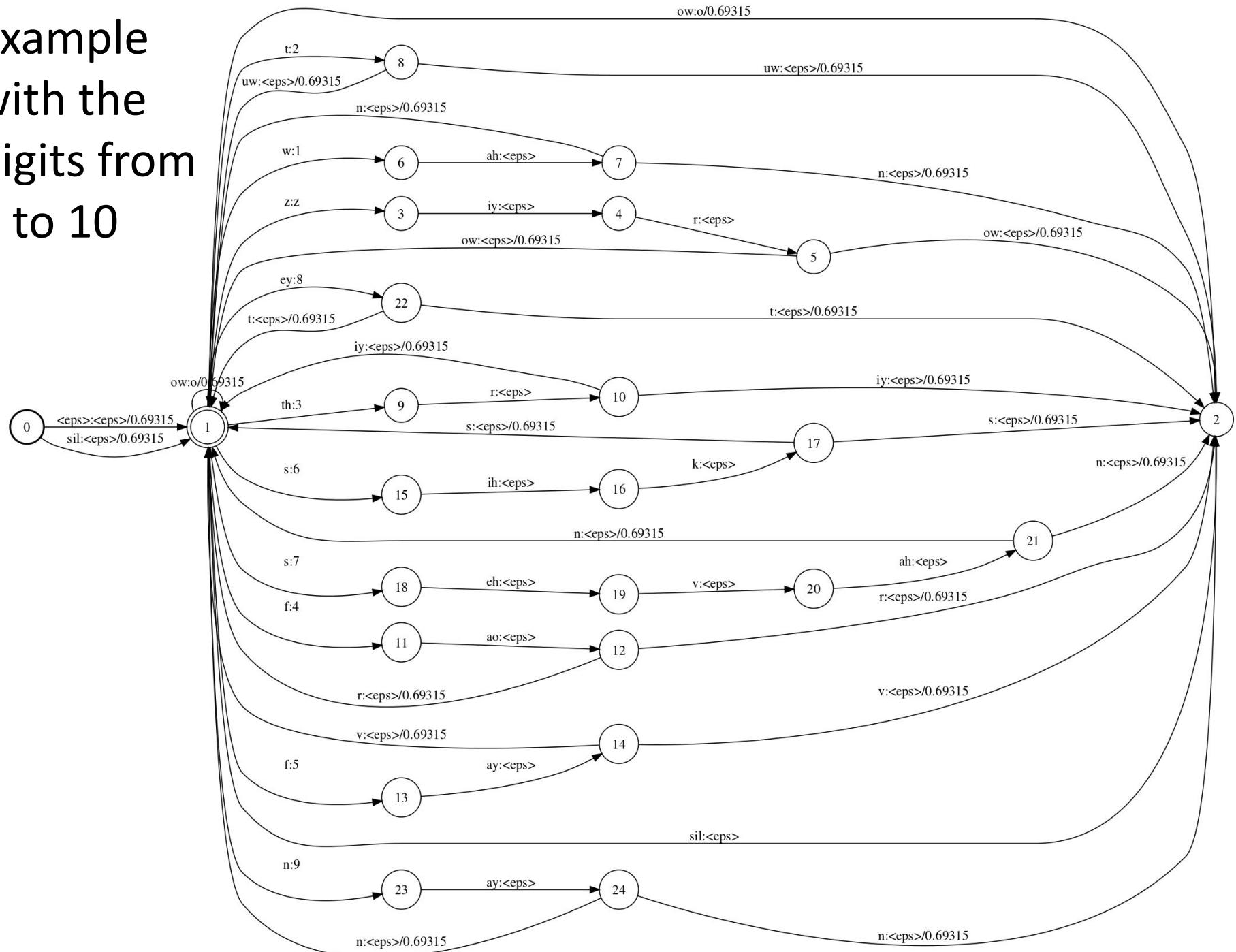
$$W^* = \operatorname{argmax}_w p_\theta(W|X)$$

$$W^* = \operatorname{argmax}_w p_\theta(X|W) \cdot p_\theta(W)$$

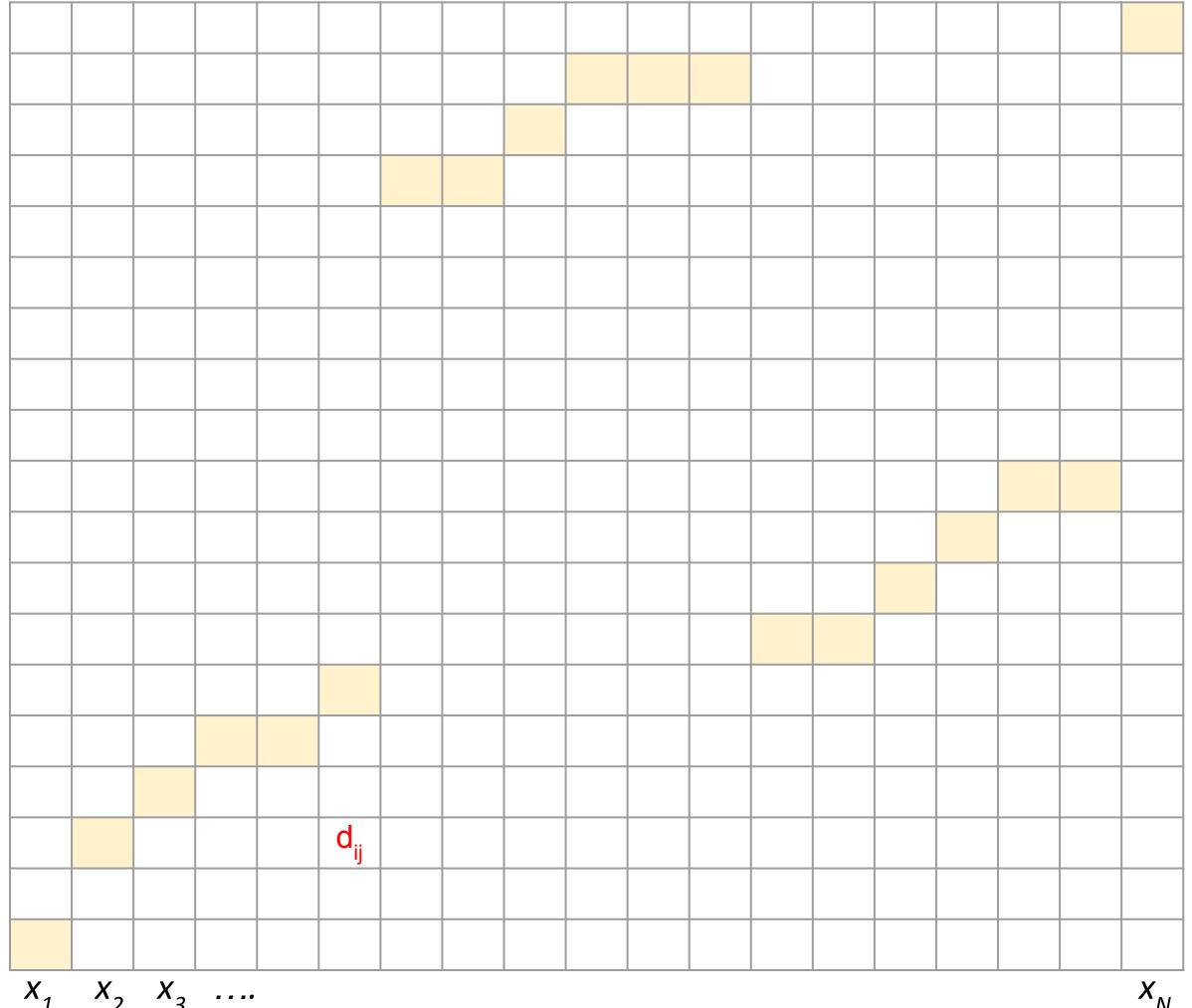
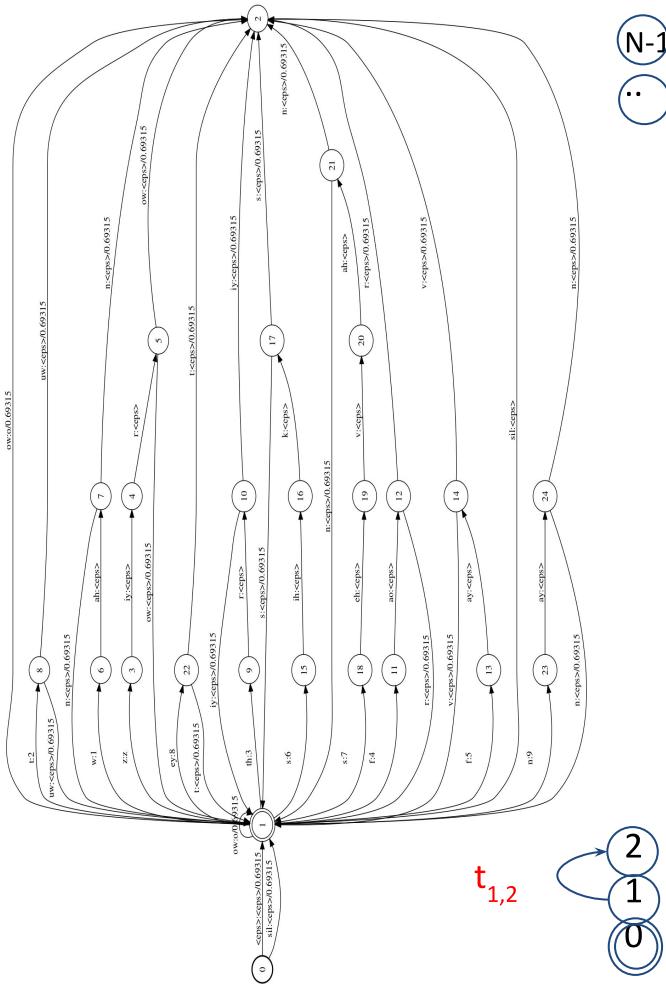
$$W^* = \operatorname{argmax}_w p_\theta(X|Q) \cdot p_\theta(Q|W) \cdot p_\theta(W)$$

- Heuristics
- Dynamic programming (a quick tour)
- Dynamic time warping
- HMM Decoding

Example with the digits from 0 to 10



K states



*minimizing $\prod_k d_{ik,jk} \cdot t_{ik-1, ik}$
number of possible paths: K^N !!*

Search heuristics

- greedy decoding
 - select the best state at each time step
 - efficient but not very good
- beam search
 - select the p best states at each time step
 - still efficient, and slightly better
- Dynamic programming/viterbi
 - efficient and exact solution !

Dynamic programming

1. define subproblems
2. identify recurrence
3. set the initial conditions

- Problem: given n , find the number of different ways to write n as the sum of 1, 3, 4
 - Example: for $n = 5$, the answer is 6

$$5 = 1+1+1+1+1$$

$$= 1+1+3$$

$$= 1+3+1$$

$$= 3+1+1$$

$$= 1+4$$

$$= 4+1$$

for $n=6$?

