

# Connectionist Temporal Classification

Peter Bell

Automatic Speech Recognition – ASR Lecture 12  
26 February 2026

# End-to-end systems

- End-to-end systems are systems which learn to directly map from an input sequence  $X$  to an output sequence  $Y$ , estimating  $P(Y|X)$ 
  - $Y$  can be a sequence of words or subwords
- ML trained HMMs are kind of end-to-end system – the HMM estimates  $P(X|Y)$ , and when combined with a language model gives an estimate of  $P(Y|X)$
- Sequence discriminative training of HMMs (using GMMs or DNNs) can be regarded as end-to-end
  - But training is quite complicated – need to estimate the denominator (total likelihood) using lattices, first train conventionally (ML for GMMs, CE for NNs) then finetune using sequence discriminative training
  - Lattice-free MMI is one way to address these issues

Approaches based purely on recurrent networks or transformers which directly map input to output sequences

- CTC – Connectionist Temporal Classification
- Encoder-decoder approaches

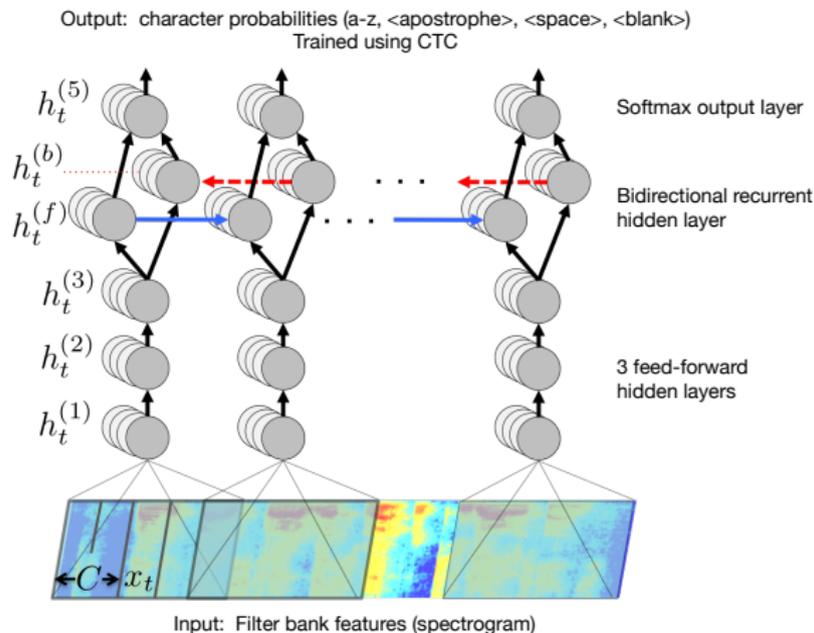
No need for specialised HMM-style decoders (although they can still be used)

Approaches based purely on recurrent networks or transformers which directly map input to output sequences

- **CTC – Connectionist Temporal Classification**
- Encoder-decoder approaches (*next lecture*)

No need for specialised HMM-style decoders (although they can still be used)

# Example: Deep Speech



Hannun et al (2014), "Deep Speech: Scaling up end-to-end speech recognition",

<https://arxiv.org/abs/1412.5567>.

# Deep Speech: Results

Model	SWB	CH	Full
Vesely et al. (GMM-HMM BMMI) [44]	18.6	33.0	25.8
Vesely et al. (DNN-HMM sMBR) [44]	12.6	24.1	18.4
Maas et al. (DNN-HMM SWB) [28]	14.6	26.3	20.5
Maas et al. (DNN-HMM FSH) [28]	16.0	23.7	19.9
Seide et al. (CD-DNN) [39]	16.1	n/a	n/a
Kingsbury et al. (DNN-HMM sMBR HF) [22]	13.3	n/a	n/a
Sainath et al. (CNN-HMM) [36]	11.5	n/a	n/a
Soltau et al. (MLP/CNN+I-Vector) [40]	<b>10.4</b>	n/a	n/a
<b>Deep Speech SWB</b>	20.0	31.8	25.9
<b>Deep Speech SWB + FSH</b>	12.6	<b>19.3</b>	<b>16.0</b>

Table 3: Published error rates (%WER) on Switchboard dataset splits. The columns labeled “SWB” and “CH” are respectively the easy and hard subsets of Hub5’00.

