



SCHOOL OF DATA ANALYSIS

## ASR II: Discriminative State-Space Models I

Andrey Malinin

14th March 2022

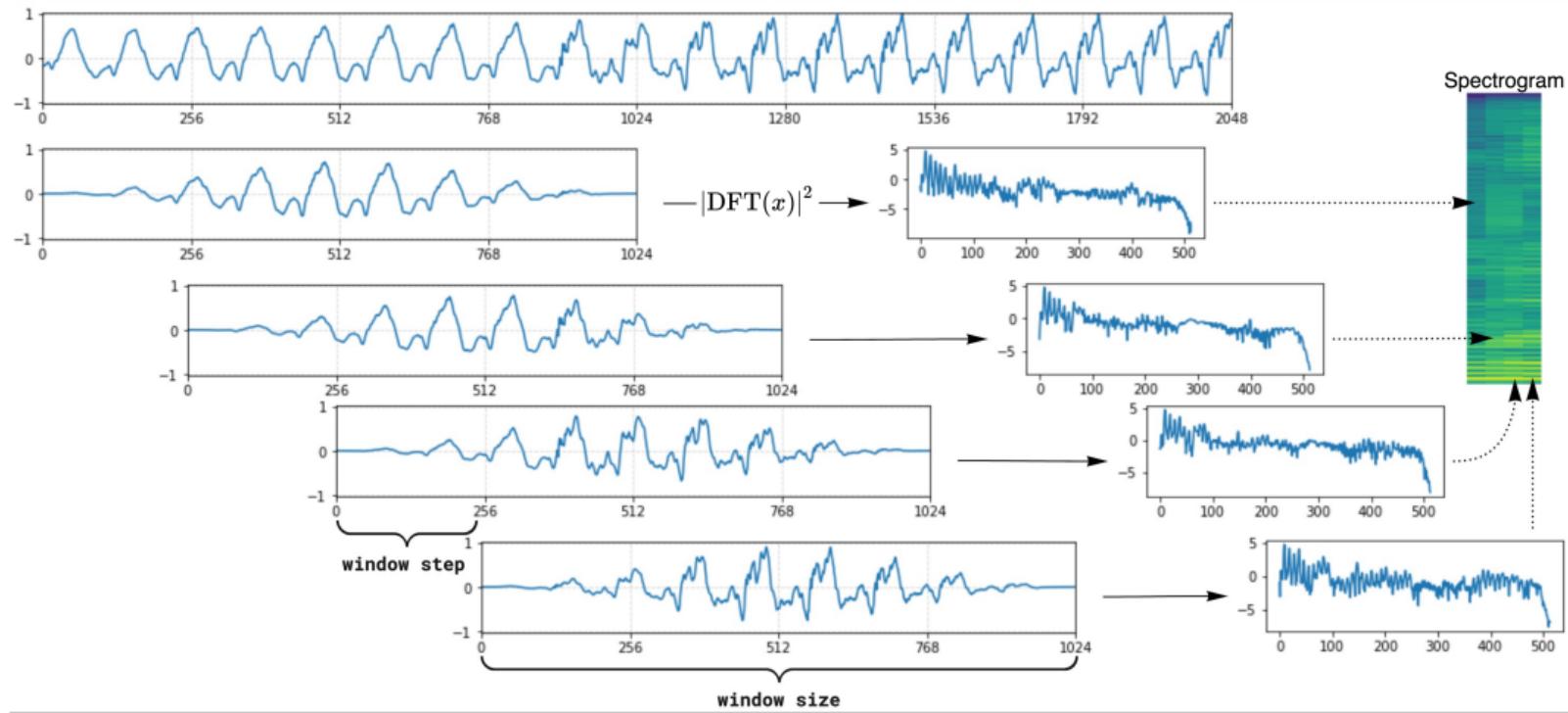
## Story so far

In this previous episode...

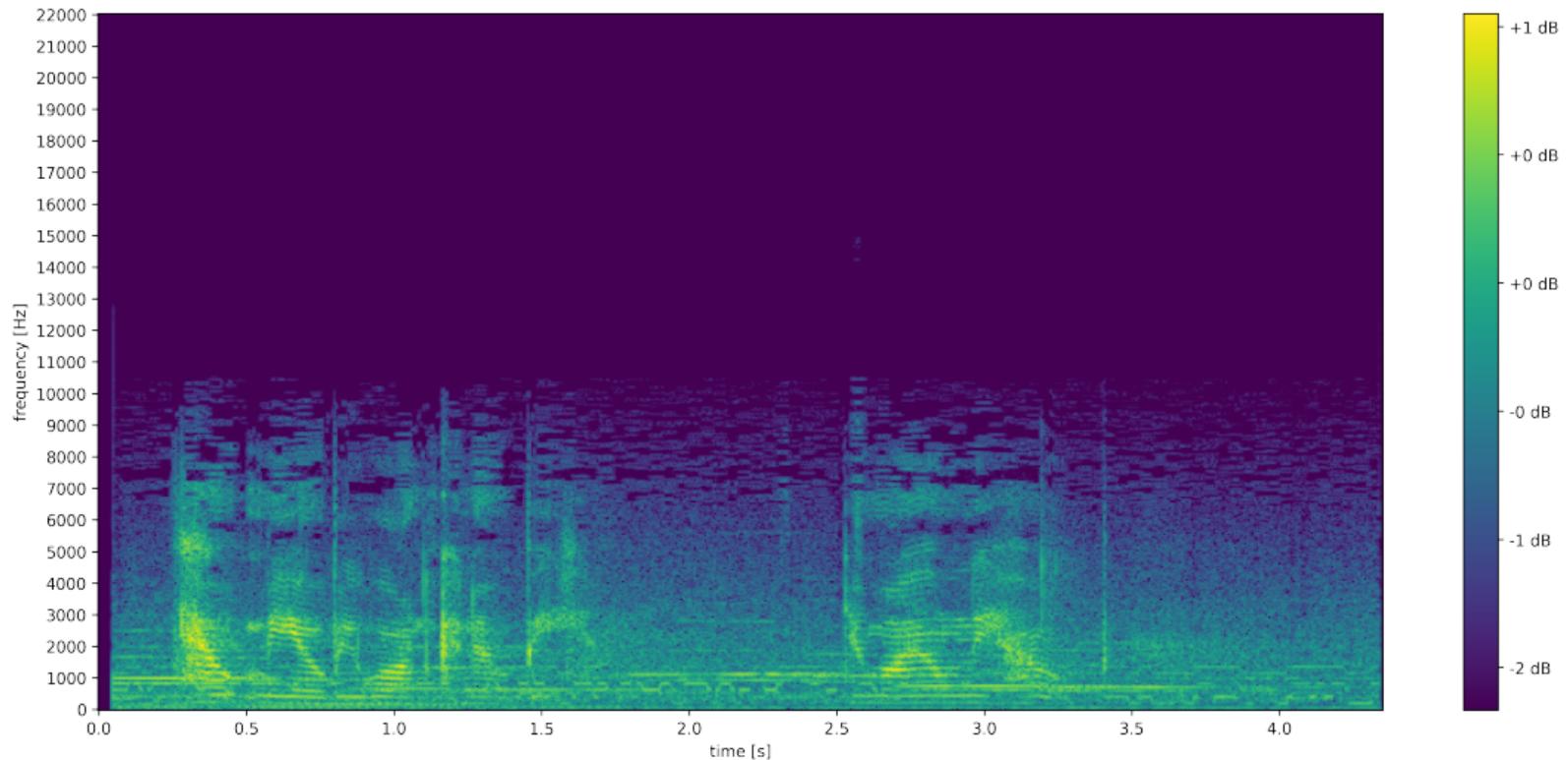
## Recap - Data Processing

- Transform audio into Mel-Spectrogram
  - Yields sequence of acoustic features or **frames**  $\mathbf{X}_{1:T}$
  - This decreases the time-resolution (downsamples) of the audio data
  - Length of each frame depends on audio sample-rate and hop-length
- Transform text into a sequence of **speech units**  $\omega_{1:L}$ 
  - Speech units can be phonemes, graphemes, syllables, etc..
  - Speech units must be compact and easy to back to / from words

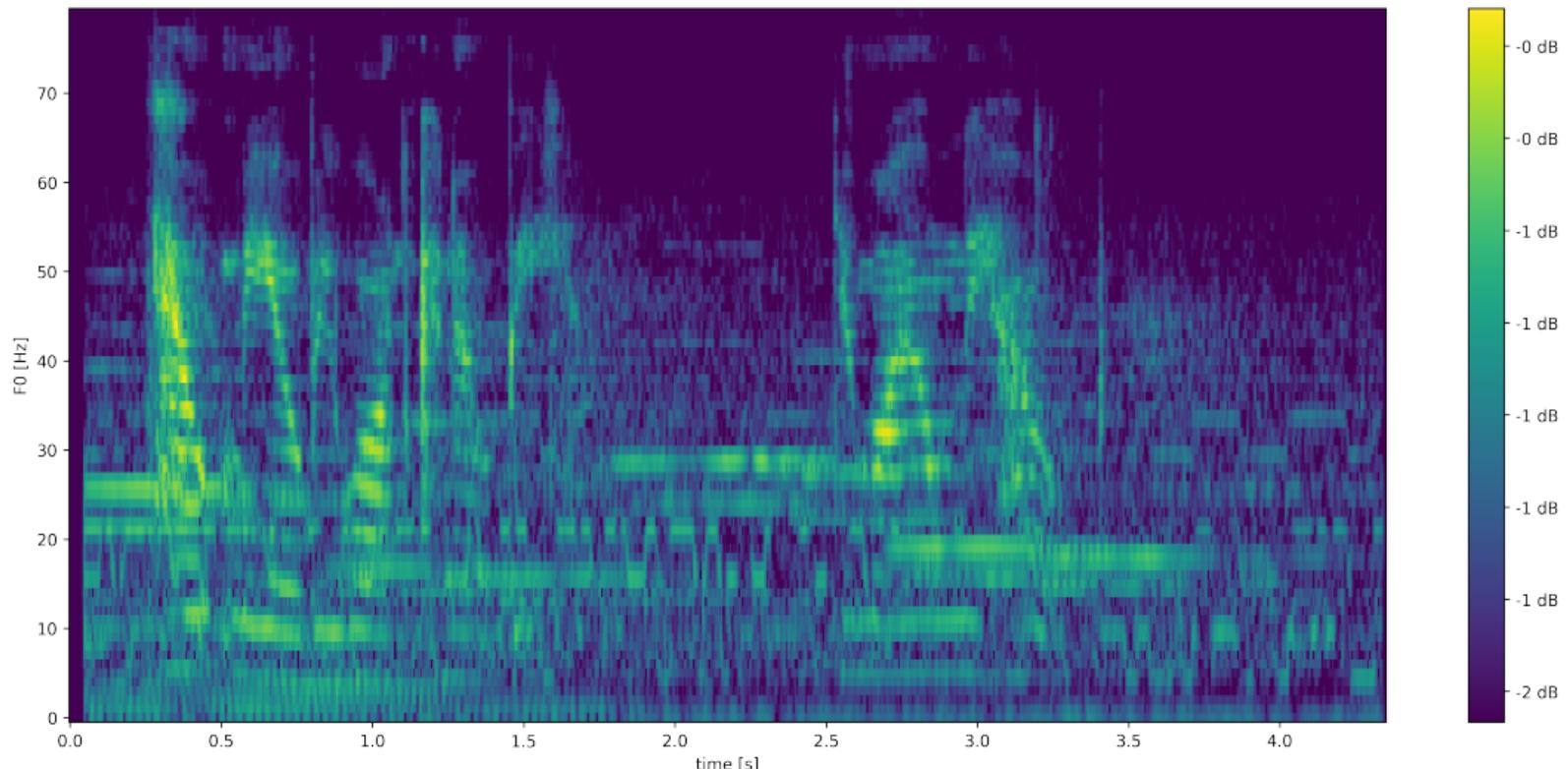
# DSP Reminder - Audio to Spectrogram via STFT



