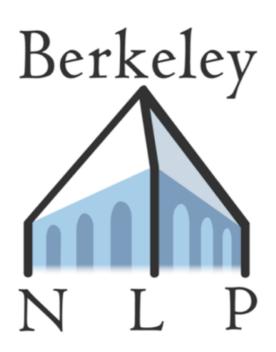
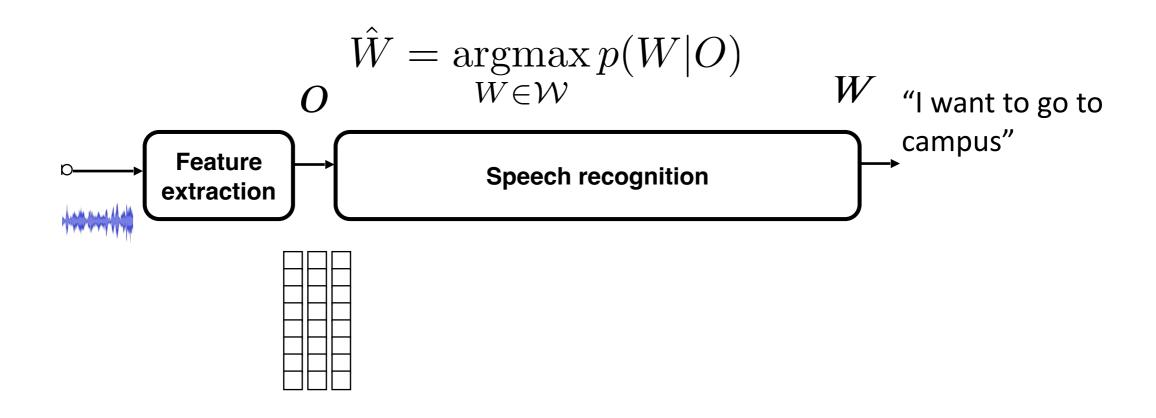
Classical and Modern formulations of ASR



EECS 183/283a: Natural Language Processing

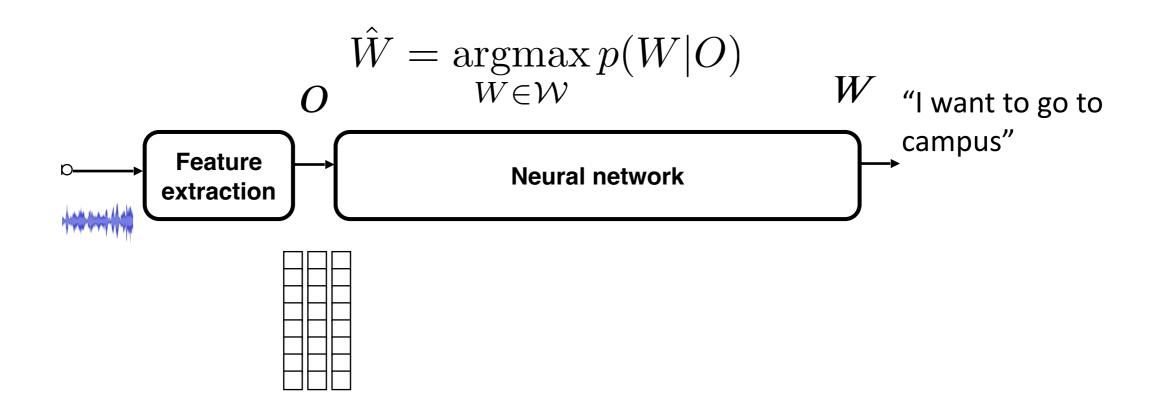
Speech recognition





End-to-end speech recognition





How to obtain the posterior p(W|O)



• We just replace it with a neural network-based function $f^{nn}(\cdot)$

$$\operatorname*{argmax}_{W} p(W|O) = \operatorname*{argmax}_{W} f^{\mathrm{nn}}(W|O)$$

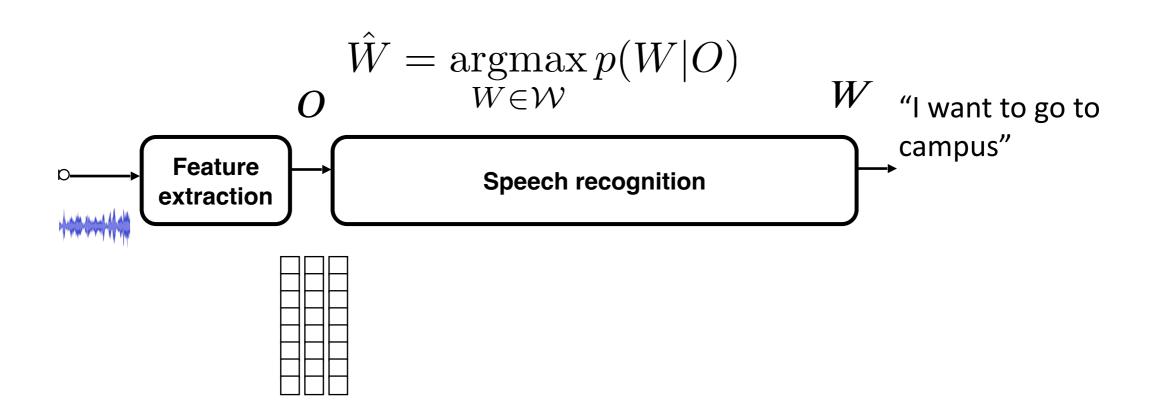
- Easy and simple, no math (in this level), however
- W is a sequence! $W = (w_n \in \mathcal{V} | n = 1, \dots, N)$
 - Very difficult to deal with it
 - Say N = 10, $|\mathcal{V}| = 100$, we have to deal with 100^{10} possible sequences
 - Also, the length N is variable
 - We have to use a special neural network (e.g., attention, CTC, and RNN-transducer)



- Classical speech recognition
 - Pipeline

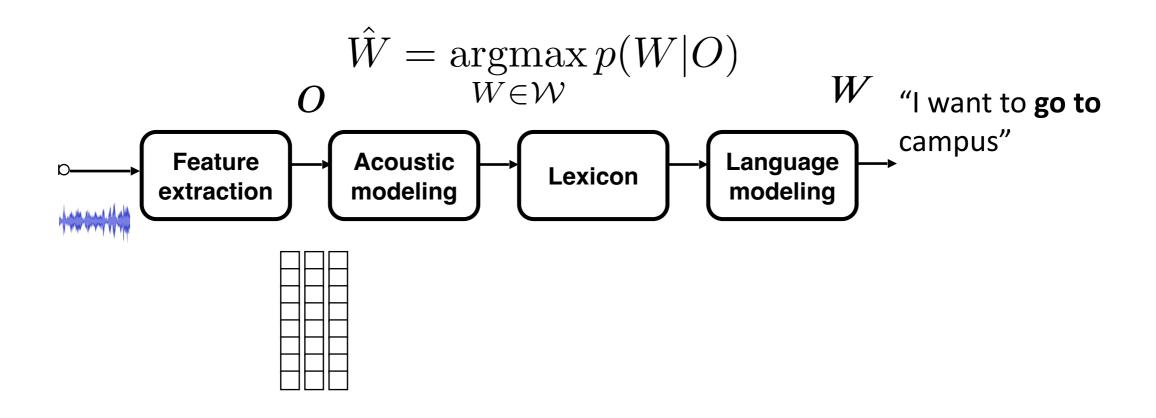
Speech recognition





Classical (non-end-to-end) speech recognition





Speech -> Text





Text W: I want to go to the campus

Speech -> Phoneme -> Text



Speech *o*:



Phoneme L: AY W AA N T T UW G OW T UW K AE M P AH S



Text W: I want to go to campus

How to obtain the posterior p(W|O)



- Factorize the model with phoneme
 - Let $L=(l_i\in\{/\mathrm{AA}/,\,/\mathrm{AE}/,\cdots\}|i=1,\cdots,J)$ be a phoneme sequence

$$\arg\max_{W} p(W|O) = \arg\max_{W} \sum_{L} p(W,L|O) \qquad \text{Sum rule}$$

$$= \arg\max_{W} \sum_{L} \frac{p(O|W,L)p(L|W)p(W)}{p(O)} \qquad \text{Bayes+ Product rule}$$

$$= \arg\max_{W} \sum_{L} p(O|W,L)p(L|W)p(W) \qquad \text{does not depend}$$
 on W
$$= \arg\max_{W} \sum_{L} p(O|L)p(L|W)p(W) \qquad \text{Conditional independence}$$
 assumption

Note: the right hand side is not the probability as it lacks a sum to one constraint

Noisy channel model



$$\underset{\mathbf{W}}{\operatorname{argmax}} p(W \mid O) = \underset{W}{\operatorname{argmax}} p(O \mid W) p(W) \approx \underset{W}{\operatorname{argmax}} \sum_{L} p(O \mid L) p(L \mid W) p(W)$$

Speech recognition

• p(O|L): Acoustic model (Hidden Markov model)

• p(L|W): Lexicon

• p(W): Language model (n-gram)

Noisy channel model



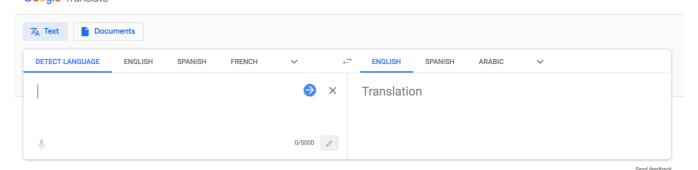
$$\underset{\mathbf{W}}{\operatorname{argmax}} p(W | Y) = \underset{\mathbf{W}}{\operatorname{argmax}} p(Y | W) p(W)$$

