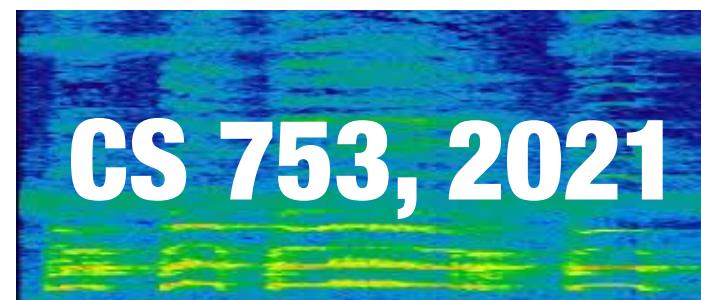


End-to-End Neural ASR Systems (II)

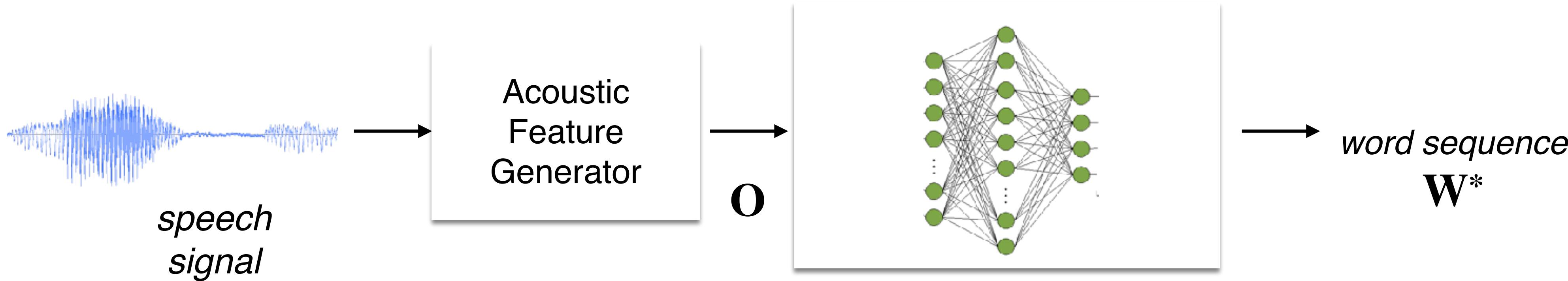
Lecture 8a



Instructor: Preethi Jyothi, IITB

Recap: Cascaded ASR \Rightarrow End-to-end ASR

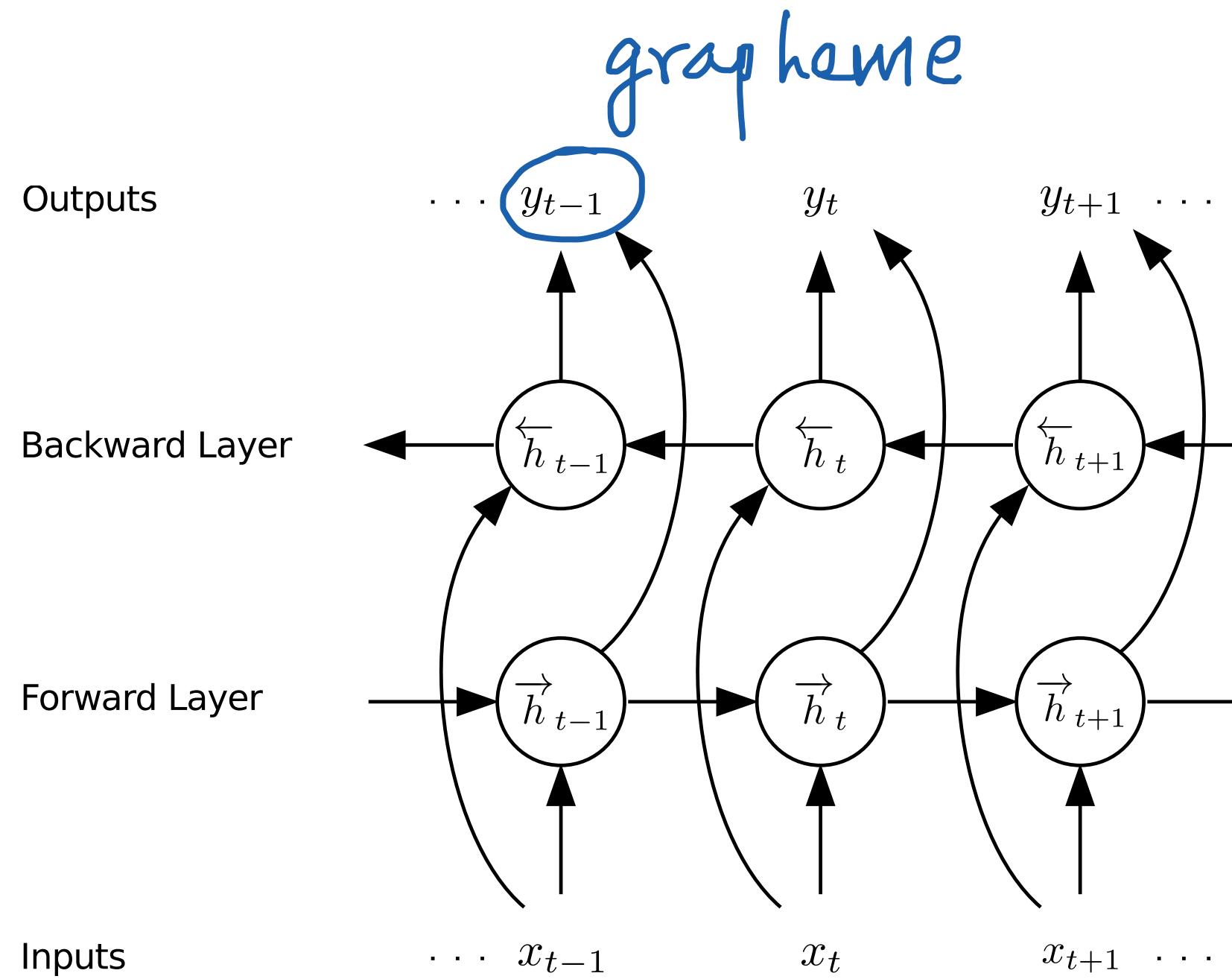
$$\mathbf{W}^* = \arg \max_W \Pr(\mathbf{W} | \mathbf{O})$$



Single end-to-end model that directly learns a mapping from speech to text

1. Encoder-decoder models with attention
2. **CTC-based models**

Network Architecture



forward

$$\vec{h}_t = \mathcal{H} \left(W_{x\vec{h}} \vec{x}_t + W_{\vec{h}\vec{h}} \vec{h}_{t-1} + b_{\vec{h}} \right)$$

backward

$$\overleftarrow{h}_t = \mathcal{H} \left(W_{x\overleftarrow{h}} \overleftarrow{x}_t + W_{\overleftarrow{h}\overleftarrow{h}} \overleftarrow{h}_{t+1} + b_{\overleftarrow{h}} \right)$$

$$y_t = W_{\vec{h}y} \vec{h}_t + W_{\overleftarrow{h}y} \overleftarrow{h}_t + b_o$$

$x_1, \dots, x_t, \dots, x_T$

- Input: Acoustic feature vectors. Output: Characters
- Long Short-Term Memory (LSTM) units (with in-built memory cells) are used to implement \mathcal{H} (in eqns above)
- Deep bidirectional LSTMs: Stack multiple bidirectional LSTM layers

Connectionist Temporal Classification (CTC)

- RNNs in ASR, if trained at the frame-level, will typically require alignments between the acoustics and the word sequence during training telling you which label (e.g. phone or character) should be output at each timestep.

Connectionist Temporal Classification (CTC)

- RNNs in ASR, if trained at the frame-level, will typically require alignments between the acoustics and the word sequence during training telling you which label (e.g. phone or character) should be output at each timestep.
- A new loss function, Connectionist Temporal Classification (CTC) tries to get around this.

Connectionist Temporal Classification (CTC)

- RNNs in ASR, if trained at the frame-level, will typically require alignments between the acoustics and the word sequence during training telling you which label (e.g. phone or character) should be output at each timestep.
- A new loss function, Connectionist Temporal Classification (CTC) tries to get around this.
- This is an objective function that allows RNN training without an explicit alignment step: CTC considers all possible alignments.

CTC Objective Function

- CTC objective function is the probability of an output label sequence y given an utterance $x \rightarrow$ Sequence of acoustic features

$$\text{CTC}(x, y) = \underbrace{\Pr(y|x)}_{\substack{a \in B^{-1}(y) \\ \text{alignment}}} = \sum_{a \in B^{-1}(y)} \Pr(a|x)$$

sequence of
chars

CTC Objective Function

- CTC objective function is the probability of an output label sequence y given an utterance x

$$\text{CTC}(x, y) = \Pr(y|x) = \sum_{a \in B^{-1}(y)} \Pr(a|x)$$

- Here, we sum over all possible alignments for y , enumerated by $B^{-1}(y)$

Is the underlying word BET or BEET?

i/p acoustic feature sequence of length 10 corresponding to the word BET

size 10

BBBEEEETTT

One possible alignment

BEEEEEETTT

BBAABEETTT

!

