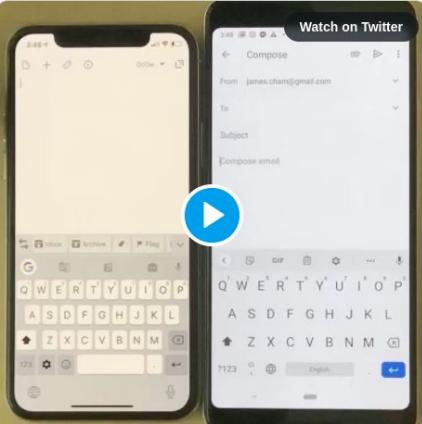


Automatic Speech Recognition (ASR)

Tal Rosenwein

 James Cham
@jamescham

I don't think that people appreciate how different the voice to text experience on a Pixel is from an iPhone. So here is a little head to head example. The Pixel is so responsive it feels like it is reading my mind!



8:19 AM · May 27, 2020

 [Read the full conversation on Twitter](#)

 17.5K  Reply  Copy link to Tweet

[Read 455 replies](#)

<https://twitter.com/i/status/1265512829806927873>

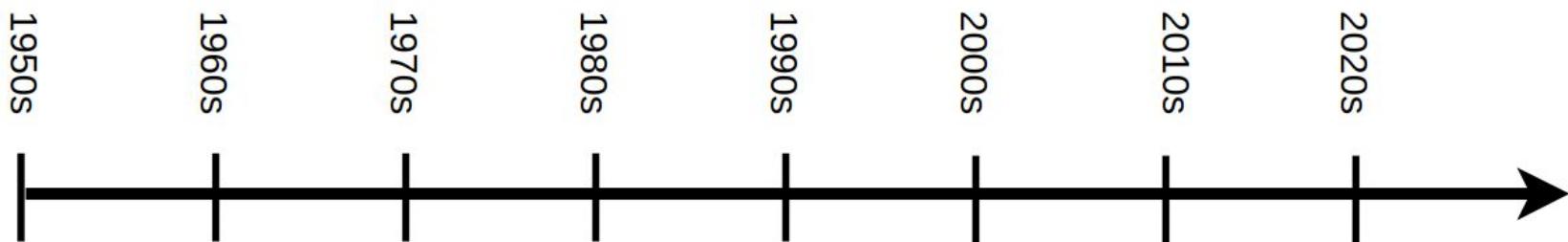
Outline

- History
- ASR as a system
- Datasets & Evaluation metrics
- ASR levels of difficulty
- Dynamic time warping (DTW)
- ASR problem formulation
- Inference pipeline
 - Decoding
 - Language models
 - Pronunciation models
 - Feature extraction
 - Acoustic models
 - HMM

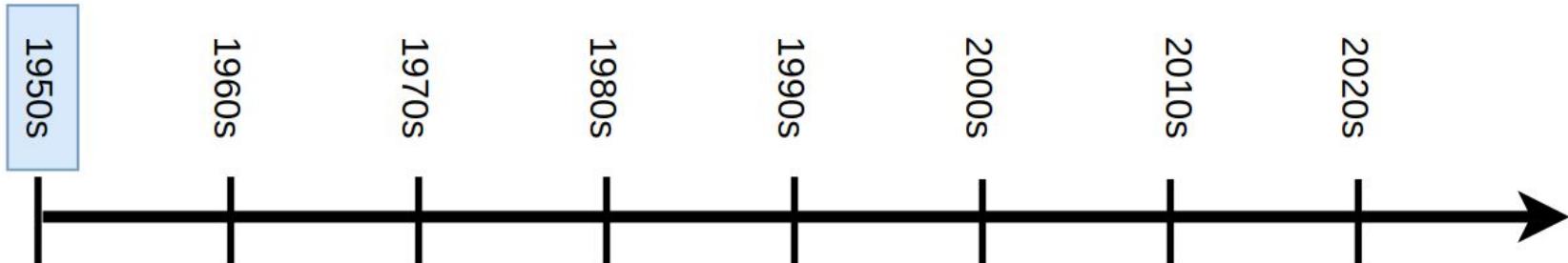
Outline

- **History**
- ASR as a system
- Datasets & Evaluation metrics
- ASR levels of difficulty
- Dynamic time warping (DTW)
- ASR problem formulation
- Inference pipeline
 - Decoding
 - Language models
 - Pronunciation models
 - Feature extraction
 - Acoustic models
 - HMM