W: Targetlanguage textY: Source languagetext

Machine translation

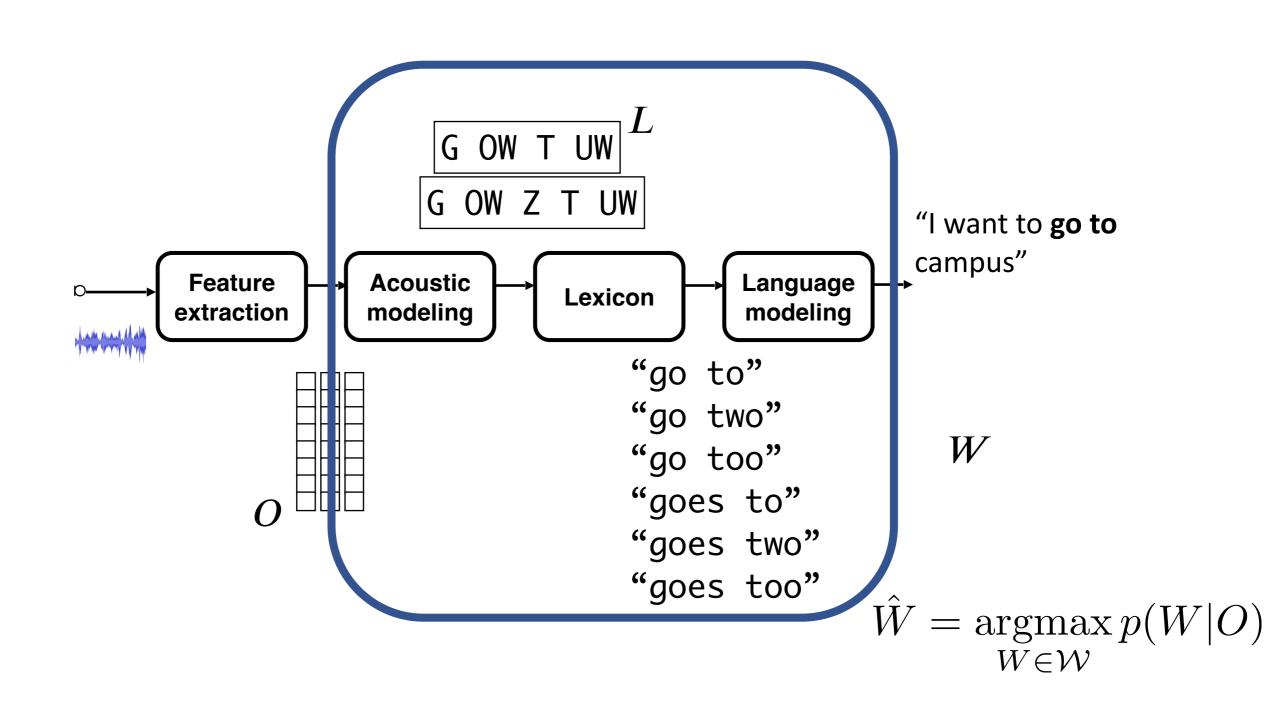
• p(Y|W): Translation model

• p(W): Language model



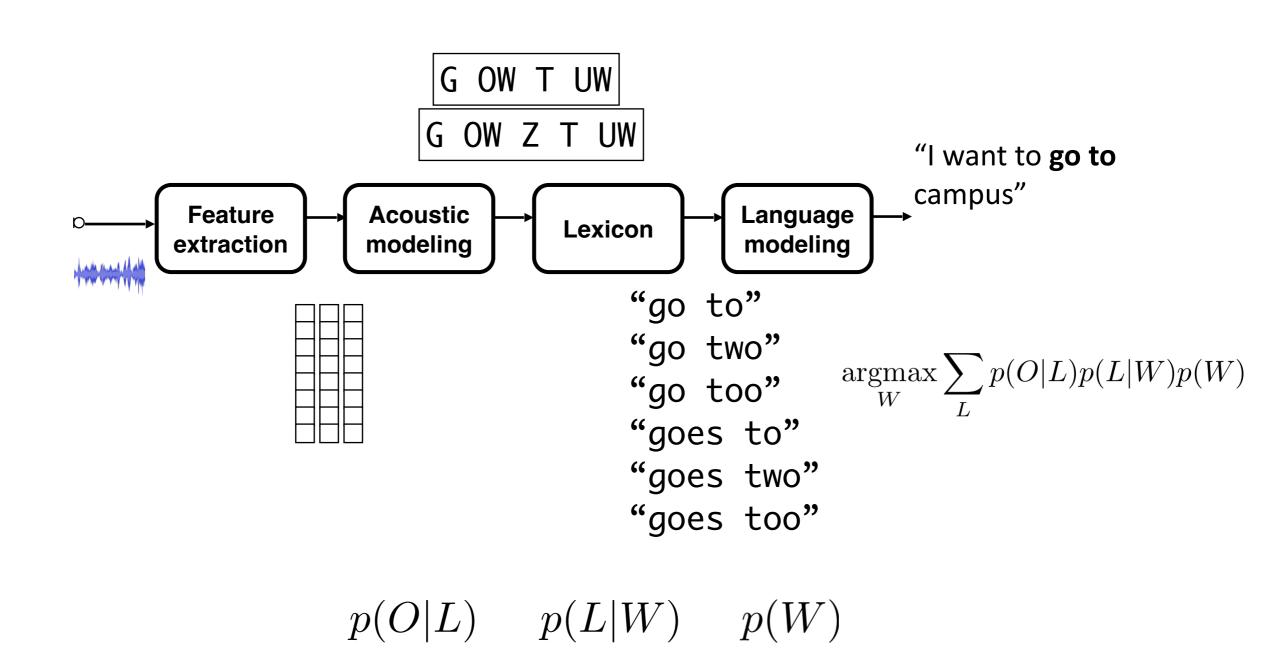
Speech recognition pipeline





Speech recognition pipeline





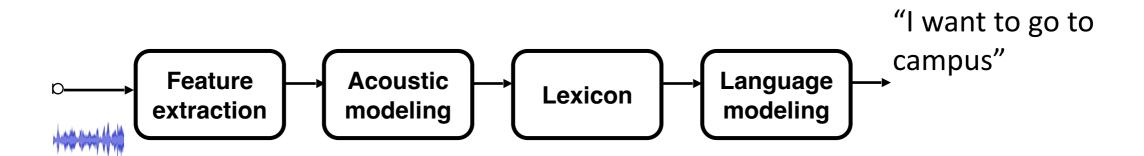
Please remember the noisy channel model

- Factorization
- Conditional independence (Markov) assumptions

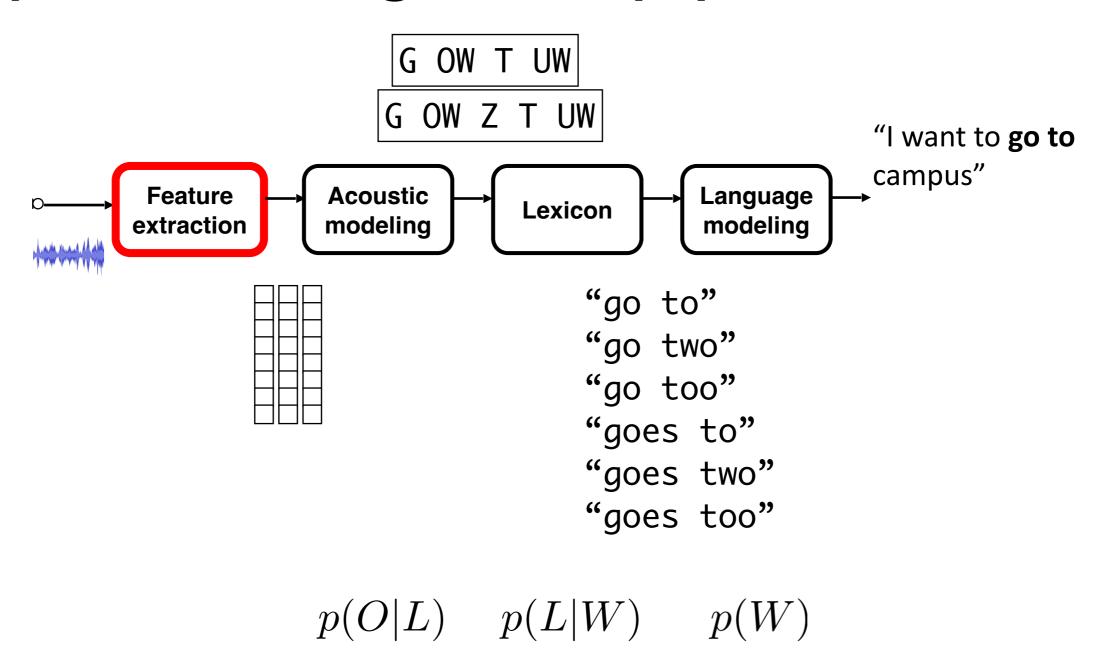
We can elegantly factorize the speech recognition problem with a tractable subproblem

Main blocks of Classical ASR



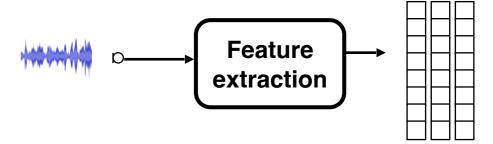


Speech recognition pipeline



Waveform to speech feature

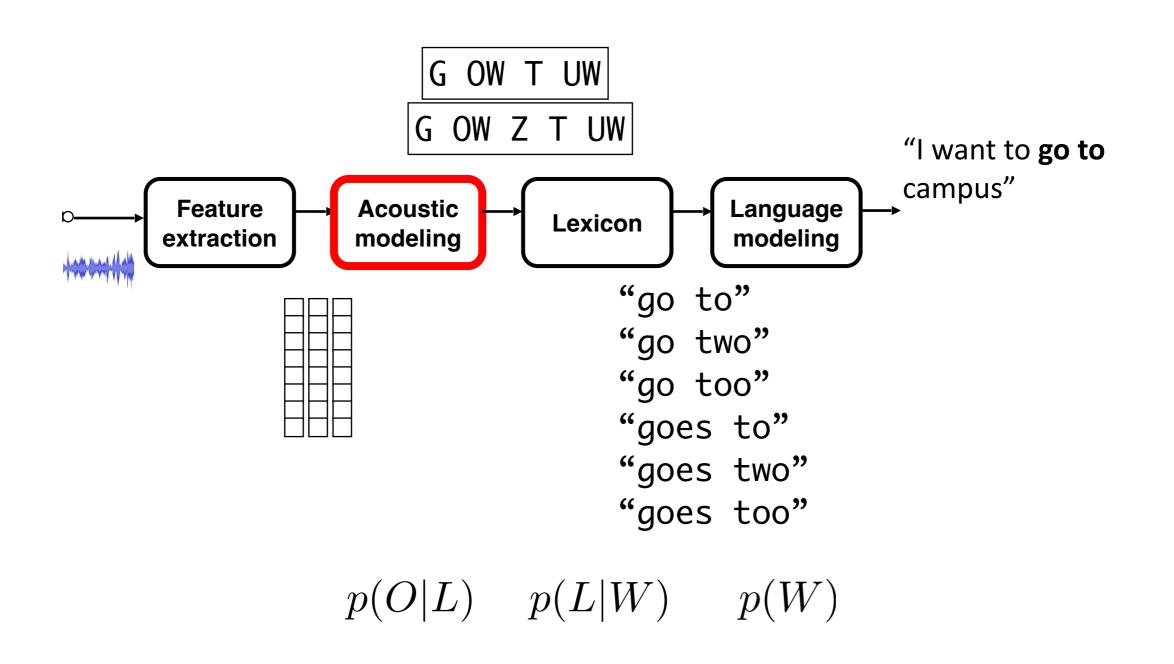




- Performed by so-called feature extraction module
 - Mel-frequency cepstral coefficient (MFCC), Perceptual Linear Prediction (PLP) used for Gaussian mixture model (GMM)
 - Log Mel filterbank used for deep neural network (DNN)
- Time scale
 - 0.0625 milliseconds (16kHz) to 10 milliseconds
- Type of values
 - Scalar (or discrete) to 12—40 dimensional vector

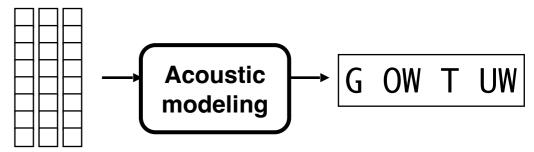
Speech recognition pipeline





Speech feature to phoneme





- Performed by so-called acoustic modeling module
 - Hidden Markov model (HMM) with GMM as an emission probability function
 - Hidden Markov model (HMM) with DNN as an emission probability function
- Time scale
 - 10 milliseconds to ~100 milliseconds (depending on a phoneme)
- Type of values
 - 12-dimensional continuous vector to 50 categorical value (~6bit)
- The most critical component to get the ASR performance
- It can be a probability of possible phoneme sequences, e.g.,

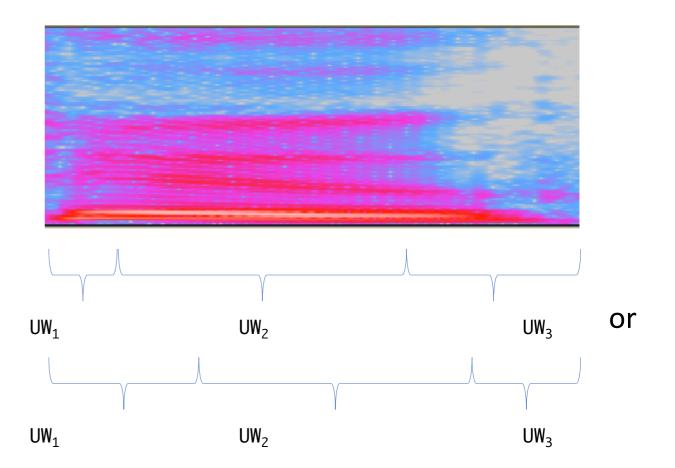
G OW T UW or G OW Z T UW with some scores

Acoustic model p(O|L)



- O and L are different lengths
- Align speech features and phoneme sequences by using HMM

- Provide p(O|L) based on this alignment and model
- The most important problem in speech recognition



How to formulate an acoustic model



- Acoustic model:
 - Again, O and L are different lengths $O = (\mathbf{o}_t \in \mathbb{R}^D | t = 1, \dots, T)$

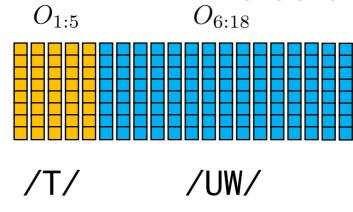
$$L = (l_i \in \{/AA/, /AE/, \dots\} | i = 1, \dots, J)$$

/T/ /UW/

Hard alignment problem similar to CTC and RNN-T

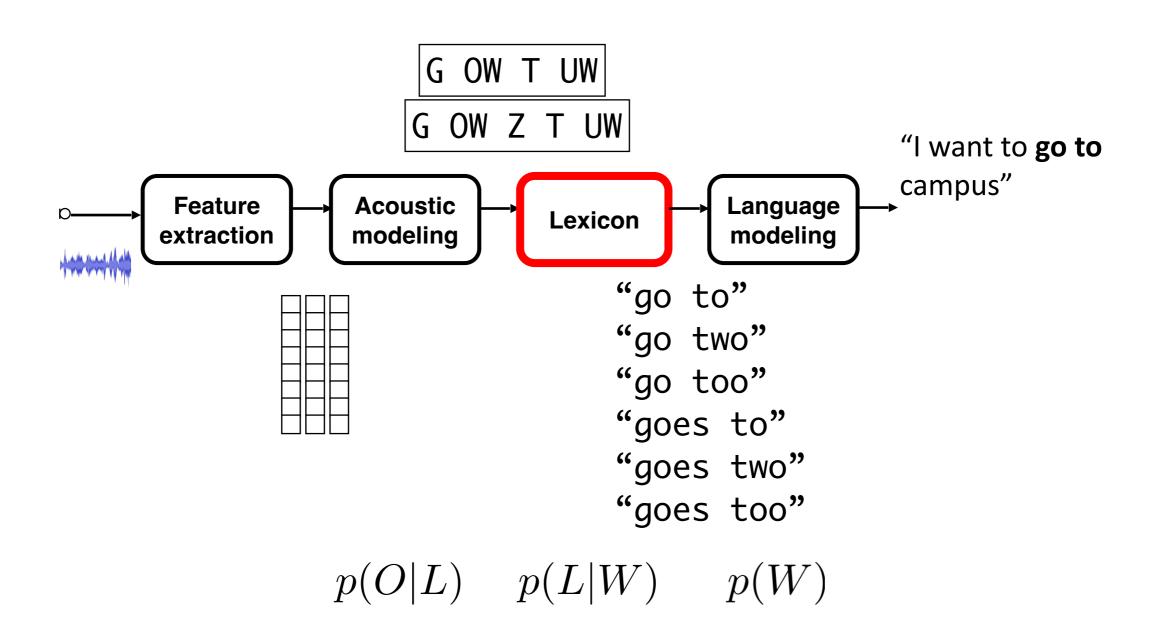
 If we assume that the alignment information is given, the problem becomes easy

$$p(O|L) = p(O_{1:T_1}, O_{T_1+1:T_2}, \cdots | l_1, l_2, \cdots)$$



Speech recognition pipeline





Phoneme to word





- Performed by lexicon module
 - American English: CMU dictionary
- Time scale
 - 100 milliseconds (depending on a phoneme) to 1 second (depending on a word and also language)
- Type of values
 - 50 categorical value (~6bit) to 100K categorical value (~2Byte)
- We need a pronunciation dictionary
- It can be multiple word sequences (one to many)

Lexicon p(L | W)

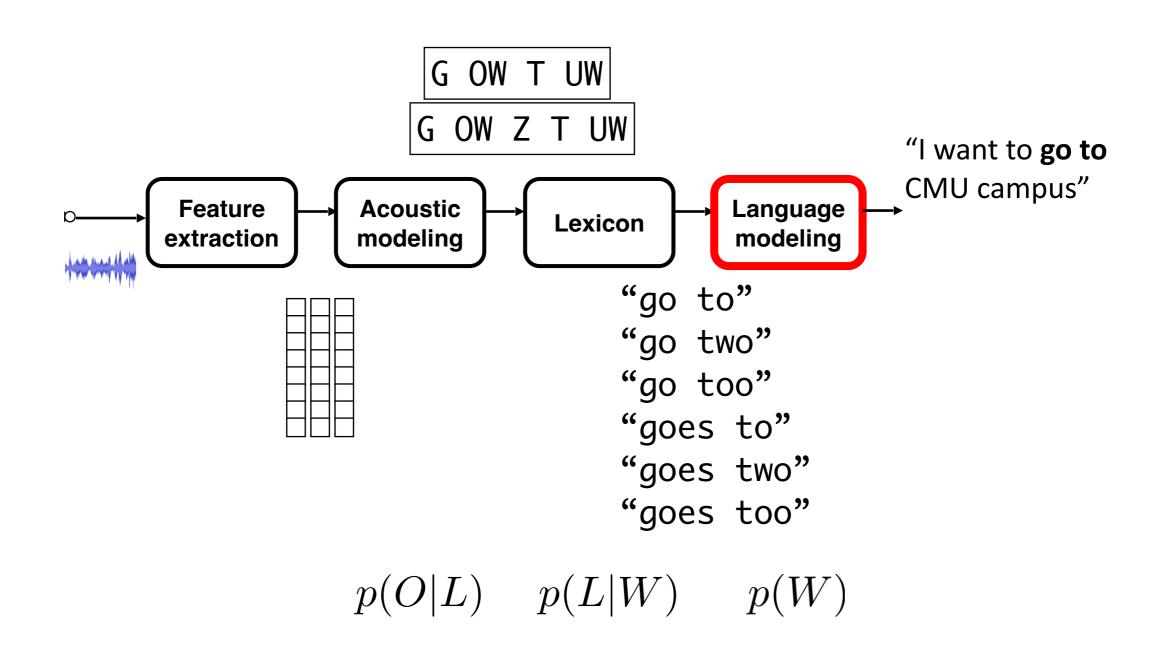


- Basically use a pronunciation dictionary, and map a word to the corresponding phoneme sequence
 - with the probability = 1.0 when single pronunciation
 - with the probability = 1.0/J when multiple (J) pronunciations

$$p(L|W) = p(/T/, /OW/|"two") = 1.0$$

Speech recognition pipeline







```
"go to"

"go two"

"go too"

Language modeling
```

- Performed by language modeling module p(W)
 - N-gram
 - Neural language model (recurrent neural network or transformer)
- From training data, we can basically find how possibly "to", "two", and "too" will be appeared after "go"
 - Part of WSJ training data, 37,416 utterances
 - "go to": **51** times
 - "go two":
 - "go too":

THE WALL STREET JOURNAL.