1. Define subproblems

- Let D_n be the number of ways to write n as the sum of 1,3,4

2. Find the recurrence

- Consider one possible solution $n=x_1+x_2+\dots+x_m$
 - If $x_m=1$, the rest of the terms must sum to $n-1$
 - Thus, the number of sums that end with $x_m=1$ is equal to D_{n-1}
 - if $x_m=3$, it is equal to D_{n-3}
 - if $x_m=4$, it is equal to D_{n-4}
- $D_n = D_{n-1} + D_{n-3} + D_{n-4}$

3. Set the initial conditions

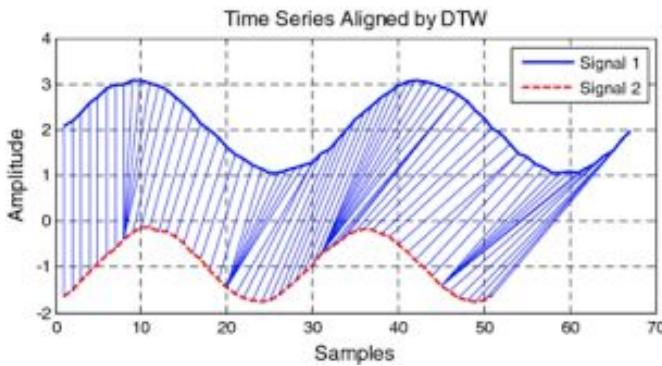
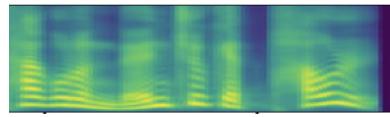
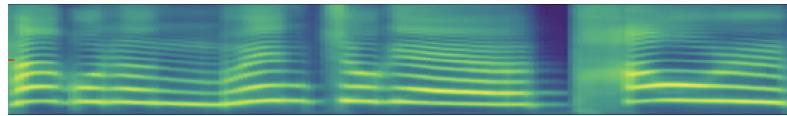
$$\begin{aligned}D[0] &= D[1] = D[2] = 1 \\D[3] &= 2\end{aligned}$$

4. Compute the solution

```
for(i = 4; i <= n; i++)  
    D[i] = D[i-1] + D[i-3] + D[i-4]
```

Done! (complexity: $O(n)$)

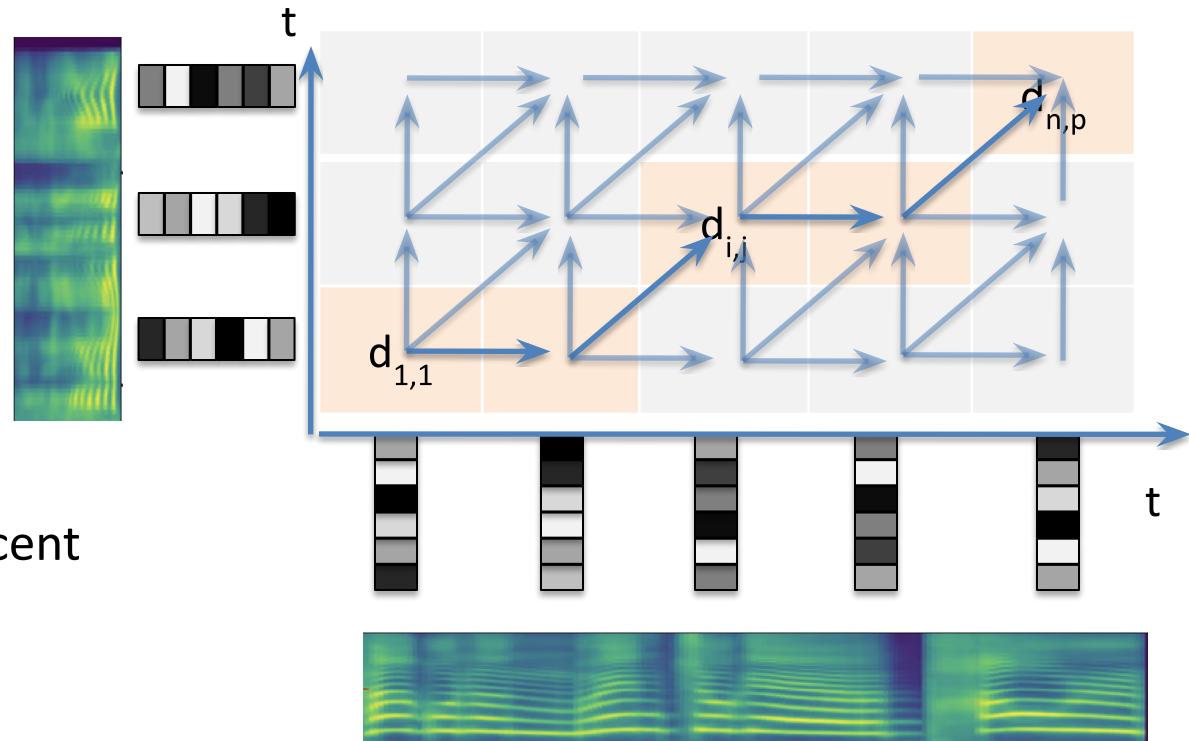
Application: Dynamic Time Warping



t

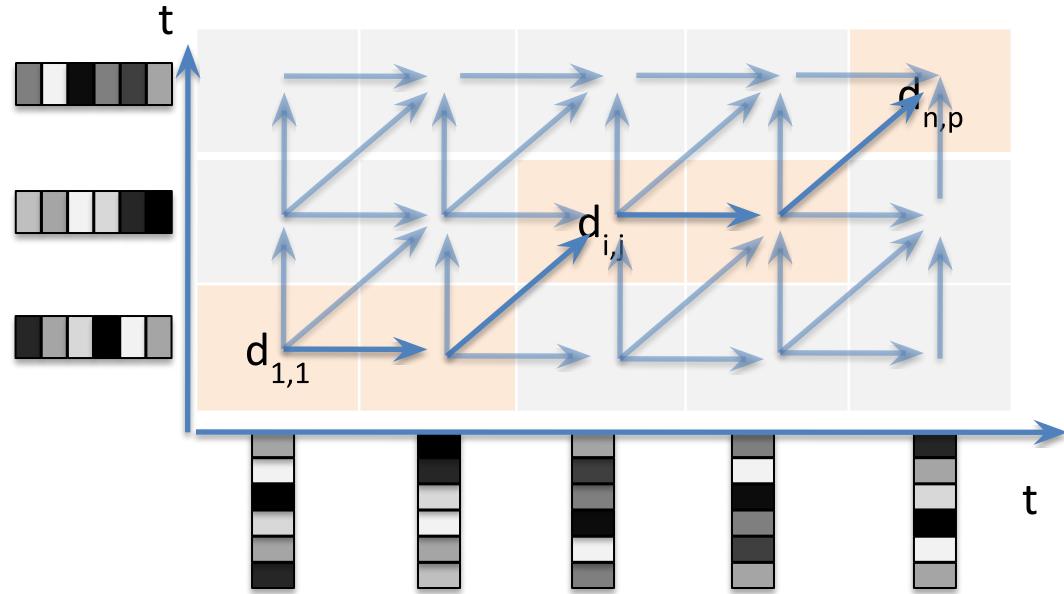
Application: Dynamic Time Warping

Find the *path* which minimizes cost: $\sum_{\text{path}} d_{i,j}$
(where the path can only go through increasing and adjacent cells)



Application: Dynamic Time Warping

Find the *path* which minimizes cost: $\sum_{\text{path}} d_{i,j}$ (where the path can only go through increasing and adjacent cells)



Dynamic programming analysis

subp: C_{ij} : cost minimal path i,j

recur: $C_{i,j} = d_{i,j} + \min(C_{i,j-1}, C_{i-1,j-1}, C_{i-1,j})$

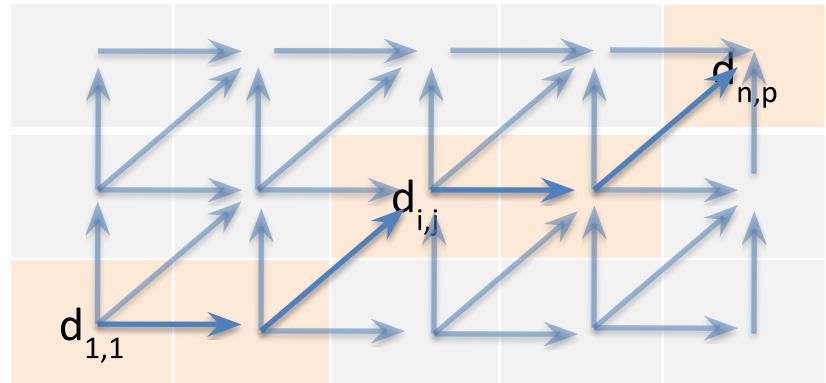
init: $C_{0,0} = 0$, $C_{i,0} = \infty$, $C_{0,j} = \infty$

($i > 0, j > 0$)

∞						
∞						
∞						
0	∞	∞	∞	∞	∞	∞

Diagram illustrating the dynamic programming step for calculating the cost $C_{i,j}$. The cost $C_{i-1,j-1}$ is added to the cost $d_{i,j}$ to get the cost $C_{i,j}$. The costs $C_{i,j-1}$ and $C_{i-1,j}$ are also shown as potential candidates for the minimum calculation.

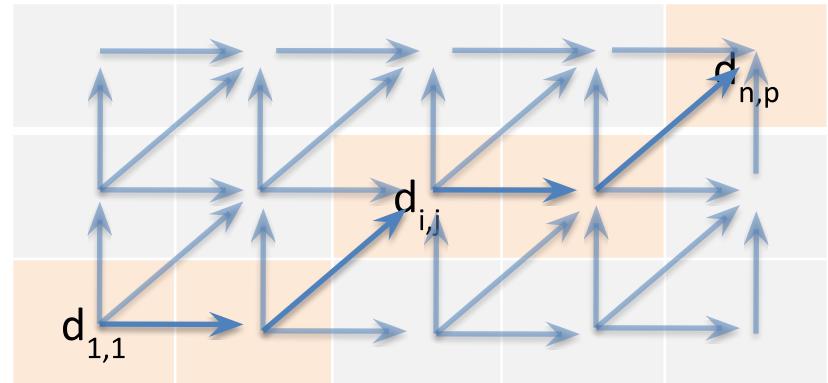
DTW algorithm



- Initialize

∞	0	0	0	0	0
∞	0	0	0	0	0
∞	0	0	0	0	0
0	∞	∞	∞	∞	∞

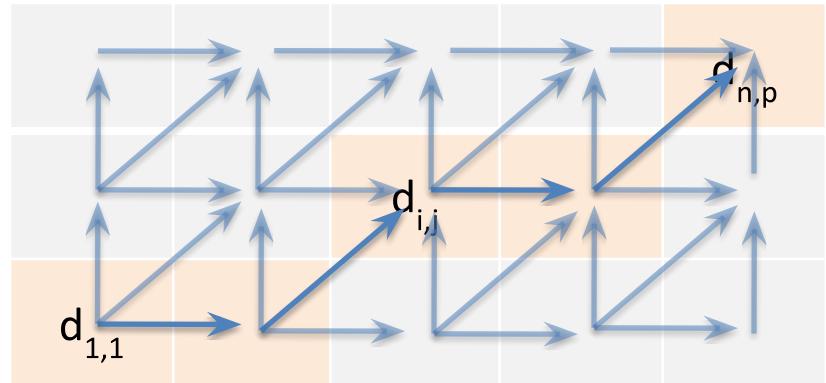
DTW algorithm



- Apply the recurrence formula
 $C_{i,j} = d_{i,j} + \min(C_{i,j-1}, C_{i-1,j-1}, C_{i-1,j})$
- Leave a pointer to which cell you are coming from

∞	0	0	0	0	0
∞	0	0	0	0	0
∞	$d_{1,1}$	0	0	0	0
0	∞	∞	∞	∞	∞

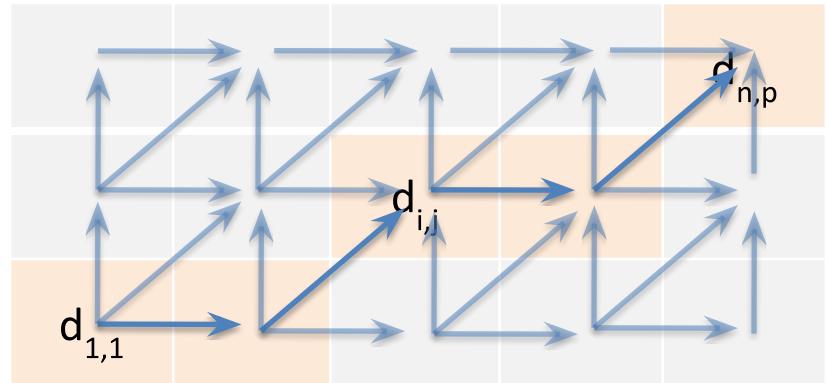
DTW algorithm



- Apply the recurrence formula
 $C_{i,j} = d_{i,j} + \min(C_{i,j-1}, C_{i-1,j-1}, C_{i-1,j})$
- Leave a pointer to which cell you are coming from

∞	0	0	0	0	0
∞	$d_{1,1}$ + $d_{2,1}$	0	0	0	0
∞	$d_{1,1}$	0	0	0	0
0	∞	∞	∞	∞	∞

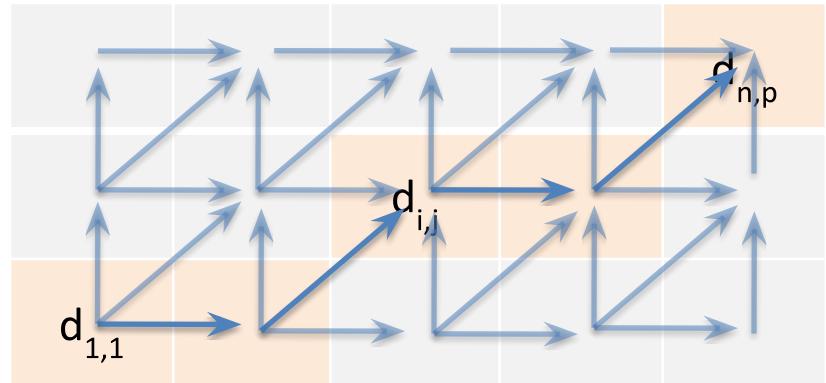
DTW algorithm



- Apply the recurrence formula
 $C_{i,j} = d_{i,j} + \min(C_{i,j-1}, C_{i-1,j-1}, C_{i-1,j})$
- Leave a pointer to which cell you are coming from