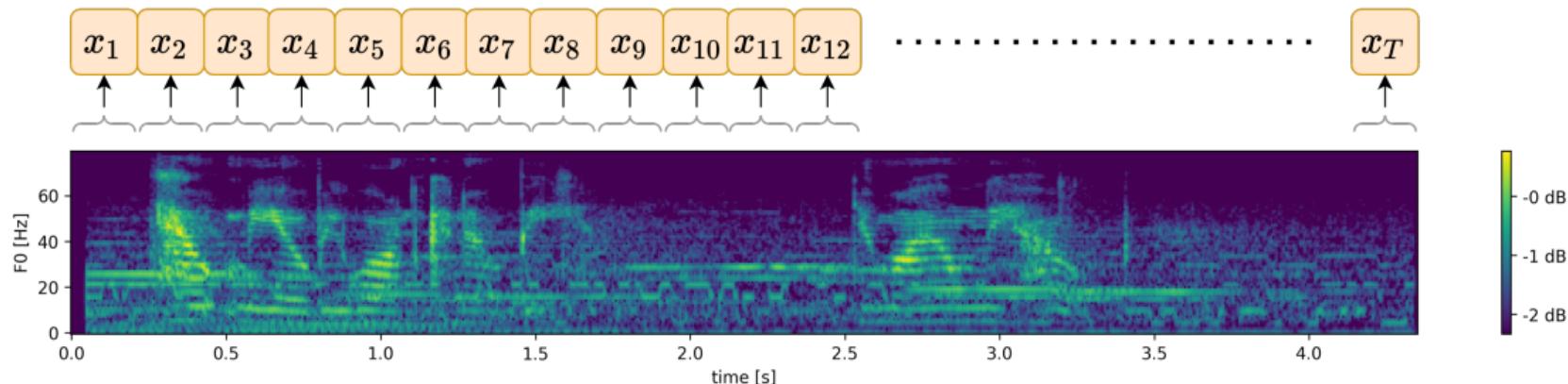
# DSP Reminder - Spectrogram



# DSP Reminder - Mel-Spectrogram



## DSP Reminder - Acoustic Features



Extract **acoustic features**, **observations** or **frames**  $\mathbf{X}_{1:T} = \{x_1, \dots, x_T\}$

- At a 22.05KHz sample-rate each second is 22050 samples long!
- DFT down-samples → At 512 sample window-step each  $x_t$  is  $\sim 25$ ms of audio.
- Acoustic sequence  $\mathbf{X}_{1:T}$  is long - 43 frames / second @22KHz & 512 frame overlap

Can use either each time-step or further down-sample via convolutions, for example.

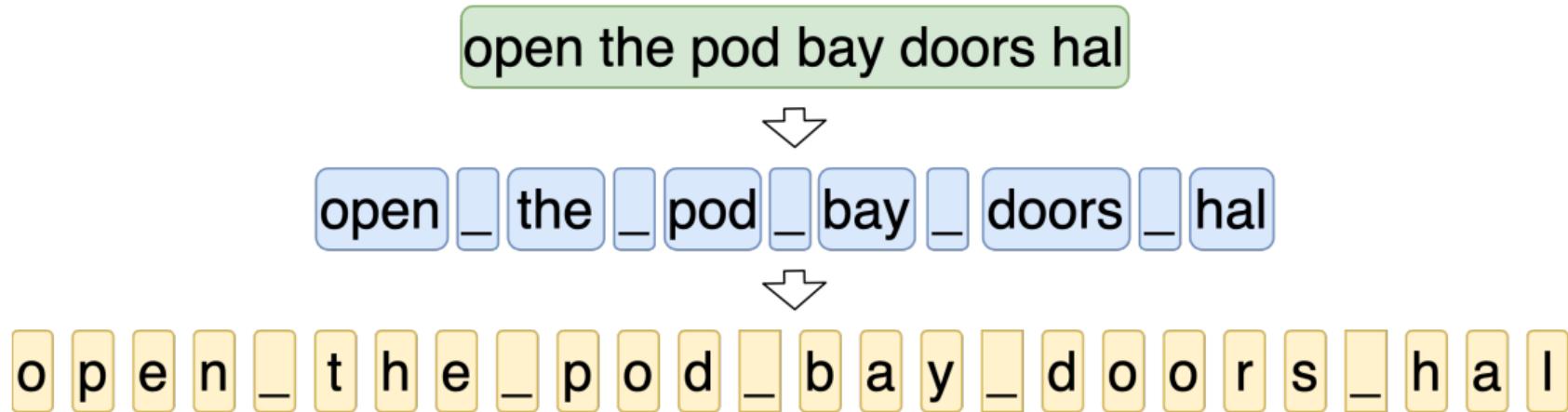
## Text Processing - sub-word speech units

- Define a mapping  $\mathcal{M}$  between words  $\mathbf{w}$  and speech units  $\omega_{1:L}$ .

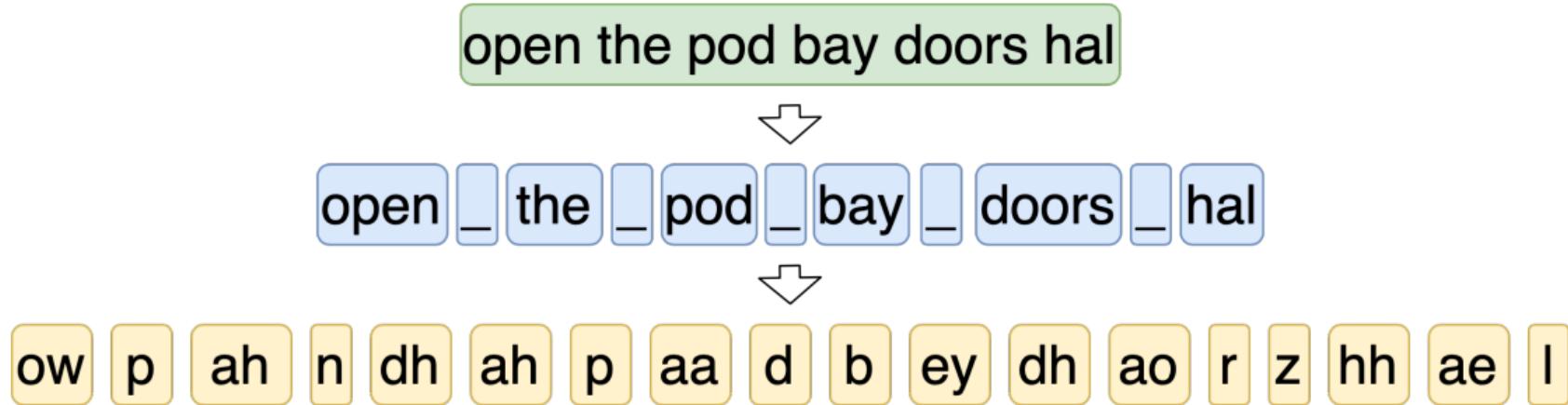
$$\{\omega_{1:L_q}^{(q)}\}_{q=1}^Q = \mathcal{M}(\mathbf{w}), \quad \{\mathbf{w}^{(p)}\}_{p=1}^P = \mathcal{M}^{-1}(\omega_{1:L})$$

- For some choices of speech units this mapping is not 1-to-1 ( $Q>1$ ,  $P>1$ )
  - Different word pronunciations → different phoneme sequences
  - Different graphemic spellings (centre vs center)
- Desirable properties of speech units →
  - Well defined, simple mapping between words and speech units
  - Compact, even for large vocabulary, and extensible to OOV words.
  - Account for phonetic variability
- Choice depends on language properties, size of training data, type of ASR system

## Mapping between words and speech units



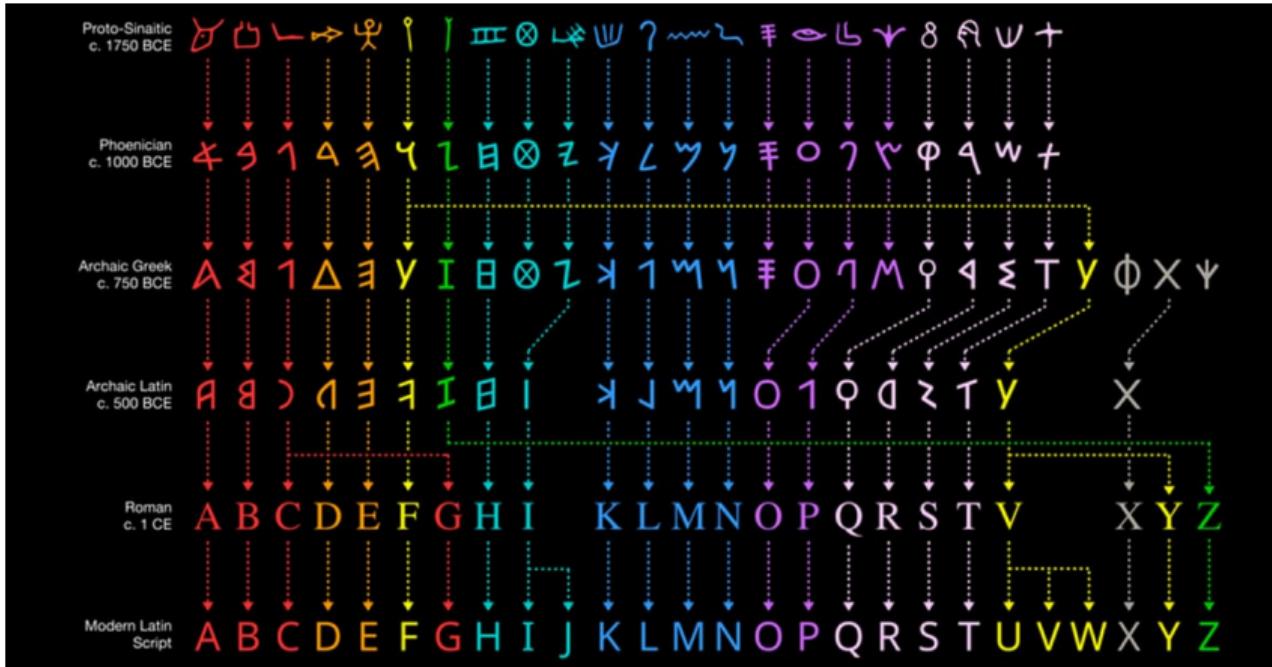
## Mapping between words and speech units