CTC: Mapping alignments to an output sequence

- Augment the output vocabulary with an additional “blank” (_) label



CTC: Mapping alignments to an output sequence

- Augment the output vocabulary with an additional “blank” (_) label
- For a given label sequence, there can be multiple alignments: (x, y, z) could correspond to (x, _, y, _, _, z) or (_, x, x, _, y, z)

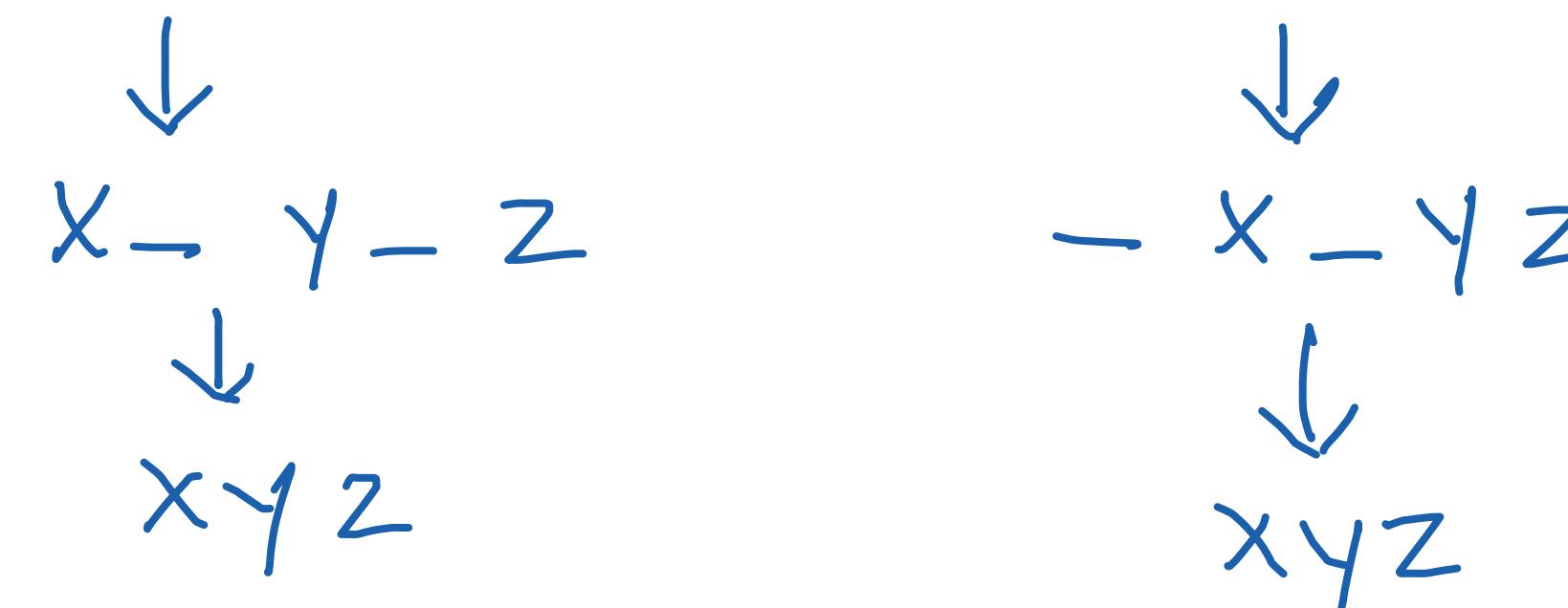
repeating chars have to be separated by a
blank _

CTC: Mapping alignments to an output sequence

$$\mathcal{B}^{-1}(y) = \left\{ \begin{matrix} a_1 & a_2 \\ _, _ & \dots, _ \\ a_M \end{matrix} \right\}$$

$$B(a_1) = B(a_2) = \dots = B(a_M) = y$$

- Augment the output vocabulary with an additional “blank” ($_$) label
- For a given label sequence, there can be multiple alignments: (x, y, z) could correspond to $(x, _, y, _, _, z)$ or $(_, x, x, _, y, z)$
- Define a 2-step operator B that reduces a label sequence by: first, removing repeating labels and second, removing blanks.
 $B("x, _, y, _, _, z") = B(" _, x, x, _, y, z") = "x, y, z"$



CTC Objective Function

- CTC objective function is the probability of an output label sequence y given an utterance x

$$\text{CTC}(x, y) = \Pr(y|x) = \sum_{a \in B^{-1}(y)} \Pr(a|x)$$

- Here, we sum over all possible alignments for y , enumerated by $B^{-1}(y)$

- CTC assumes that $\Pr(a|x)$ can be computed as $\prod_{t=1}^T \Pr(a_t|x)$

$$\Pr(a|x) = \prod_t \Pr(a_t|x)$$

CTC Objective Function

- CTC objective function is the probability of an output label sequence y given an utterance x

$$\text{CTC}(x, y) = \Pr(y|x) = \sum_{a \in B^{-1}(y)} \Pr(a|x)$$

- Here, we sum over all possible alignments for y , enumerated by $B^{-1}(y)$
- CTC assumes that $\Pr(a|x)$ can be computed as $\prod_{t=1}^T \Pr(a_t|x)$
 - i.e. CTC assumes that outputs at each time-step are conditionally independent given the input

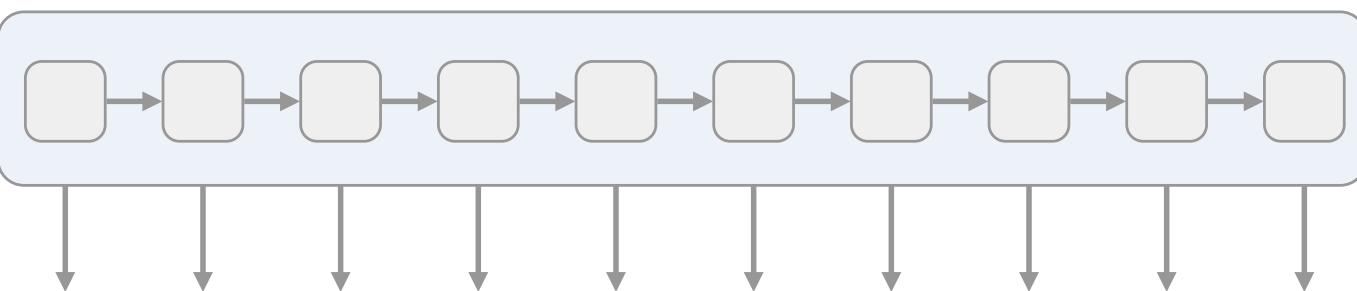
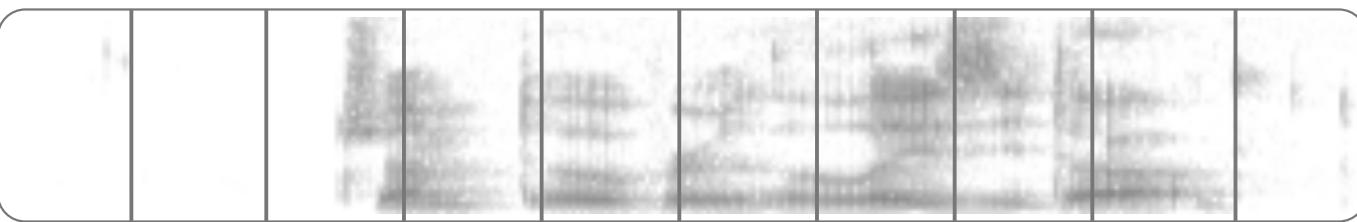
CTC Objective Function

- CTC objective function is the probability of an output label sequence y given an utterance x

$$\text{CTC}(x, y) = \Pr(y|x) = \sum_{a \in B^{-1}(y)} \Pr(a|x)$$

- Here, we sum over all possible alignments for y , enumerated by $B^{-1}(y)$
- CTC assumes that $\Pr(a|x)$ can be computed as $\prod_{t=1}^T \Pr(a_t|x)$
 - i.e. CTC assumes that outputs at each time-step are conditionally independent given the input
 - Efficient dynamic programming algorithm to compute this loss function and its gradients

Connectionist Temporal Classification: Overview



h	h	h	h	h	h	h	h	h	h	h
e	e	e	e	e	e	e	e	e	e	e
o	o	o	o	o	o	o	o	o	o	o
€	€	€	€	€	€	€	€	€	€	€