History

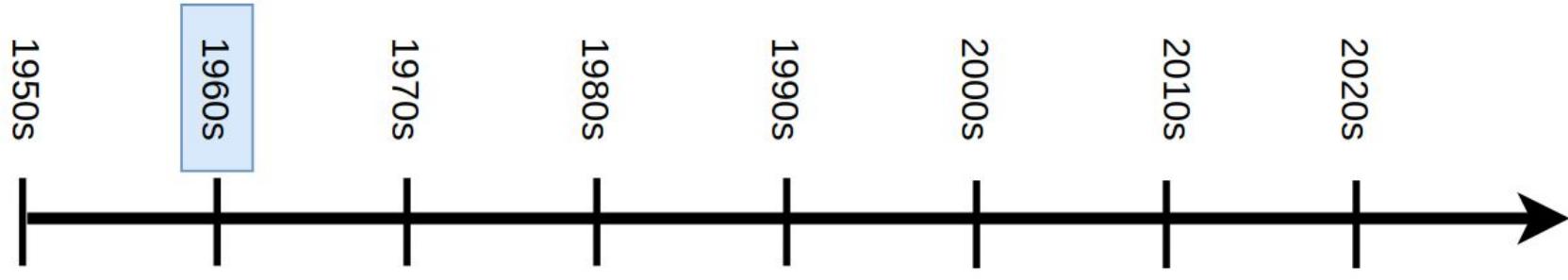


History



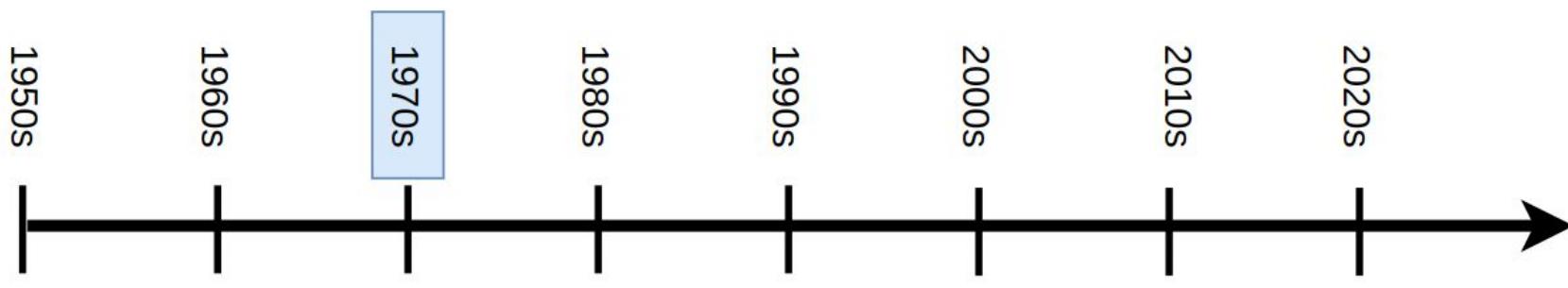
1952 Bell Labs Audrey. Not shown is six foot high rack of supporting electronics.