- Maps from acoustic frames  $X$  to subword sequences  $Y$ , where  $Y$  is a sequence of characters (in some other CTC approaches,  $Y$  can be a sequence of phones)
- CTC loss function
- Makes good use of large training data
  - Synthetic additional training data by jittering the signal and adding noise
- Many computational optimisations
- n-gram language model to impose word-level constraints
- Competitive results on standard tasks

- Maps from acoustic frames  $X$  to subword sequences  $Y$ , where  $Y$  is a sequence of characters (in some other CTC approaches,  $Y$  can be a sequence of phones)
- **CTC loss function**
- Makes good use of large training data
  - Synthetic additional training data by jittering the signal and adding noise
- Many computational optimisations
- n-gram language model to impose word-level constraints
- Competitive results on standard tasks

# Connectionist Temporal Classification (CTC)

- Train a recurrent network (or transformer) to map from input sequence  $X$  to output sequence  $Y$ 
  - sequences can be different lengths – for speech, input sequence  $X$  (acoustic frames) is much longer than output sequence  $Y$  (characters or phonemes)
  - CTC does not require frame-level alignment (matching each input frame to an output token)
- CTC sums over all possible alignments (similar to forward-backward algorithm) – “alignment free”
- Possible to back-propagate gradients through CTC loss function

Good overview of CTC: Awni Hannun, “Sequence Modeling with CTC”, *Distill*. <https://distill.pub/2017/ctc>

# CTC: Alignment

- Imagine mapping  $(x_1, x_2, x_3, x_4, x_5, x_6)$  to  $[a, b, c]$ 
  - Possible alignments:  $aaabbc, aabbcc, abbbbc, \dots$
- However
  - Don't always want to map every input frame to an output symbol (e.g. if there is "inter-symbol silence")
  - Want to be able to have two identical symbols adjacent to each other – keep the difference between
- Solve this using an additional *blank* symbol ( $\epsilon$ )
- CTC output compression
  - 1 Merge repeating characters
  - 2 Remove blanks

Thus to model the same character successively, separate with a blank

- Some possible alignments for  $[h, e, l, l, o]$  and  $[h, e, l, o]$  given a 10-element input sequence
  - $[h, e, l, l, o]: h\epsilon\epsilon\epsilon\ell\ell\epsilon\ell o; h\epsilon\ell\ell\epsilon\ell\epsilon o o$
  - $[h, e, l, o]: h\epsilon\epsilon\epsilon\ell\ell\ell\ell o; h h\epsilon\epsilon\epsilon\ell\epsilon\epsilon o\epsilon$

# CTC: Alignment example

h h e  $\epsilon$   $\epsilon$  l l l  $\epsilon$  l l o

h e  $\epsilon$  l  $\epsilon$  l o

h e l l o

h e l l o

First, merge repeat characters.

Then, remove any  $\epsilon$  tokens.

The remaining characters are the output.

# CTC: Valid and invalid alignments

Consider an output [c, a, t] with an input of length six

## Valid Alignments

€ c c € a t

c c a a t t

c a € € € t

## Invalid Alignments

c € c € a t

c c a a t   

c € € € | t t

corresponds to  
 $Y = [c, c, a, t]$

has length 5

missing the 'a'

# CTC: Alignment properties

- Monotonic – Alignments are monotonic (left-to-right model); no re-ordering (unlike neural machine translation)
- Many-to-one – Alignments are many-to-one; many inputs can map to the same output
- But a single input cannot map to many outputs – could be a problem for sounds like “th”
- CTC doesn't find a single alignment: it sums over all possible alignments

# CTC: Loss function (1)

- Let  $C$  be an output label sequence, including blanks and repetitions – same length as input sequence  $X$
- Posterior probability of output labels  $C = (c_1, \dots, c_t, \dots, c_T)$  given the input sequence  $X = (x_1, \dots, x_t, \dots, x_T)$ :

$$P(C|X) = \prod_{t=1}^T P_t(c_t|X)$$

where  $P_t(c_t|X)$  is the probability of outputting label  $c_t$  at time  $t$

- This is the probability of a single alignment – we need to sum over all alignments consistent with  $Y$

## CTC: Loss function (2)

- Let  $Y$  be the compressed target output sequence
- Compute the posterior probability of the target sequence  $Y = (s_1, \dots, s_m, \dots, s_M)$  ( $M \leq T$ ) given  $X$  by summing over the possible CTC alignments:

$$P(Y|X) = \sum_{C \in A(Y)} P(C|X)$$

where  $A$  is the set of possible output label sequences  $C$  that can be mapped to  $Y$  using the CTC compression rules (merge repeated labels, then remove blanks)

- The CTC loss function  $\mathcal{L}_{CTC}$  is given by the negative log likelihood of the sum of CTC alignments:

$$\mathcal{L}_{CTC} = -\log P(Y|X)$$

- Various NN architectures can be used for CTC – traditionally used a deep bidirectional LSTM RNN

## CTC: Dynamic programming

Perform the sum over alignments,  $A(Y)$ , using dynamic programming – very similar to the forward algorithm for classic HMMs.