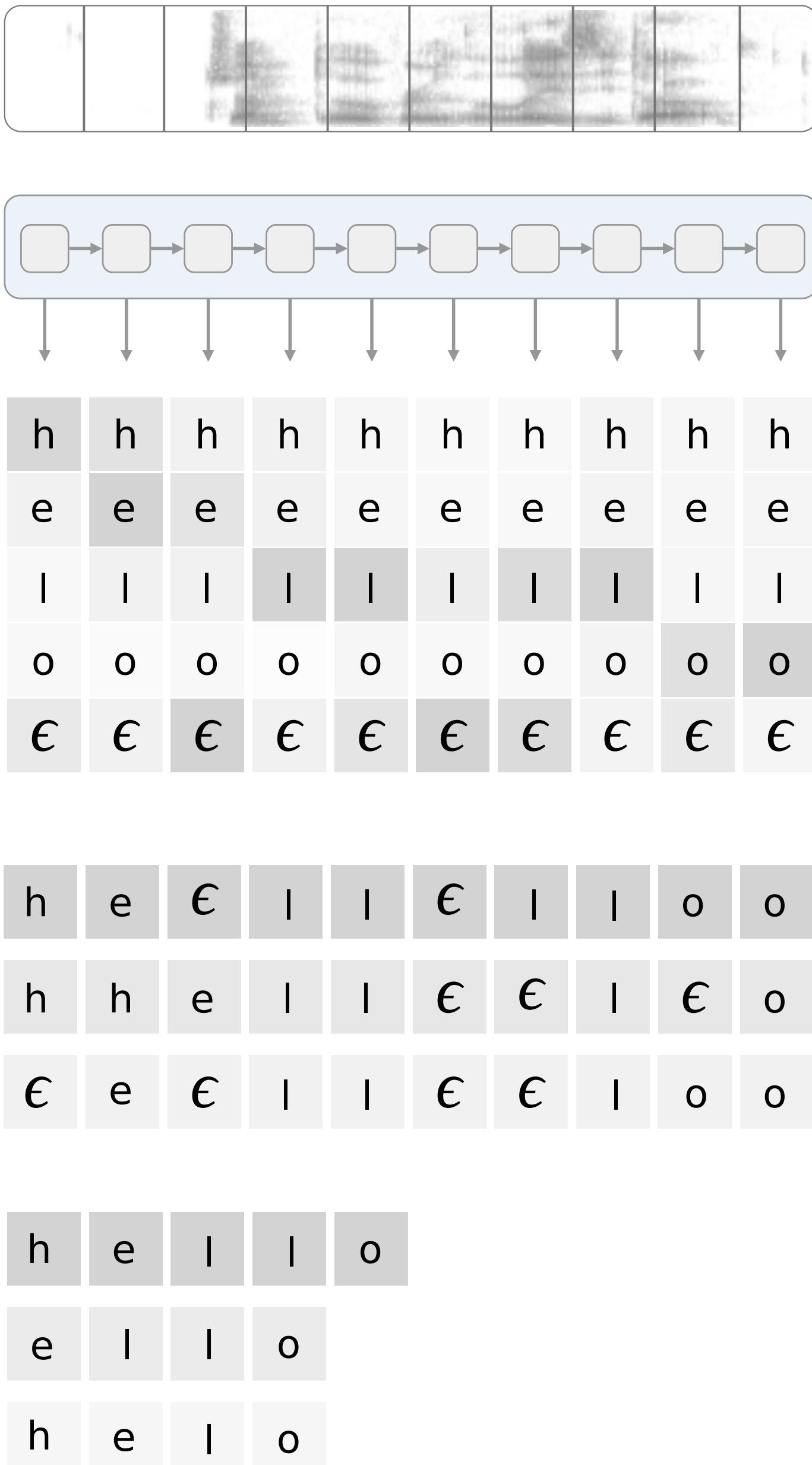
h	e	€			€			o	o
h	h	e			€	€		€	o
€	e	€			€	€		o	o

h	e			o
e			o	
h	e		o	

- CTC objective function is the probability of an output label sequence y given an utterance x (by summing over all possible alignments for y provided by $B^{-1}(y)$):

$$\begin{aligned} \text{CTC}(x, y) &= \underbrace{\Pr(y | x)}_{\sum_{a \in B^{-1}(y)} \Pr(a | x)} = \sum_{a \in B^{-1}(y)} \prod_{t=1}^T \Pr(a_t | \underbrace{x}_{a}) \\ &= \sum_{a \in B^{-1}(y)} \prod_{t=1}^T \Pr(a_t | \underbrace{x}_{a}) \end{aligned}$$

Connectionist Temporal Classification: Overview



- CTC objective function is the probability of an output label sequence y given an utterance x (by summing over all possible alignments for y provided by $B^{-1}(y)$):

$$\begin{aligned} \text{CTC}(x, y) &= \Pr(y | x) = \sum_{a \in B^{-1}(y)} \Pr(a | x) \\ &= \sum_{a \in B^{-1}(y)} \prod_{t=1}^T \Pr(a_t | x) \end{aligned}$$

- Efficient forward+backward algorithm to compute this loss function and its gradients [GJ14]

[GJ14] Towards End-to-End Speech Recognition with Recurrent Neural Networks, ICML 14

Illustration: Forward Algorithm to compute $\alpha_t(j)$

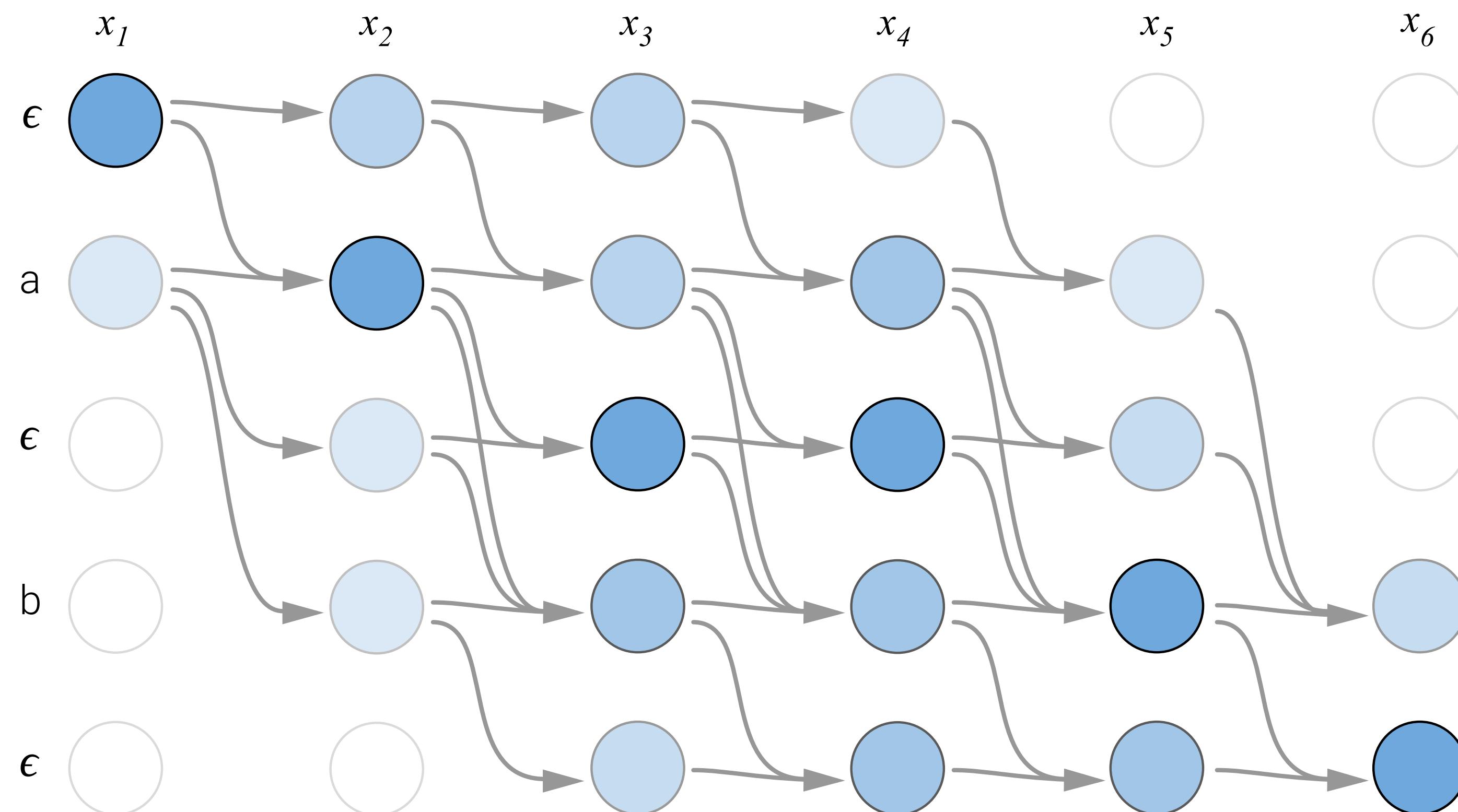
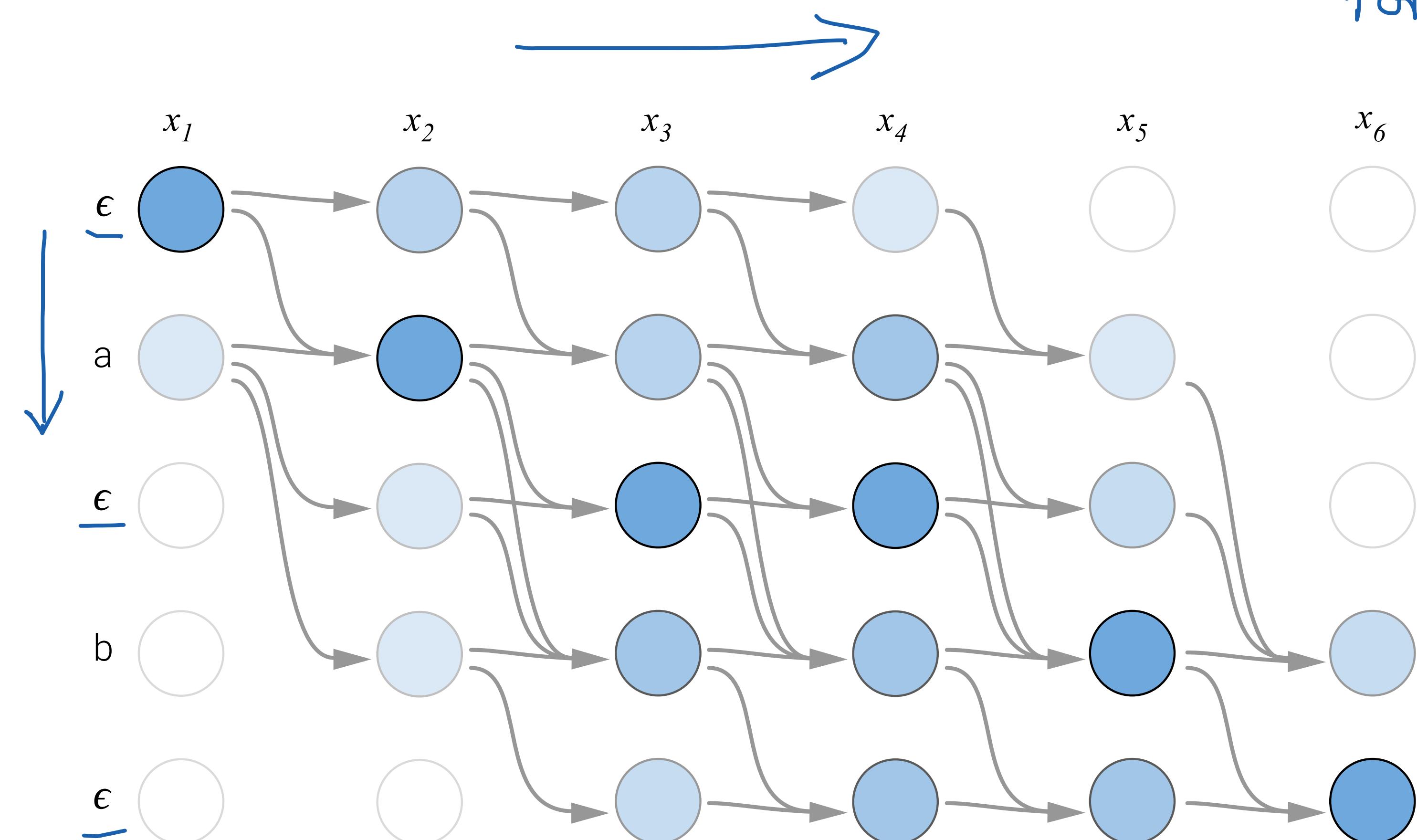