History



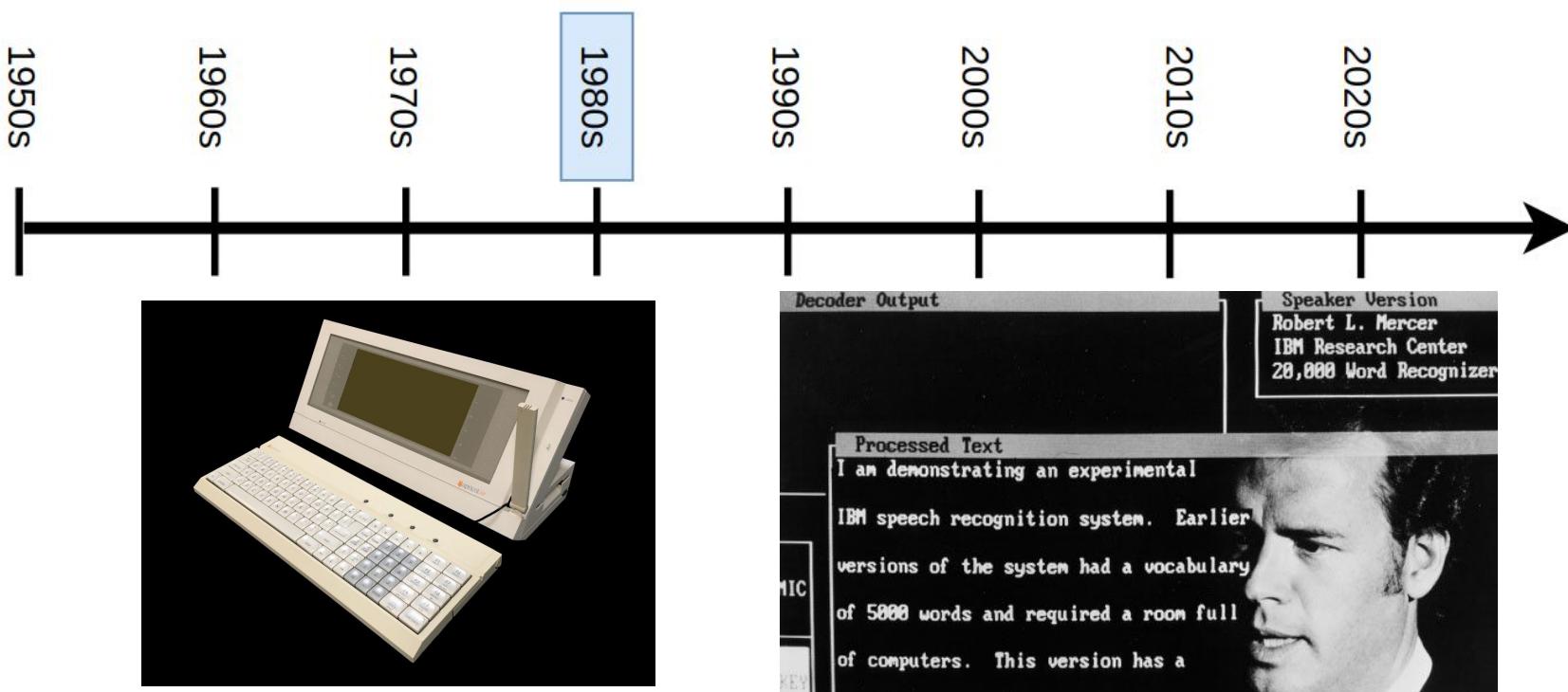
“Shoebox” by IBM

History



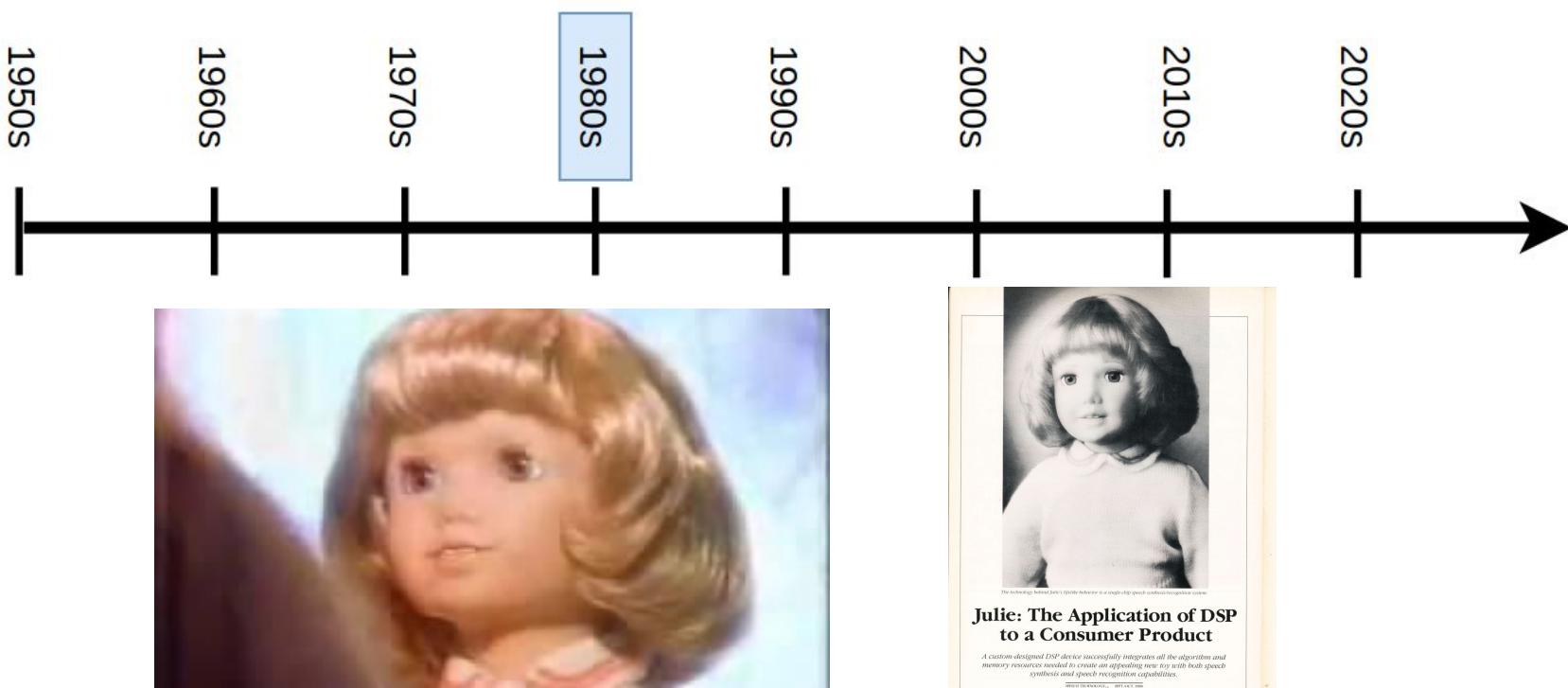