# Syllabic (basic Katakana)

	a	i	u	e	o
Vowels	ア a	イ i	ウ u	エ e	オ o
k	カ ka	キ ki	ク ku	ケ ke	コ ko
s	サ sa	シ shi	ス su	セ se	ソ so
t	タ ta	チ chi	ツ tsu	テ te	ト to
n	ナ na	ニ ni	ヌ nu	ネ ne	ノ no
h	ハ ha	ヒ hi	フ fu	ヘ he	ホ ho
m	マ ma	ミ mi	ム mu	メ me	モ mo
y	ヤ ya		ユ yu		ヨ yo
r	ラ ra	リ ri	ル ru	レ re	ロ ro
w	ワ wa				ヲ wo
		ン n			

## Graphemic (Alphabet)



Matt Baker. Usefulcharts.com

# Phonetic (IPA)

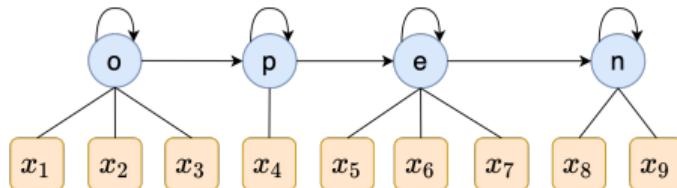
VOWELS	monophthongs				diphthongs		<b>Phonemic Chart</b> voiced unvoiced
	i: <u>sheep</u>	I <u>ship</u>	ʊ <u>good</u>	u: <u>shoot</u>	ɪə <u>here</u>	eɪ <u>wait</u>	
	e <u>bed</u>	ə <u>teacher</u>	ɜ: <u>bird</u>	ɔ: <u>door</u>	ʊə <u>tourist</u>	ɪɔ <u>boy</u>	
	æ <u>cat</u>	ʌ <u>up</u>	ɑ: <u>far</u>	ɒ <u>on</u>	eə <u>hair</u>	aɪ <u>my</u>	au <u>cow</u>
CONSONANTS	p <u>pea</u>	b <u>boat</u>	t <u>tea</u>	d <u>dog</u>	tʃ <u>cheese</u>	dʒ <u>June</u>	k <u>car</u>
	f <u>fly</u>	v <u>video</u>	θ <u>think</u>	ð <u>this</u>	s <u>see</u>	z <u>zoo</u>	ʃ <u>shall</u>
	m <u>man</u>	n <u>now</u>	ŋ <u>sing</u>	h <u>hat</u>	l <u>love</u>	r <u>red</u>	w <u>wet</u>
							j <u>yes</u>

## Recap - Data Processing

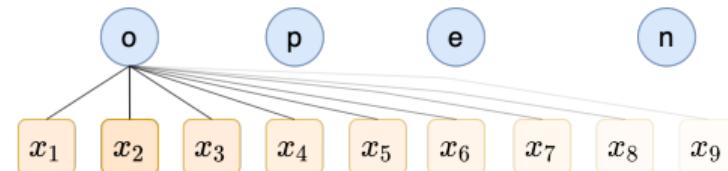
- Transform audio into Mel-Spectrogram
  - Yields sequence of acoustic features or **frames**  $\mathbf{X}_{1:T}$
  - This decreases the time-resolution (downsamples) of the audio data
  - Length of each frame depends on audio sample-rate and hop-length
- Transform text into a sequence of **speech units**  $\omega_{1:L}$ 
  - Speech units can be phonemes, graphemes, syllables, etc..
  - Speech units must be compact and easy to back to / from words

# Problem of Alignment

- Which feature vectors  $X_{1:T}$  and speech units  $\omega_{1:L}$  relate or **align** to each other?
  - Models need to be able to dynamically align features vectors to words.
- Two common approaches to constructing models which can align
  - State-Space models
  - **Neural Attention Mechanisms**