Illustration: Forward Algorithm to compute $\alpha_t(j)$

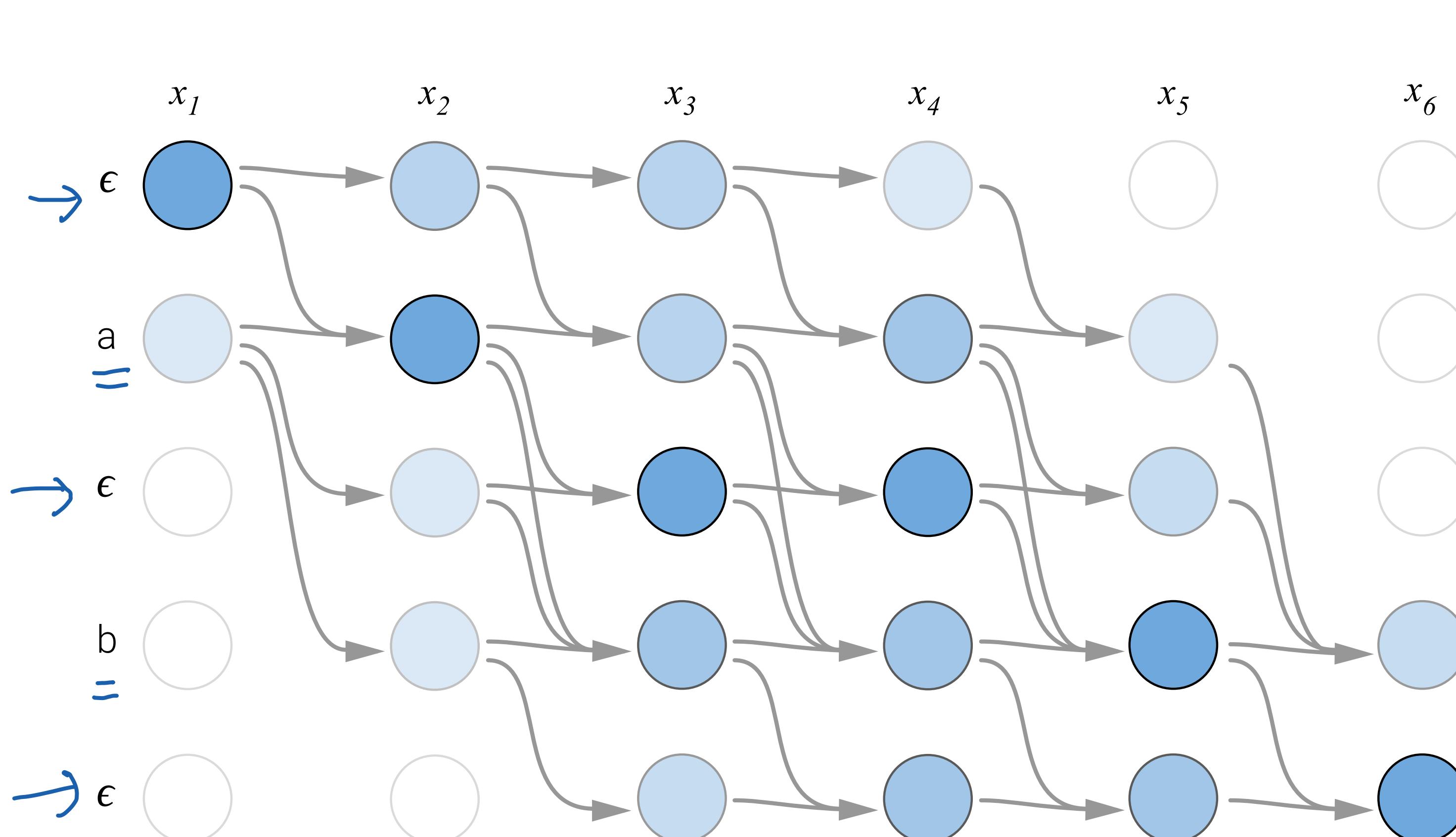
$\alpha_t(j)$ = Prob of $y_{1:j}$ at time t



forward probability $\alpha_t(j)$

$$\alpha_t(j) = \sum_{i=j-2}^j \alpha_{t-1}(i) a_{ij} b_t(y'_j)$$

Illustration: Forward Algorithm to compute $\alpha_t(j)$



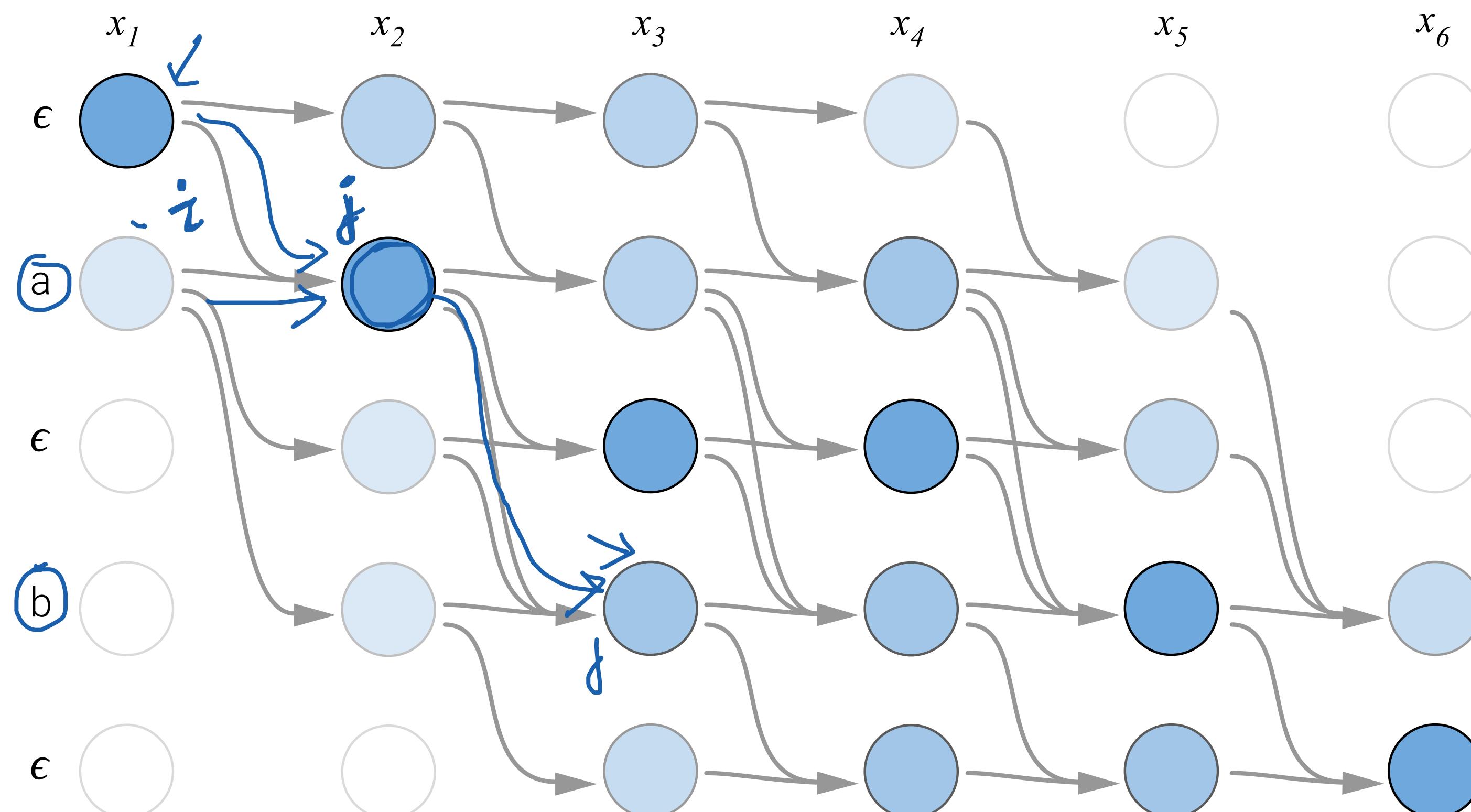
$$\alpha_t(j) = \sum_{i=j-2}^j \alpha_{t-1}(i) a_{ij} b_t(y'_j)$$

where

$b_t(y'_j)$ is the probability given by NN to the symbol y'_j for $t = 1 \dots T$, when $|x| = T$

$$y'_j = \begin{cases} y_{j/2} & \text{if } j \text{ is even} \\ \epsilon & \text{otherwise} \end{cases} \quad (j = 1 \dots 2l + 1 \text{ when } |y| = l)$$