- Performed by language modeling module p(W)
 - N-gram
 - Neural language model (recurrent neural network or transformer)
- From training data, we can basically find how possibly "to", "two", and "too" will be appeared after "go"
 - Part of WSJ training data, 37,416 utterances
 - "go to": **51** times
 - "go two": **0** times
 - "go too": **0** times



```
"go to"

"go two"

"go too"

Language modeling
```

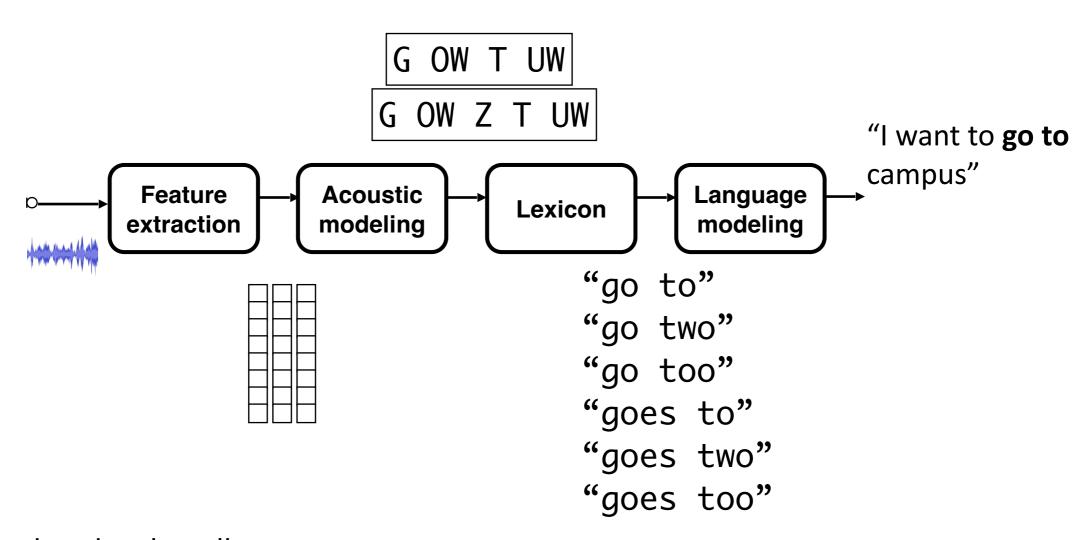
- Performed by language modeling module p(W)
 - N-gram
 - Neural language model (recurrent neural network or transformer)
- From training data, we can basically find how possibly "to", "two", and "too" will be appeared after "go"
 - WSJ all text data, 6,375,622 sentences
 - "go to": **2710** times
 - "go two":
 - "go too":



- Performed by language modeling mod "go too"
 - N-gram
 - Neural language model (recurrent neural network or transformer)
- From training data, we can basically find how possibly "to", "two", and "too" will be appeared after "go"
 - WSJ all text data, 6,375,622 sentences
 - "go to": **2710** times
 - "go two": 2 times, e.g., "those serving shore plants often go two hundred miles or more"
 - "go too":

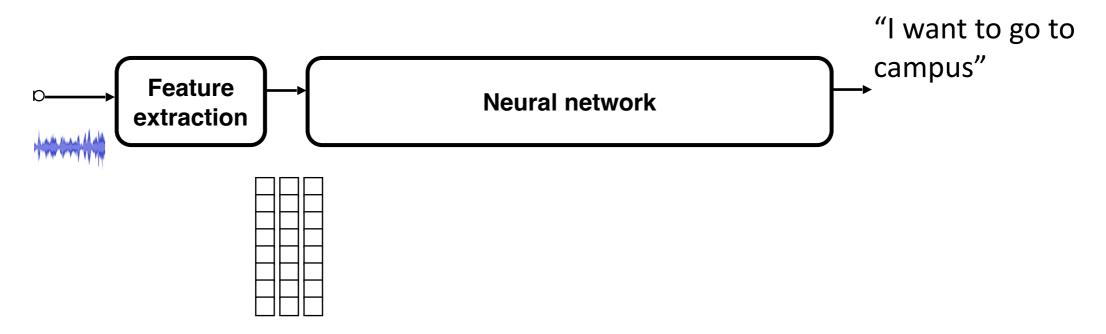
Building speech recognition was really difficult...





- We need to develop all components
- Each component requires a lot of background knowledge
- We need to tune hyper-parameters in each module

Next: End-to-end speech recognition



- We can simply the complicated models
- Optimize all components by using back propagation
- We still need some formulations to make a problem tractable

Today's agenda

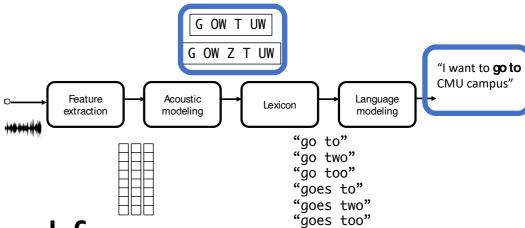


- Output
- Alignments

Output unit



• Our final goal: output a text



- Text can be represented by several forms
 - Word, Character, phoneme, etc.
 - We will discuss the characteristics of each output unit

Which unit? End-to-End ASR case



How to describe the phrase "go to"?

Word

- "go" and "to"
- More semantic/syntactic
- Very large vocabulary size, e.g., $|\mathcal{Y}|$ would be 100K
- Out of vocabulary issue ☺
- The length (N=2) is very short. Less computational cost, but larger mismatch between the input and output lengths

Character

- "g" "o" "_" "t" and "o" ("_" means the space)
- The vocabulary size is not large in general. ~30 in the Roman script, ~10K in the Chinese script
- No out of vocabulary issue (rarely happens) ©
- The length (N=5) becomes longer. More computational cost, but relaxing the mismatch between the input and output lengths
- BPE: Byte Pair Encoding (sentence piece)
 - "go to" \rightarrow "__g" "o" "_to" (N = 3)
 - Something between, we can also control the vocabulary size

In general, we do not specify which unit we use in our lecture since this is one of the model configurations.

Which unit? HMM-based (classical) system



How to describe the phrase "go to"?

Phoneme

- "go to" → "G OW T UW"
- More acoustic
- The phoneme vocabulary size is not very large in general
- The length (N = 5) becomes long (like the discussion in the (Roman) character).
- We need a dictionary
- If we use a phone, we can make this part language-independent

State (hidden Markov model state)

- We further decompose a phoneme into several states
- Furthermore acoustic
- Classically, we use this representation a lot (e.g., 3-state HMM)
- · It makes the unit further longer and the mismatch between the input and output lengths is further relaxed
- It will be introduced later in HMM

Which unit practically?

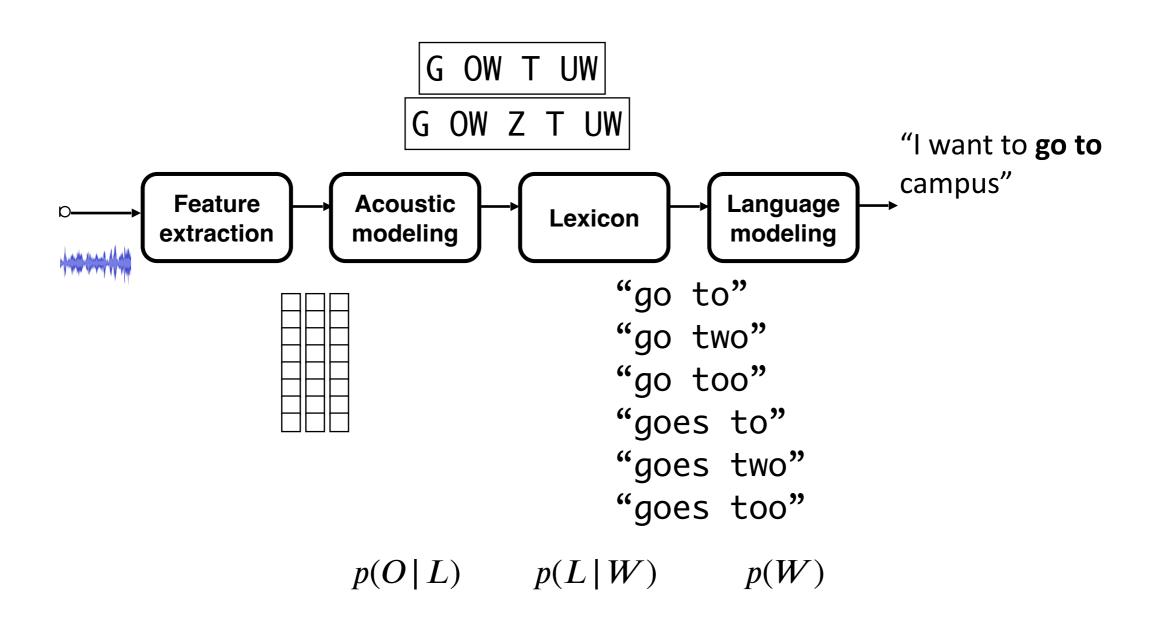


- Mostly, we use BPE (or sentence piece)
 - We need to set the maximum number depending on the training data
- Character for the low-resource case or Chinese and Japanese
 - Chinese/Japanese has ~10,000 characters
- Some languages do not have scripts
 - Phoneme, phone
 - Translation to the other languages (not ASR but speech translation)

Today's agenda

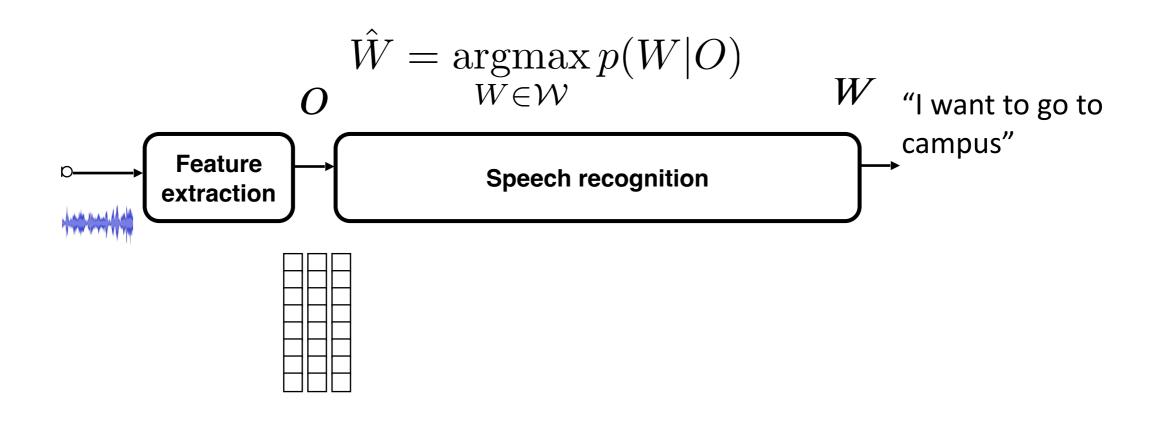
- Output
- Alignments

HMM-based (classical) speech recognition pipeline



Speech recognition





End-to-End ASR

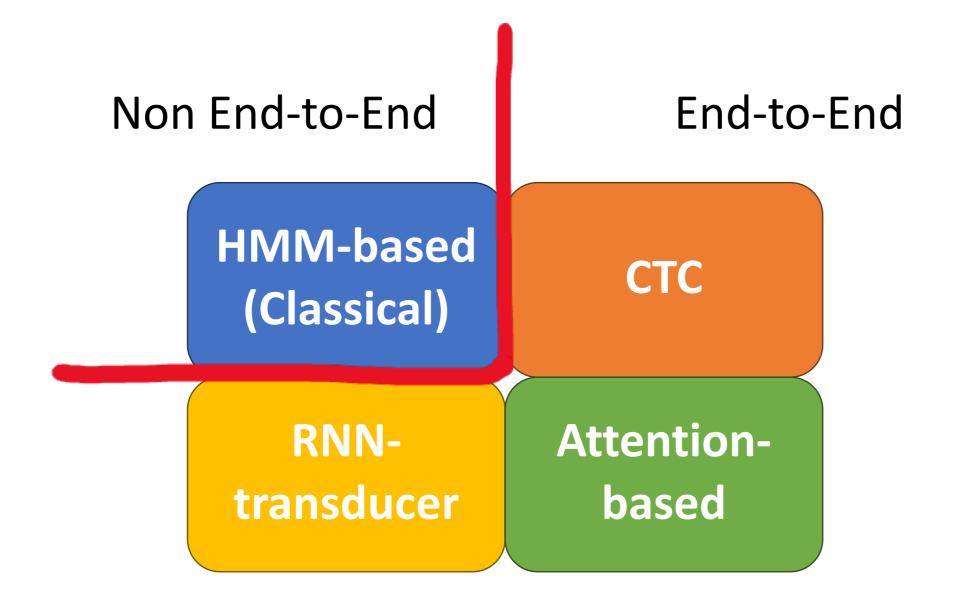


- Direct mapping function from speech feature sequence o
 to text w
- Usually, it does not deal with the phoneme-based intermediate representation
- Mainly three architectures
 - Attention-based
 - Decoder-only is also included here (not fully end-to-end)
 - Connectionist temporal classification (CTC)
 - Recurrent neural network transducer (RNN-T)

HMM-based (Classical)

CTC

RNNtransducer Attentionbased



Attention-based ASR



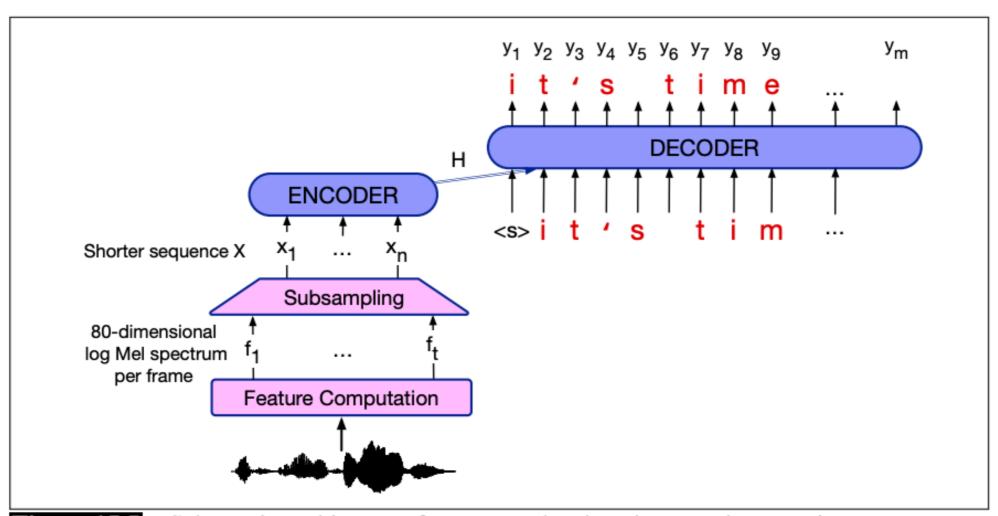


Figure 15.5 Schematic architecture for an encoder-decoder speech recognizer.