We first define the expanded symbol sequence,

$$Z = (z_1, \dots, z_i, \dots, z_J) = (\epsilon, s_1, \epsilon, s_2, \epsilon, \dots, \epsilon, s_M, \epsilon)$$

(where  $J = 2M + 1$ )

The forward probability is:

$$\begin{aligned}\alpha_j(t) &= P(z_1, \dots, z_j | X) \\ &= \sum_{(c_1, \dots, c_t) \in A(z_1, \dots, z_j)} P(c_1, \dots, c_t | X)\end{aligned}$$

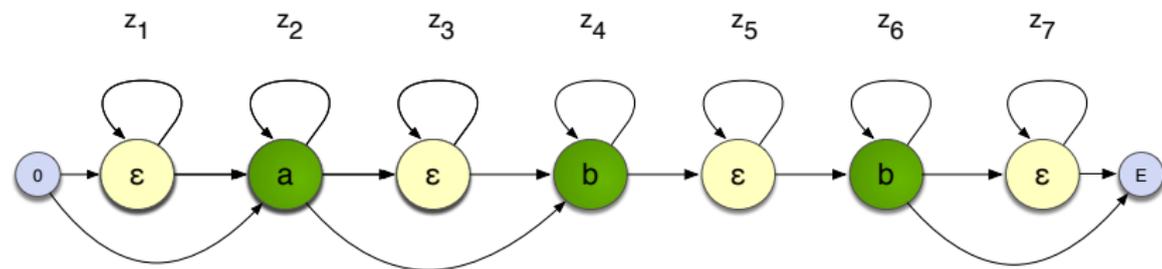
This computes the probability over all label sequences up to time  $t$  that are consistent with  $(z_1, \dots, z_j)$ .

# CTC: HMM topology

We can encode the valid transitions of  $Z$  over time using an HMM.

This is a standard left-to-right HMM topology, with the addition of a skip  $z_{i-2} \rightarrow z_i$  if  $z_i \neq \epsilon$  and  $z_i \neq z_{i-2}$

Example for original sequence  $Y = [a, b, b]$ :



- **Initialisation:**

$$\begin{aligned}\alpha_i(0) &= 1 & i &= 1 \\ &= 0 & & \textit{otherwise}\end{aligned}$$

- **Recursion:**

If  $z_i = \epsilon$  or  $z_i = z_{i-2}$ :

$$\alpha_i(t) = [\alpha_{i-1}(t-1) + \alpha_i(t-1)]p_t(z_i|X)$$

Otherwise:

$$\alpha_i(t) = [\alpha_{i-2}(t-1) + \alpha_{i-1}(t-1) + \alpha_i(t-1)]p_t(z_i|X)$$

- **Termination:**

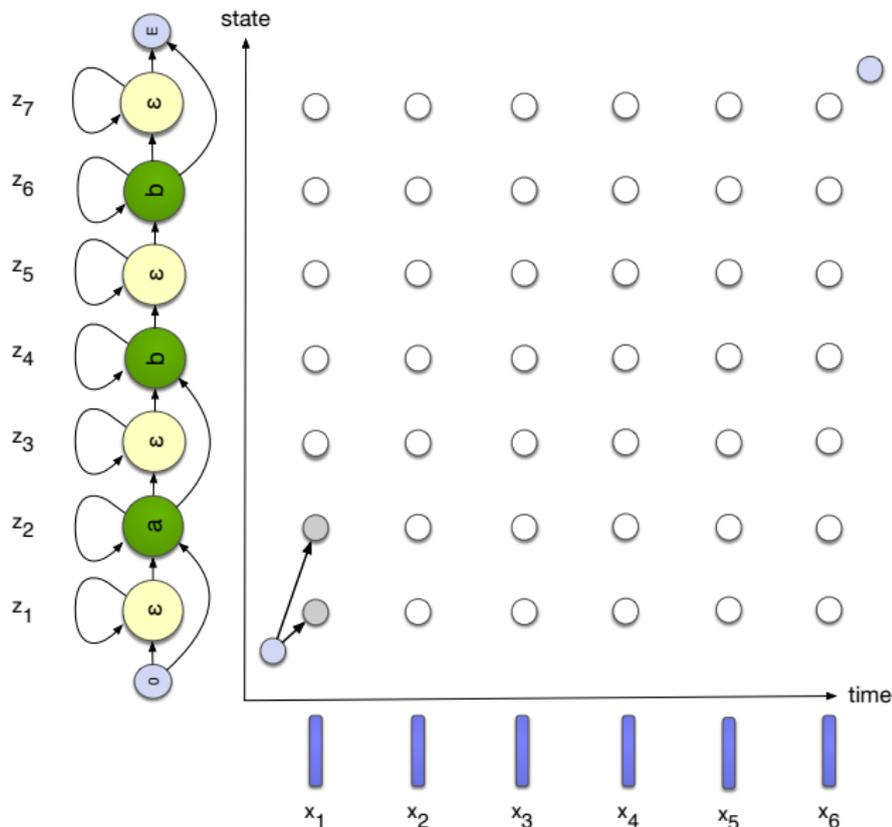
$$P(Z|X) = \alpha_{J-1}(t) + \alpha_J(t)$$