Illustration: Forward Algorithm to compute $\alpha_t(j)$



$$\alpha_t(j) = \sum_{i=j-2}^j \alpha_{t-1}(i) a_{ij} b_t(y'_j)$$

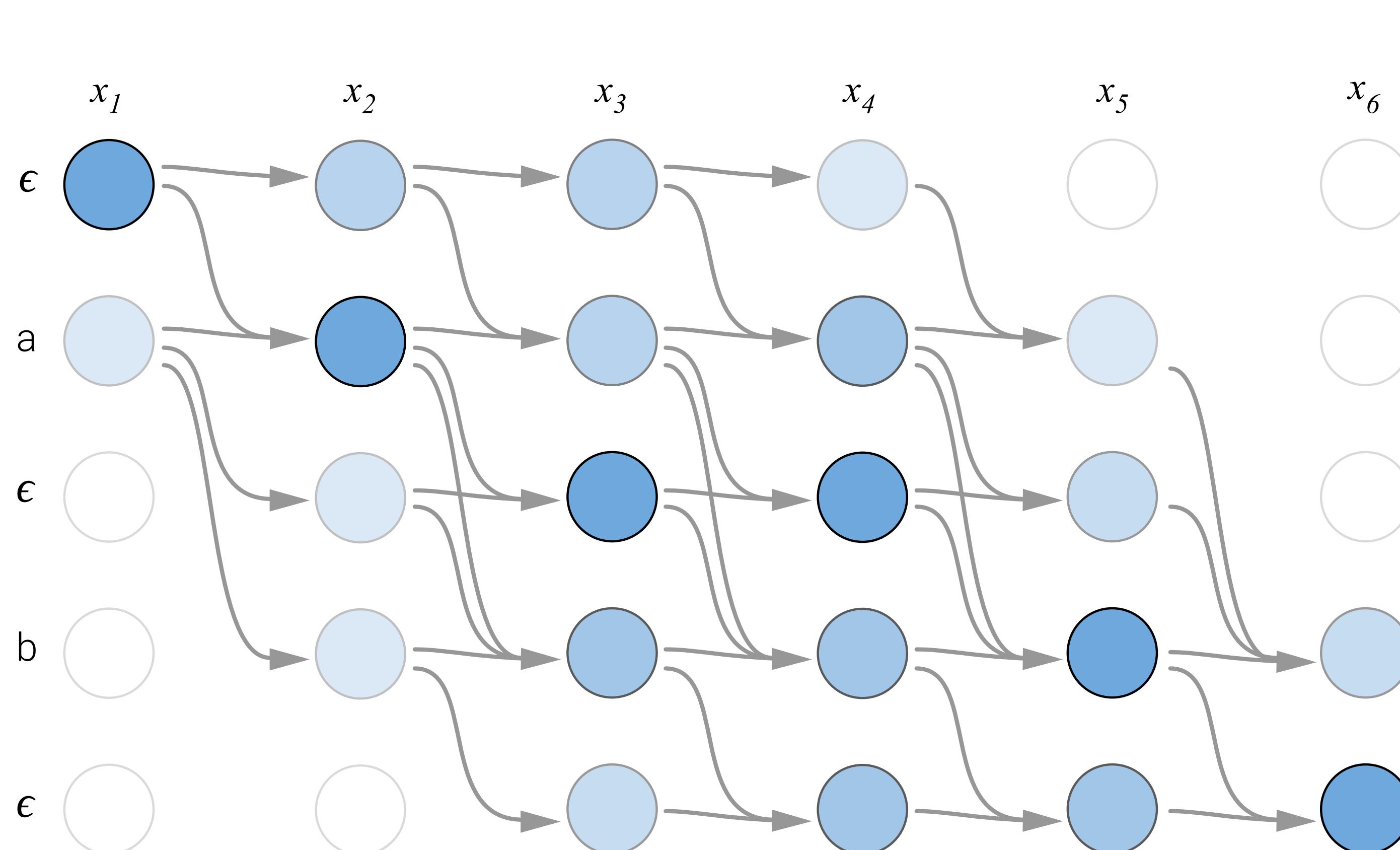
where

$b_t(y'_j)$ is the probability given by NN to the symbol y'_j for $t = 1 \dots T$, when $|x| = T$

$$y'_j = \begin{cases} y_{j/2} & \text{if } j \text{ is even} \\ \epsilon & \text{otherwise} \end{cases} \quad (j = 1 \dots 2l + 1 \text{ when } |y| = l)$$

$$a_{ij} = \begin{cases} 1 & \text{if } i = j \text{ or } i = j - 1 \\ 1 & \text{if } i = j - 2 \text{ and } y'_j \neq y'_{j-2} \\ 0 & \text{otherwise} \end{cases}$$

Illustration: Forward Algorithm to compute $\alpha_t(j)$



$$\alpha_t(j) = \sum_{i=j-2}^j \alpha_{t-1}(i)a_{ij}b_t(y'_j)$$

where

$b_t(y'_j)$ is the probability given by NN to the symbol y'_j for $t = 1 \dots T$, when $|x| = T$

$$y'_j = \begin{cases} y_{j/2} & \text{if } j \text{ is even} \\ \epsilon & \text{otherwise} \end{cases} \quad (j = 1 \dots 2l + 1 \text{ when } |y| = l)$$

$$a_{ij} = \begin{cases} 1 & \text{if } i = j \text{ or } i = j - 1 \\ 1 & \text{if } i = j - 2 \text{ and } y'_j \neq y'_{j-2} \\ 0 & \text{otherwise} \end{cases}$$

$$CTC(x, y) = \sum_{a \in B^{-1}(y)} \Pr(a \mid x) = \alpha_T(2l) + \alpha_T(2l + 1)$$

$\Pr(y \mid x)$

Decoding

Greedy decoding

- Pick the single most probable output at every time step

$$\arg \max_y \Pr(y|x) \approx B(\arg \max_a Pr(a|x))$$

Decoding

- Pick the single most probable output at every time step

$$\arg \max_y \Pr(y|x) \approx B(\arg \max_a Pr(a|x))$$

- Use a beam search algorithm to integrate a dictionary and a language model

Experimental Results

Table 1. Wall Street Journal Results. All scores are word error rate/character error rate (where known) on the evaluation set. ‘LM’ is the Language model used for decoding. ‘14 Hr’ and ‘81 Hr’ refer to the amount of data used for training.

SYSTEM	LM	14 HR	81 HR
RNN-CTC	NONE	74.2/30.9	30.1/9.2
RNN-CTC	DICTIONARY	69.2/30.0	24.0/8.0
RNN-CTC	MONOGRAM	25.8	15.8
RNN-CTC	BIGRAM	15.5	10.4
RNN-CTC	TRIGRAM	13.5	8.7
BASELINE	NONE	—	—
BASELINE	DICTIONARY	56.1	51.1
BASELINE	MONOGRAM	23.4	19.9
BASELINE	BIGRAM	11.6	9.4
BASELINE	TRIGRAM	9.4	7.8
COMBINATION	TRIGRAM	—	6.7

Sample Character-level Transcripts

target: *TO ILLUSTRATE THE POINT A PROMINENT MIDDLE EAST ANALYST IN WASHINGTON RECOUNTS A CALL FROM ONE CAMPAIGN*

output: *TWO ALSTRAIT THE POINT A PROMINENT MIDILLE EAST ANALYST IM WASHINGTON RECOUNCACALL FROM ONE CAMPAIGN*

target: *T. W. A. ALSO PLANS TO HANG ITS BOUTIQUE SHINGLE IN AIRPORTS AT LAMBERT SAINT*

output: *T. W. A. ALSO PLANS TOHING ITS BOOTIK SINGLE IN AIRPORTS AT LAMBERT SAINT*

target: *ALL THE EQUITY RAISING IN MILAN GAVE THAT STOCK MARKET INDIGESTION LAST YEAR*

output: *ALL THE EQUITY RAISING IN MULONG GAVE THAT STACRK MARKET IN TO JUSTIAN LAST YEAR*

target: *THERE'S UNREST BUT WE'RE NOT GOING TO LOSE THEM TO DUKAKIS*

output: *THERE'S UNREST BUT WERE NOT GOING TO LOSE THEM TO DEKAKIS*

Another end-to-end system

- Decoding is at the word level if we use a dictionary+LM.
Out-of-vocabulary (OOV) words cannot be handled.