Listen Attend & Spell (2016)

Attention-based ASR



- Our staring point p(W|O)
 - Input: $O = (\mathbf{o}_t | t = 1,...,T)$
 - It is difficult to deal with $W = (w_i \in \mathcal{V} | i = 1,..., N)$
 - $T \neq N$
- Instead, we factorize p(W|O) as follows based on a probabilistic chain rule

$$p(W|O) = \prod_{i=1}^{N} p(w_i|w_{1:i-1}, O)$$

Attention-based ASR



- Our staring point p(W|O)
 - Input: $O = (\mathbf{o}_t | t = 1,...,T)$
 - It is difficult to deal with $W = (w_i \in \mathcal{V} | i = 1,..., N)$
 - $T \neq N$
- Instead, we factorize p(W|O) as follows based on a probabilistic chain rule

$$p(W|O) = \prod_{i=1}^{N} p(w_i|w_{1:i-1}, O)$$

- This neural network is handled by an attention-based method to align the input and output (soft alignments)
- We usually do not use the phoneme

Whisper (OpenAI)



 $p(W|O) = \prod_{i=1}^{N} p(w_i|w_{1:i-1}, O)$

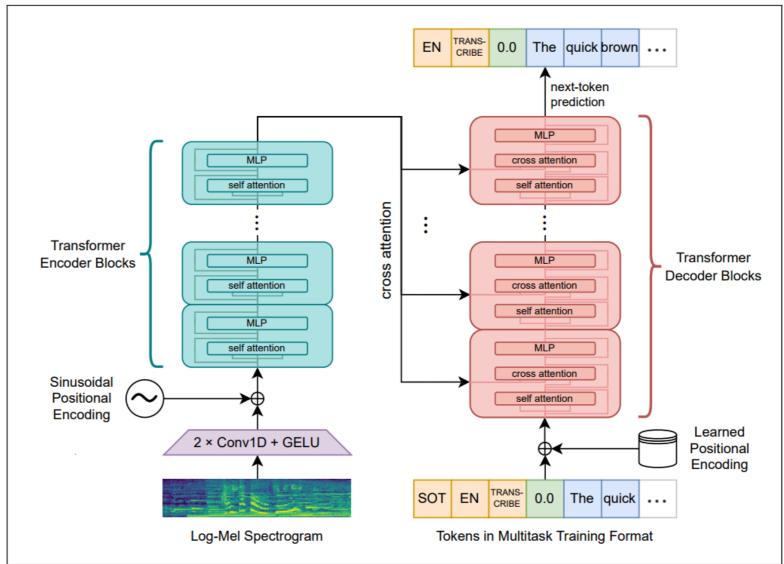
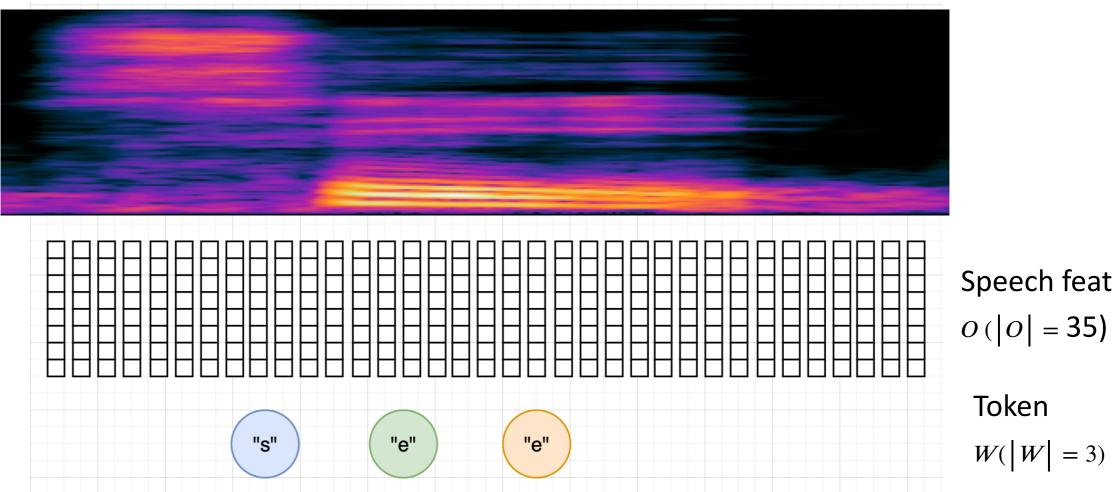


Figure 15.6 A sketch of the Whisper architecture from Radford et al. (2023). Because Whisper is a multitask system that also does translation and diarization (we'll discuss these non-ASR tasks in the following chapter), Whisper's transcription format has a Start of Transcript (SOT) token, a language tag, and then an instruction token for whether to transcribe or translate.

Alignments



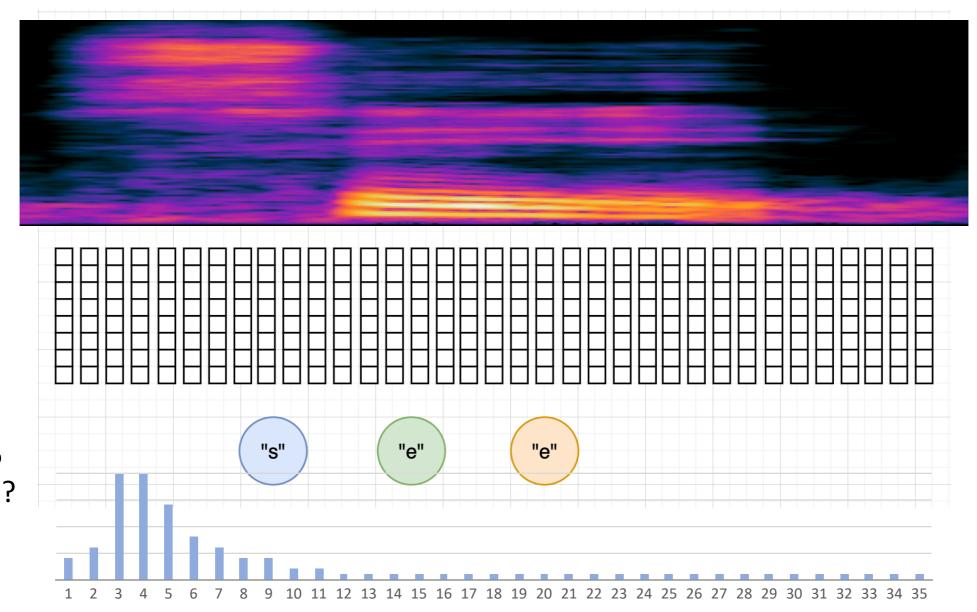


Speech features

$$W(|W|=3)$$

Soft alignments

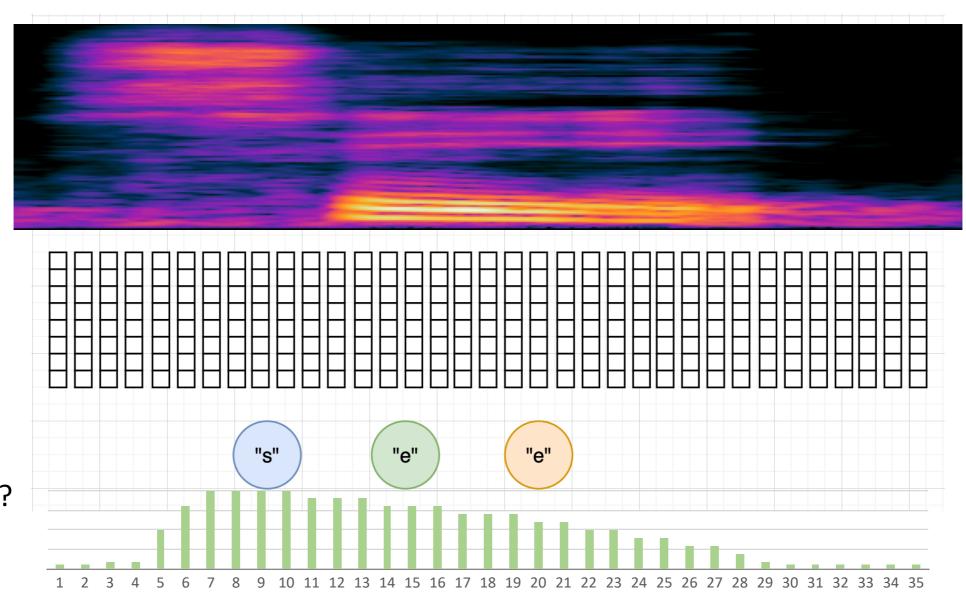




How many % "s" is aligned?

Soft alignments

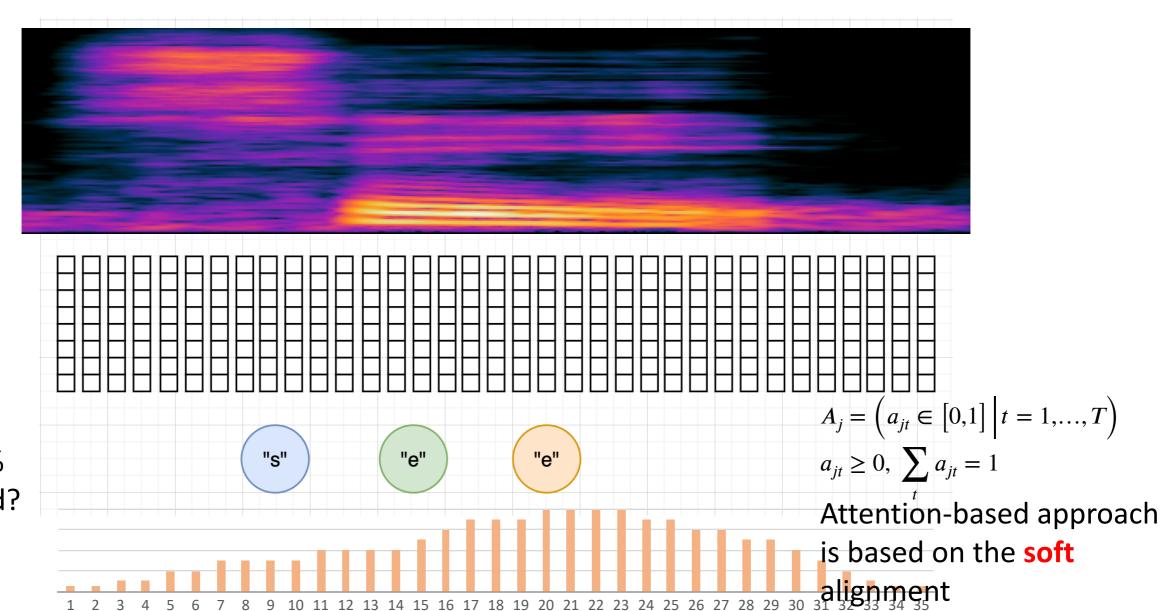




How many % "e" is aligned?

Soft alignments





How many % "e" is aligned?

Connectionist Temporal Classification



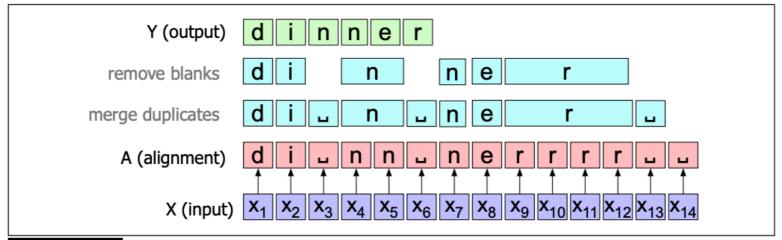


Figure 15.13 The CTC collapsing function B, showing the space blank character \Box ; repeated (consecutive) characters in an alignment A are removed to form the output Y.

The CTC collapsing function is many-to-one; lots of different alignments map to the same output string. For example, the alignment shown in Fig. 15.13 is not the only alignment that results in the string *dinner*. Fig. 15.14 shows some other alignments that would produce the same output.

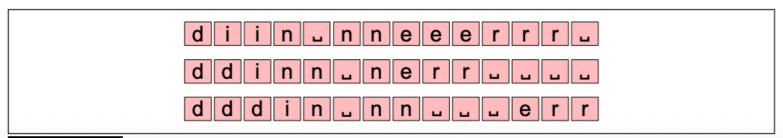


Figure 15.14 Three other legitimate alignments producing the transcript *dinner*.

RNN-Transducer



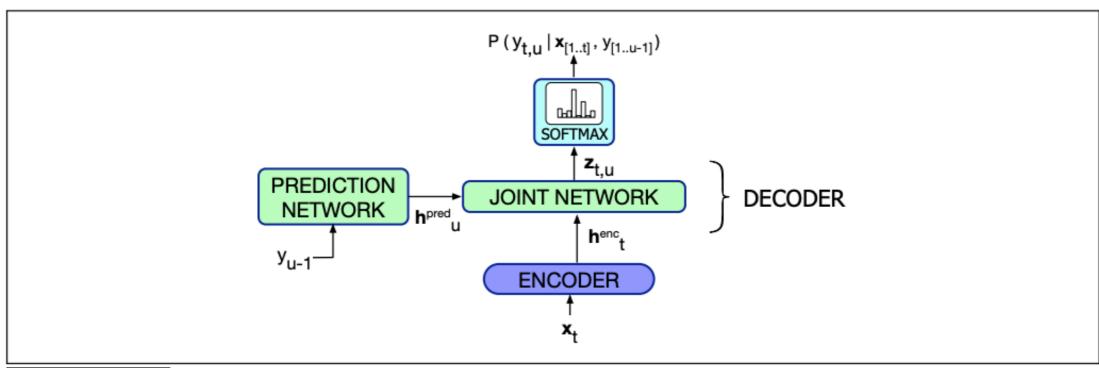


Figure 15.17 The RNN-T model computing the output token distribution at time t by integrating the output of a CTC acoustic encoder and a separate 'predictor' language model.

CTC or RNN transducer cases



- Again, our staring point p(W|O)
- We introduce a hard alignment path as a random variable:

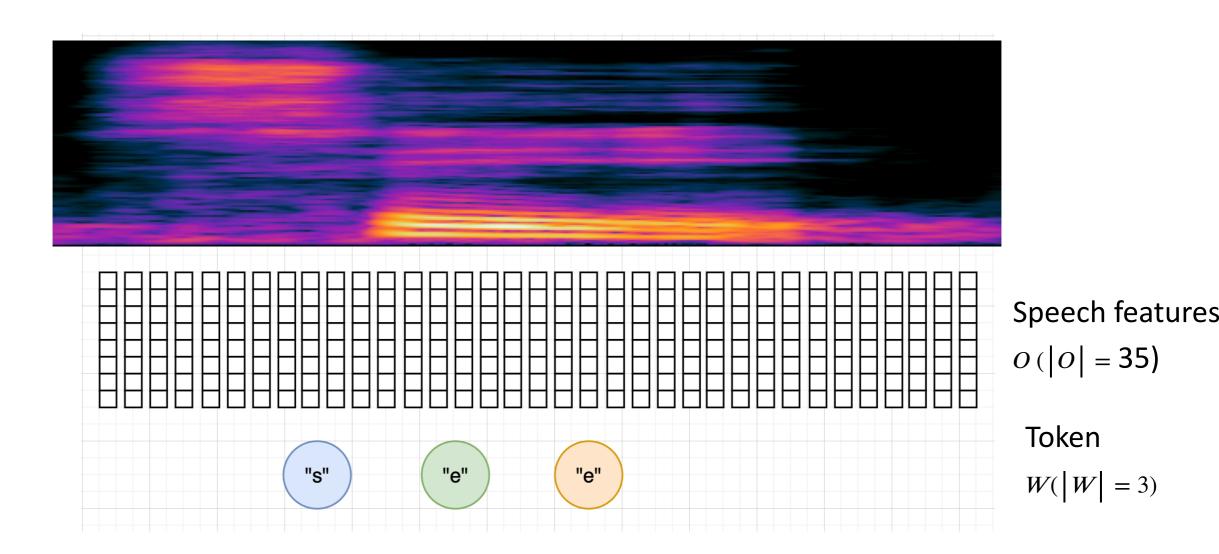
$$Z = (z_t \in \mathcal{V}' | t = 1, ..., T)$$