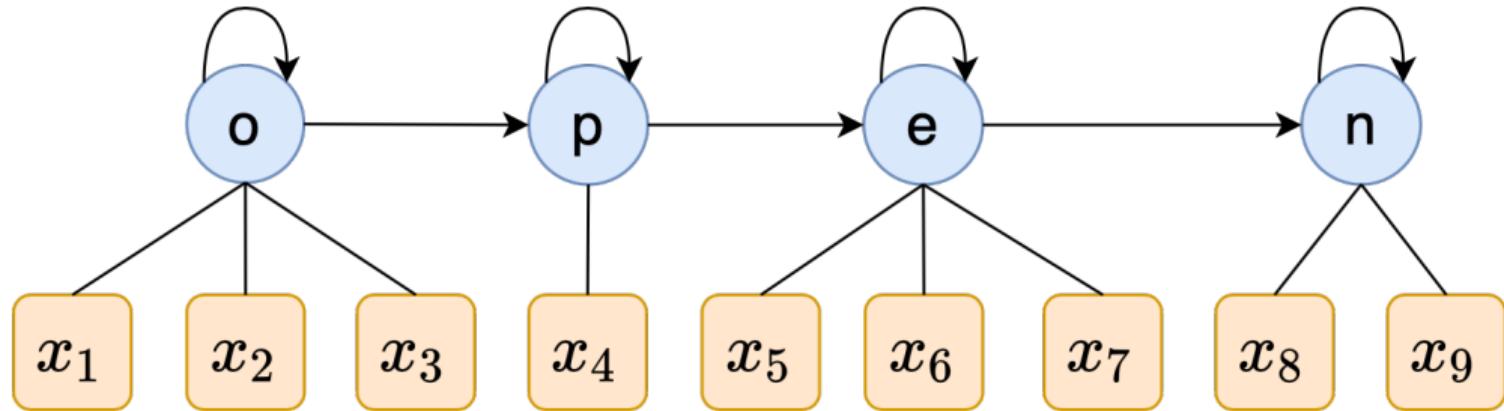


(a) State-Space

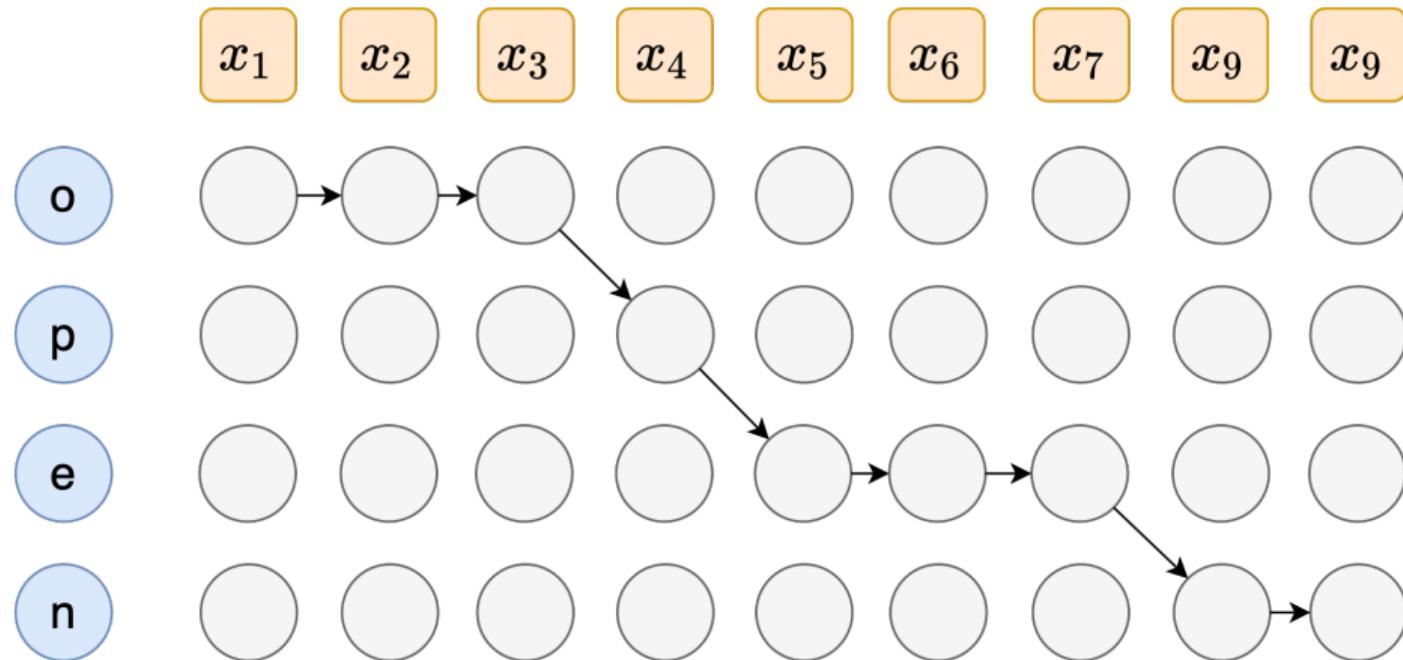


(b) Attention Mechanism

# State-Space Alignment Models

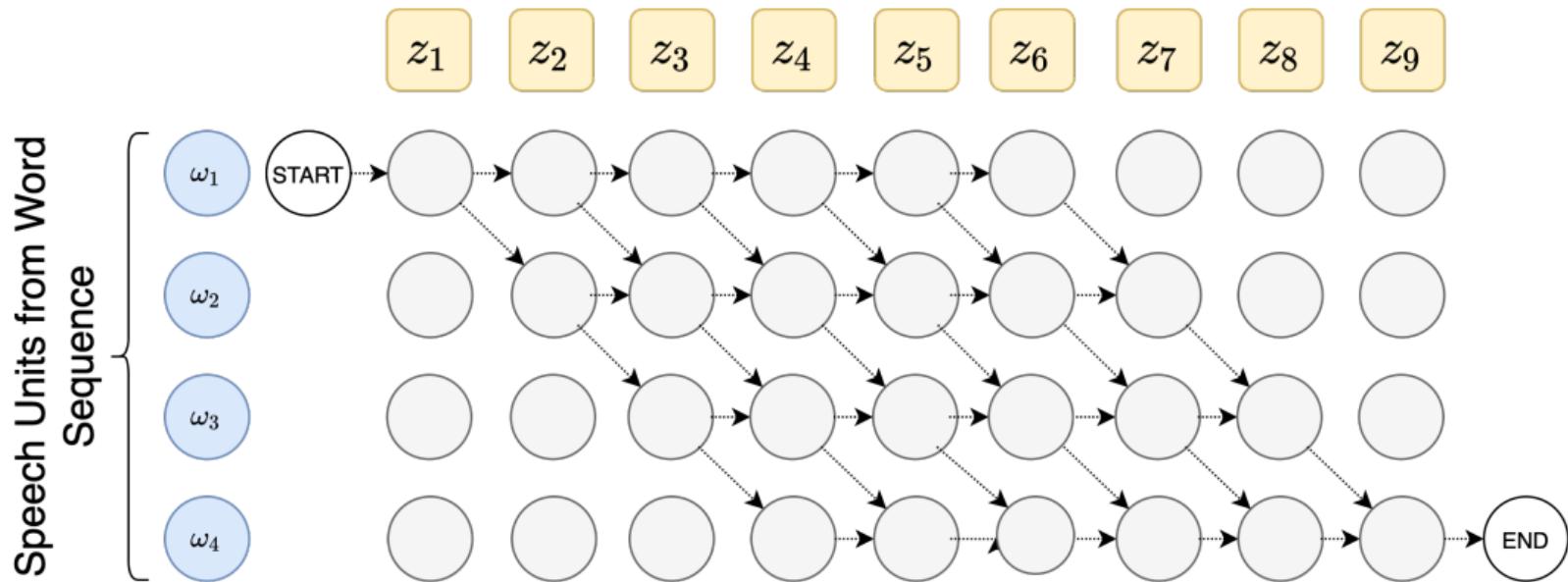


## State-Space Trellis

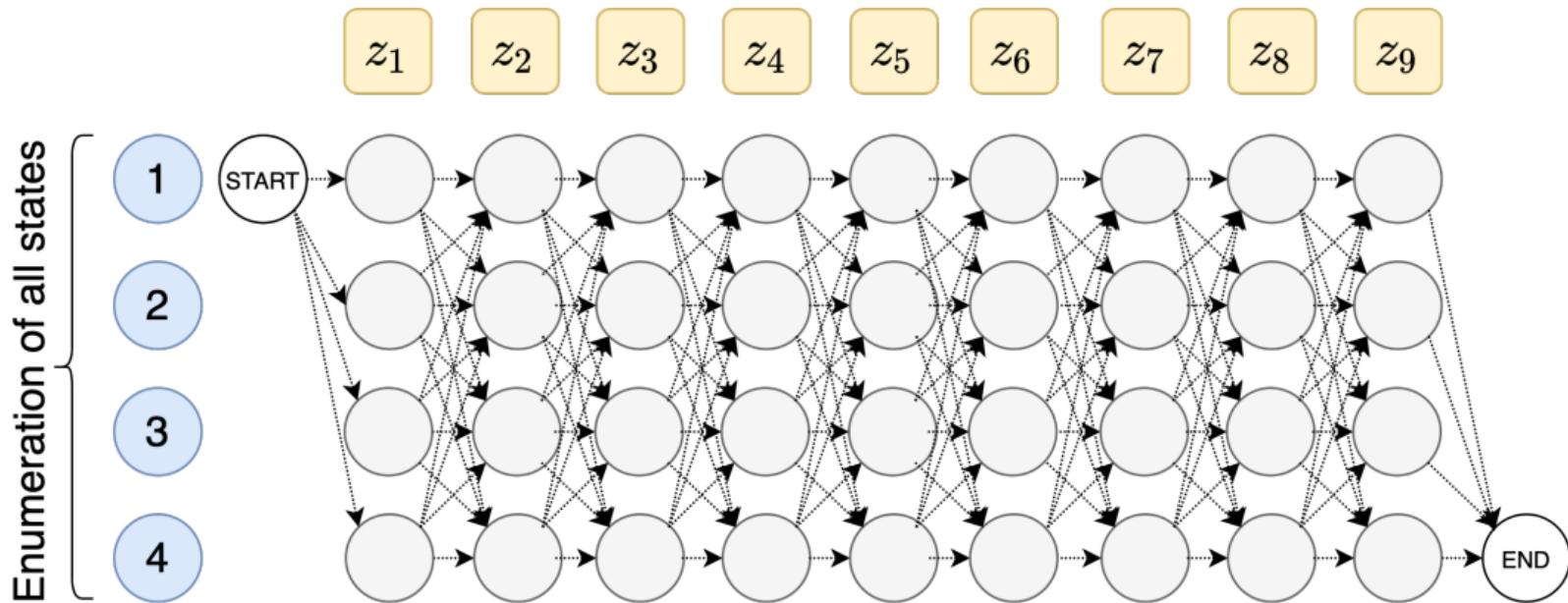


- Inference and Training / Alignment Trellises have different structure!
  - Inference Trellis - Any path is valid
  - Training / Alignment Trellis - speech unit sequence defines set of valid paths

## HMM Trellis - Training / Alignment Trellis



## HMM Trellis - Inference Trellis



## Discriminative State-Space ASR Models

# Automatic Speech Recognition (ASR) - Discriminative vs Generative Approach

- **Discriminative** ASR → "Which sequence  $\hat{\mathbf{w}}$  is likely given the audio?"

$$\hat{\mathbf{w}} = \mathcal{M}^{-1}(\hat{\omega}_{1:L}), \quad \hat{\omega}_{1:L} = \arg \max_{\hat{\omega}_{1:L}} P(\hat{\omega}_{1:L} | \mathbf{X}_{1:T}; \theta)$$

- **Generative** ASR → "Given which sequence  $\hat{\mathbf{w}}$  is the audio mostly likely?"

$$\hat{\mathbf{w}} = \mathcal{M}^{-1}(\hat{\omega}_{1:L}), \quad \hat{\omega}_{1:L} = \arg \max_{\hat{\omega}_{1:L}} P(\mathbf{X}_{1:T} | \hat{\omega}_{1:L}; \theta) P(\hat{\omega}_{1:L}; \theta)$$

# Automatic Speech Recognition (ASR) - Discriminative Approach

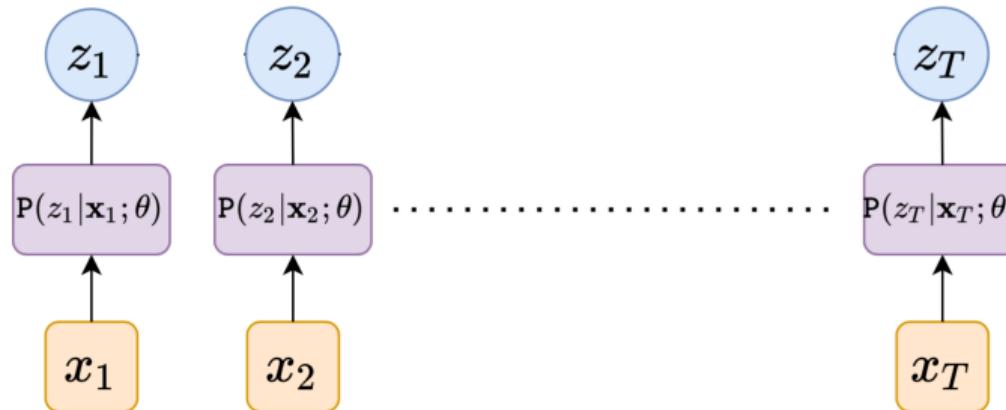
- ASR is a **discriminative** task → "Which sequence  $\hat{\mathbf{w}}$  is likely given the audio?"

$$\hat{\mathbf{w}} = \mathcal{M}^{-1}(\hat{\omega}_{1:L}), \quad \hat{\omega}_{1:L} = \arg \max_{\hat{\omega}_{1:L}} P(\hat{\omega}_{1:L} | \mathbf{X}_{1:T}; \theta)$$

- Discriminative State-Space ASR systems:

$$P(\omega_{1:L} | \mathbf{X}_{1:T}; \theta) = \sum_{\pi_{1:T} \in \mathcal{A}(\omega_{1:L}, T)} P(\pi_{1:T} | \mathbf{X}_{1:T}; \theta)$$

# Connectionist Temporal Classification (CTC)



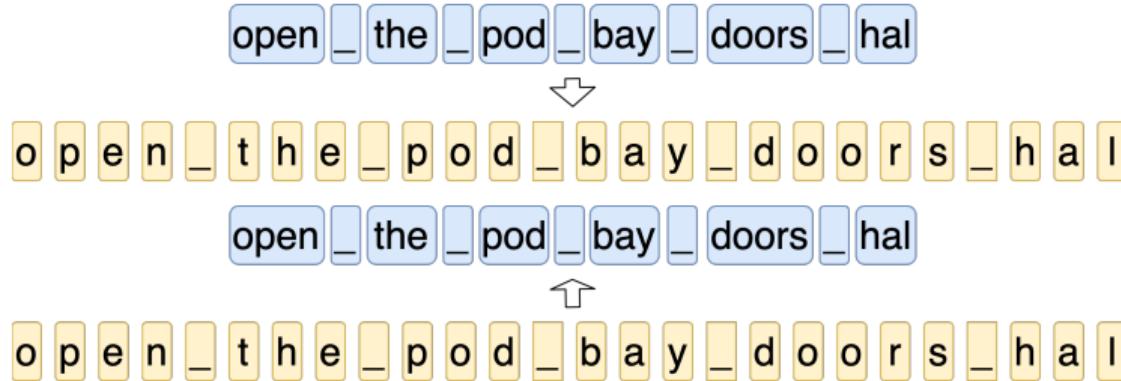
- CTC is a **discriminative** state-space ASR model defined as:

$$p(\mathbf{w}|\mathbf{X}_{1:T}; \theta) = \sum_{\pi_{1:T} \in \mathcal{A}(\omega_{1:L}, T)} \prod_{t=1}^T p(z_t|x_t; \theta), \quad \omega_{1:L} = \mathcal{M}(\mathbf{w}), \quad \pi_{1:T} = \{z_t\}_{t=1}^T$$

## Connectionist Temporal Classification (CTC) - Properties

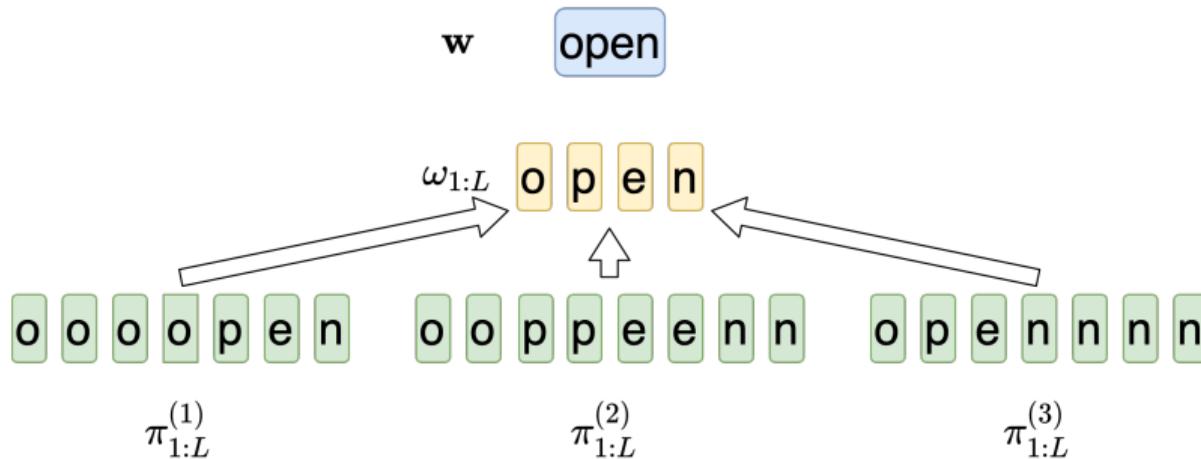
- CTC assumes all states are conditionally independent
  - Generally poor assumption, but context modelled implicitly by Neural Network.
  - State duration can also be poorly modelled
  - However, allows model to be more easily adapted by using external LM
- Alignment free - Do not need prior alignment for training
  - Train by maximizing probability of all valid alignments
- CTC models are typically graphemic with an extra black symbol  $\epsilon$ .
  - Allows simple mapping between words and speech units
- CTC, like HMM allows only monotonic alignment
  - Good assumption for ASR
  - Cannot ignore any frames