Another end-to-end system

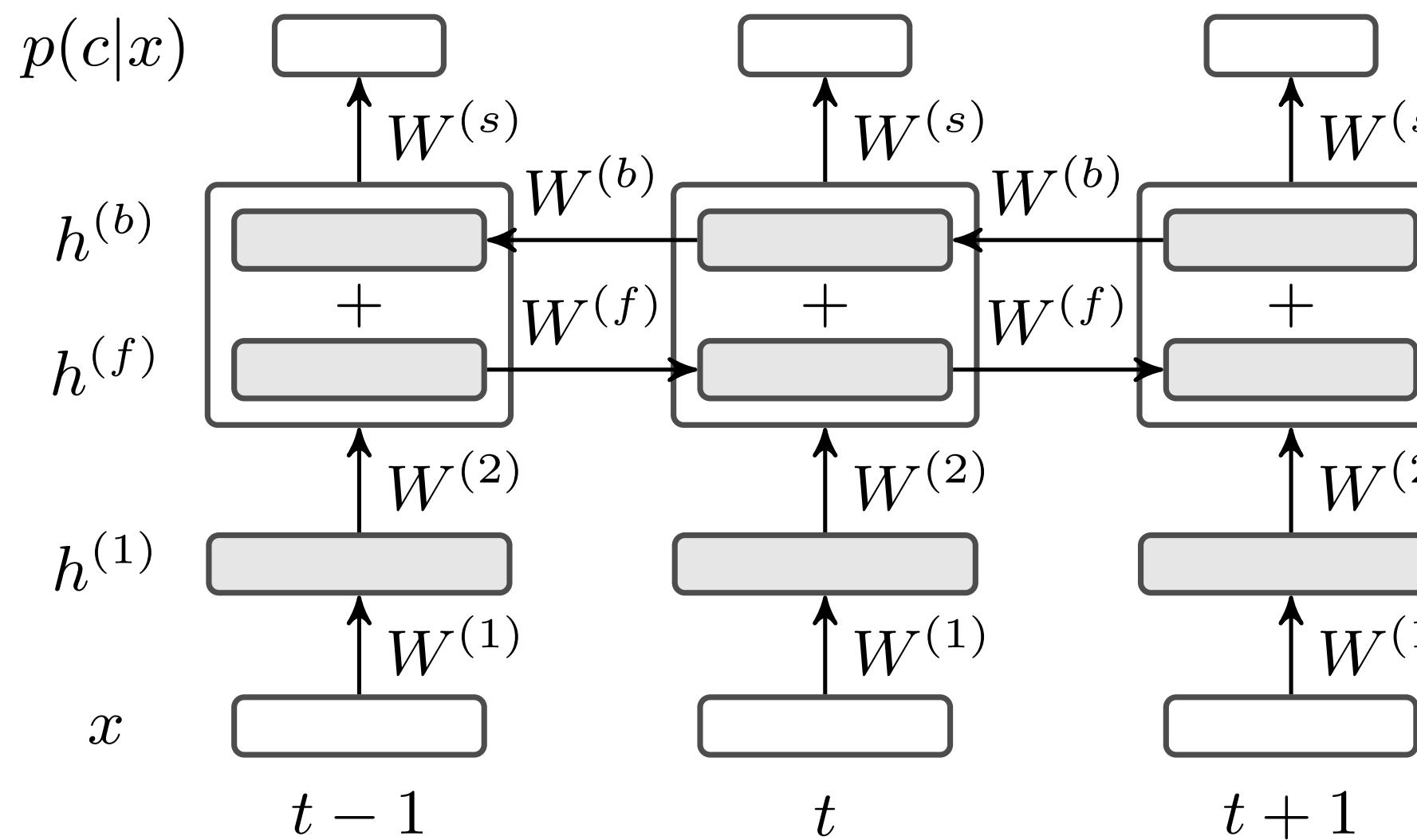
- Decoding is at the word level if we use a dictionary+LM.
Out-of-vocabulary (OOV) words cannot be handled.
- Build a system that is trained and decoded entirely at the character-level [Maas et al.]
- This would enable the transcription of OOV words, disfluencies, etc.

Another end-to-end system

- Decoding is at the word level if we use a dictionary+LM.
Out-of-vocabulary (OOV) words cannot be handled.
- Build a system that is trained and decoded entirely at the character-level [Maas et al.]
- This would enable the transcription of OOV words, disfluencies, etc.
- Shows results on the Switchboard task. Matches a GMM-HMM baseline system but underperforms compared to an HMM-DNN baseline.

Model Specifics

- Approach consists of two neural models:
 - A deep bidirectional RNN (DBRNN) mapping acoustic features to character sequences (Trained using CTC.)
 - A neural network character language model



Decoding

- Simplest form: Decode without any language model

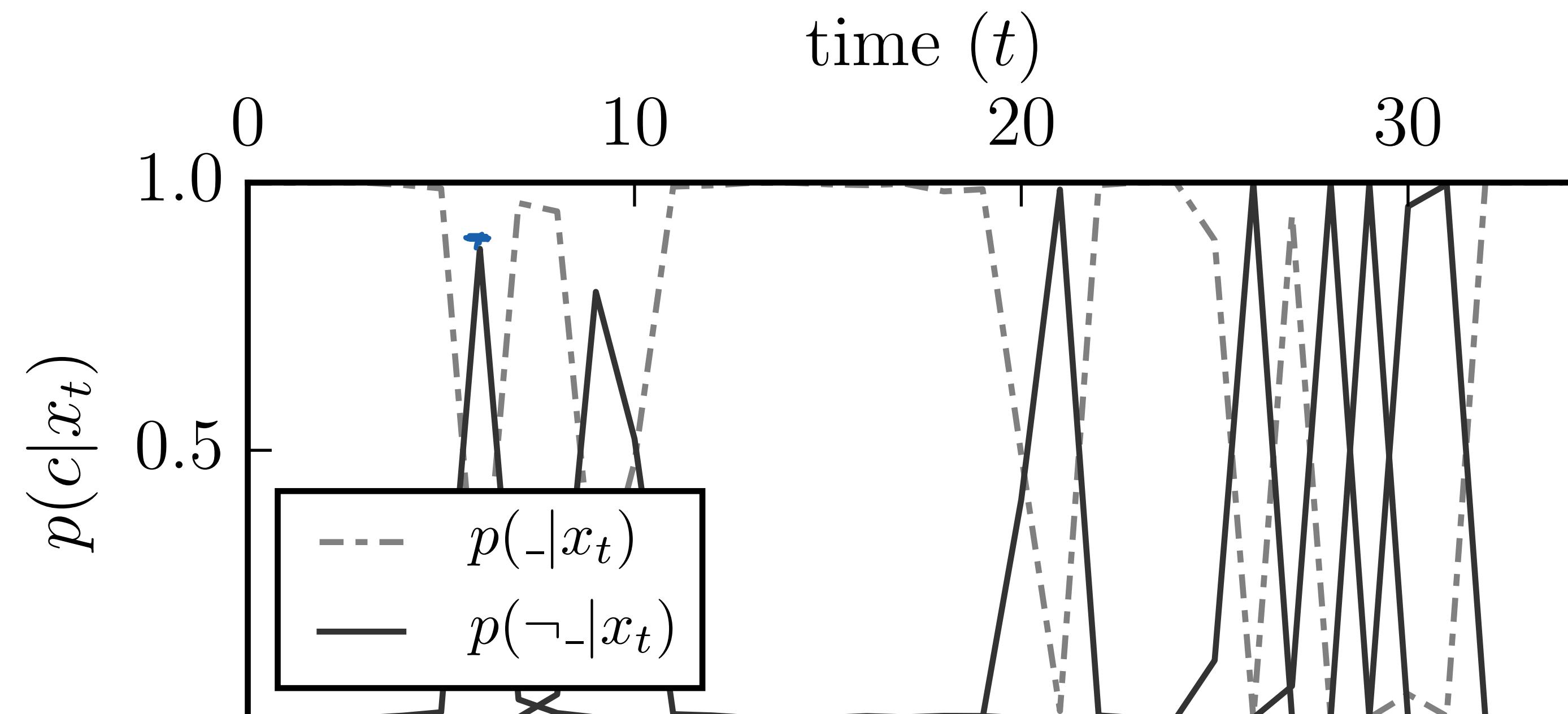
Decoding

- Simplest form: Decode without any language model
- Beam Search decoding:
 - Combine DBRNN outputs with a char-level language model
 - Char-level language model applied at every time step (unlike word models)
 - Circumvents the issue of handling OOV words during decoding

Sample Test Utterances

#	Method	Transcription
(1)	Truth	yeah i went into the i do not know what you think of <i>fidelity</i> but
	HMM-GMM	yeah when the i don't know what you think of fidel it even them
	CTC+CLM	yeah i went to i don't know what you think of <u>fidelity</u> but um
(2)	Truth	no no speaking of weather do you carry a altimeter slash <i>barometer</i>
	HMM-GMM	no i'm not all being the weather do you uh carry a uh helped emitters last brahms her
	CTC+CLM	no no beating of whether do you uh carry a uh a time or less <u>barometer</u>
(3)	Truth	i would ima- well yeah it is i know you are able to stay home with them
	HMM-GMM	i would amount well yeah it is i know um you're able to stay home with them
	CTC+CLM	i would ima- well yeah it is i know uh you're able to stay home with them