Growing vocabulary- DARPA funding

History

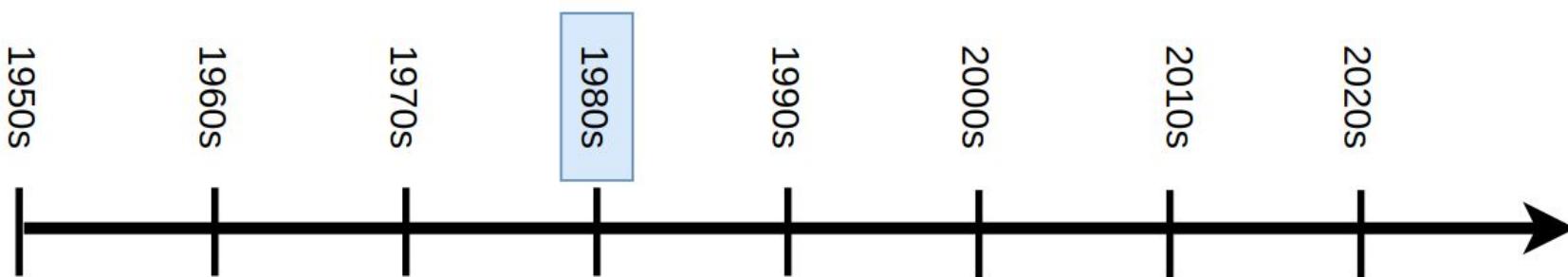


First speech recognition capability for a personal computer.