∞	$d_{1,1}$ + $d_{2,1}$ + $d_{3,1}$	0	0	0	0
∞	$d_{1,1}$ + $d_{2,1}$	0	0	0	0
∞	$d_{1,1}$	0	0	0	0
0	∞	∞	∞	∞	∞

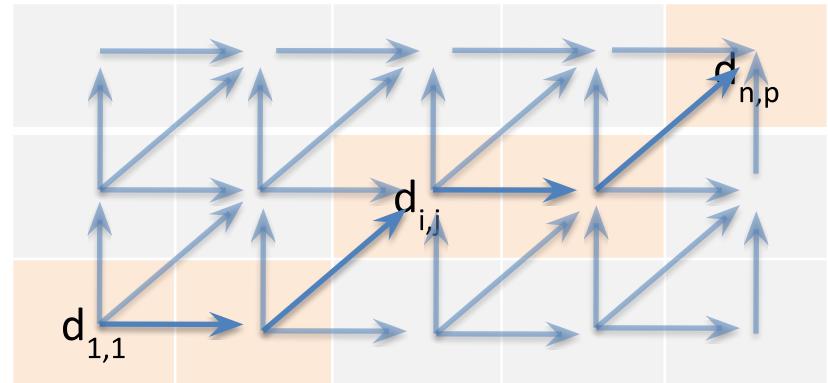
DTW algorithm



- Apply the recurrence formula
 $C_{i,j} = d_{i,j} + \min(C_{i,j-1}, C_{i-1,j-1}, C_{i-1,j})$
- Leave a link to which cell you are coming from

∞	$d_{1,1}$ + $d_{2,1}$ + $d_{3,1}$	0	0	0	0
∞	$d_{1,1}$ + $d_{2,1}$	0	0	0	0
∞	$d_{1,1}$	$d_{1,1}$ + $d_{1,2}$	0	0	0
0	∞	∞	∞	∞	∞

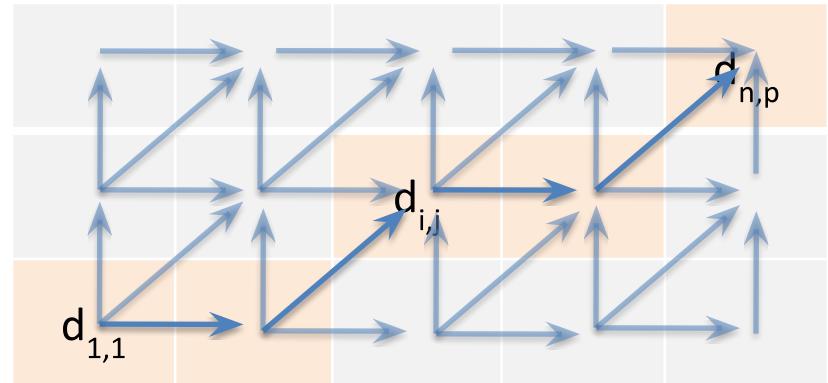
DTW algorithm



etc...

∞	$d_{1,1} + d_{2,1} + d_{3,1}$	$d_{1,1} + d_{2,1} + d_{3,2}$	$d_{1,1} + d_{2,1} + d_{2,2} + d_{3,3}$	$d_{1,1} + d_{2,1} + d_{2,2} + d_{3,3} + d_{3,4}$	$d_{1,1} + d_{1,2} + d_{2,3} + d_{2,4} + d_{3,5}$
∞	$d_{1,1} + d_{2,1}$	$d_{1,1} + d_{2,1} + d_{2,2}$	$d_{1,1} + d_{1,2} + d_{2,3}$	$d_{1,1} + d_{1,2} + d_{2,3} + d_{2,4}$	$d_{1,1} + d_{1,2} + d_{2,3} + d_{2,4} + d_{2,5}$
∞	$d_{1,1}$	$d_{1,1} + d_{1,2}$	$d_{1,1} + d_{1,2} + d_{1,3}$	$d_{1,1} + d_{1,2} + d_{1,3} + d_{1,4}$	$d_{1,1} + d_{1,2} + d_{1,3} + d_{1,4} + d_{1,5}$
0	∞	∞	∞	∞	∞

DTW algorithm

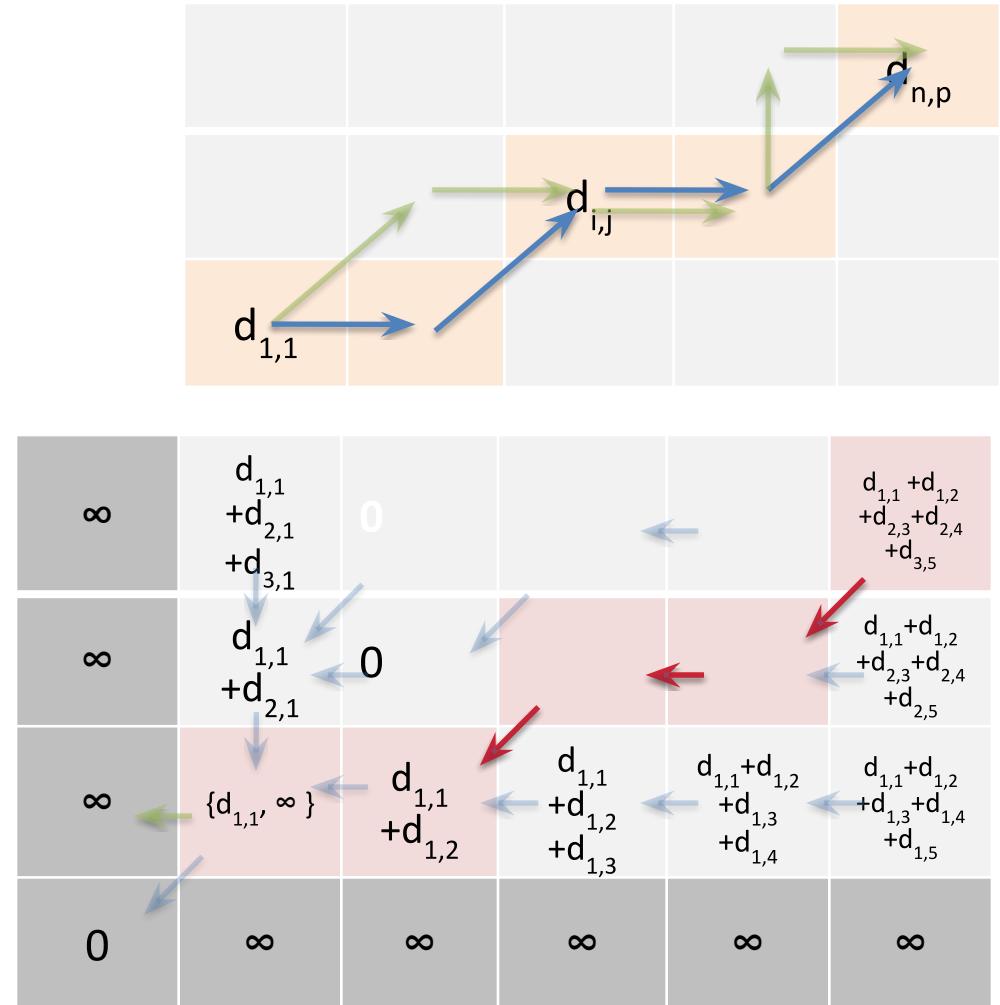


∞	$d_{1,1} + d_{2,1} + d_{3,1}$	$d_{1,1} + d_{2,1} + d_{3,2}$	$d_{1,1} + d_{2,1} + d_{2,2} + d_{3,3}$	$d_{1,1} + d_{2,1} + d_{2,2} + d_{3,3} + d_{3,4}$	$d_{1,1} + d_{2,1} + d_{2,2} + d_{3,3} + d_{3,4} + d_{3,5}$
∞	$d_{1,1} + d_{2,1}$	$d_{1,1} + d_{2,1} + d_{2,2}$	$d_{1,1} + d_{2,1} + d_{2,2} + d_{2,3}$	$d_{1,1} + d_{2,1} + d_{2,2} + d_{2,3} + d_{2,4}$	$d_{1,1} + d_{2,1} + d_{2,2} + d_{2,3} + d_{2,4} + d_{2,5}$
∞	$d_{1,1}$	$d_{1,1} + d_{1,2}$	$d_{1,1} + d_{1,2} + d_{1,3}$	$d_{1,1} + d_{1,2} + d_{1,3} + d_{1,4}$	$d_{1,1} + d_{1,2} + d_{1,3} + d_{1,4} + d_{1,5}$
0	∞	∞	∞	∞	∞

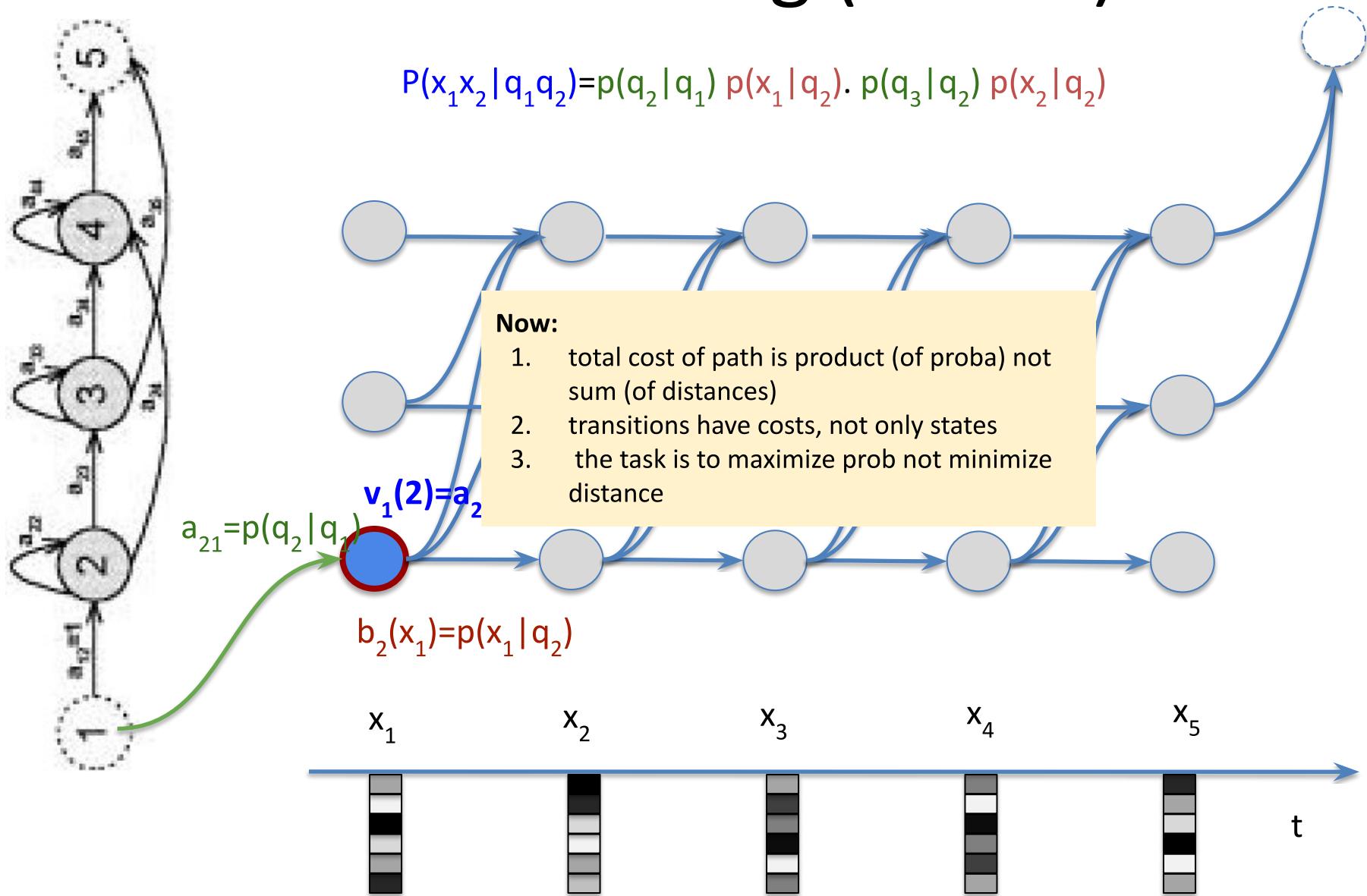
- Start from the final point
- Follow the path backward