# Example

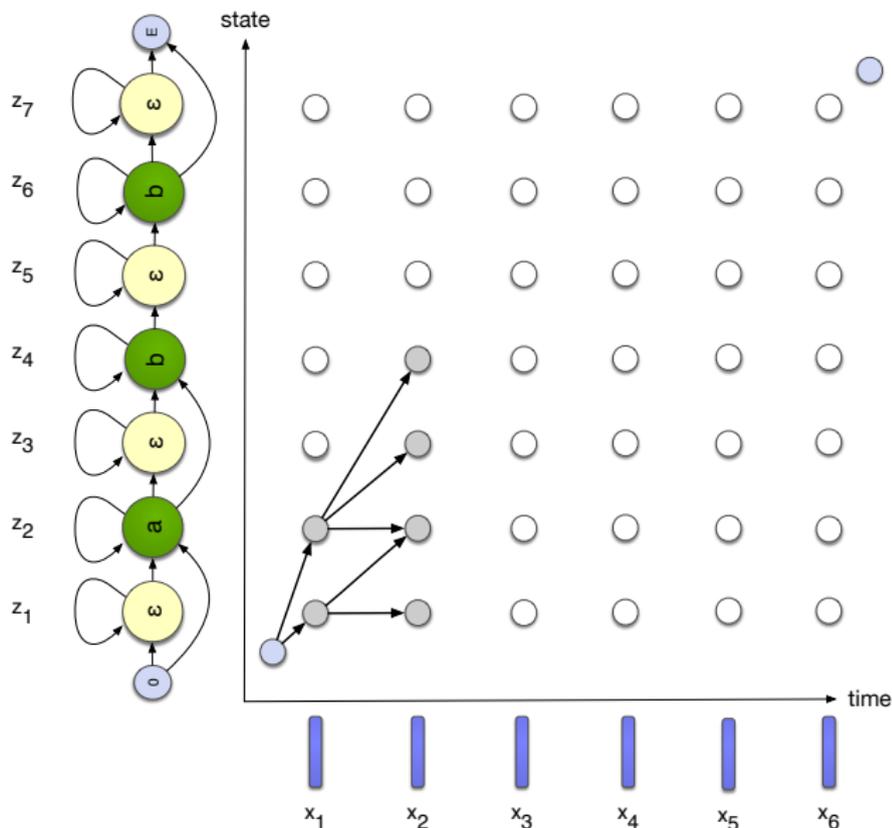
Example alignments for  $[a, b, b]$  to an utterance of six frames:

a	$\epsilon$	b	$\epsilon$	b	b
a	a	$\epsilon$	b	$\epsilon$	b
$\epsilon$	a	a	b	$\epsilon$	b
$\epsilon$	a	b	$\epsilon$	b	$\epsilon$
a	b	$\epsilon$	$\epsilon$	b	b
$\vdots$	$\vdots$	$\vdots$			