# Connectionist Temporal Classification (CTC)



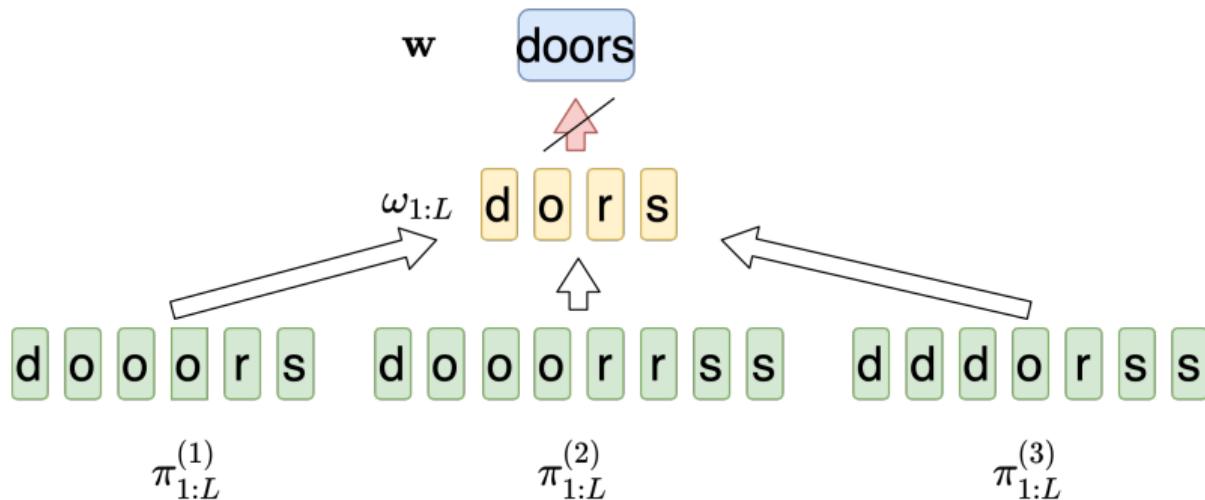
- Mapping word sequence to grapheme sequences is simple
  - We also need a procedure to map back from graphemes to words!
- A grapheme could be repeated for many states. Ex: oooopen → open
  - Could simply merge repeated graphemes!

# Connectionist Temporal Classification (CTC)



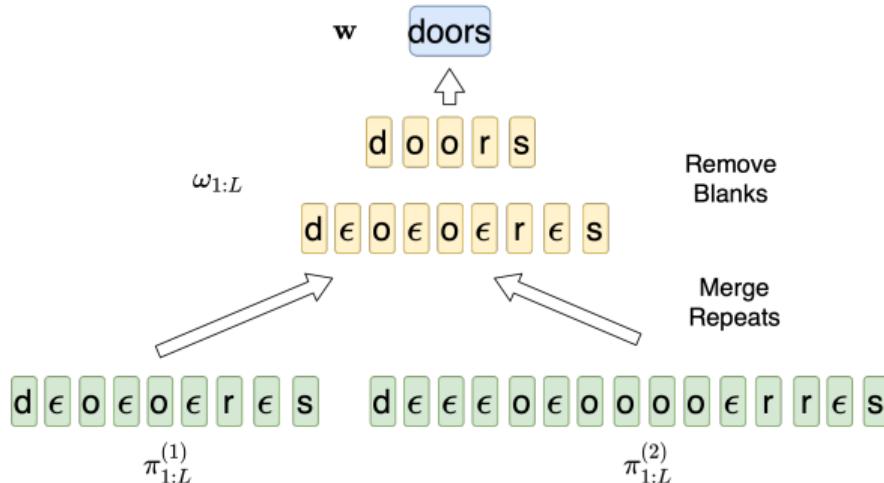
- Merging repeats is an easy operation
  - Can map many state sequences back to original grapheme sequence
  - But does this always work?

# Connectionist Temporal Classification (CTC)



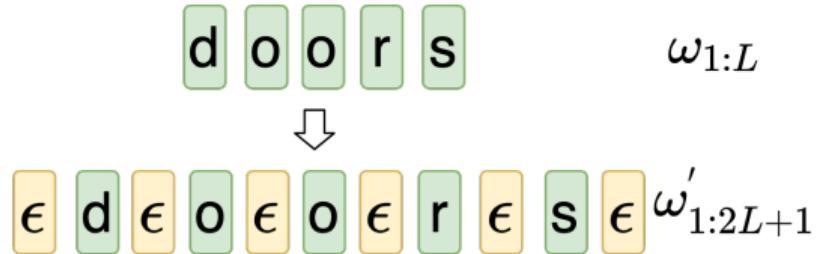
- Cannot properly restore words with repeated letters!

# Connectionist Temporal Classification (CTC)



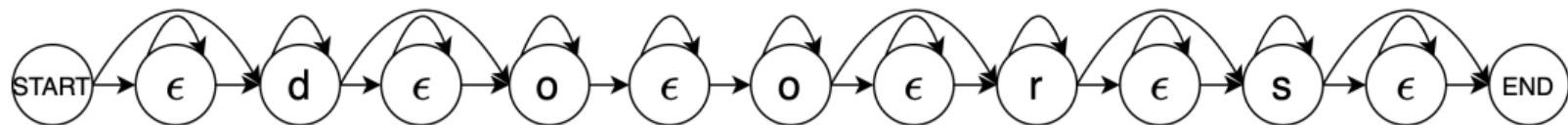
- Solve by introducing a **blank** symbol  $\epsilon$ !
  - Restore works by first merging repeats, then removing blanks.
- Blanks follow a few simple rules
  - States sequence can start and end with blank
  - Blanks can come after any character, but can also be omitted - ex: do $\epsilon$ ors

# Connectionist Temporal Classification (CTC)



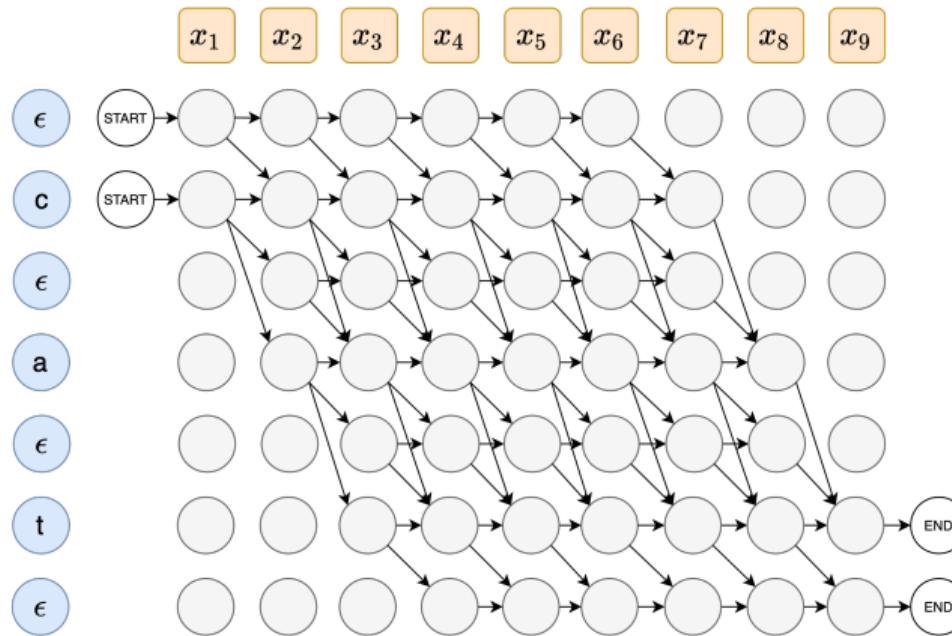
- For training use a sequence  $\omega'_{1:2L+1}$  with a blank before and after each character

# Connectionist Temporal Classification (CTC)



- Following state transitions in are  $\omega'_{1:2L+1}$  valid:
  - Self-transitions
  - Transitions between blank and non-blank
  - Transitions between distinct non-blank labels
- Allows sequence  $\omega'_{1:2L+1}$  contains all valid alignments.

# Connectionist Temporal Classification (CTC) - Alignment and Training Trellis



- CTC trellis has 2 state and end states.

Ok, so how can we predict?

## Greedy Decoding

- Predicting with CTC is simple as it does not model inter-state dependencies:

$$P(z_{1:T} | \mathbf{X}_{1:T}) = \prod_{t=1}^T P(z_t | \mathbf{X}_{1:T})$$

$$\pi_{1:T}^* = \arg \max_{\pi_{1:T}} \prod_{t=1}^T P(z_t = \pi_t | \mathbf{X}_{1:T})$$

- This yields most probable state sequence - then merge repeats and remove blanks.
- Greedy Decoding works well as CTC yield sharp distributions
- GD can still fail to find the most likely  $\omega_{1:L}^*$ 
  - Most likely path  $\omega_{1:L}^*$  may be composed of multiple, less likely state sequences.
- Can do better with Prefix-Decoding, but need LM for big gains.

Ok, so how can we train CTC models?

## Training - CTC-Loss Function

- CTC loss is alignment free → do not need to know a-prior alignment
  - Maximize probability of all valid alignments  $\pi_{1:T} \in \mathcal{A}(\omega_{1:L})$

$$\begin{aligned}\mathcal{L}(\theta) &= \frac{1}{N} \sum_{n=1}^N -\ln P(\omega_{1:L_n}^{(n)} | \mathbf{X}_{1:T_n}^{(n)}, \theta) \\ &= \frac{1}{N} \sum_{n=1}^N -\ln \sum_{\pi_{1:T_n} \in \mathcal{A}(\omega_{1:L_n}^{(n)})} P(\pi_{1:T_n} | \mathbf{X}_{1:T_n}^{(n)}, \theta) \\ &= \frac{1}{N} \sum_{n=1}^N -\ln \sum_{\pi_{1:T_n} \in \mathcal{A}(\omega_{1:L_n}^{(n)})} \prod_{t=1}^{T_n} P(\pi_t | \mathbf{X}_{1:T_n}^{(n)}, \theta)\end{aligned}$$

- Number of possible paths is  $O(K^{2L+1}) \rightarrow$  use Dynamic Programming!