Variation: n-best

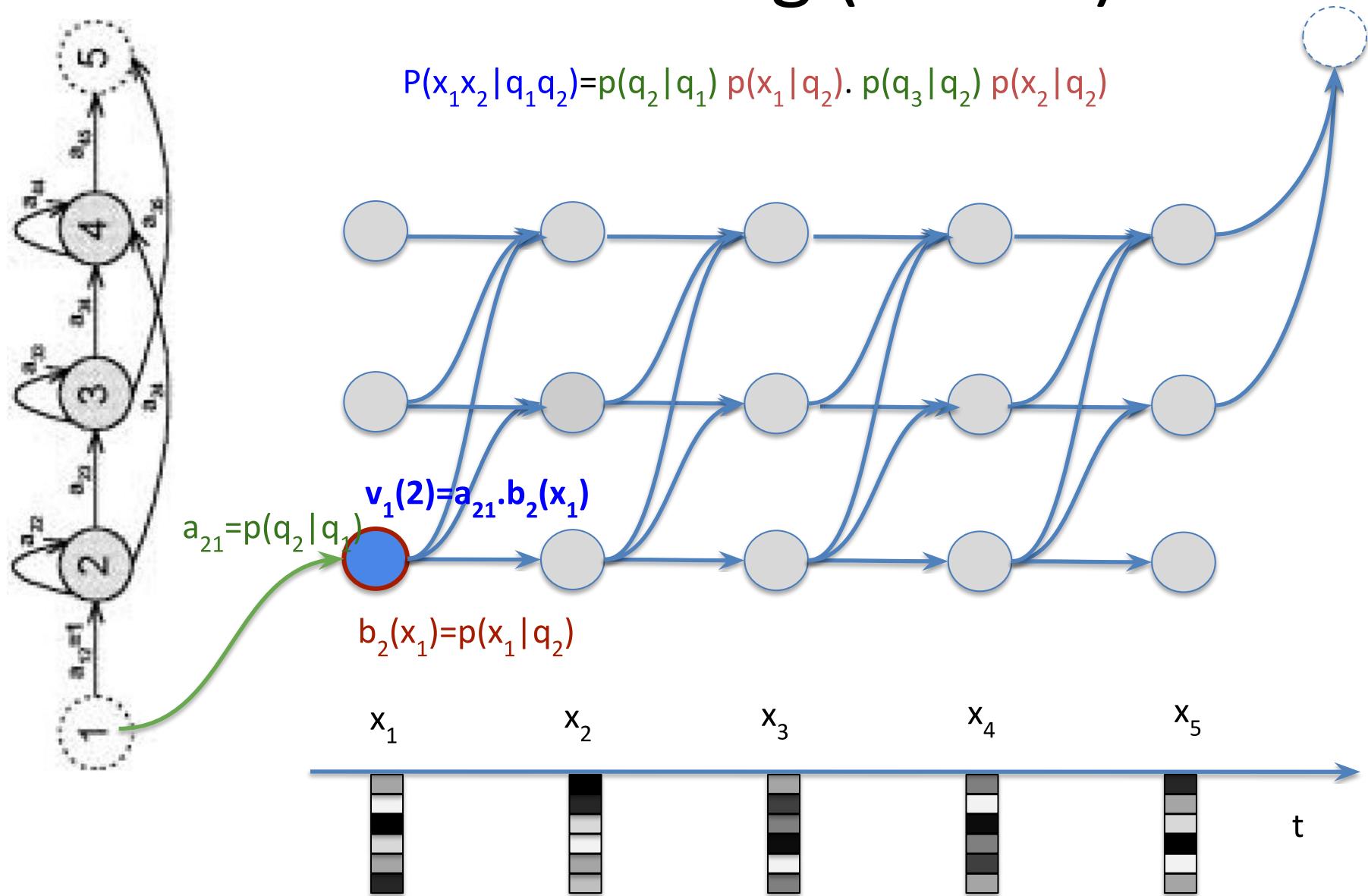
- Change the recurrence formula
 $C_{i,j} = d_{i,j} + \text{best}(n, C_{i,j-1} \cup C_{i-1,j-1} \cup C_{i-1,j})$
 where $\text{best}(n, S)$ =the n smallests elements in set S.
- Memorize n links backward



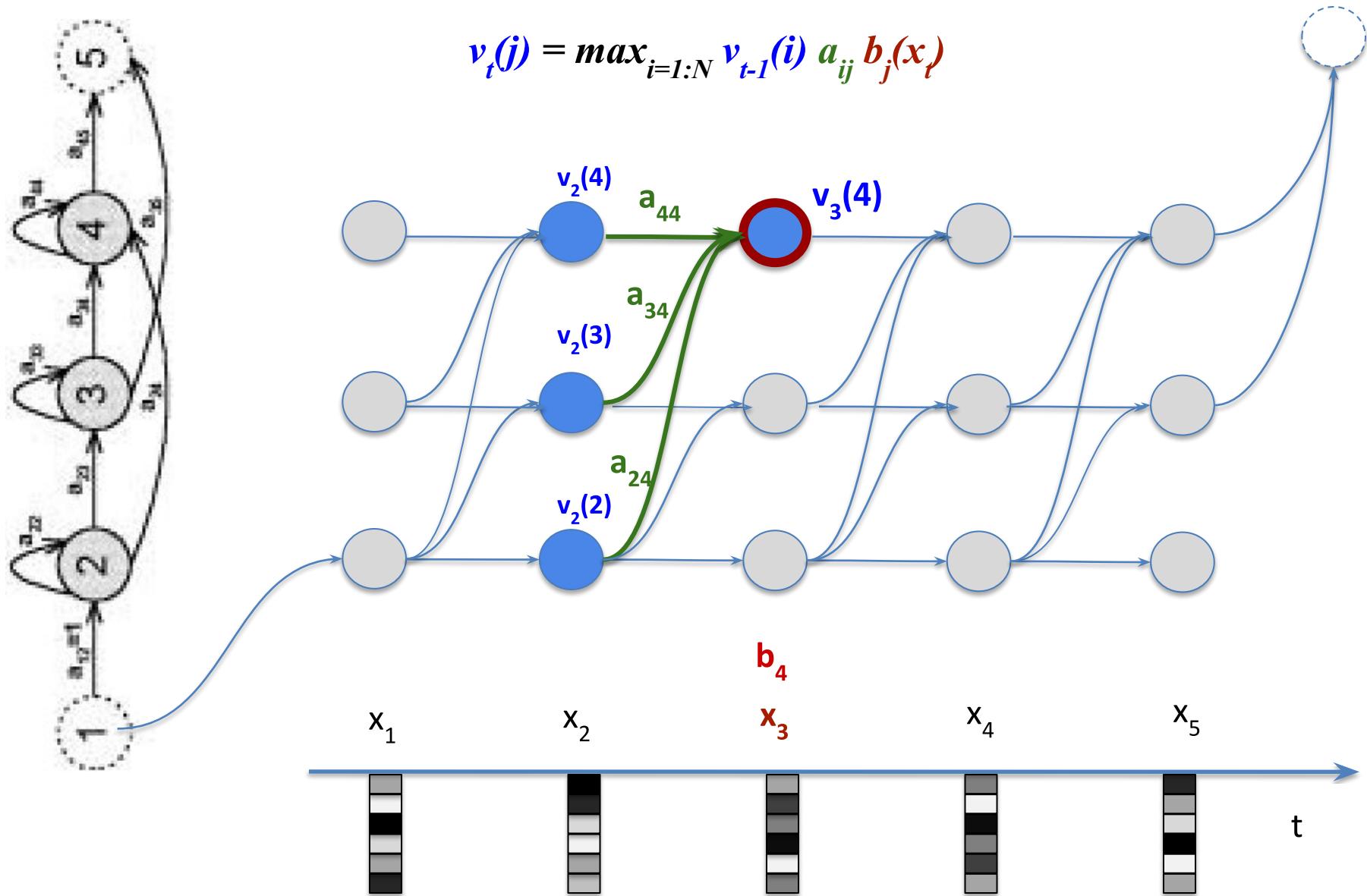
HMM decoding (viterbi)



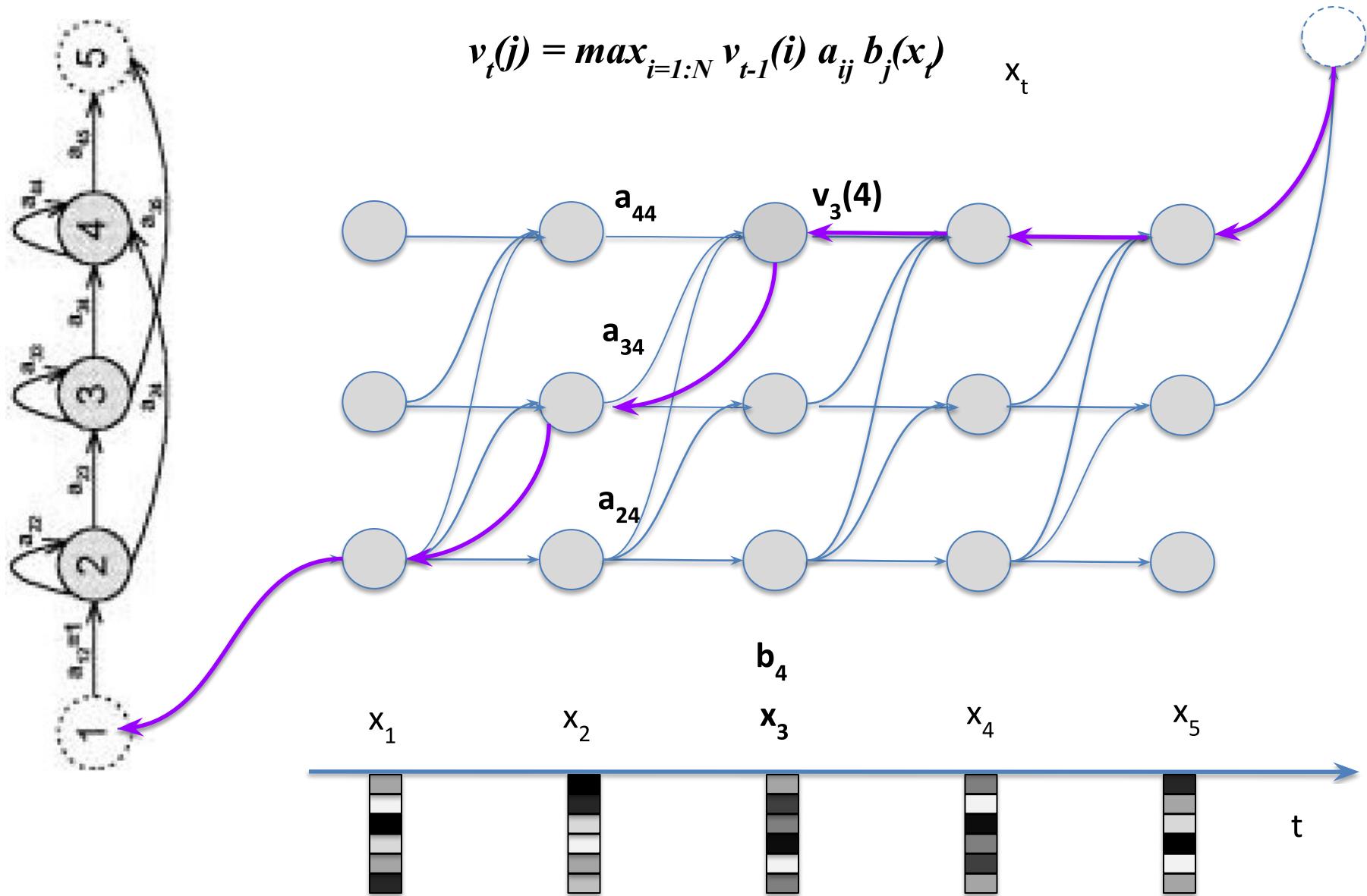
HMM decoding (viterbi)