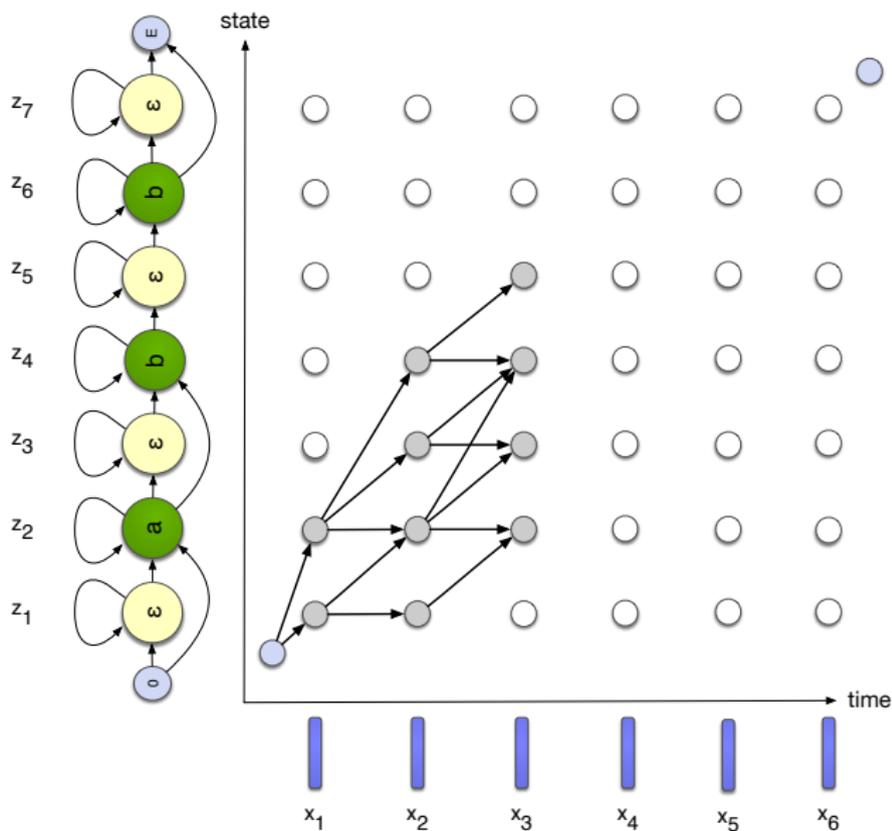
# Forward recursion



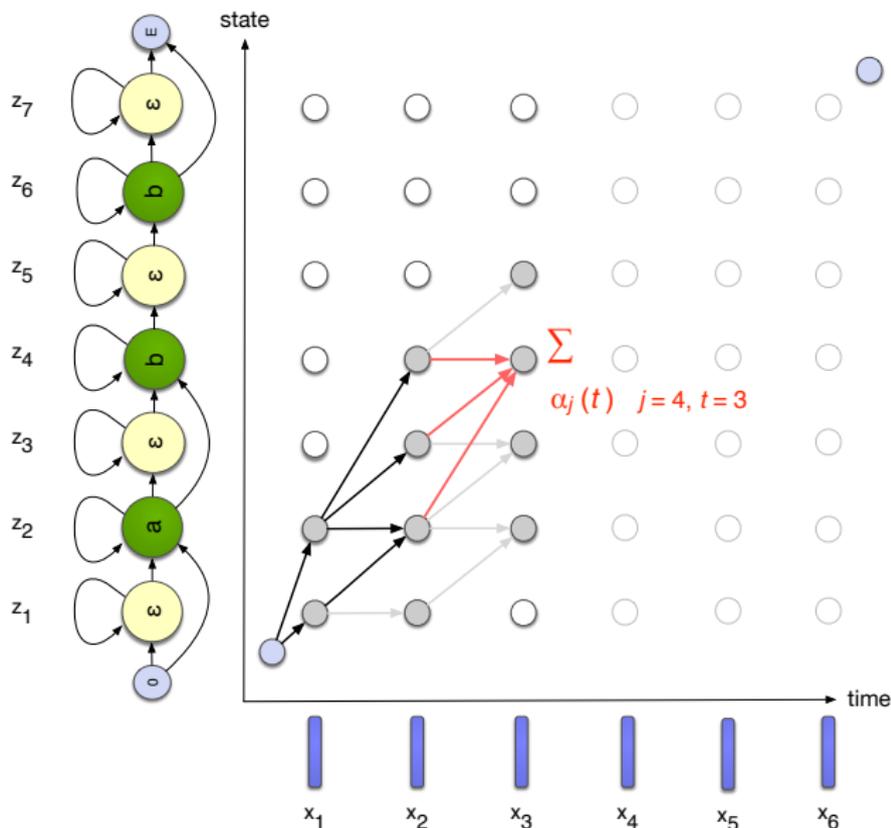
# Forward recursion



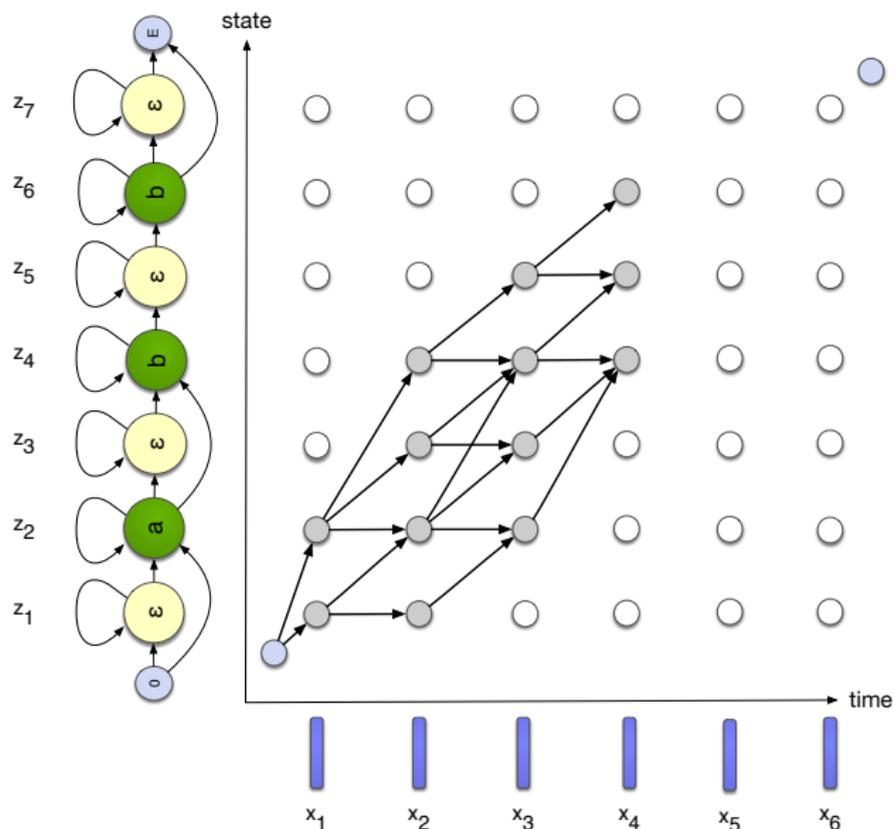