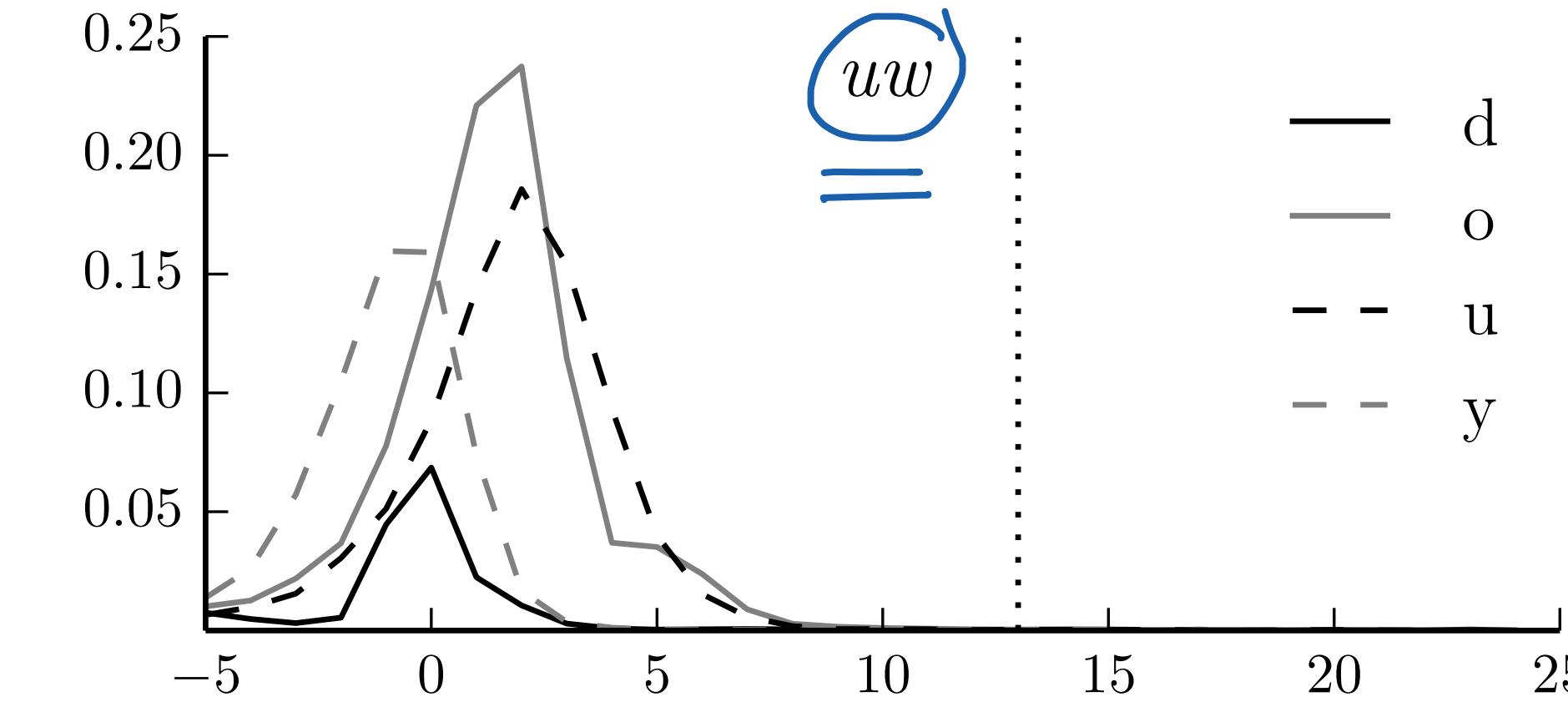
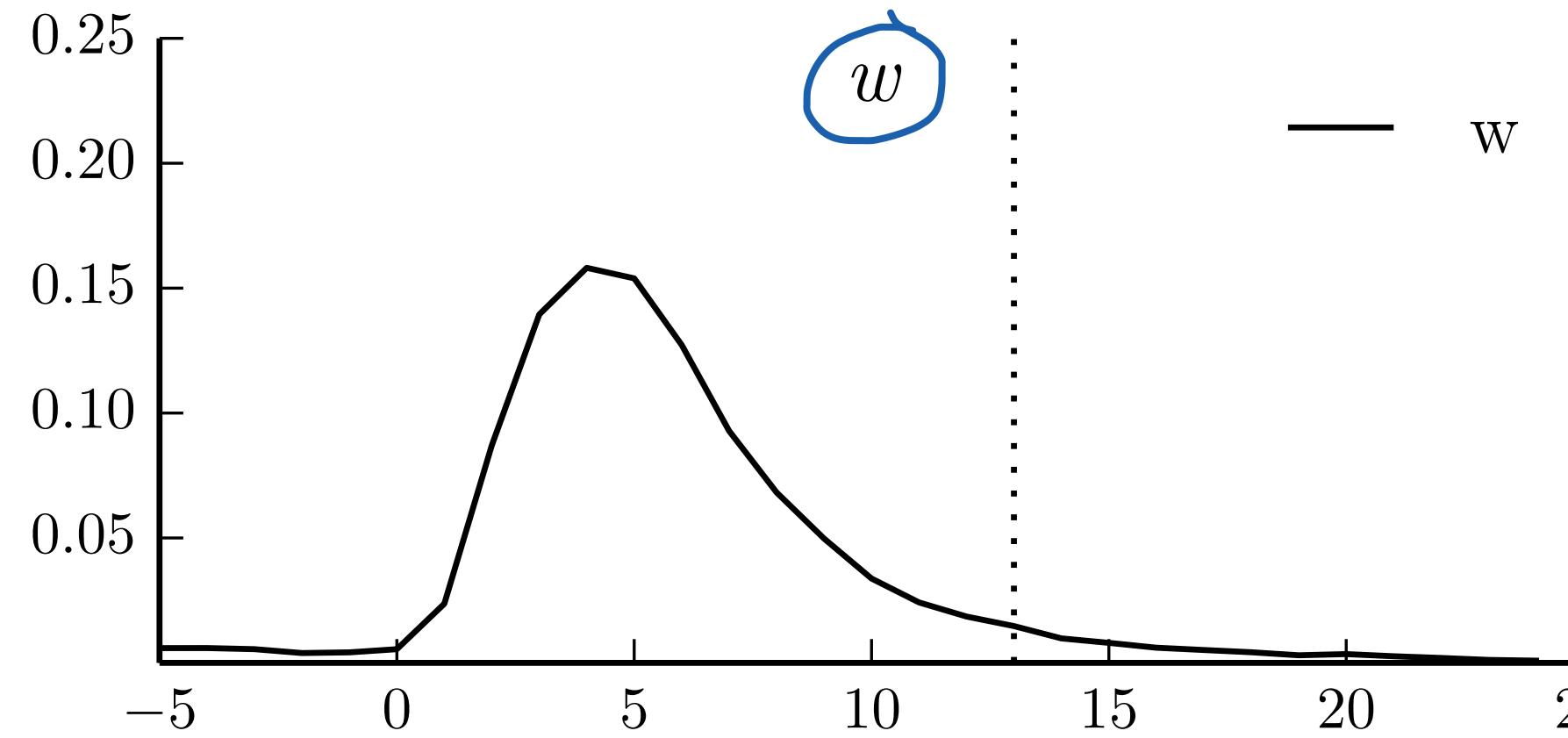
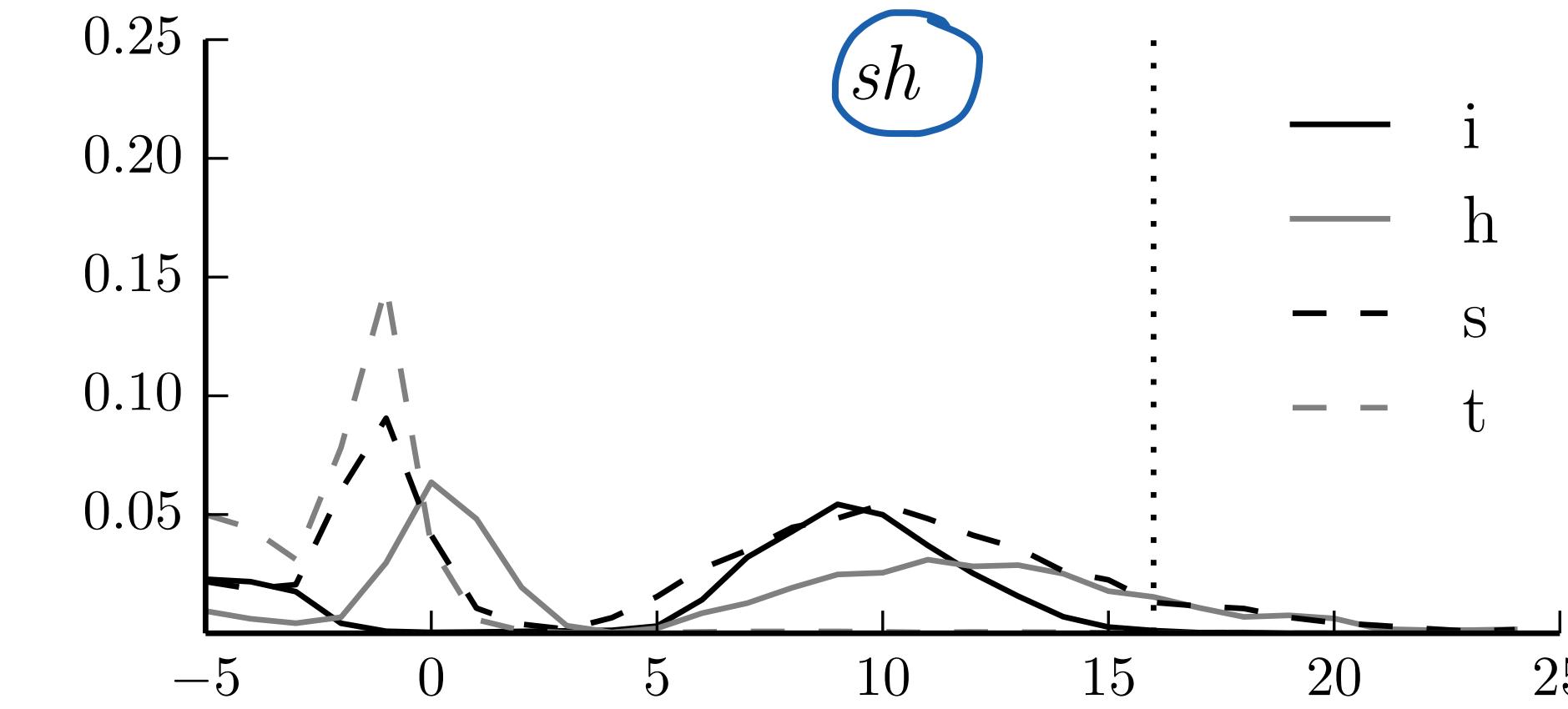
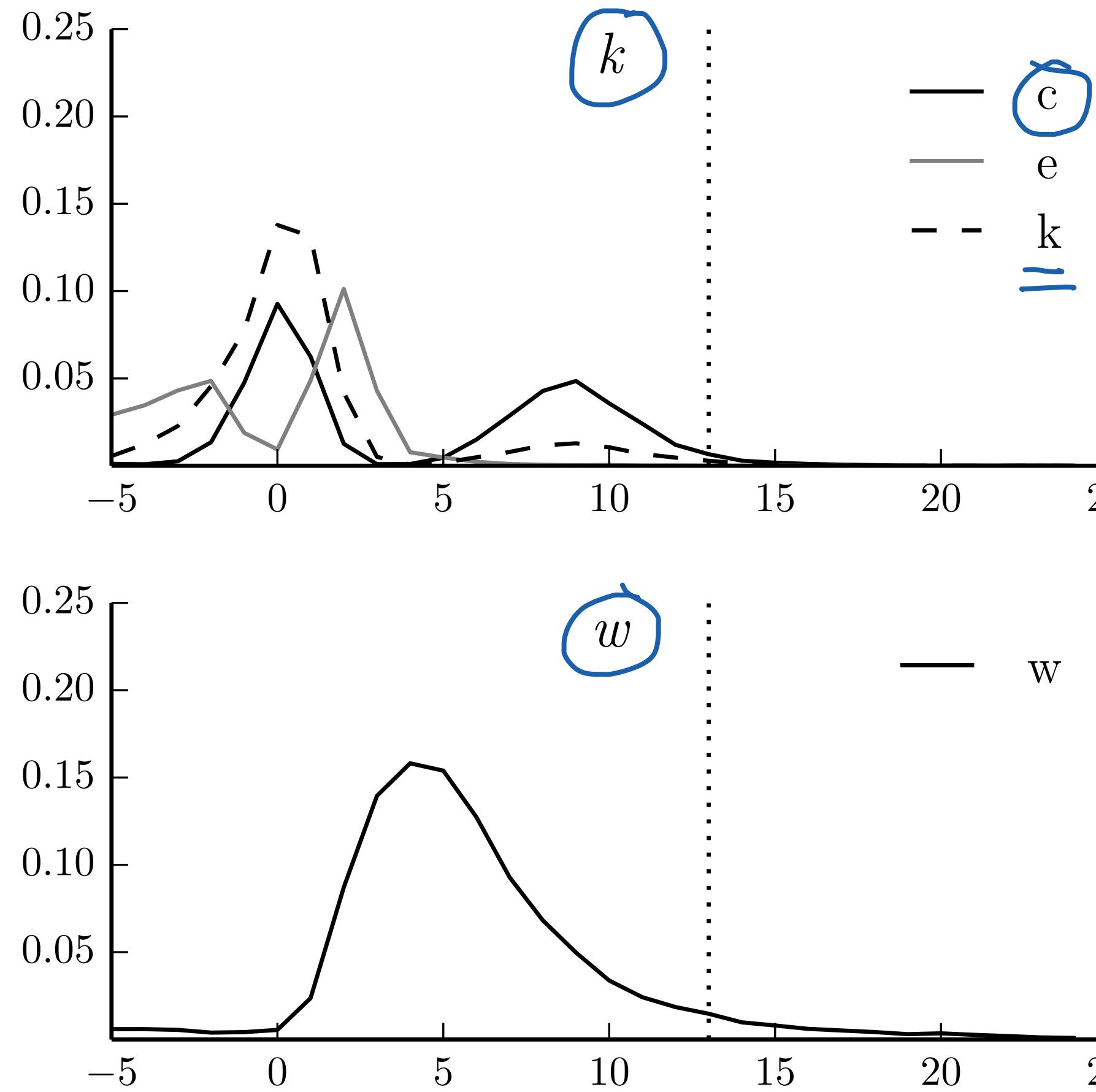
Analysing character probabilities