decoding: $Q^* = \operatorname{argmax}_Q P(X|Q)$



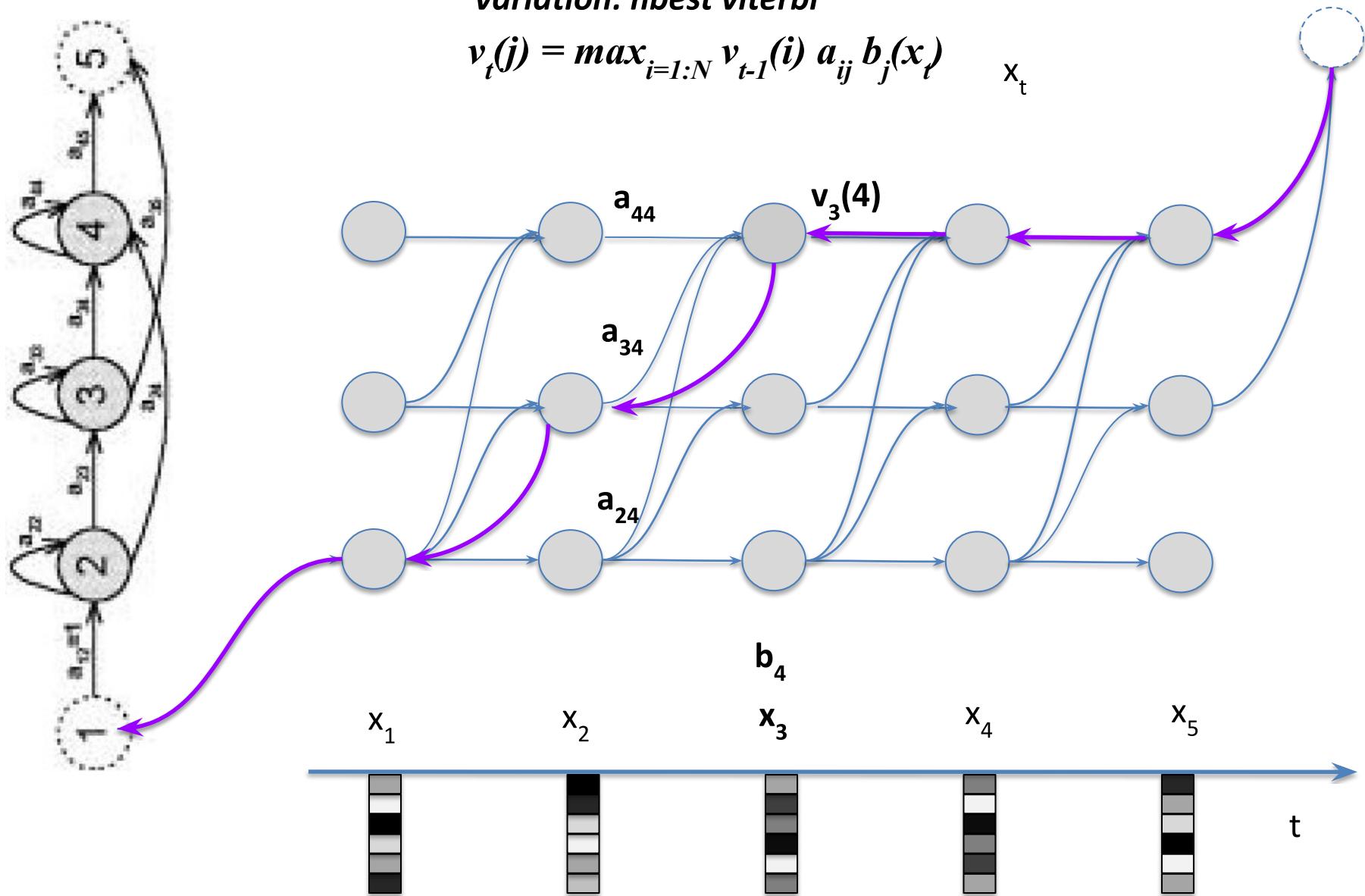
decoding: $Q^* = \operatorname{argmax}_Q P(X|Q)$



decoding: $Q^* = \operatorname{argmax}_Q P(X|Q)$

variation: nbest viterbi

$$v_t(j) = \max_{i=1:N} v_{t-1}(i) a_{ij} b_j(x_t)$$

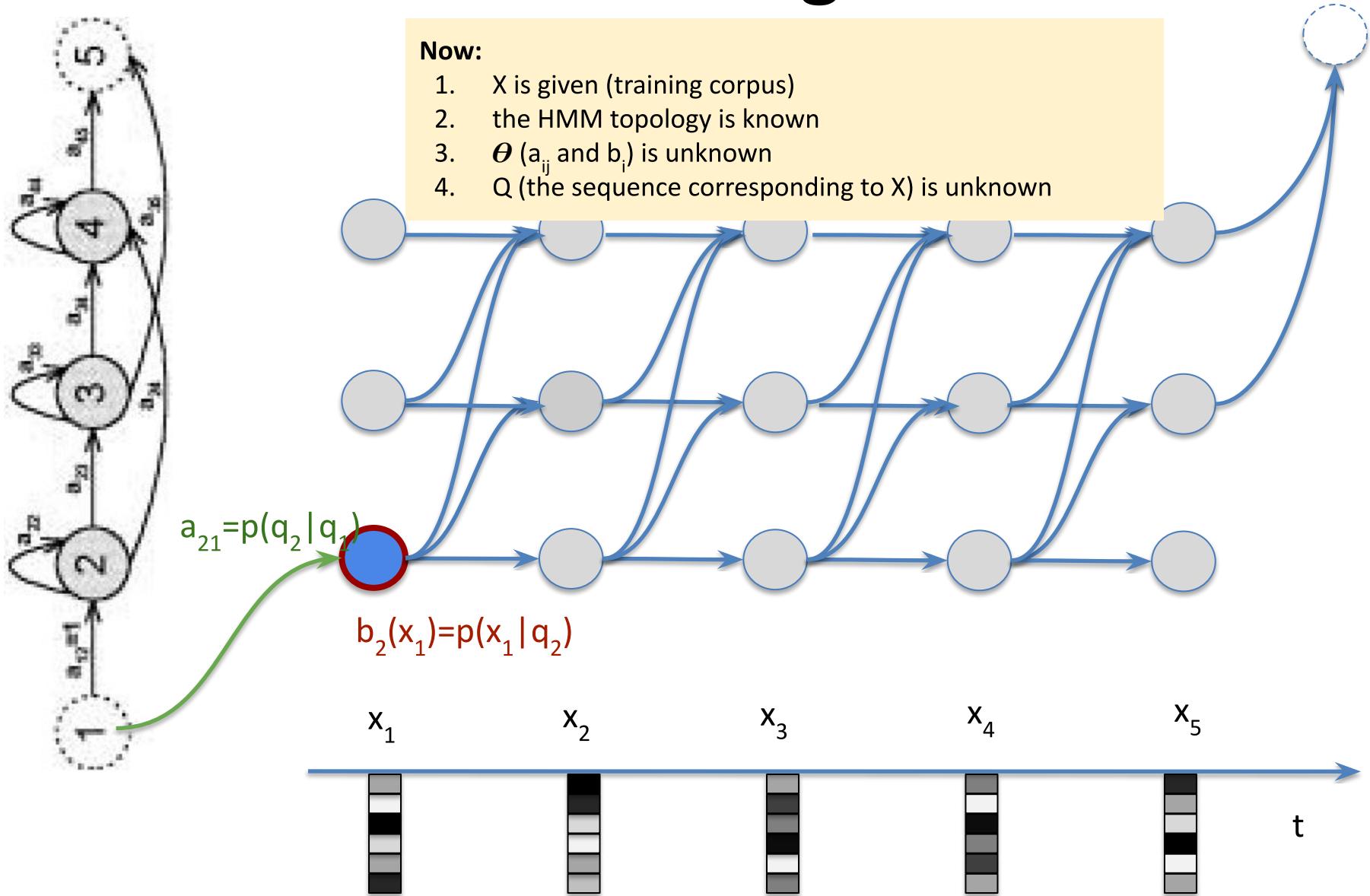


Learning

Learning

Now:

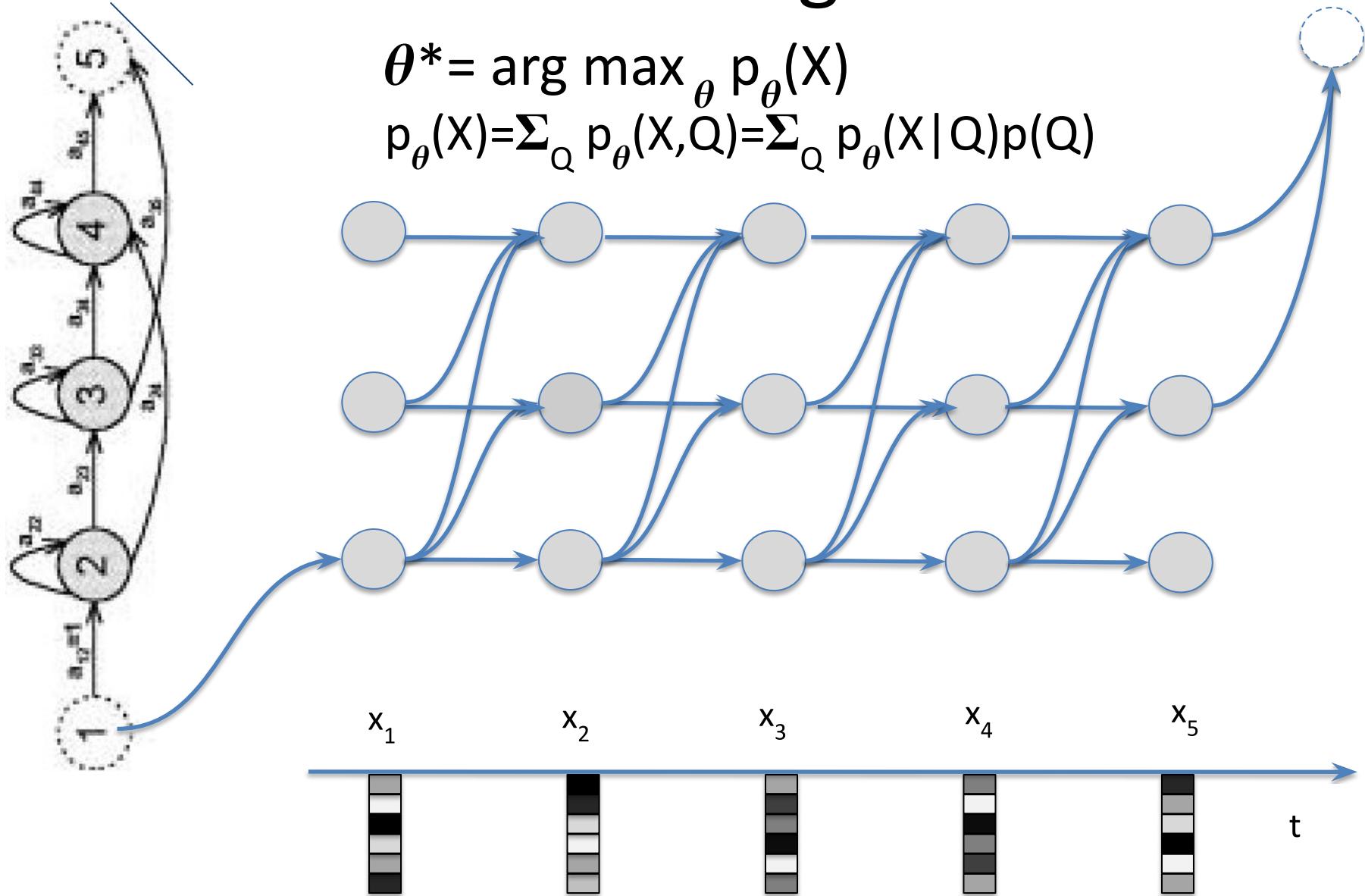
1. X is given (training corpus)
2. the HMM topology is known
3. θ (a_{ij} and b_i) is unknown
4. Q (the sequence corresponding to X) is unknown