We consider the following equations

$$p(W|O) = \sum_{Z} p(Z, W|O)$$
$$= \sum_{Z} p(Z|O)p(W|Z, O)$$

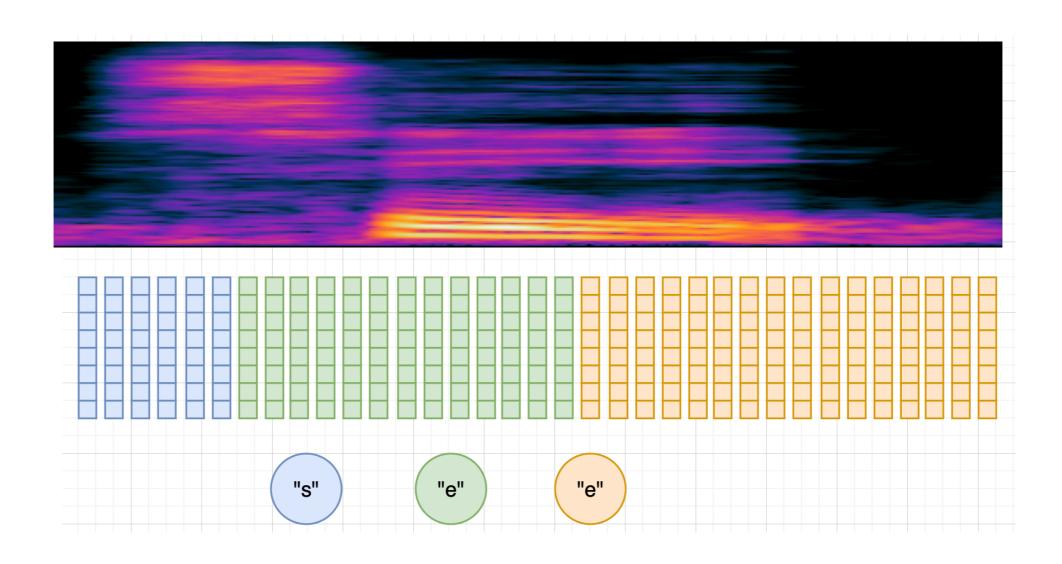
Alignments





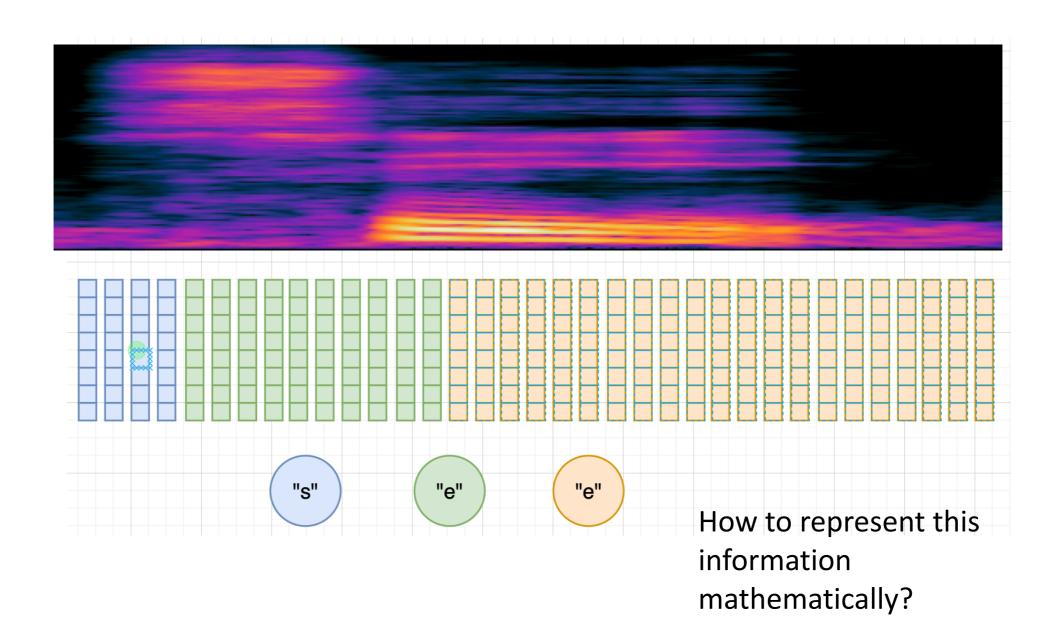
Hard Alignments





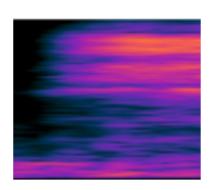
Hard Alignments

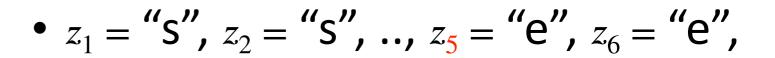




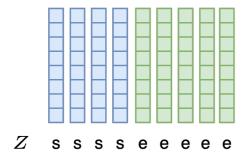
Alignment variable







One possible realization of alignments

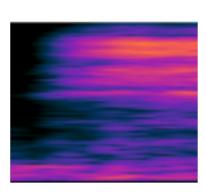






Alignment variable





• $z_1 =$ "s", $z_2 =$ "s", .., $z_7 =$ "e", $z_8 =$ "e",

The other possibility

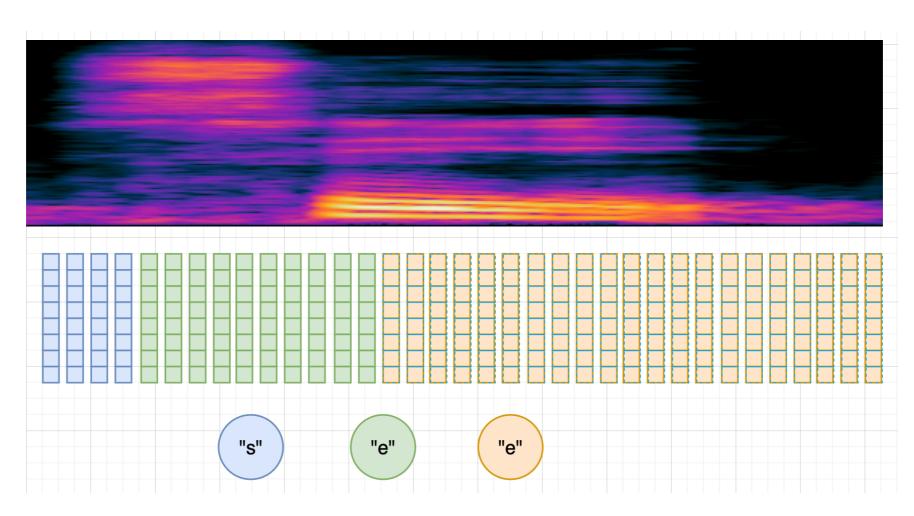






Hard alignments





We can represent a path z_1 = "s", z_2 = "s", .., z_5 = "e", ..., z_T = "e" Transducer are based on as a sequence variable This is called a hard alignment $z_t = (z_t \in \{ s, e^n, e^n, ... \} | t = 1,...,T)$

HMM, CTC, and RNNthe **hard** alignment

CTC case



- Again, our staring point p(W|O)
- We introduce a hard alignment path as a random variable: $Z = (z_t \in \mathcal{V}' | t = 1,...,T)$
- We consider a set of all possible alignment paths representing W. We call this set as $\mathcal{Z}(W)$
- We incorporate this random variable to p(W|O) as follows:

$$p(W|O) = \sum_{Z \in \mathcal{Z}(W)} p(Z|O)$$

• z is a sequence and it is difficult to deal with. Similar to the attention case, we use a chain rule to further factorize this equation:

$$p(W|O) = \sum_{Z \in \mathcal{Z}(W)} \prod_{t=1}^{T} p(z_t|z_{1:t}, O)$$

- This is not very difficult compared with an attention-based method, since z and o are based on the same length T
- RNN transducer can also be formulated in a similar manner

Difference between E2E and HMM-based

• p(W|O) (end-to-end)

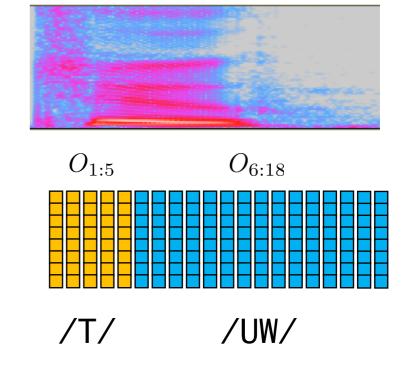
• p(O|L) (HMM based)

• In p(W|O) and p(O|L), there are two differences

Simple case



• We know that phoneme /T/ is aligned with $o_{1:5}$, and phoneme / UW/ is aligned with $o_{6:18}$



$$p(O_{1:18}|/T/, /UW/)$$

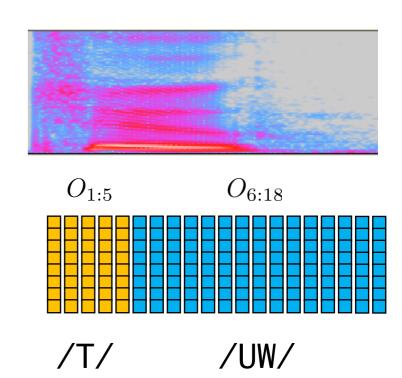
$$= p(O_{1:5}|/T/)p(O_{6:18}|/UW/)$$

This factorization is very natural if they are aligned

Simple case



• We know that phoneme /T/ is aligned with $o_{1:5}$, and phoneme / UW/ is aligned with $o_{6:18}$



$$p(O_{1:18}|/T/,/UW/) = p(O_{1:5}, O_{6:18}|/T/,/UW/)$$

$$= p(O_{1:5}|Q_{6:18},/T/,/UW/)p(O_{6:18}|/T/,/UW/)$$

$$= p(O_{1:5}|/T/)p(O_{6:18}|/UW/)$$

Conditional independence assumptions

- Note that the conditional independence assumptions are "approximation or modeling"
- Not based on the derivation.

General case



• Factorize for each phoneme:

```
\begin{split} p(O|L) &= p(O_{1:T_1}, O_{T_1+1:T_2}, \cdots | l_1, l_2, \cdots) \\ &= p(O_{1:T_1} | O_{T_1+1:T_2}, \cdots, l_1, l_2, \cdots) p(O_{T_1+1:T_2}, \cdots, | l_1, l_2, \cdots) \\ &= p(O_{1:T_1} | l_1) p(O_{T_1+1:T_2}, \cdots, | l_1, l_2, \cdots) \\ &\vdots \\ &= p(O_{1:T_1} | l_1) p(O_{T_1+1:T_2} | l_2) \dots \\ &= \prod_{i=1}^{J} p(O_{T_{j-1}+1:T_j} | l_j) \end{split} \qquad \text{ information, we can factor as a factor of the first of the fi
```

- If we know the alignment information, we can factorize the acoustic model probability for each phoneme
- We use conditional independence assumptions very aggressively

General case



(chain rule)

Conditional independe

assumption

Factorize for each phoneme:

$$\begin{split} p(O|L) &= p(O_{1:T_1}, O_{T_1+1:T_2}, \cdots | l_1, l_2, \cdots) \\ &= p(O_{1:T_1}|O_{T_1+1:T_2}, \cdots, l_1, l_2, \cdots) p(O_{T_1+1:T_2}, \cdots, | l_1, l_2, \cdots) \\ &= p(O_{1:T_1}|l_1) p(O_{T_1+1:T_2}, \cdots, | l_1, l_2, \cdots) \\ &\vdots \\ &= p(O_{1:T_1}|l_1) p(O_{T_1+1:T_2}|l_2) \dots \\ &= p(O_{1:T_1}|l_1) p(O_{T_1+1:T_2}|l_2) \dots \\ &= \prod_{i=1}^J p(O_{T_{j-1}+1:T_j}|l_j) \end{split} \qquad \text{ information, we can factorize the accustic model probability for accustic model probability for a solution of the probability for accustic model probability for accustic model probability for a solution of the probability for accustic model probabi$$

- If we know the alignment information, we can factorize the acoustic model probability for each phoneme
- We use conditional independence assumptions very aggressively

Word -> Phoneme -> State

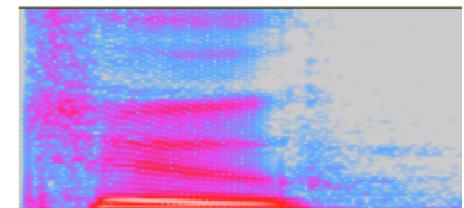


We need a more precise unit to present

an acoustic dynamics



Introduce phoneme representation



/SILB/ /T/ /UW/ /SILE/

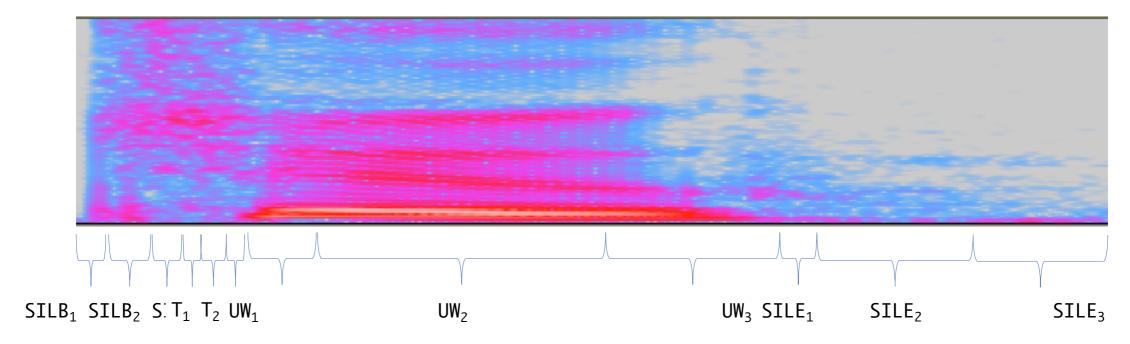
SILB₁ SILB₂ SILB₃ T₁ T₂ T₃ UW₁ UW₂ UW₃ SILE₁ SILE₂ SILE₃

Introduce further precise representation (we call it state)

How to align states to features?

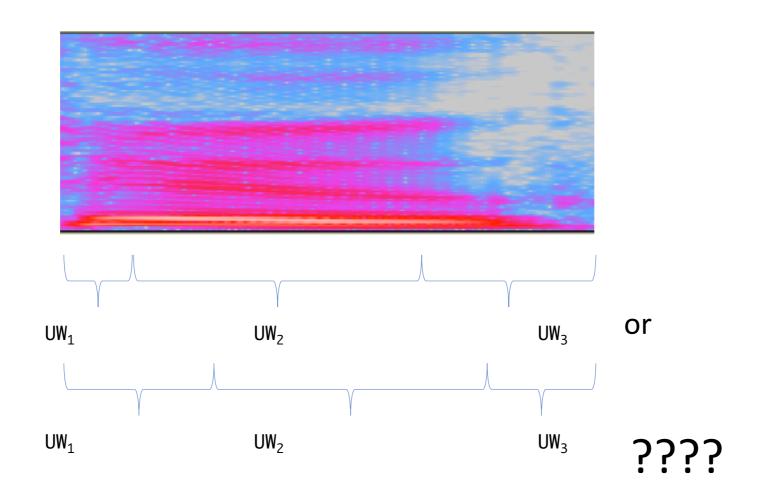


Alignment problem



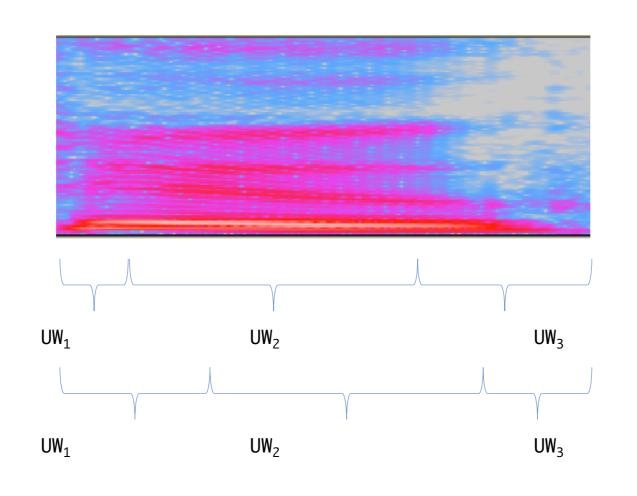
We don't know the alignment information...

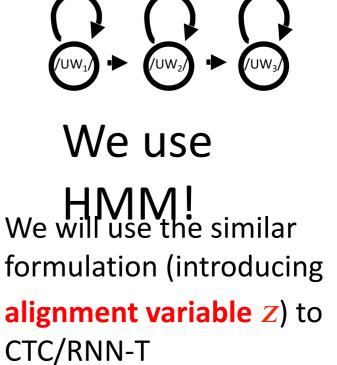




We don't know the alignment information...