## Training - CTC Forward Algorithm

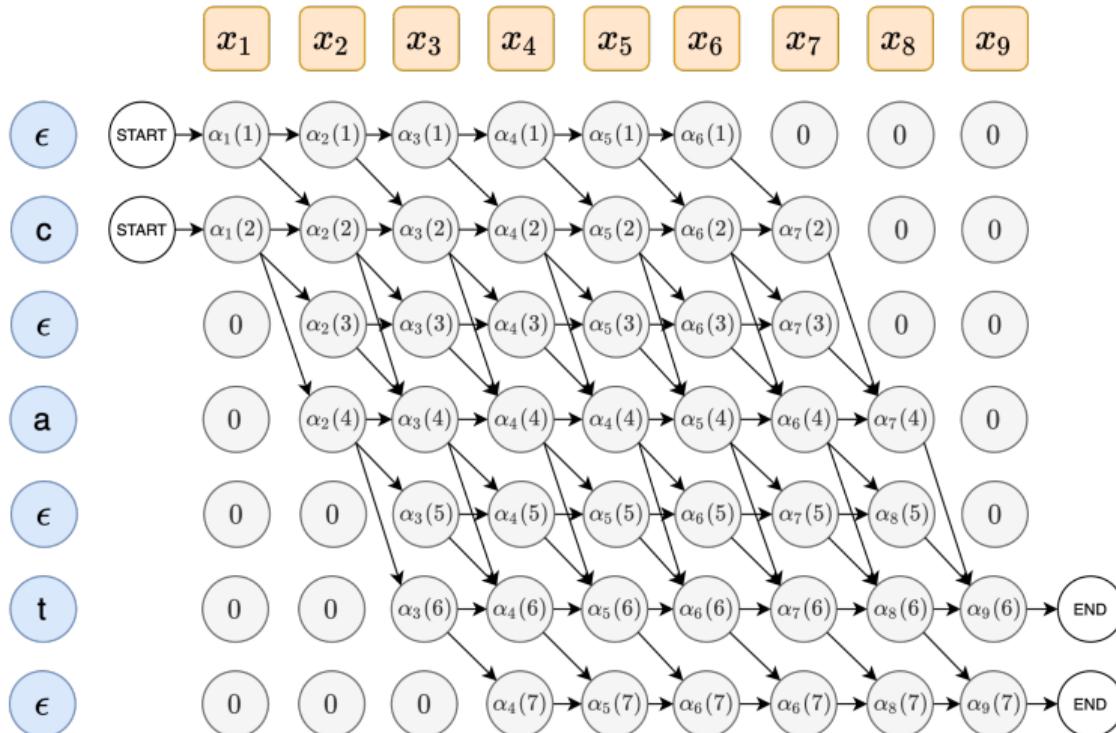
- The CTC forward algorithm recursively computes the forward variable  $\alpha_t(s)$

$$\begin{aligned}\alpha_t(s) &= P(\omega_{1:s/2}, \pi_t = \omega'_s | \mathbf{X}_{1:T}, \theta) \\ &= \sum_{\pi_{1:t-1} \in \mathcal{A}(\omega_{1:s/2}), \pi_t = \omega'_s} P(\pi_{1:t} | \mathbf{X}_{1:T}, \theta)\end{aligned}$$

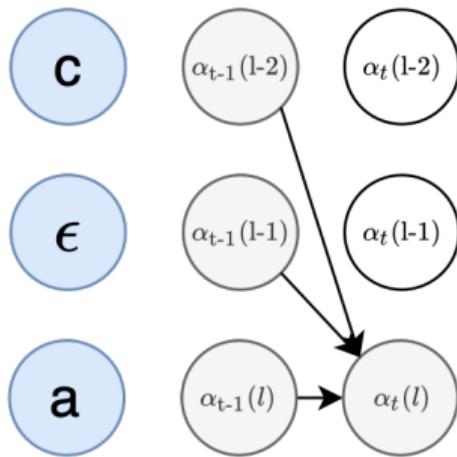
- $\alpha_t(s)$  is the probability of all paths of length  $t$  which go through state  $\omega'_s$ 
  - Note we consider alignments of **subsequence**  $\omega_{1:s/2}$ !
- Final loss is given by the sum of alphas of two end states:

$$P(\omega_{1:L} | \mathbf{X}_{1:T}, \theta) = \alpha_T(2L) + \alpha_T(2L + 1)$$

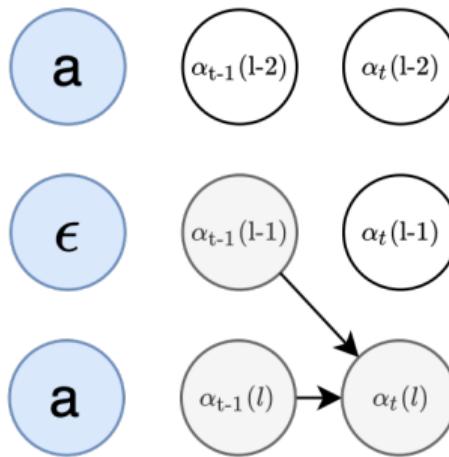
# Training - CTC Forward Algorithm



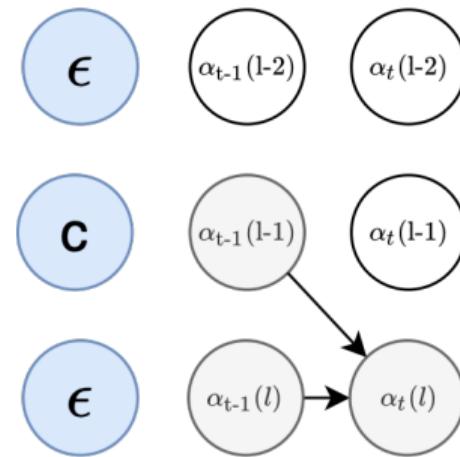
## Training - CTC Forward Algorithm - Edge Cases



(a) Can skip blank



(b) Can't skip blank



(c) Go to blank

## Training - CTC Forward Algorithm

$$\alpha_t(0) = 0, \forall t$$

$$\alpha_1(1) = P(z_1 = \epsilon | \mathbf{X}_{1:T}),$$

$$\alpha_1(2) = P(z_1 = \omega'_2 | \mathbf{X}_{1:T}),$$

$$\alpha_1(s) = 0, \forall s > 2$$

$$\alpha_t(s) = 0, \forall s < (2L + 1) - 2(T - t) - 1$$

$$\alpha_t(s) = \begin{cases} (\alpha_{t-1}(s) + \alpha_{t-1}(s-1))P(z_t = \omega'_s | \mathbf{X}_{1:T}) & \text{if } \omega'_s = \epsilon \text{ or } \omega'_s = \omega'_{s-2} \\ (\alpha_{t-1}(s) + \alpha_{t-1}(s-1) + \alpha_{t-1}(s-2))P(z_t = \omega'_s | \mathbf{X}_{1:T}) & \text{otherwise} \end{cases}$$

## Training - CTC Backward Algorithm

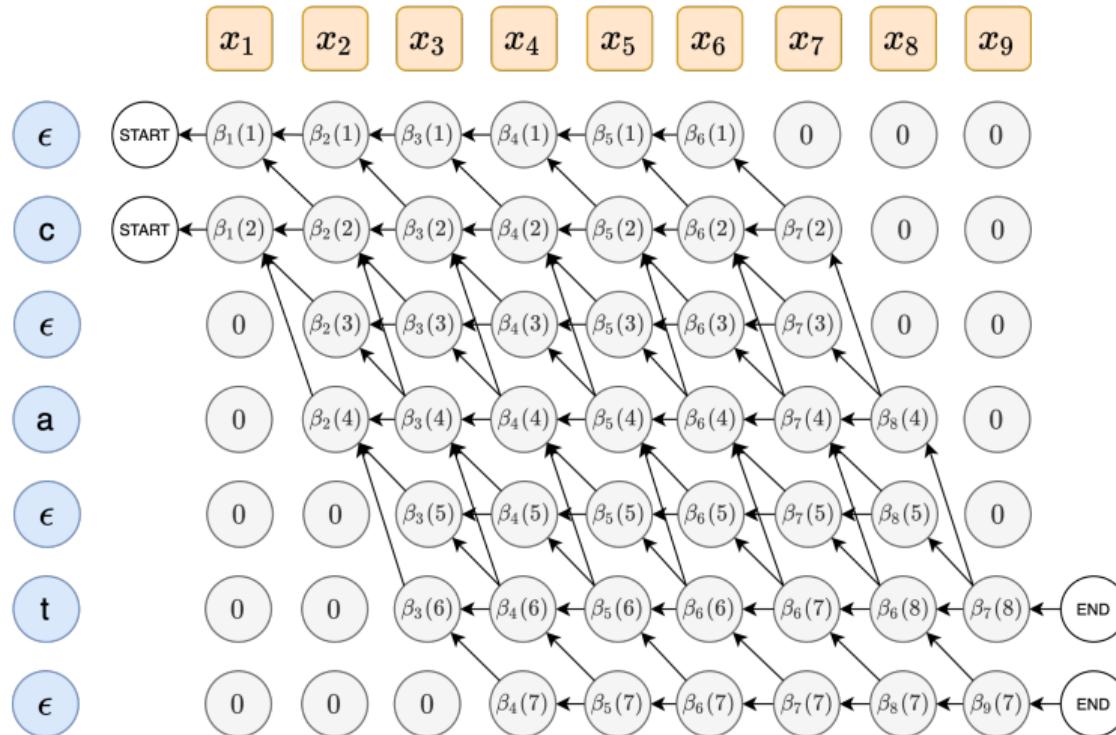
- The CTC backward algorithm recursively computes the backward variable  $\beta_t(s)$

$$\begin{aligned}\beta_t(s) &= P(\omega_{s/2:L}, \pi_t = \omega'_s | \mathbf{X}_{1:T}, \theta) \\ &= \sum_{\pi_{t+1:T} \in \mathcal{A}(\omega_{s/2:L}), \pi_t = \omega'_s} P(\pi_{t+1:T} | \mathbf{X}_{1:T}, \theta)\end{aligned}$$

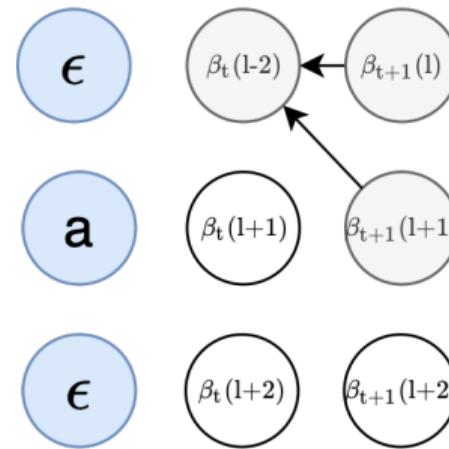
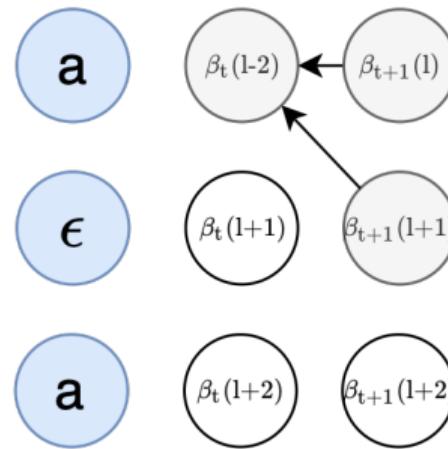
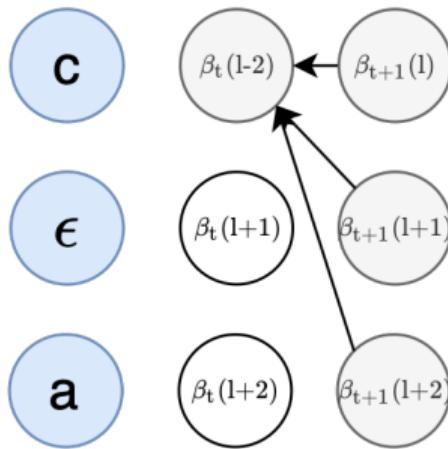
- $\beta_t(s)$  is the probability of all valid alignments  $\omega_{s/2:L}$  starting in state  $\omega'_s$
- Final loss also given by the sum of betas of two start states:

$$P(\omega_{1:L} | \mathbf{X}_{1:T}, \theta) = \beta_1(1) + \beta_1(2)$$

# Training - CTC Backward Algorithm



## Training - CTC Backward Algorithm - Edge Cases



(a) Can skip blank

(b) Can't skip blank

(c) Go to blank

## Training - CTC Backward Algorithm

$$\beta_T(2L + 1) = 1$$

$$\beta_T(2L) = 1$$

$$\beta_T(s) = 0, \forall s < 2L$$

$$\beta_t(s) = 0, \forall s > 2t$$

$$\beta_t(2L + 2) = 0, \forall t$$

$$\beta_t(s) = \begin{cases} (\beta_{t+1}(s) + \beta_{t+1}(s + 1))P(z_t = \omega'_s | \mathbf{X}_{1:T}) & \text{if } \omega'_s = \epsilon \text{ or } \omega'_s = \omega'_{s+2} \\ (\beta_{t+1}(s) + \beta_{t+1}(s + 1) + \beta_{t+1}(s + 2))P(z_t = \omega'_s | \mathbf{X}_{1:T}) & \text{otherwise} \end{cases}$$