Learning

$$\theta^* = \arg \max_{\theta} p_{\theta}(X)$$

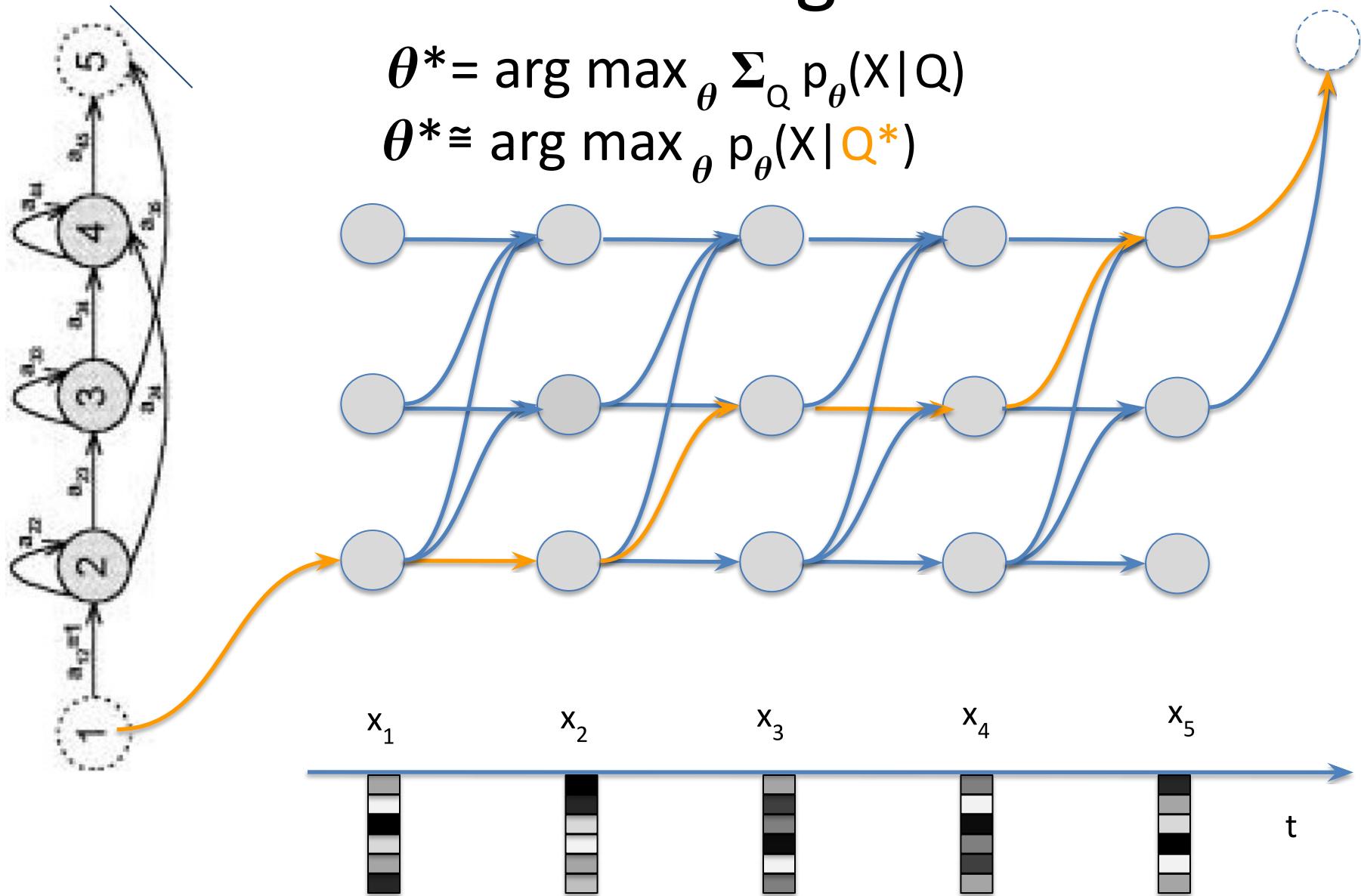
$$p_{\theta}(X) = \sum_Q p_{\theta}(X, Q) = \sum_Q p_{\theta}(X | Q)p(Q)$$



Learning

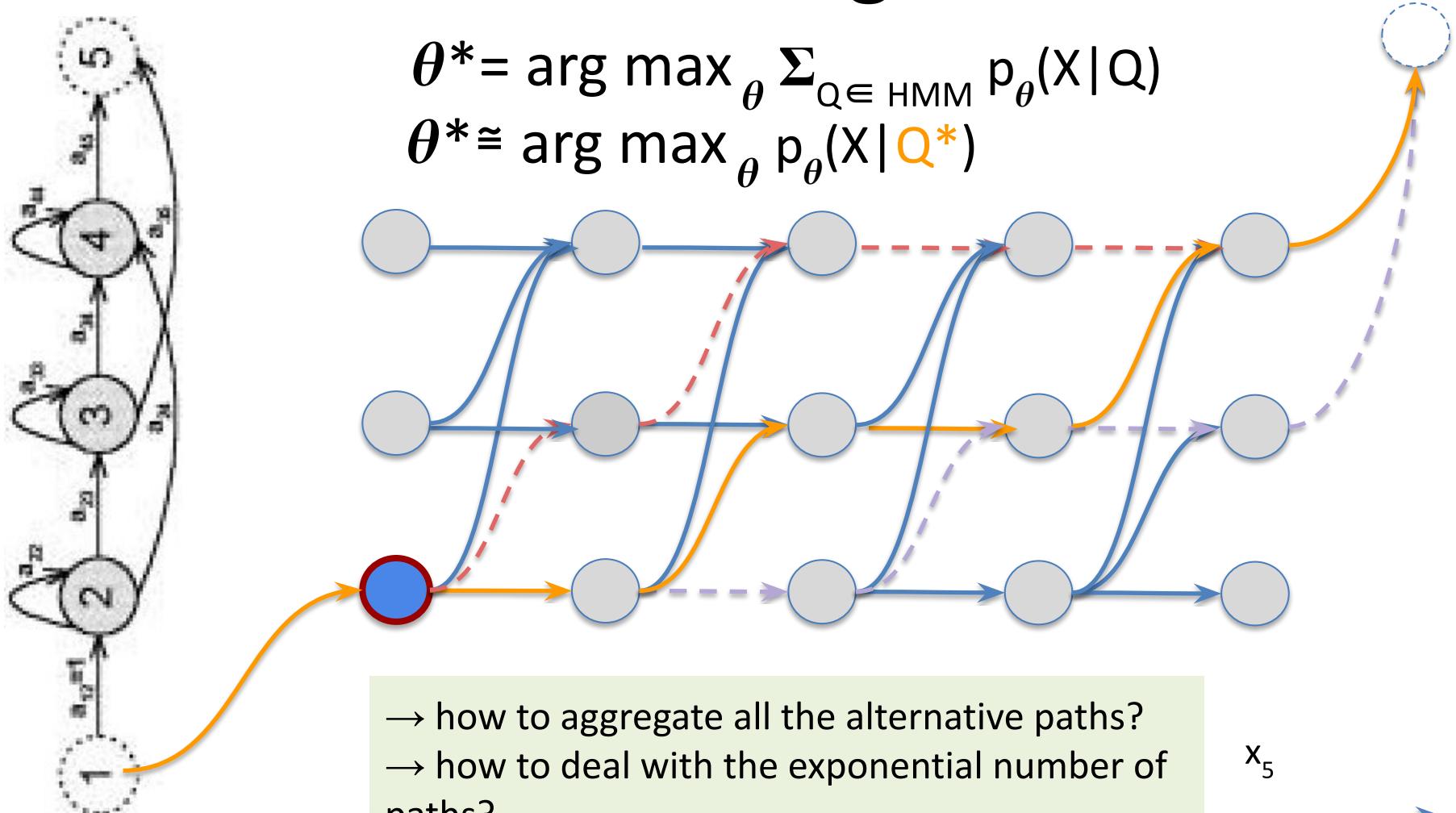
$$\theta^* = \arg \max_{\theta} \sum_Q p_{\theta}(X|Q)$$

$$\theta^* \cong \arg \max_{\theta} p_{\theta}(X|Q^*)$$



Learning

$$\theta^* = \arg \max_{\theta} \sum_{Q \in \text{HMM}} p_{\theta}(X|Q)$$
$$\theta^* \approx \arg \max_{\theta} p_{\theta}(X|Q^*)$$



Learning

- $\theta^* = \arg \max_{\theta} p(X|\theta)$
- EM (Baum Welch, forward backward) algorithm:
 - E step: estimate distribution of hidden variable Q (given X and θ): $p(Q|X, \theta)$
 - M step: maximize θ given the estimate of Q:
 $\theta^* = \arg \max_{\theta} p(X|Q, \theta)$
 - Iterate

an example: GMMs

- mixture model

$$p(x) = \sum_{j=1}^K w_j \cdot N(x | \mu_j, \Sigma_j)$$

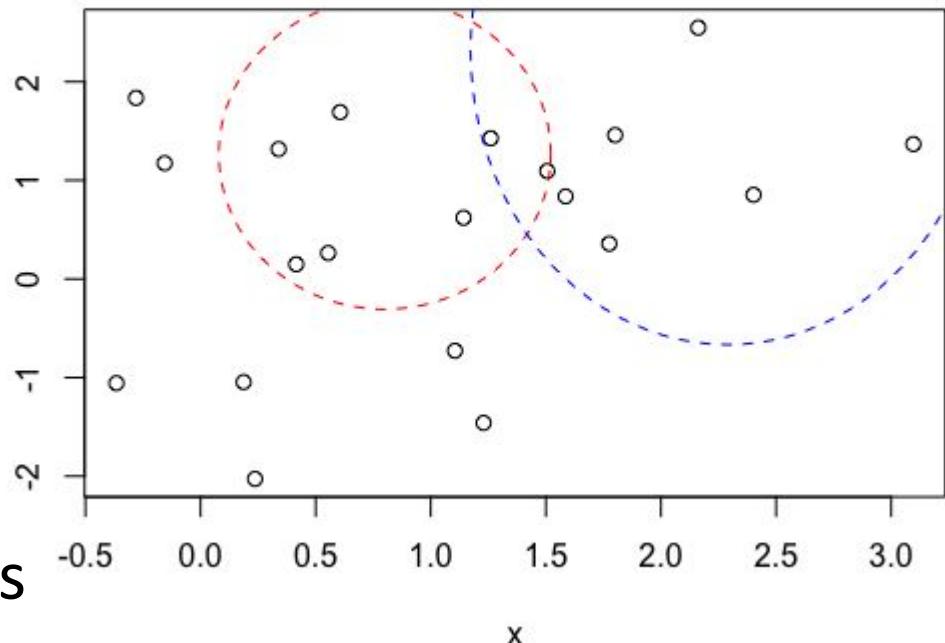
$$\theta = \{\mu_i, \Sigma_i\}$$

- problem: given observations
 $X = \{x_1 \dots x_N\}$

$$\theta^* = \arg \max_{\theta} \prod_i p(x_i | \theta)$$

- introduce hidden cluster membership $q_i = 1..k$

$$\theta^* = \arg \max_{\theta} \prod_i p(x_i, q_i | \theta)$$



K-means solution:

- initialize random values for θ
- alternate between
 - finding membership q
 - update the parameter θ

an example: GMMs

- mixture model

$$p(x) = \sum_{j=1}^K w_j \cdot N(x | \mu_j, \Sigma_j)$$

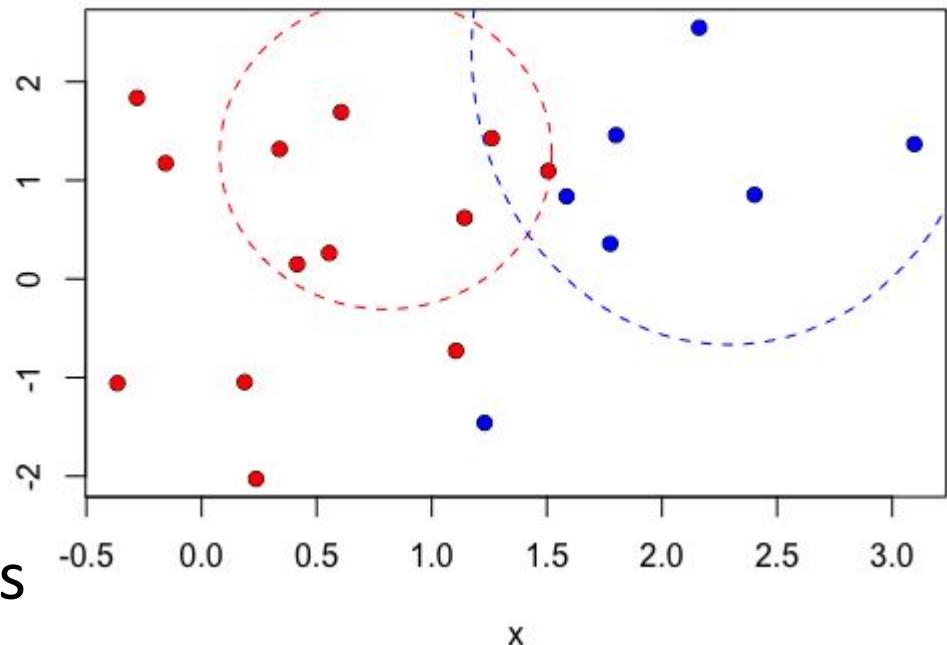
$$\theta = \{\mu_i, \Sigma_i\}$$

- problem: given observations
 $X = \{x_1 \dots x_N\}$

$$\theta^* = \arg \max_{\theta} \prod_i p(x_i | \theta)$$

- introduce hidden cluster membership $q_i = 1..k$

$$\theta^* = \arg \max_{\theta} \prod_i p(x_i, q_i | \theta)$$



K-means solution:

- initialize random values for θ
- alternate between
 - **finding membership q**
 - update the parameter θ

an example: GMMs

- mixture model

$$p(x) = \sum_{j=1}^K w_j \cdot N(x | \mu_j, \Sigma_j)$$

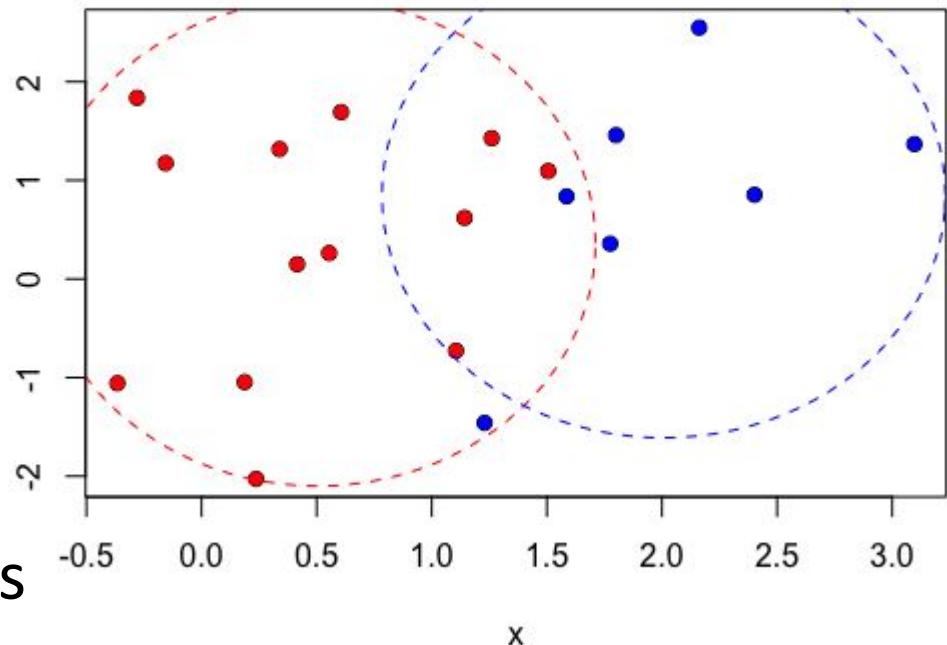
$$\theta = \{\mu_i, \Sigma_i\}$$

- problem: given observations
 $X = \{x_1 \dots x_N\}$

$$\theta^* = \arg \max_{\theta} \prod_i p(x_i | \theta)$$

- introduce hidden cluster membership $q_i = 1..k$

$$\theta^* = \arg \max_{\theta} \prod_i p(x_i, q_i | \theta)$$



K-means solution:

- initialize random values for θ
- alternate between
 - finding membership q
 - update the parameter θ

an example: GMMs

- mixture model

$$p(x) = \sum_{j=1}^K w_j \cdot N(x | \mu_j, \Sigma_j)$$

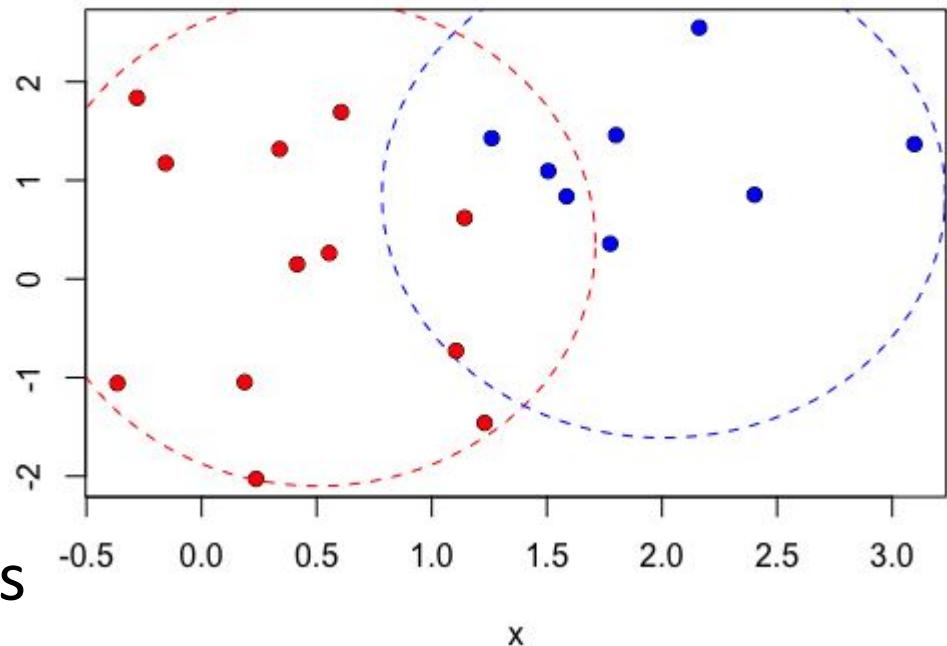
$$\theta = \{\mu_i, \Sigma_i\}$$

- problem: given observations
 $X = \{x_1 \dots x_N\}$

$$\theta^* = \arg \max_{\theta} \prod_i p(x_i | \theta)$$

- introduce hidden cluster membership $q_i = 1..k$

$$\theta^* = \arg \max_{\theta} \prod_i p(x_i, q_i | \theta)$$



K-means solution:

- initialize random values for θ
- alternate between
 - **finding membership q**
 - update the parameter θ

an example: GMMs

- mixture model

$$p(x) = \sum_{j=1}^K w_j \cdot N(x | \mu_j, \Sigma_j)$$

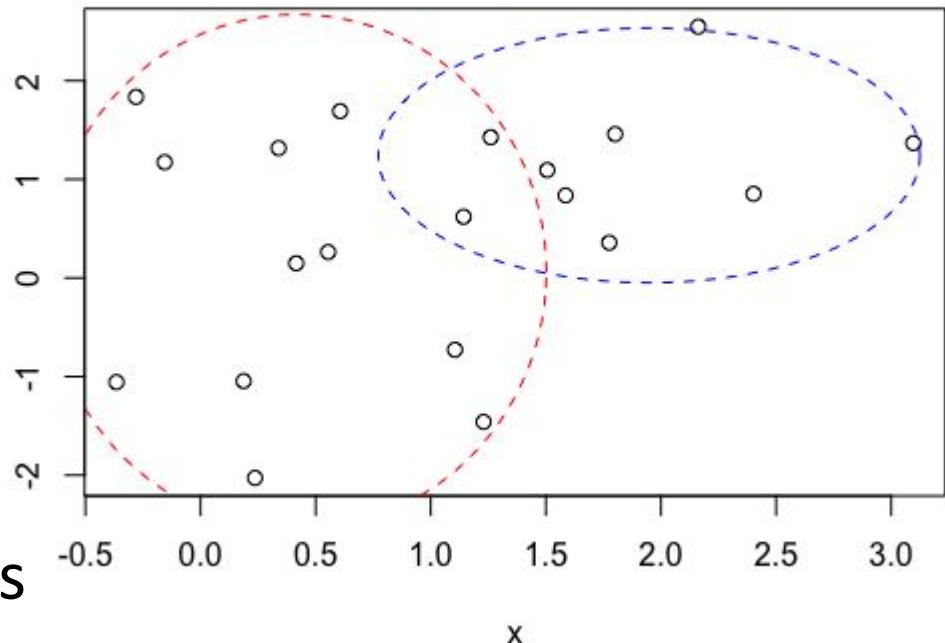
$$\theta = \{\mu_i, \Sigma_i\}$$

- problem: given observations
 $X = \{x_1 \dots x_N\}$

$$\theta^* = \arg \max_{\theta} \prod_i p(x_i | \theta)$$

- introduce hidden cluster membership $q_i = 1..k$

$$\theta^* = \arg \max_{\theta} \prod_i p(x_i, q_i | \theta)$$



K-means solution:

- initialize random values for θ
- alternate between
 - finding membership q
 - update the parameter θ

→ guaranteed to converge to a local max

an example: GMMs

- mixture model

$$p(x) = \sum_{j=1}^K w_j \cdot N(x | \mu_j, \Sigma_j)$$

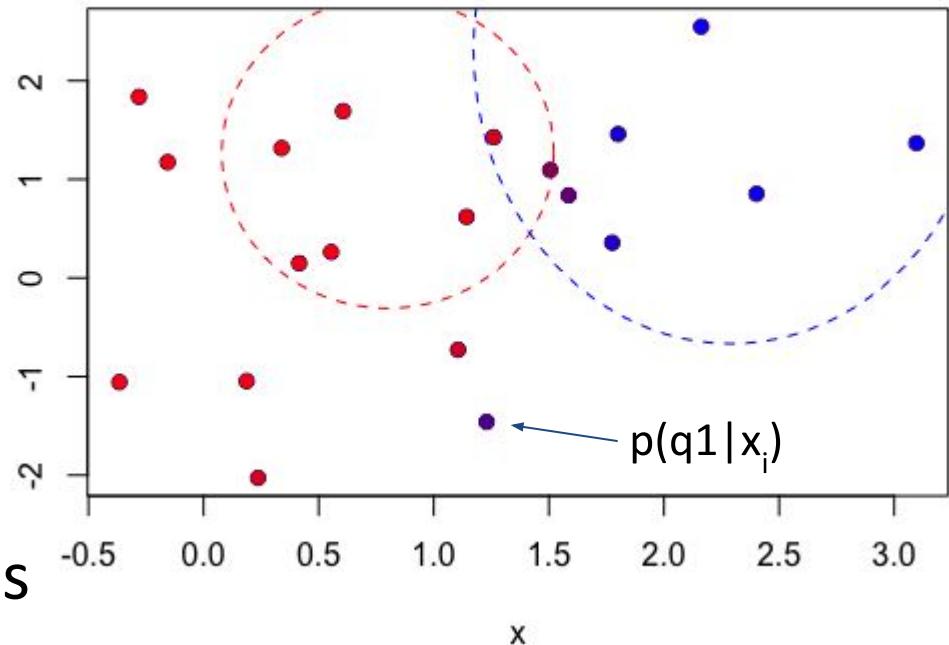
$$\theta = \{\mu_i, \Sigma_i\}$$

- problem: given observations
 $X = \{x_1 \dots x_N\}$

$$\theta^* = \arg \max_{\theta} \prod_i p(x_i | \theta)$$

- introduce hidden cluster membership $q_i = 1..k$

$$\theta^* = \arg \max_{\theta} \prod_i p(x_i, q_i | \theta)$$



The EM solution

- same thing, but instead of a discrete membership, use a graded one:
 $p(q_j | x_i)$
- now the update is weighed

$$\mu_j = \sum_i p(q_j | x_i) x_i / \sum_i p(q_j | x_i)$$

→ guaranteed to converge to a local max

going back to HMMs

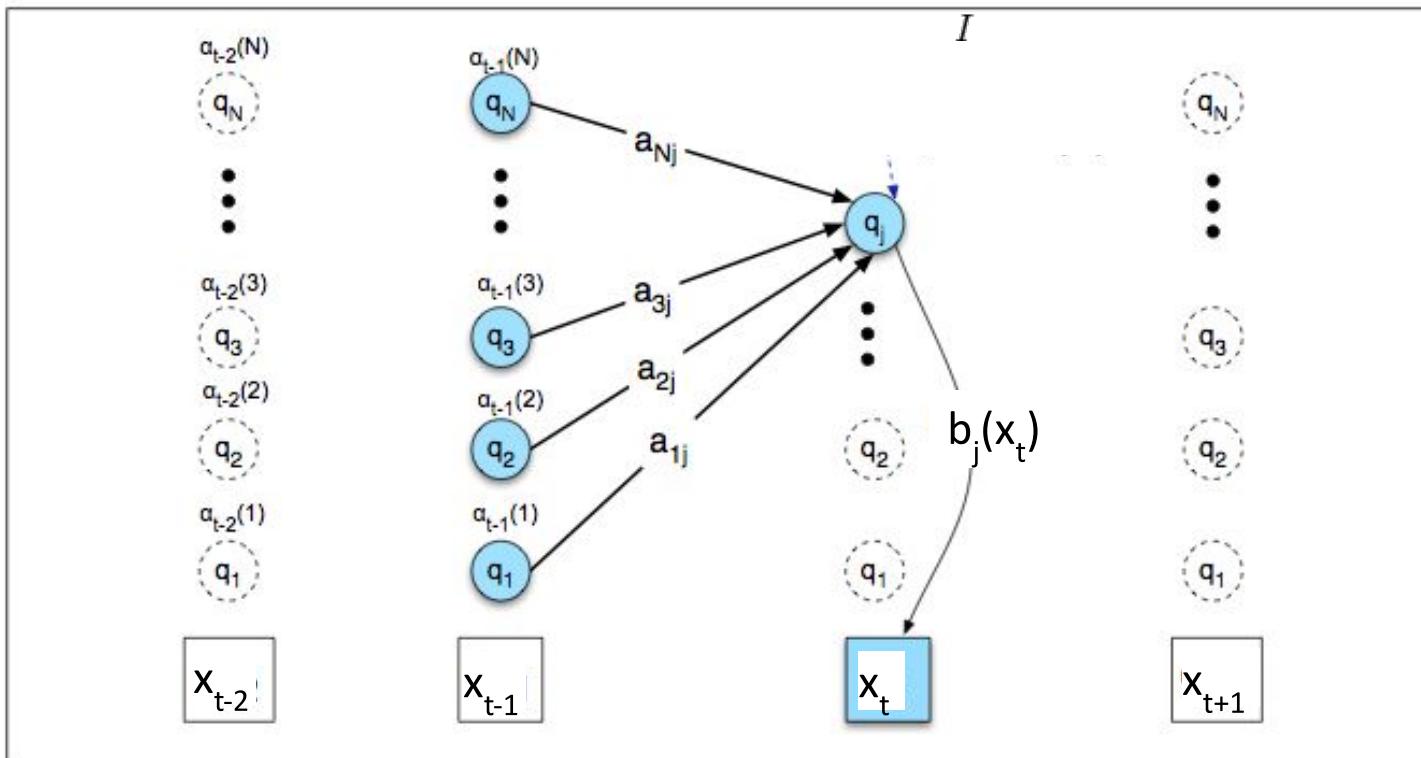
- intuition
 - in an HMM: iteratively estimate the counts
 - start with an estimate a_{ij} and b_j , use them to decode and get better estimates
 - compute these estimates on all paths, weighted by the *probability of the paths*