Interim Summary



- The most difficult issue in speech recognition: the input and output lengths are different
 - We need some alignments
- The soft alignment case
 - We use an attention-based neural network (this will be introduced later).
- The hard alignment cases
 - We explicitly introduce an alignment path z.
 - We can use it for CTC, RNN-transducer, and HMM-based approach



HMM-based (Classical)

CTC

RNNtransducer Attentionbased



Non End-to-End

End-to-End

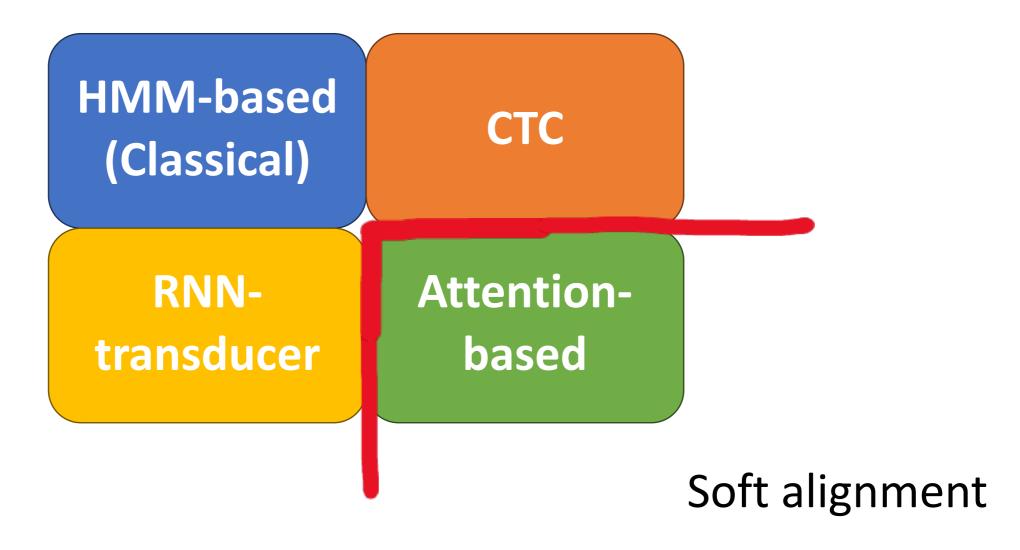
HMM-based (Classical)

CTC

RNNtransducer Attentionbased



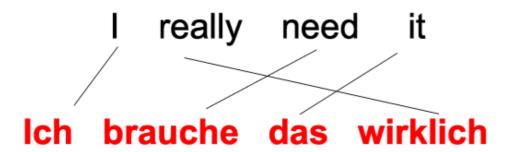
Hard alignment



Monotonic vs. non-monotonic



Machine translation case (non monotonic)



• Speec ition case (monotonic)

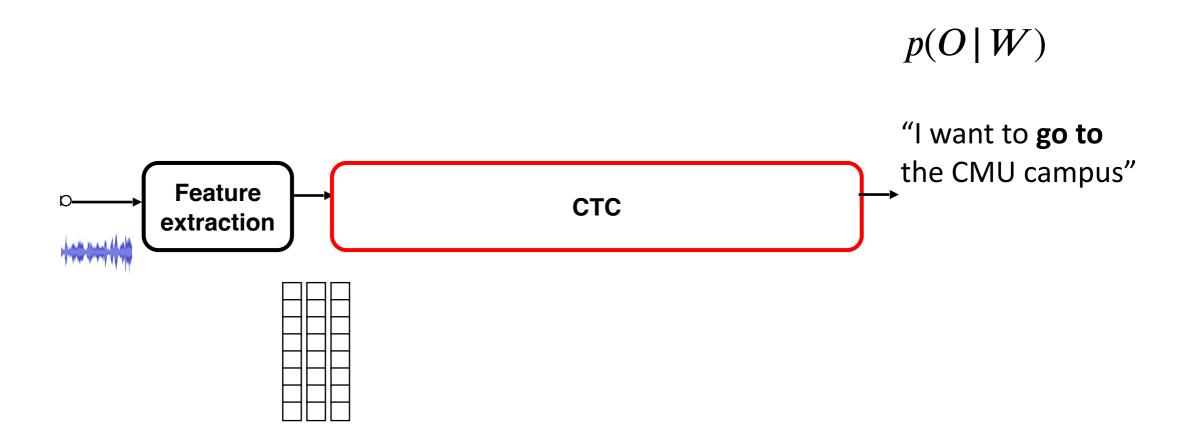
Agenda



- Alignment paths
 - CTC
 - RNN-Transducer
 - HMM (3-state left-to-right)

Speech recognition pipeline

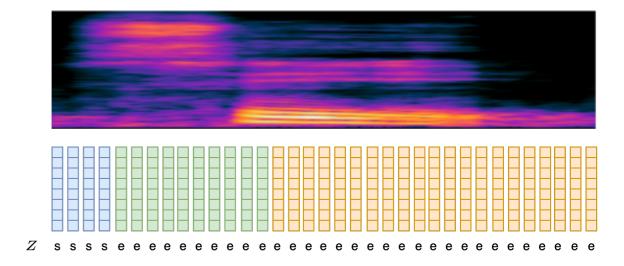




Hard alignments for repeated tokens



- Hard alignment examples
- How to distinguish them?
 - ("s", "e", "e")
 vs.
 - ("s", "e")
 - Both are written as





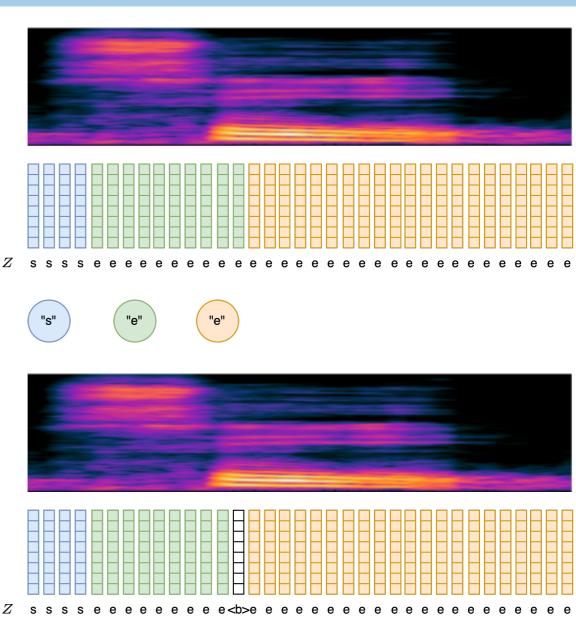




Hard alignments for repeated tokens



- Hard alignment examples
- How to distinguish them?
 - ("s", "e", "e")
 vs.
 - ("s", "e")
 - Both are written as
- We introduce the blank symbol!
 - Z = (ssseeeeeeeee < b > eee...)









Introduction of blank symbol



First, we insert to the character sequence "see"

$$\rightarrow W = (\text{"s", "e", "e"}), \text{ where } |W| = J$$

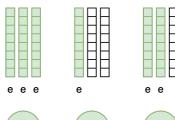
- Then, expand W' to the frame length T to form Z
 - Assuming that T > J in general: We cannot use CTC if it is not satisfied
 - All tokens and can be repeated to adjust the length
 - For example, if "e" is repeated three times

•
$$W = \text{"e"} \rightarrow Z = (\text{"e"}, \text{"e"}, \text{"e"})$$

•
$$W = \text{"e"} \rightarrow Z = (\text{"e"}, \text{"e"}, \text{""})$$

•
$$W = \text{"e"} \rightarrow Z = (\text{"e"}, \text{""}, \text{""})$$

- must be inserted between repeated character
 - W = ("e", "e"), then Z = (..."e", "", "e"...): we cannot skip
 - W = ("s", "e"), then, Z = (..."s", "e"...): we can skip







Example of z

•
$$W = (\text{"s", "e", "e"}), T = 5$$

Then

$$Z = (\text{``s''}, \text{``s''}, \text{``e''}, \text{``e''}), \text{ or }$$
 $Z = (\text{``s''}, \text{``s''}, \text{``e''}, \text{``e''}, \text{``e''}), \text{ or }$
 $Z = (\text{``s''}, \text{``s''}, \text{``e''}, \text{``e''}), \text{ or }$

This is an alignment problem

"see" (only 5 frames e s or

Example of z

•
$$W = (\text{"s", "e", "e"}), T = 5$$

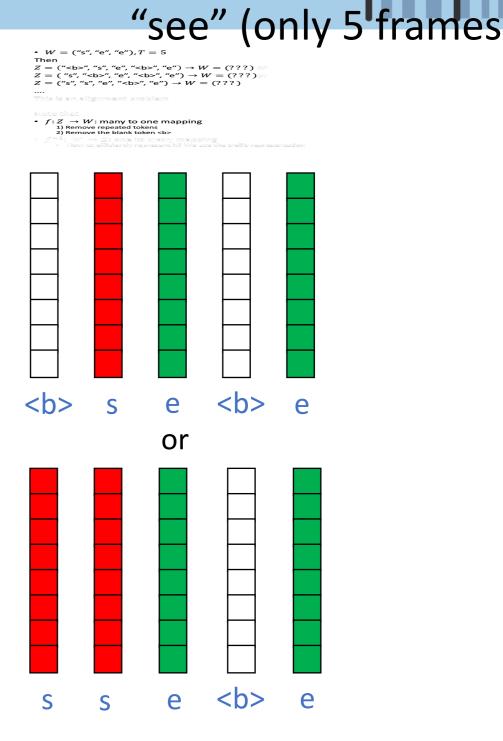
Then

$$Z = (\text{``s''}, \text{``s''}, \text{``e''}, \text{``e''}) \rightarrow W = (???)$$

$$Z = (\text{``s''}, \text{``e''}, \text{``e''}, \text{``e''}) \rightarrow W = (???)$$

$$Z = (\text{``s''}, \text{``s''}, \text{``e''}, \text{``e''}) \rightarrow W = (???)$$

- $f: Z \to W$: many to one mapping
 - 1) Remove repeated tokens
 - 2) Remove the blank token



Example of z

•
$$W = (\text{"s", "e", "e"}), T = 5$$

Then

$$Z = (\text{``} < \text{b}>\text{''}, \text{``} < \text{b}>\text{''},$$

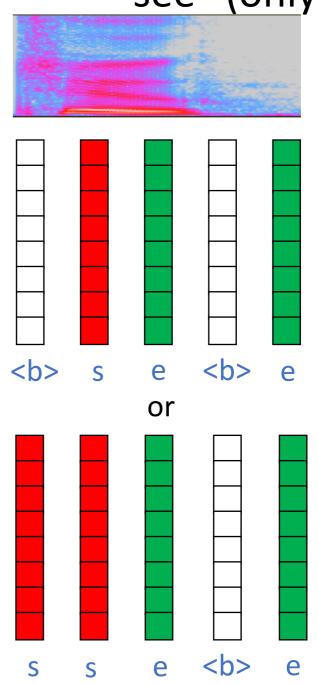
...

This is an alignment problem

Note that

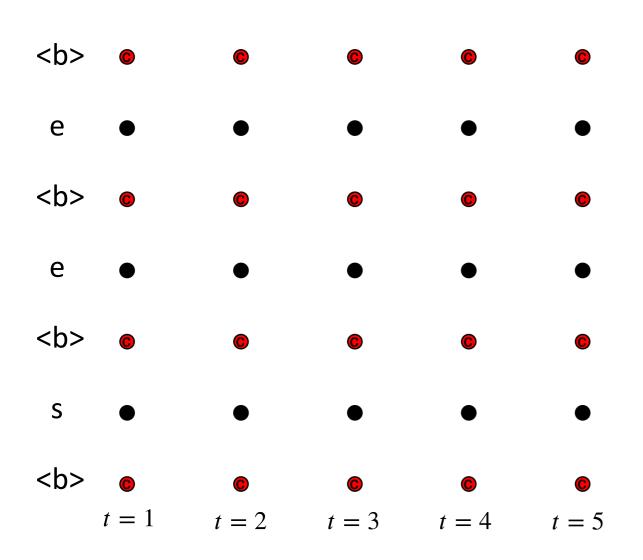
- $f: Z \to W$: many to one mapping
 - 1) Remove repeated tokens
 - 2) Remove the blank token
- f^{-1} : $W \to Z$: one to many mapping
 - How to efficiently represent it? We use the trellis representation

"see" (only 5 frames)





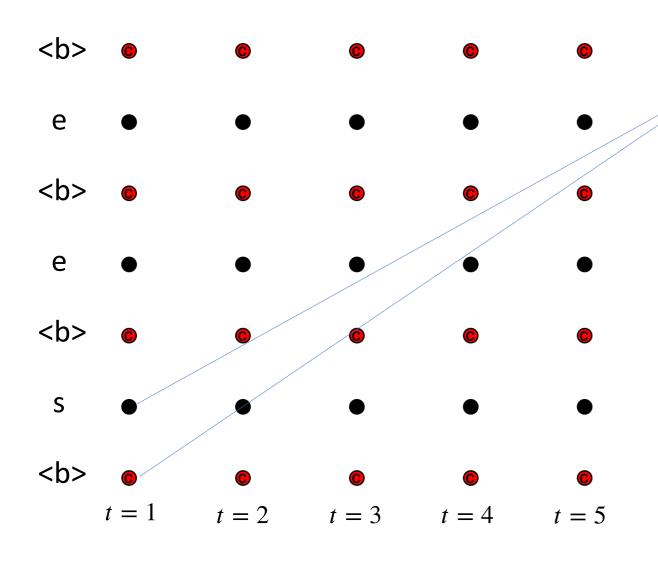




- 1. Must start at the first or s
- 2. Must end at the last or e
- 3. All characters can be repeated
- can be skipped except when it is inserted between repeated character
 - "s", "", "e": we can skip
 - "e", "", "e": we cannot skip

85

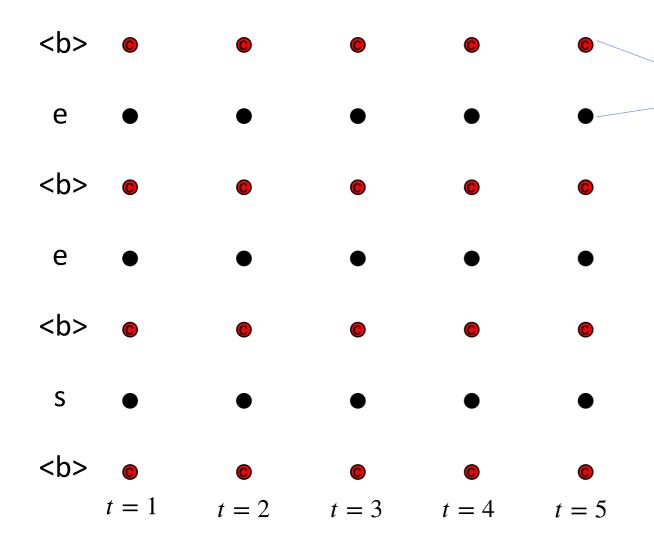




- 1. Must start at the first or s
- 2. Must end at the last or e
- 3. All characters can be repeated
- 4. can be skipped except when it is inserted between repeated character
 - "s", "", "e": we can skip
 - "e", "", "e": we cannot skip