## CTC Forward-Backward Algorithm, Loss Computation and Soft-Alignment

- The probability of all paths passing through a state  $\pi_t = \omega'_s$  is the product of forward and backward variables
  - This is important for the calculation of the gradient (not shown)

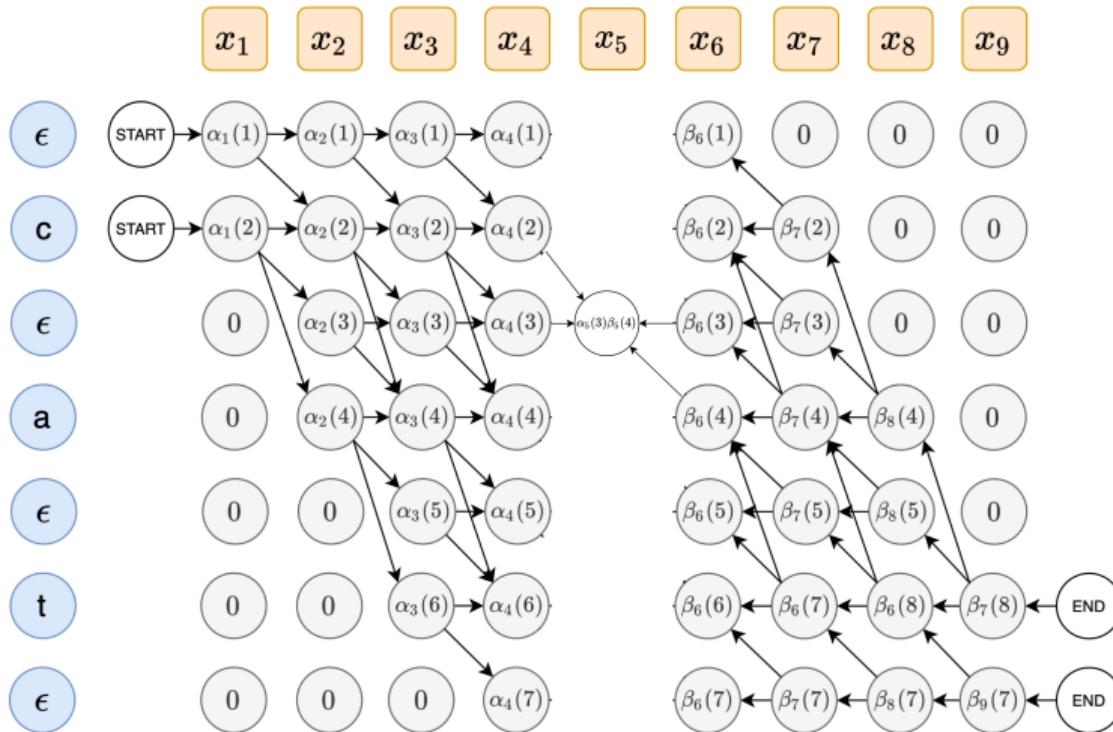
$$\alpha_t(s)\beta_t(s) = \sum_{\pi_{1:T} \in \mathcal{A}(\omega_{1:L}), \pi_t = \omega'_s} P(\pi_{1:T} | \mathbf{X}_{1:T}, \theta)$$

- Sum of all such products yields total probability:

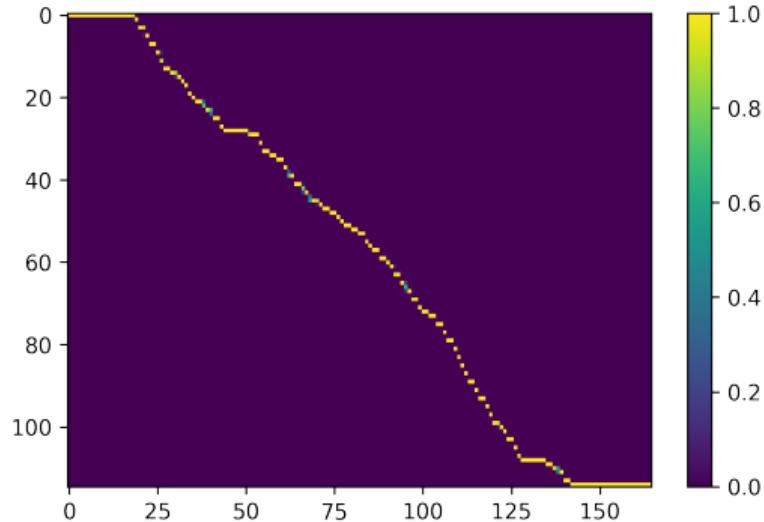
$$\sum_{s=1}^{2L+1} \alpha_t(s)\beta_t(s) = P(\omega_{1:L} | \mathbf{X}_{1:T}, \theta)$$

- We can also use  $\alpha_t(s)\beta_t(s)$  as a measure of soft-alignment
  - Often normalize to align =  $\frac{\alpha_t(s)\beta_t(s)}{\sum_{s=1}^{2L+1} \alpha_t(s)\beta_t(s)}$

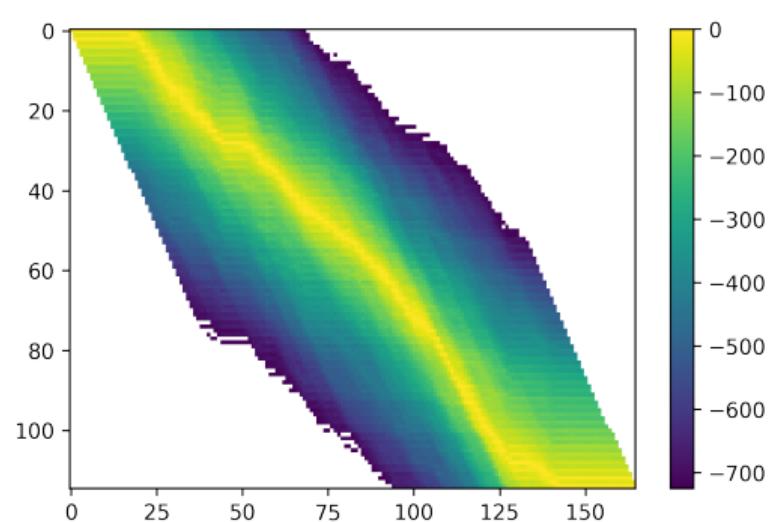
# Training - CTC Forward-Backward Algorithm



# CTC Soft-Alignment



(a) Soft-Alignment



(b) Log Soft-Alignment

# CTC Forward-Backward Algorithm and Loss Computation

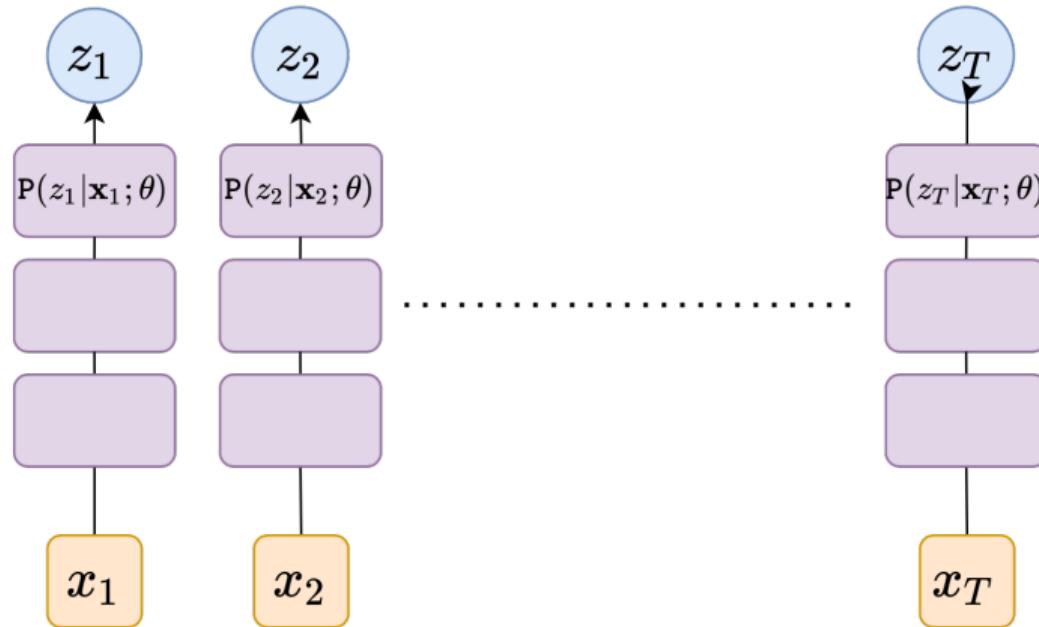
- Why use the Fw-Bw algorithm instead of just the Fw or Bw algorithms??

$$P(\omega_{1:L} | \mathbf{X}_{1:T}, \theta) = \underbrace{\sum_{s=1}^{2L+1} \alpha_t(s) \beta_t(s)}_{\text{Forward-Backward}} = \underbrace{\alpha_T(2L) + \alpha_T(2L+1)}_{\text{Forward}} = \underbrace{\beta_1(1) + \beta_1(2)}_{\text{Backward}}$$

- For gradient computation during training and alignment.