# Forward recursion



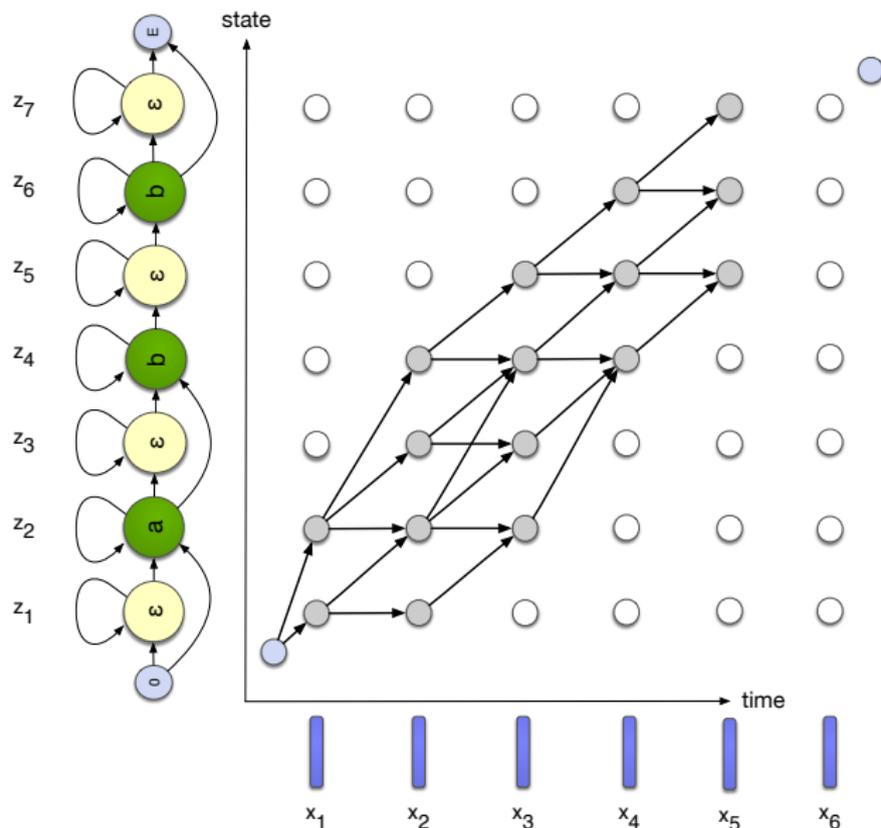
# Forward recursion



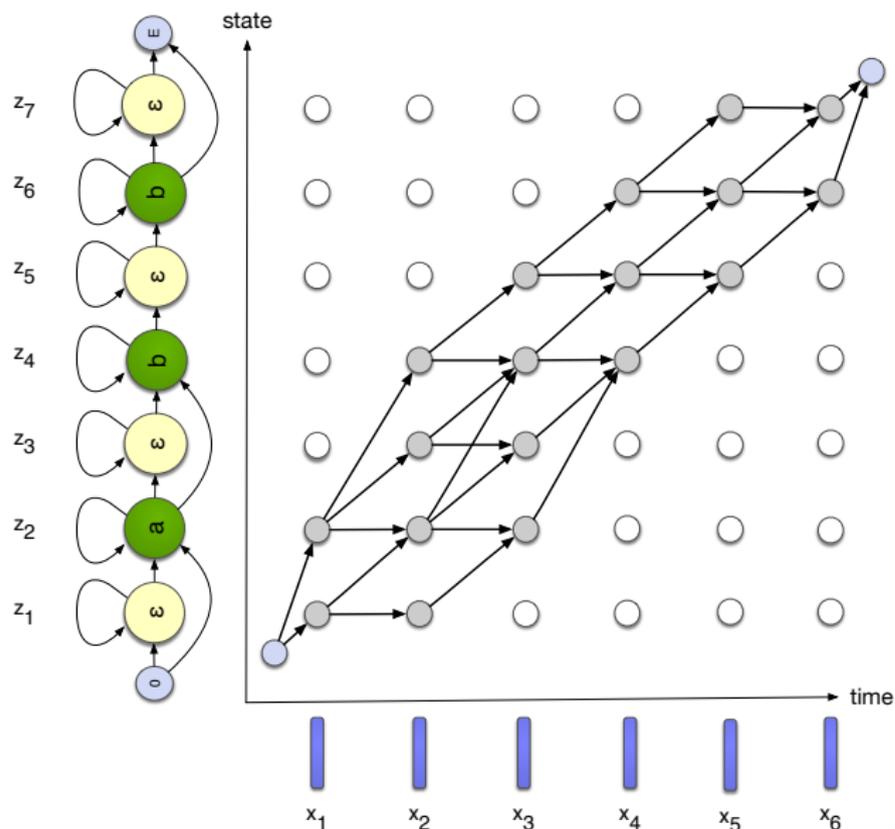
# Forward recursion



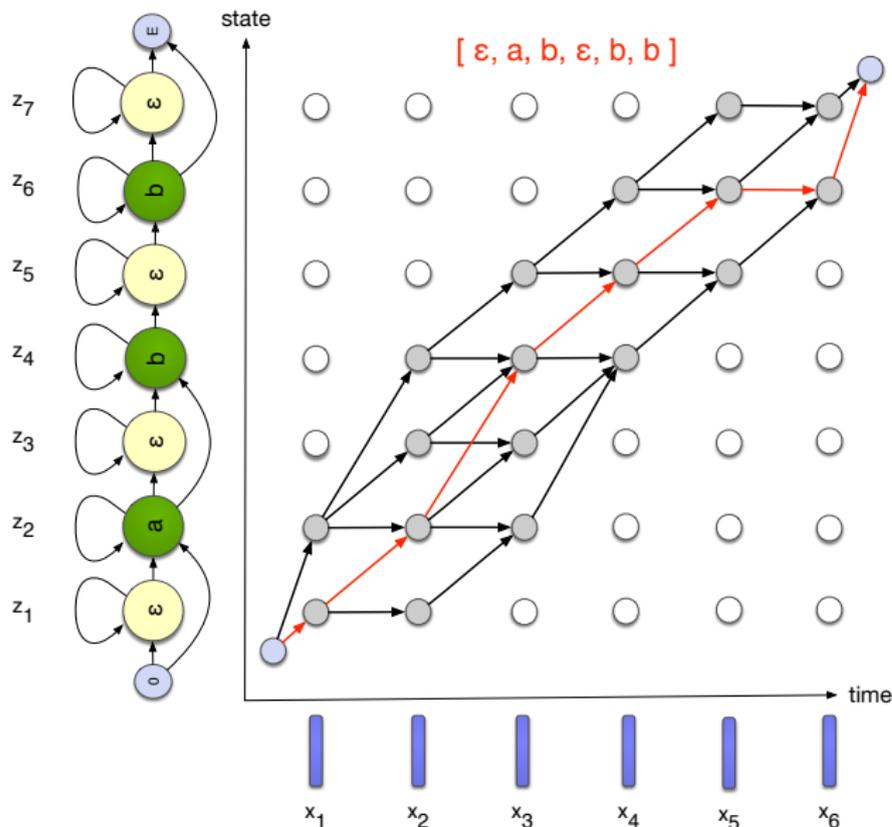