86

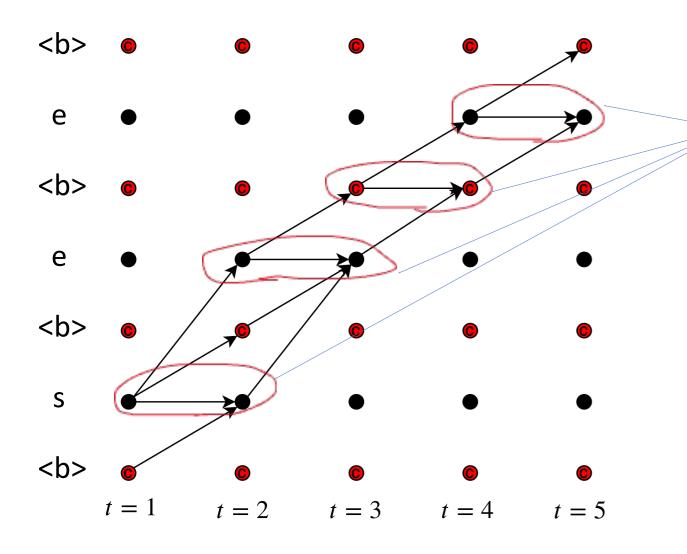




- 1. Must start at the first or s
- 2. Must end at the last or e
- 3. All characters can be repeated
- 4. can be skipped except when it is inserted between repeated character
 - "s", "", "e": we can skip
 - "e", "", "e": we cannot skip

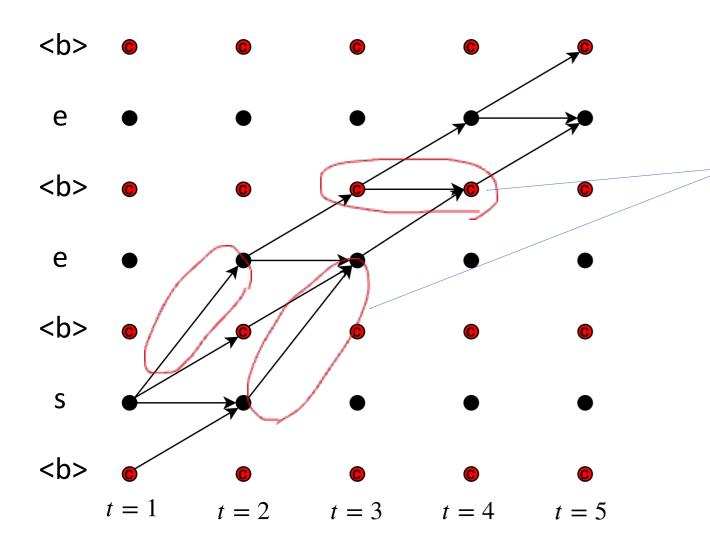
87





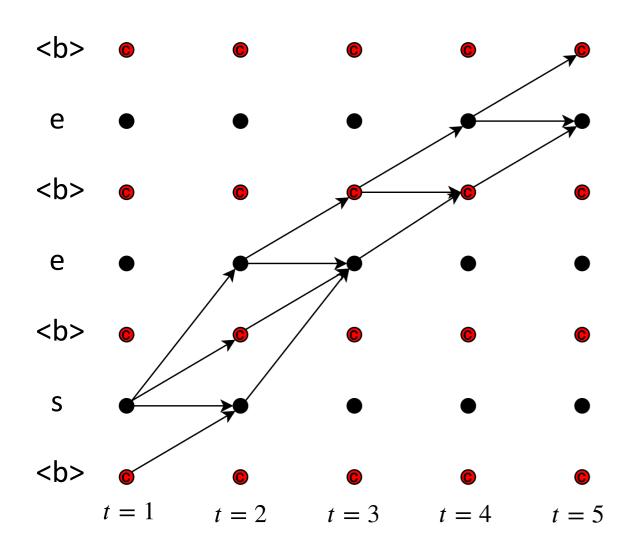
- 1. Must start at the first or s
- 2. Must end at the last or e
- 3. All characters can be repeated
- can be skipped except when it is inserted between repeated character
 - "s", "", "e": we can skip
 - "e", "", "e": we cannot skip





- 1. Must start at the first or s
- 2. Must end at the last or e
- 3. All characters can be repeated
- can be skipped except when it is inserted between repeated character
 - "s", "", "e": we can skip
 - "e", "", "e": we cannot skip





- 1. Must start at the first or s
- 2. Must end at the last or e
- 3. All characters can be repeated
- can be skipped except when it is inserted between repeated character
 - "s", "", "e": we can skip
 - "e", "", "e": we cannot skip



	•	•	•	•	•	•
е	•	•	•	•	•	•
	•	•	•	•	•	•
е	•	•	•	•	•	•
	•	•	•	•	•	•
S	•	•	•	•	•	•
	•	©	©	•	©	•
	t = 1	t = 2	t = 3	t = 4	t = 5	t = 6

How to represent it as an actual function?

- We define an alignment variable as $Z = (z_t \in \mathcal{V} \cup \langle b \rangle | t = 1,...,T)$
 - |Z| = |O| = T
 - The vocabulary set is augmented by the blank symbol
- We incorporate this random variable to $p(W \mid O)$ as follows:

$$p(W|O) = \sum_{Z \in f^{-1}(W)} p(Z|O)$$

$$= \sum_{Z \in f^{-1}(W)} \prod_{t=1}^{T} p(z_t|\underline{z_{1:t-1}}, O)$$

$$= \sum_{Z \in f^{-1}(W)} \prod_{t=1}^{T} p(z_t|\underline{z_{1:t-1}}, O)$$

$$= \sum_{Z \in f^{-1}(W)} \prod_{t=1}^{T} p(z_t|O)$$
• $f^{-1}(W)$: all possible Z representing W
• Use a chain rule to further factorize Z
• Perform conditional independence assumptions and ignore the dependency of previous alignments
• (this is an approximation)

- $f^{-1}(W)$: all possible Z representing W
 - - assumptions and ignore the dependency of previous alignments $z_{1:t-1}$ (this is an approximation)
- $p(z_t|O)$ is represented by a neural network (Bidirectional LSTM or self-attention) over all possible paths \sum

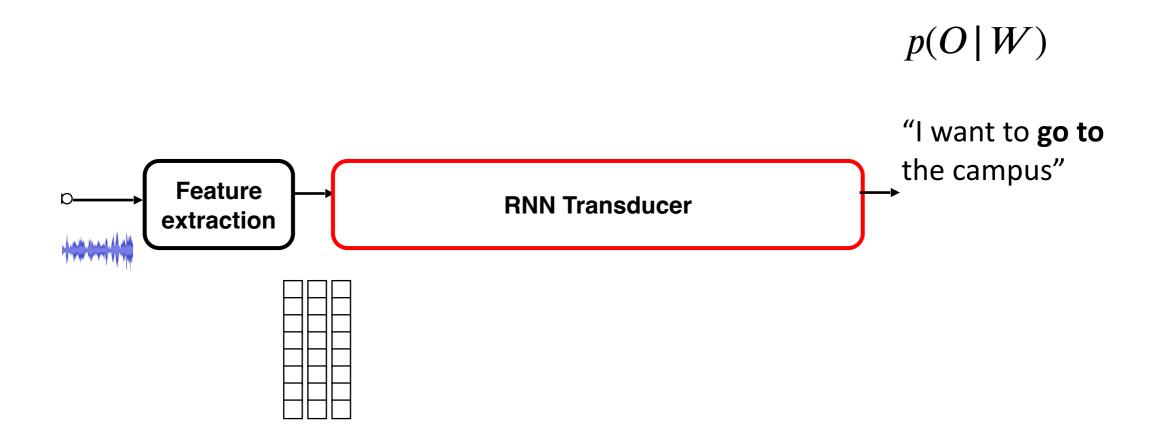
←Exponential

Agenda

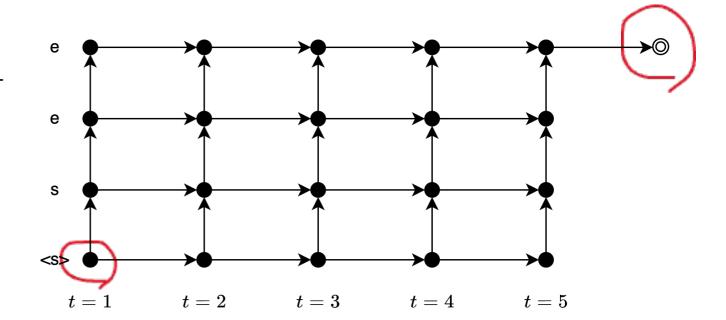
- Alignment paths
 - CTC
 - RNN-Transducer
 - HMM (3-state left-to-right)

Speech recognition pipeline

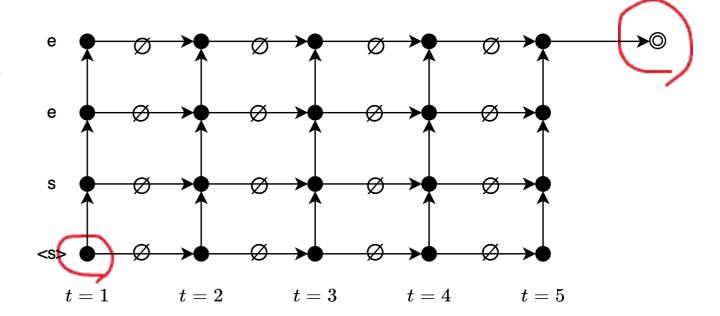




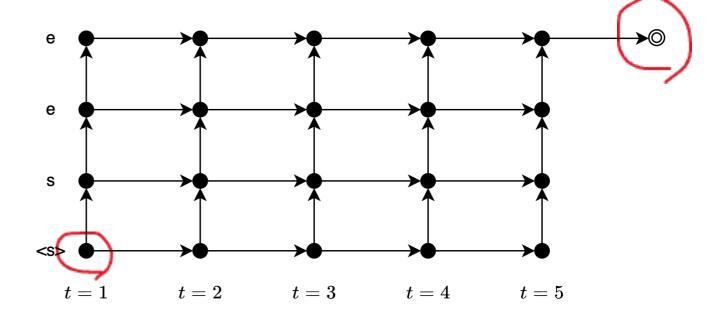
- Alternative alignment paths
 - $W = (w_1, ..., w_J)$: *J*-length label sequence
 - $O = (o_1, ..., o_T)$: T-length input sequence
- Similar to CTC, consider all possible paths in the rightside trellis
 - Start: **left bottom corner**, i.e., <s> (start of sentence) at *t* = 1
 - End: top right corner
 - ↑: output token (it does not consume the input frame, unlike CTC)!



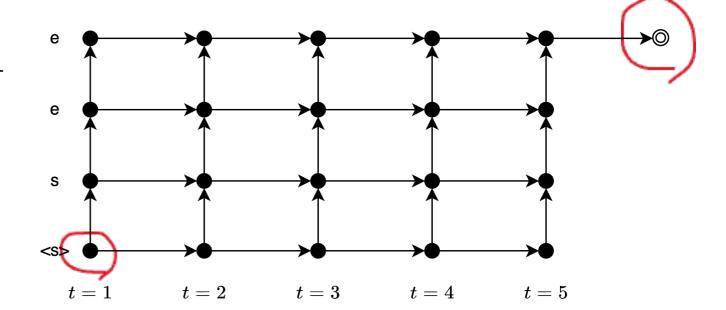
- Alternative alignment paths
 - $W = (w_1, ..., w_J)$: *J*-length label sequence
 - $O = (o_1, ..., o_T)$: T-length input sequence
- Similar to CTC, consider all possible paths in the rightside trellis
 - Start: **left bottom corner**, i.e., <s> (start of sentence) at *t* = 1
 - End: top right corner
 - ↑: output token (it does not consume the input frame, unlike CTC)!
 - →: no output (∅)



- Alternative alignment paths
 - $W = (w_1, ..., w_J)$: *J*-length label sequence
 - $O = (o_1, ..., o_T)$: T-length input sequence
- Similar to CTC, consider all possible paths in the rightside trellis
 - Start: left bottom corner, i.e., <s> (start of sentence)
 at t = 1
 - End: top right corner
 - ↑: output token (it does not consume the input frame, unlike CTC)!
 - →: no output (∅)
 - f^{-1} : W=('s', 'e', 'e') \rightarrow Z=(\varnothing , 's', \varnothing , \varnothing , 'e', 'e', \varnothing , \varnothing), ('s', \varnothing , 'e', \varnothing , \varnothing , \varnothing , 'e', \varnothing) etc. if T = 5 and J = 3
 - $f: Z \to W$ is a many to one mapping (we can just remove \emptyset
- Consider all possible paths including the case that it allows the output of the token ${\it W}$

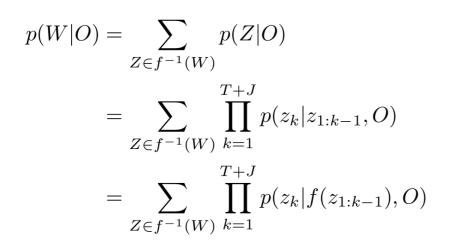


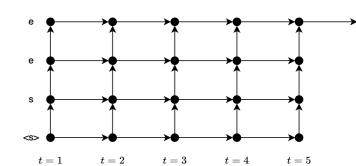
- Alternative alignment paths
 - $W = (w_1, ..., w_J)$: *J*-length label sequence
 - $O = (o_1, ..., o_T)$: T-length input sequence
- Similar to CTC, consider all possible paths in the rightside trellis
 - Start: left bottom corner, i.e., <s> (start of sentence)
 at t = 1
 - End: top right corner
 - ↑: output token (it does not consume the input frame, unlike CTC)!
 - →: no output (∅)
 - f^{-1} :W=('s', 'e', 'e') \rightarrow Z=(\varnothing , 's', \varnothing , \varnothing , 'e', 'e', \varnothing , \varnothing), ('s', \varnothing , 'e', \varnothing , \varnothing , \varnothing , 'e', \varnothing) etc. if T = 5 and J = 3
 - $f: Z \to W$ is a many to one mapping (we can just remove \emptyset
- Consider all possible paths including the case that it allows the output of the token ${\it W}$



How to represent it as an actual function?

- We define an alignment variable as $Z = (z_k \in \mathcal{V} \cup \emptyset \mid k = 1,...,T + J)$
 - Vocabulary set is augmented by the blank symbol \varnothing
 - Note that |Z| = T + J, not T
- We can introduce the following equations





- $f^{-1}(W)$: all possible Z representing W
- Use a chain rule to further factorize Z
- $f(z_{1:k-1})$ represents the partial token sequence up to k-1

How to represent it as an actual function?

 $f(z_{1:k-1}) = s, e,$

- We define an alignment variable as $Z = (z_k \in \mathcal{V} \cup \emptyset \mid k = 1,...,T + J)$
 - Vocabulary set is augmented by the blank symbol \varnothing
 - Note that |Z| = T + J, not T
- We can introduce the following equations

$$p(W|O) = \sum_{Z \in f^{-1}(W)} p(Z|O)$$

$$= \sum_{Z \in f^{-1}(W)} \prod_{k=1}^{T+J} p(z_k|z_{1:k-1}, O)$$

$$= \sum_{Z \in f^{-1}(W)} \prod_{k=1}^{T+J} p(z_k|f(z_{1:k-1}), O)$$

- $f^{-1}(W)$: all possible Z representing W_7
- Use a chain rule to further factorize $\overset{ullet}{Z}$
- $f(z_{1:k-1})$ represents the partial token sequence up to k-1
- We need to solve the summuzation over all possible paths \sum_{z}

Example of z in CTC

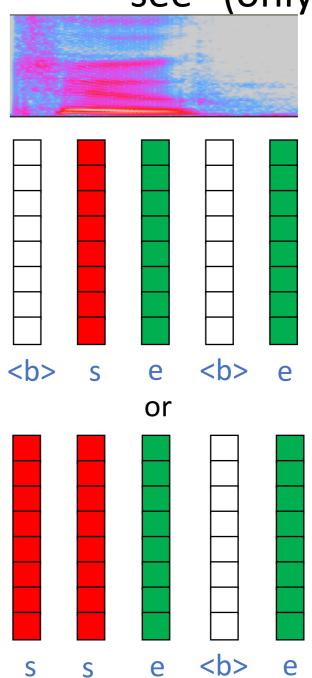
"see" (only 5 frames)

•
$$W = (\text{"s", "e", "e"}), T = 5$$

Then

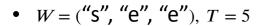
$$Z = (\text{``s''}, \text{``s''}, \text{``e''}, \text{``e''}), \text{ or }$$
 $Z = (\text{``s''}, \text{``s''}, \text{``e''}, \text{``e''}, \text{``e''}), \text{ or }$
 $Z = (\text{``s''}, \text{``s''}, \text{``e''}, \text{``e''}), \text{ or }$

This is an alignment problem



Example of z in RNN-trans.