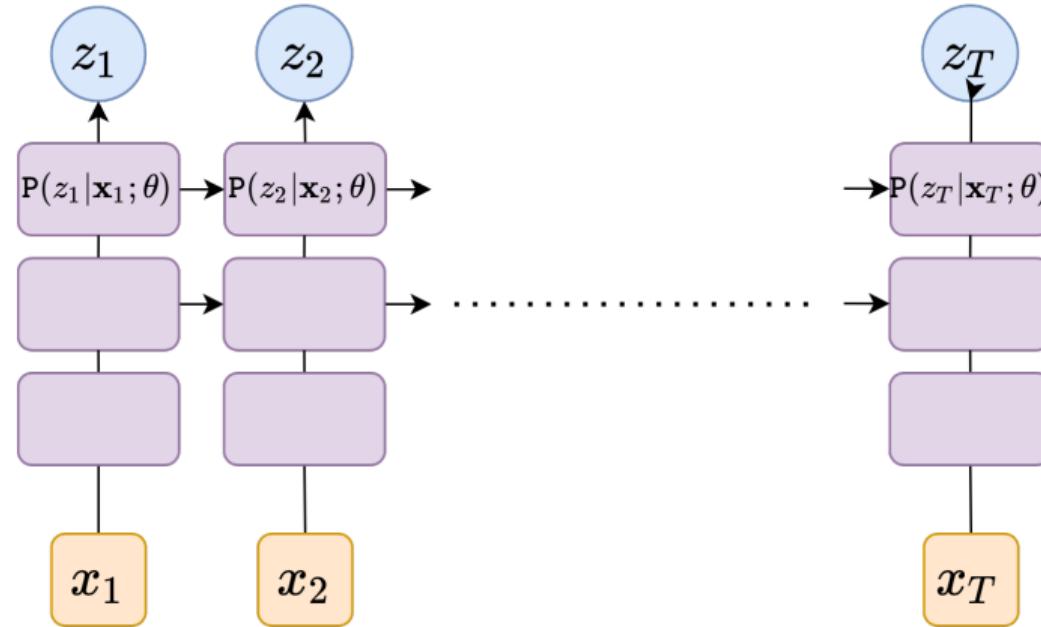
Math is cool (maybe?), but how do I build this?

# Neural Architecture for CTC ASR



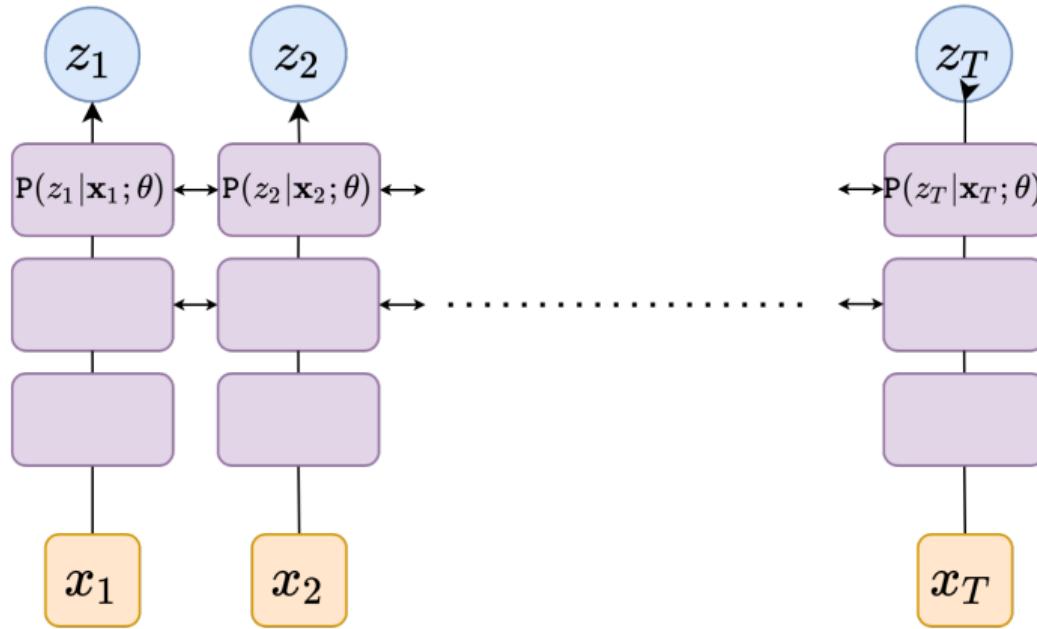
- We can use a Feed-Forward DNN, but not context modelling

# Neural Architecture for CTC ASR



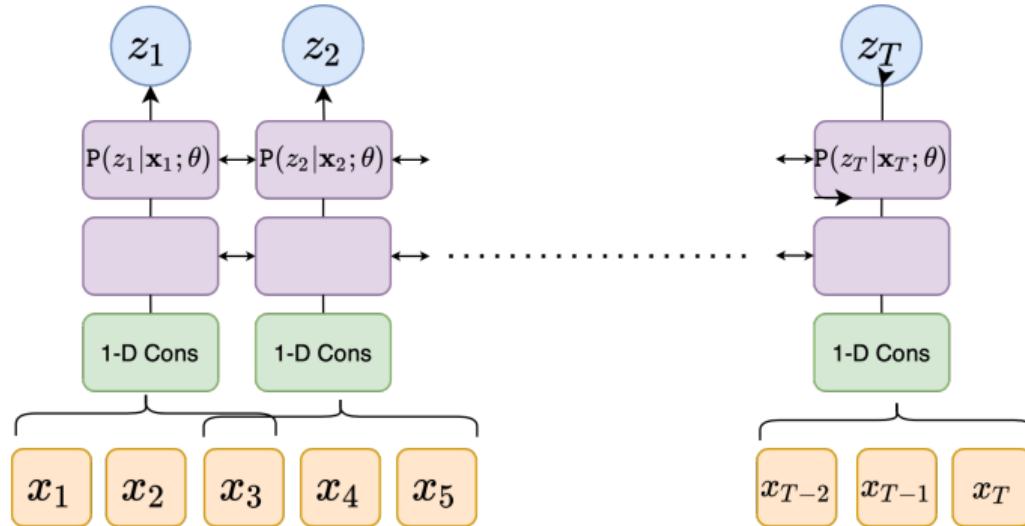
- Using Recurrent Units allows us to capture causal relations
  - Can use different units - RNN, GRU, LSTM, etc...

# Neural Architecture for CTC ASR



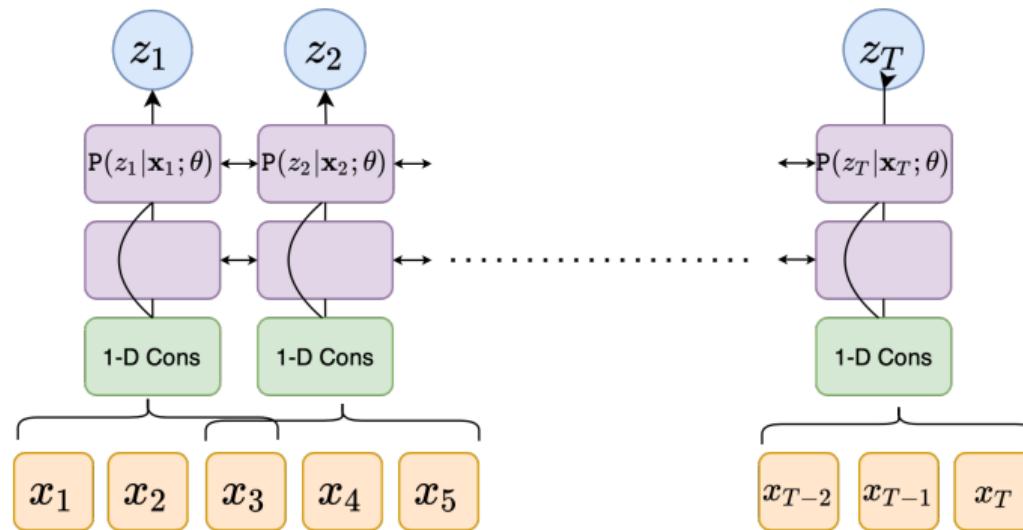
- Bidirectional Recurrent Units allows us to capture future context better

# Neural Architecture for CTC ASR



- Using convolution filters low-lever representation learning
- Strided convolution allows down-sampling!

# Neural Architecture for CTC ASR



- Adding residual layers improves convergence, adds regularization
  - Used in many large-scale architectures!

- Deep Learning is like lego - can add many architectural elements
  - BatchNorm / LayerNorm
  - Skip-connection
  - Self-Attention Layers
  - Strided convolutions for pooling
  - Dropout and other stochastic layers
- Make sure more not too large for about of training data!

# Neural Architecture for CTC ASR - Jasper (2019 SOTA)

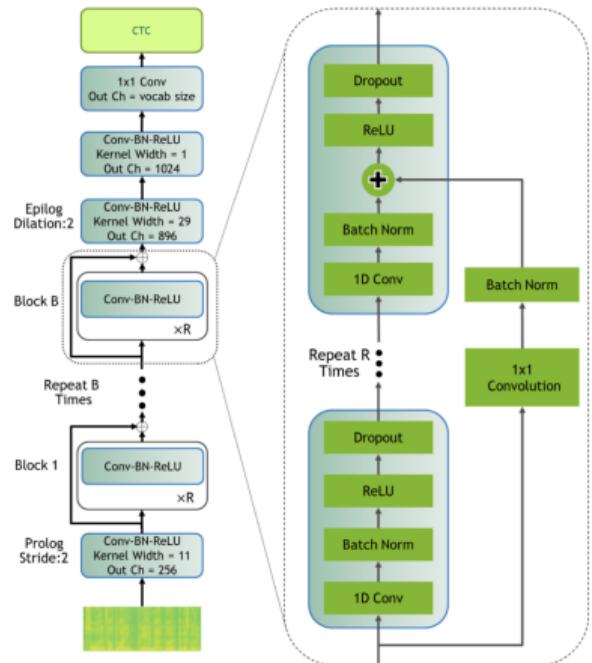


Figure 1: Jasper BxR model: B - number of blocks, R - number of sub-blocks.

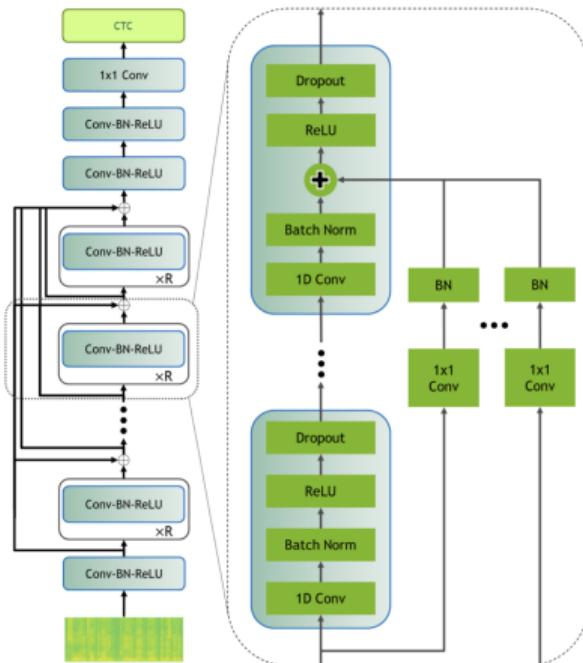


Figure 2: Jasper Dense Residual

So what have we learned?

# Conclusion

- CTC is a discriminative state-space ASR system
  - Uses graphemic states with special blank symbol
  - ‘Alignment-free’ training
  - Fully factorizes - doesn’t explicitly model context
- CTC Trellis has different structure from HMM trellis
  - Transfer between blanks and distinct non-blanks
- Can compute CTC loss using Forward-Backward Algorithm
- CTC and HMMs are closely related
  - CTC model is HMM with uniform prior and transition probs.
  - CTC training and inference criteria are matched
  - State duration and context modelling are poorer
  - CTC needs much more data
- Combining convolution, bidirectional recurrent and residual layers works well