$s:$ -----o__hh----- -----y_eahh-----

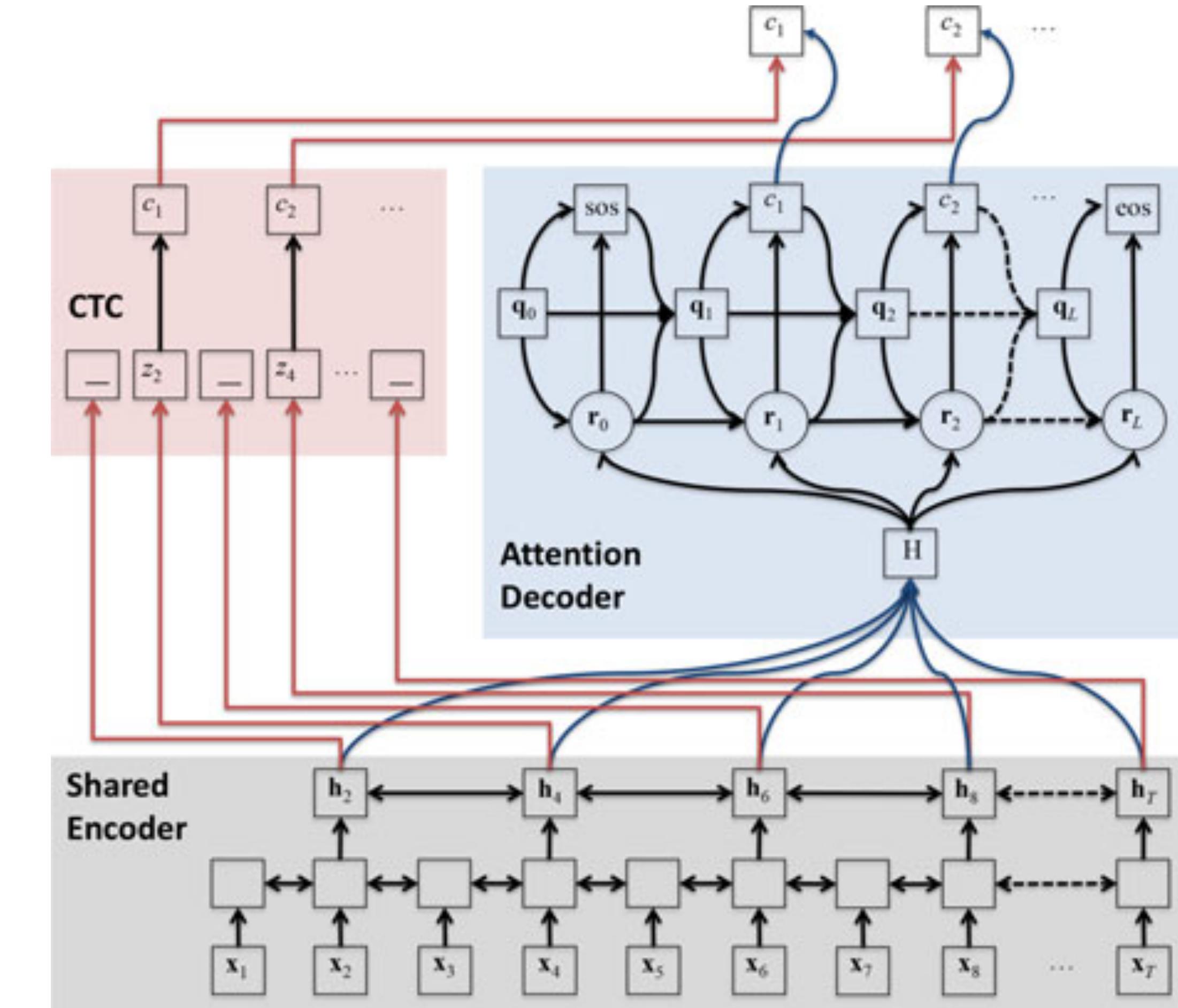
$\kappa(s):$ oh yeah

Character probabilities for various monophones



Hybrid CTC/Attention-based End-to-End ASR

- The objective to be maximised is a linear combination of the CTC and attention objective functions.



Hybrid CTC/Attention-based End-to-End ASR

- The objective to be maximised is a linear combination of the CTC and attention objective functions.

	Task1	Task2	Task3
Cascaded GMM/ HMM Model	11.2	9.2	12.1
Cascaded HMM/ DNN Model	9.0	7.2	9.6
Hybrid CTC/ Attention	8.4	6.1	6.9