"see" (only 5 frames)



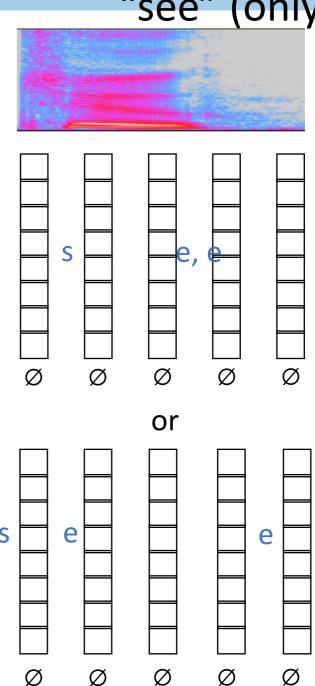
Then

$$Z = (\emptyset, \text{`s'}, \emptyset, \emptyset, \text{`e'}, \text{`e'}, \emptyset, \emptyset), \text{ Or }$$

$$Z = (\text{`s'}, \emptyset, \text{`e'}, \emptyset, \emptyset, \emptyset, \text{`e'}, \emptyset), Or$$

....

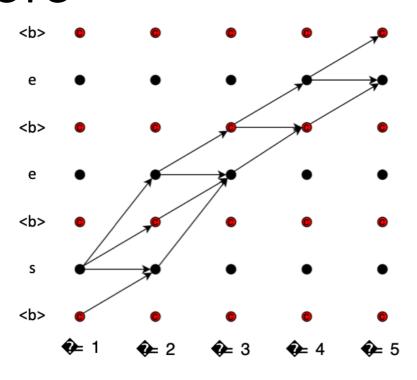
This is an alignment problem



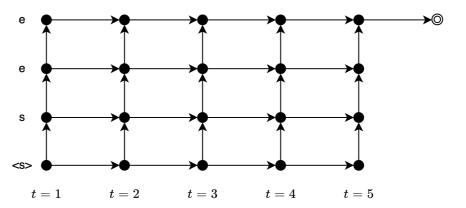
CTC vs. RNN-T



• CTC



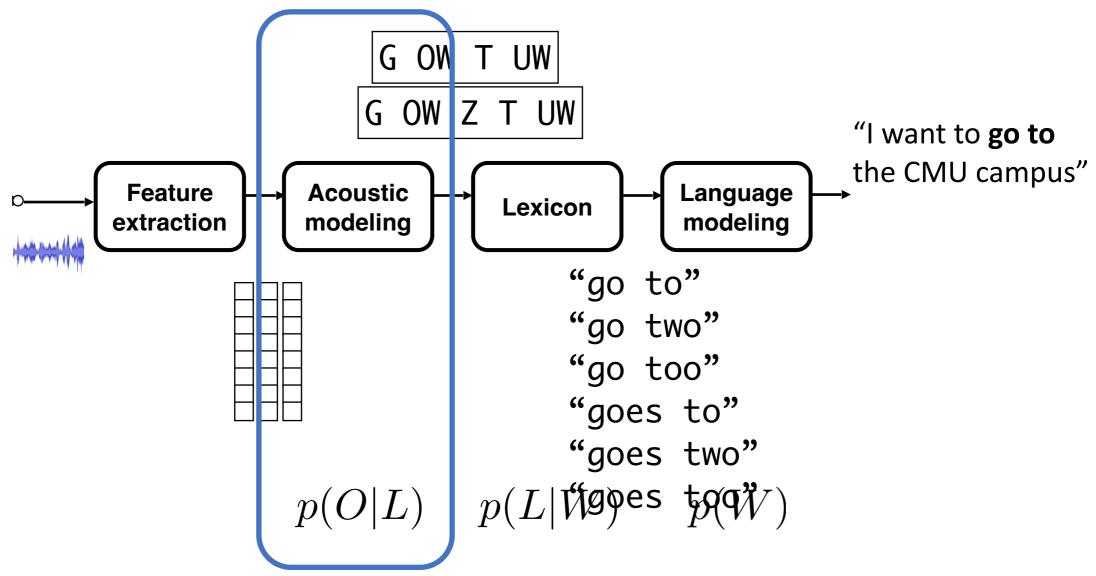
• RNN-transducer



Agenda

- Alignment paths
 - CTC
 - RNN-Transducer
 - HMM (3-state left-to-right)

Speech recognition pipeline



Introduction of HMM states

- Two differences
 - Phoneme sequence L instead of word sequence W
 - p(O|L) instead of p(L|O) or p(W|O) due to the Bayes theorem
- 1. Decompose the word sequence to the phoneme sequence by using a pronunciation dictionary

- 2. Introduce the silence related symbols
 - Silence begin (/SilB/): placed in the beginning of the sentence
 - Silence end (/SilE/): placed in the end of the sentence
 - Short pause (/SP/): placed between the words (can be skipped)
 /W/, /AH/, /N/, /T/, /UW/ → /SilB/, /W/, /AH/, /N/, /SP/, /T/, /UW/, /SilE/
- 3. Expand each phoneme with **three** states

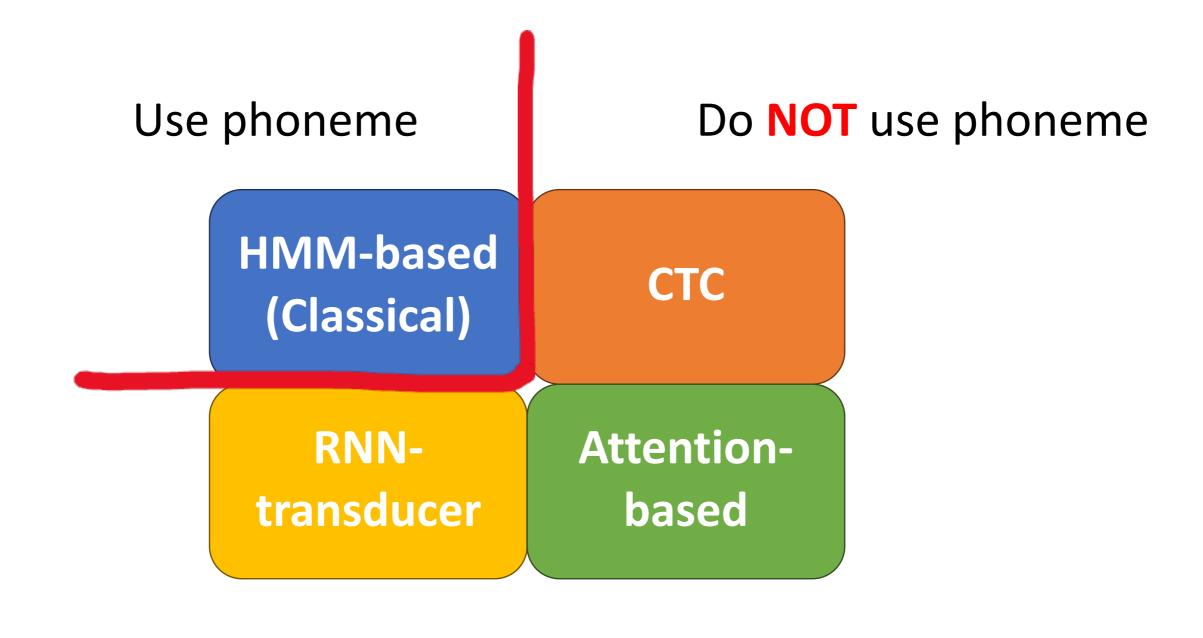
```
/SilB/, /W/, /AH/, /N/, /SP/, /T/, /UW/, /SilE/
```

- → /SilB/₁, /SilB/₂, /SilB/₃, ..., /T/₁, /T/₂, /T/₃, /UW/₁, /UW/₂, /UW/₃, ...
- This structure is widely used, but there are some variants (e.g., short pause can be one state)

HMM-based (Classical)

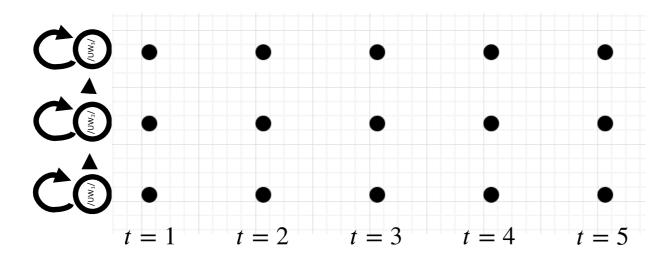
CTC

RNNtransducer Attentionbased

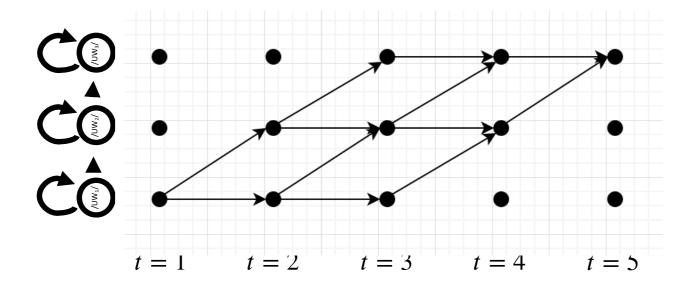




- Left-to-right HMM without skip
- (/UW₁/) **** (/UW₂/) **** (/UW₃/)
 - Each state must have at least one frame
- Example, 3 state left-to-right HMM (/UW₁/, /UW₂/, / UW₃/) $S = (s_t \in \{/\text{UW}_1/, /\text{UW}_2/, /\text{UW}_3/\}|t = 1, 2, 3, 4, 5)$ quences



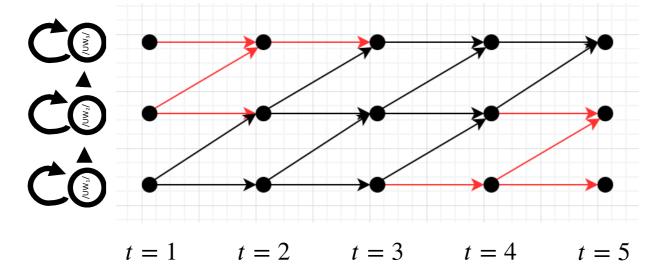




- It must start at the initial state (/UW₁/)
- It must end at the final state (/UW₃/)

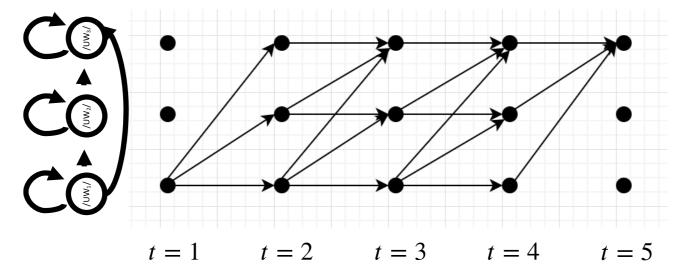


Red arrows cannot reach the final state



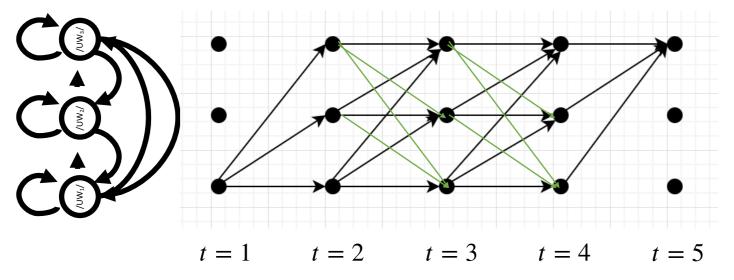


If we have some skips...





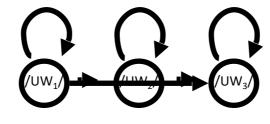
If we have some skips...



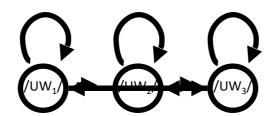
Variant of HMMs



- Left to right
 - We mainly use this
- Left to right with skip
 - 2nd order in this example



- Ergodic (fully connected)
 - No specific order anymore...



Phoneme sequence



We just connect the HMMs for each phoneme

Let's check three distributions



• CTC

$$p(W|O) = \sum_{Z \in f^{-1}(W)} p(Z|O)$$

$$= \sum_{Z \in f^{-1}(W)} \prod_{t=1}^{T} p(z_t | \underline{z_{1:t-1}}, O)$$

$$= \sum_{Z \in f^{-1}(W)} \prod_{t=1}^{T} p(z_t | O)$$

RNN transducer

$$p(W|O) = \sum_{Z \in f^{-1}(W)} p(Z|O)$$

$$= \sum_{Z \in f^{-1}(W)} \prod_{k=1}^{T+J} p(z_k|z_{1:k-1}, O)$$

$$= \sum_{Z \in f^{-1}(W)} \prod_{k=1}^{T+J} p(z_k|f(z_{1:k-1}), O)$$

• HMM

$$p(W|O) \propto \sum_{L} p(O|L)p(L|W)p(W)$$

$$= \sum_{L} \sum_{S} p(O, S|L)p(L|W)p(W)$$

$$= \sum_{L} \sum_{S} p(\mathbf{o}_{1}|s_{1}, L)p(s_{1}|L) \prod_{t=2}^{T} (\mathbf{o}_{t}|s_{t}, L)p(s_{t}|s_{t-1}, L)p(L|W)p(W)$$

Which part is similar?



Let's check three distributions

• CTC

$$p(W|O) = \sum_{Z \in f^{-1}(W)} p(Z|O)$$

$$= \sum_{Z \in f^{-1}(W)} \prod_{t=1}^{T} p(z_t | \underline{z_{1:t-1}}, O)$$

$$= \sum_{Z \in f^{-1}(W)} \prod_{t=1}^{T} p(z_t | O)$$

RNN transducer

$$p(W|O) = \sum_{Z \in f^{-1}(W)} p(Z|O)$$

$$= \sum_{Z \in f^{-1}(W)} \prod_{k=1}^{T+J} p(z_k|z_{1:k-1}, O)$$

$$= \sum_{Z \in f^{-1}(W)} \prod_{k=1}^{T+J} p(z_k|f(z_{1:k-1}), O)$$

• HMM

$$p(W|O) \propto \sum_{L} p(O|L)p(L|W)p(W)$$

$$= \sum_{L} \sum_{S} p(O, S|L)p(L|W)p(W)$$

$$= \sum_{L} \sum_{S} p(\mathbf{o}_{1}|s_{1}, L)p(s_{1}|L) \prod_{t=2}^{T} (\mathbf{o}_{t}|s_{t}, L)p(s_{t}|s_{t-1}, L)p(L|W)p(W)$$

Which part is similar?

- \rightarrow All are based on $\sum_{Z \text{ or } S} \prod_{t} f(t)$
- → We can compute them based on the dynamic programming

Summary



- CTC, RNN-transducer, and HMM are represented by a similar function form, i.e., $\sum_{Z \text{ or } S} \prod_{t} f(t)$
- All possible alignments are constrained by each method
- The most difficult issue in speech recognition: the input and output lengths are different
 - We need some alignments
- The soft alignment case
 - We use an attention-based neural network (this will be introduced later).
- The hard alignment cases
 - We explicitly introduce an alignment path z.
 - We can use it for CTC, RNN-transducer, and HMM-based approach