back to HMMs

- forward probabilities

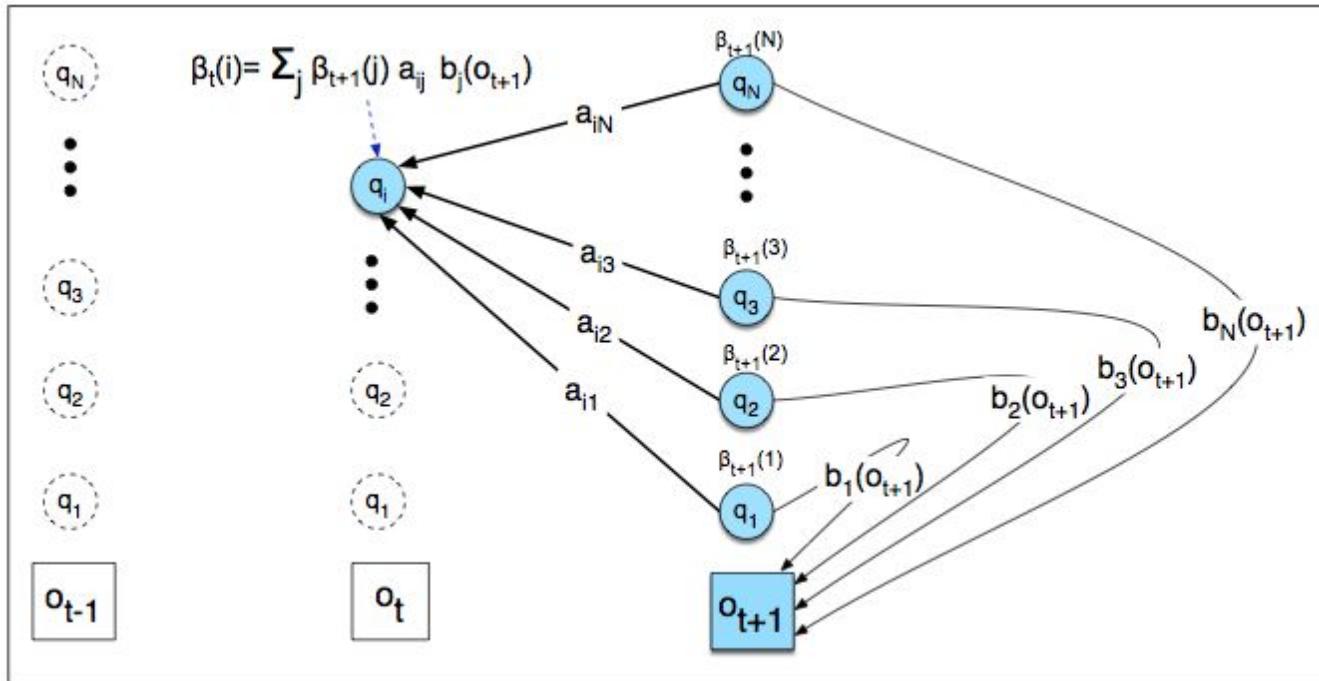
$$\alpha_t(j) = P(x_1, x_2..x_t, q_t = j | \theta)$$

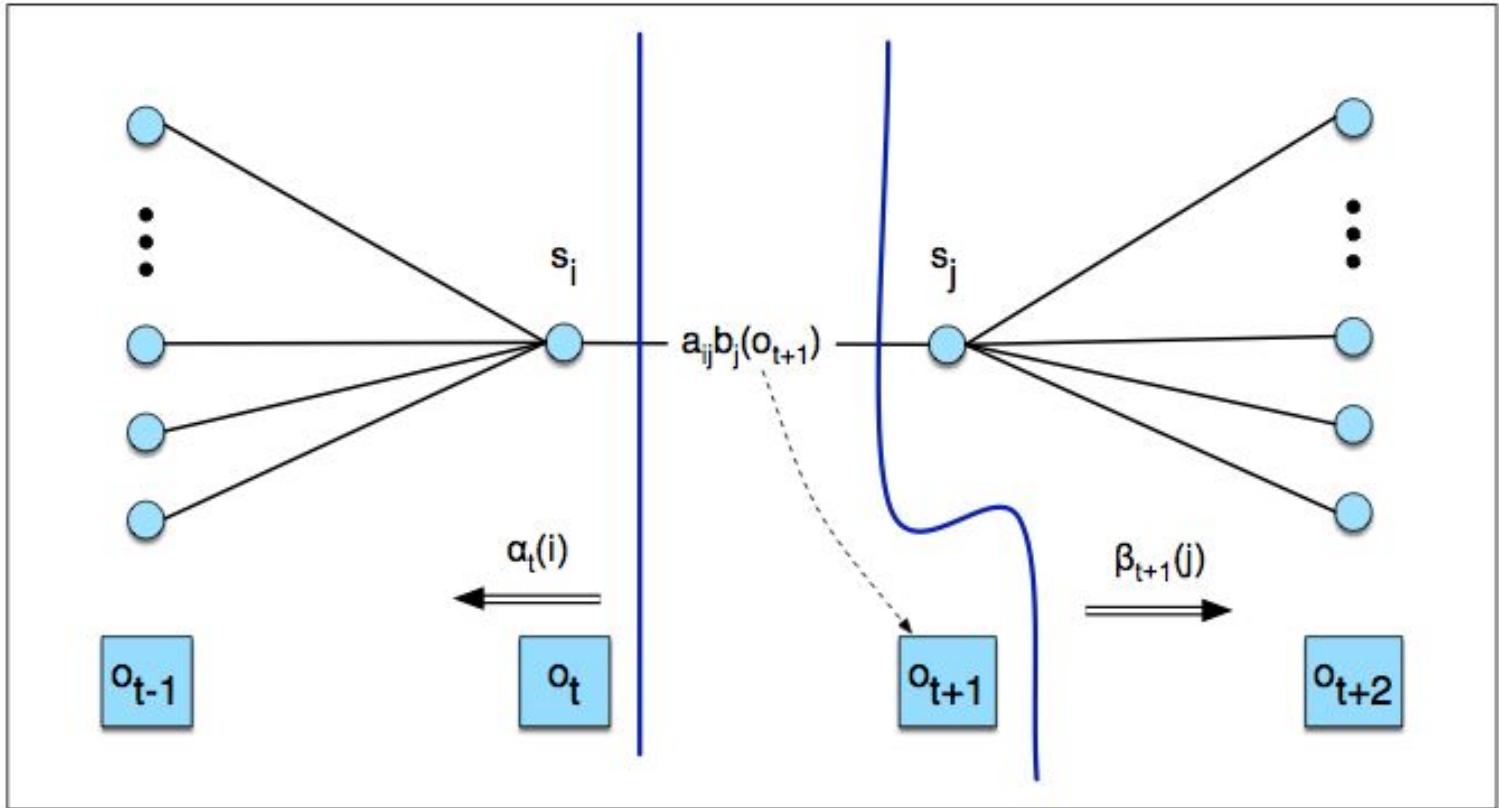
$$\alpha_t(j) = \sum_I \alpha_{t-1}(j) a_{ij} b_j(x_t))$$



- backward probabilities

$$\beta_t(i) = P(o_{t+1}, o_{t+2} \dots o_T | q_t = i, \Theta)$$





$$P(q_t = i, q_{t+1} = j, O | \Theta) = \alpha_t(i) a_{ij} b_j(o_{t+1}) \beta_{t+1}(j)$$

$$\xi_t(i, j) = P(q_t = i, q_{t+1} = j | O, \Theta) = \frac{\alpha_t(i) a_{ij} b_j(o_{t+1}) \beta_{t+1}(j)}{\alpha_T(q_F)}$$

from J&M (2017)

Baum-Welch algorithm

- E step

$$\xi_t(i, j) = \frac{\alpha_t(i) a_{ij} b_j(o_{t+1}) \beta_{t+1}(j)}{\alpha_T(q_F)}$$

$$\gamma_t(j) = \frac{\alpha_t(j) \beta_t(j)}{\alpha_T(q_F)}$$

- M Step

$$\hat{a}_{ij} = \frac{\sum_{t=1}^{T-1} \xi_t(i, j)}{\sum_{t=1}^{T-1} \sum_{k=1}^N \xi_t(i, k)}$$

$$\hat{b}_j(v_k) = \frac{\sum_{t=1}^T s.t. O_t=v_k \gamma_t(j)}{\sum_{t=1}^T \gamma_t(j)}$$

References

- Language models:
 - Jurasky & Martin (2017). Chapter 4. Language modeling with N-grams
- Dynamic programming:
 - Jaehyun Park (2015). Course materials for CS 97SI, Stanford University
- Viterbi decoding, Baum Welch
 - Jurafsky & Martin (2017). Chapter 9. Hidden Markov Models
- to know more
 - Gales & Young (2007). The application of Hidden Markow Models in Speech Recognition. *Foundations and Trends in Signal Processing*, **1(3)**, 195-304.

Baum-Welch algorithm:

- unsupervised (if you have a very generic HMM topology)
- supervised (if you have the transcription of X and only need to find the alignment)
- in between (if you have a lot of latent variables, eg, the precise lexical variant used for each W_i , etc)

MMI

Optimizing $P(X)$ does not necessarily optimizes the metric of interest (Word Error Rate)

$$\text{MMI}(\theta) = \log P_{\theta}(X|W)P(W) / \sum_w P_{\theta}(X|w)P(w)$$

discriminative loss

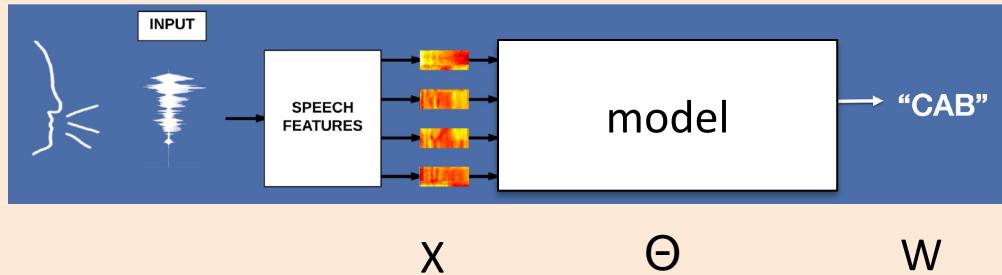
see for details <https://desh2608.github.io/2019-05-21-chain/>

Course

I. Generative Models

II. Decoding

III. Learning



Notations

θ : model parameters

$X=x_1 \dots x_t$: observed speech frames

$W=w_1..w_n$: word transcriptions

Decoding: given observed X and param. θ ,
find the most likely W .

$$W^* = \arg \max_w p(W|X, \theta)$$

Learning: given an observed X and
model, find the most likely θ

$$\theta^* = \arg \max_{\theta} p(X|\theta)$$