# Forward recursion



# Forward recursion



# One alignment...



Need to solve

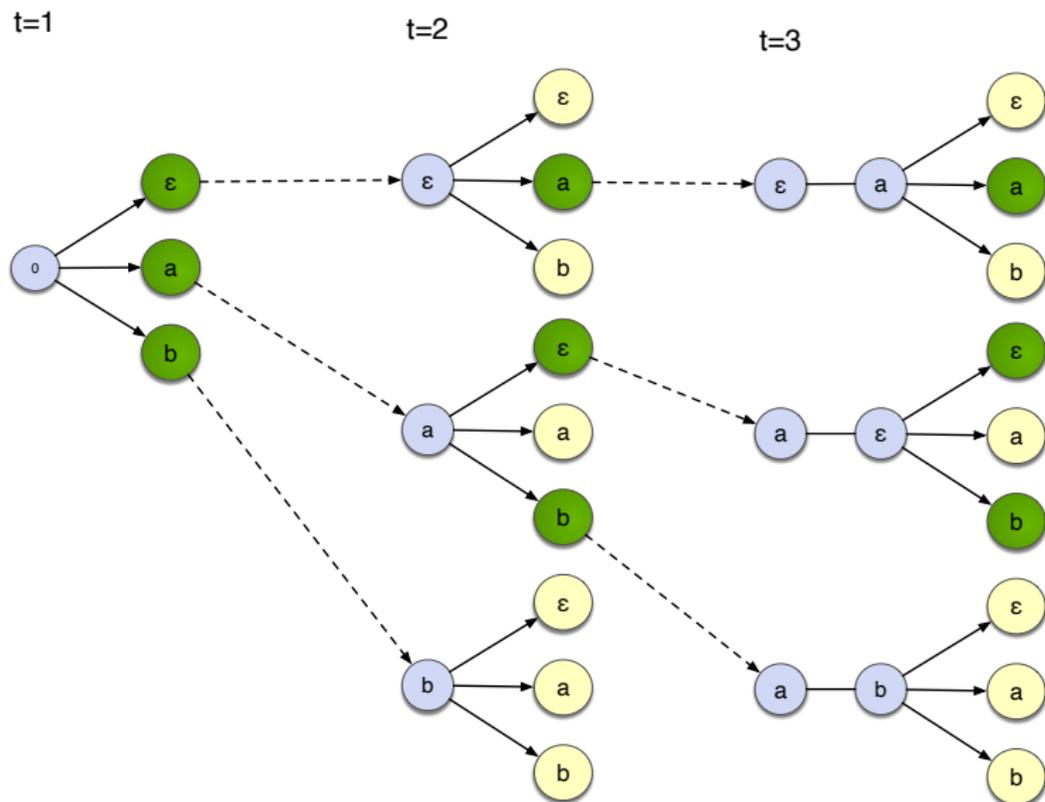
$$Y^* = \arg \max_Y P(Y|X)$$

Find best alignment:

$$C^* = \arg \max_C \prod_t^T P(c_t|X)$$

Solve using beam search

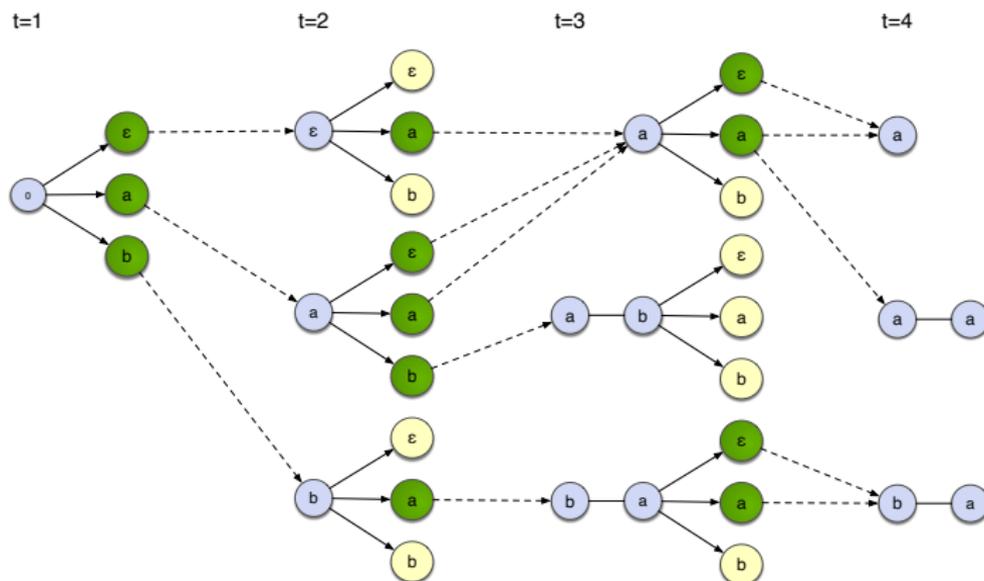
# CTC: Decoding with beam search



# Merge hypotheses

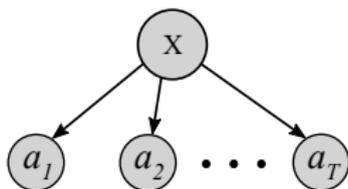
Problem: the best-scoring sequence  $C^*$  might not correspond to the overall best  $Y$

Solution: merge hypotheses with the same prefix



# Understanding CTC: Conditional independence assumption

- Each output is dependent on the entire input sequence (in Deep Speech this is achieved using a bidirectional recurrent layer)
- Given the inputs, each output is independent of the other outputs (conditional independence)
- CTC does not learn a language model over the outputs, although a language model can be applied later
- Graphical model showing dependences in CTC:



# Applying language models to CTC

- Direct interpolation of a language model with the CTC acoustic model:

$$\hat{W} = \arg \max_W (\log P(Y|X) + \lambda \log P(W)) + \eta L(W)$$

Only consider word sequences  $W$  which correspond to the sub-word sequence  $Y$  (using a lexicon)

- $\lambda, \eta$  are empirically determined scaling factor/insertion bonus
- Lexicon-free CTC: use a sub-word language model  $P(Y)$  (Maas et al, 2015)
- WFST implementation: create an FST  $T$  which transforms a framewise label sequence  $\mathbf{c}$  into the subword sequence  $Y$ , then compose with  $L$  and  $G$ :  $T \circ \min(\det(L \circ G))$  (Miao et al, 2015)

- Mozilla have released an Open Source TensorFlow implementation of the Deep Speech architecture:
- <https://hacks.mozilla.org/2017/11/a-journey-to-10-word-error-rate/>
- <https://github.com/mozilla/DeepSpeech>
- Close to state-of-the-art results on librispeech
- Mozilla Common Voice project: <https://voice.mozilla.org/en>

# Summary and reading

- CTC is an alternative approach to sequence discriminative training, typically applied to RNN systems
- Used in “Deep Speech” architecture for end-to-end speech recognition
- Reading
  - A Hannun et al (2014), “Deep Speech: Scaling up end-to-end speech recognition”, ArXiv:1412.5567.  
<https://arxiv.org/abs/1412.5567>
  - A Hannun (2017), “Sequence Modeling with CTC”, *Distill*.  
<https://distill.pub/2017/ctc>
- Background reading
  - Y Miao et al (2015), “EESSEN: End-to-end speech recognition using deep RNN models and WFST-based decoding”, ASRU-2105. <https://ieeexplore.ieee.org/abstract/document/7404790>
  - A Maas et al (2015). “Lexicon-free conversational speech recognition with neural networks”, NAACL HLT 2015, <http://www.aclweb.org/anthology/N15-1038>