History

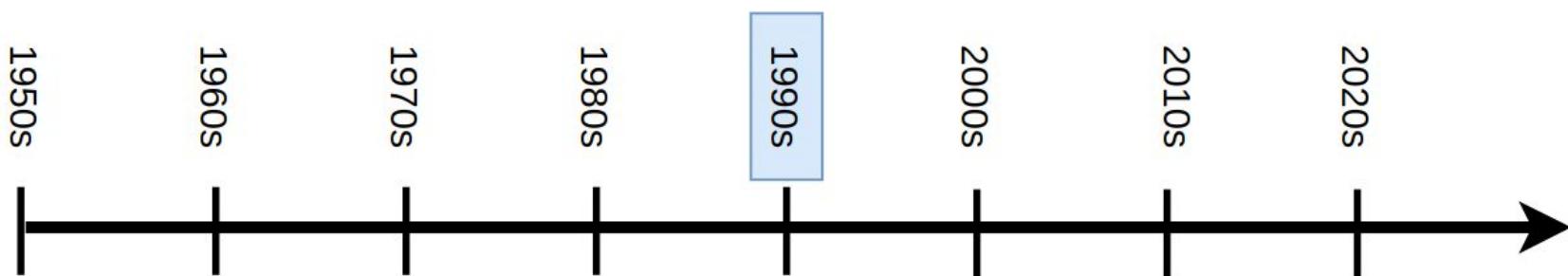


History



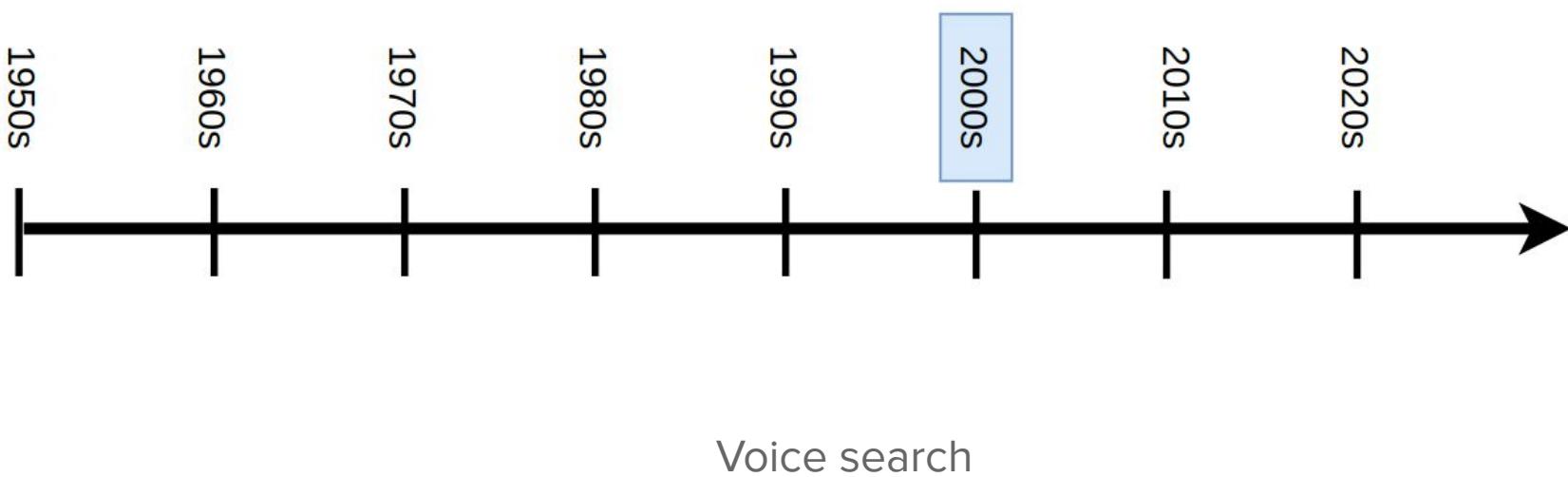
1989: Hidden Markov Model Toolkit ([HTK](#)) was introduced

History

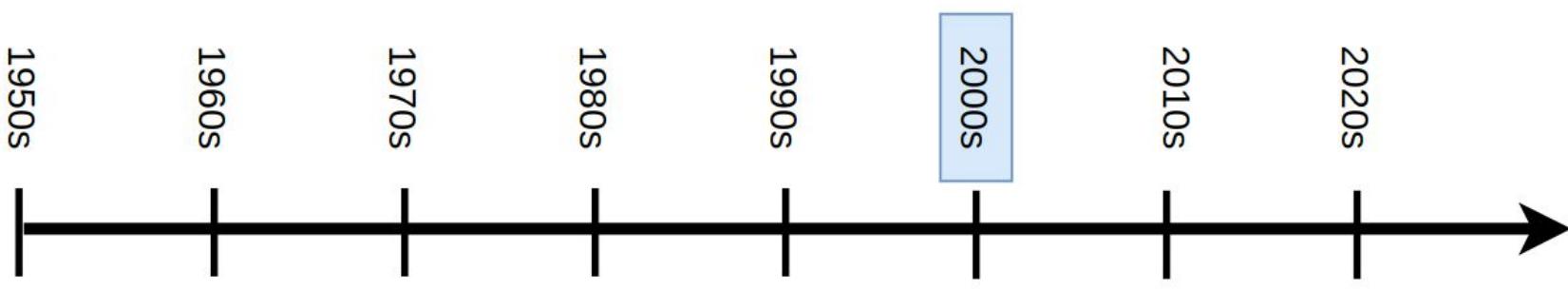


General purpose continuous speech product for consumer PCs

History

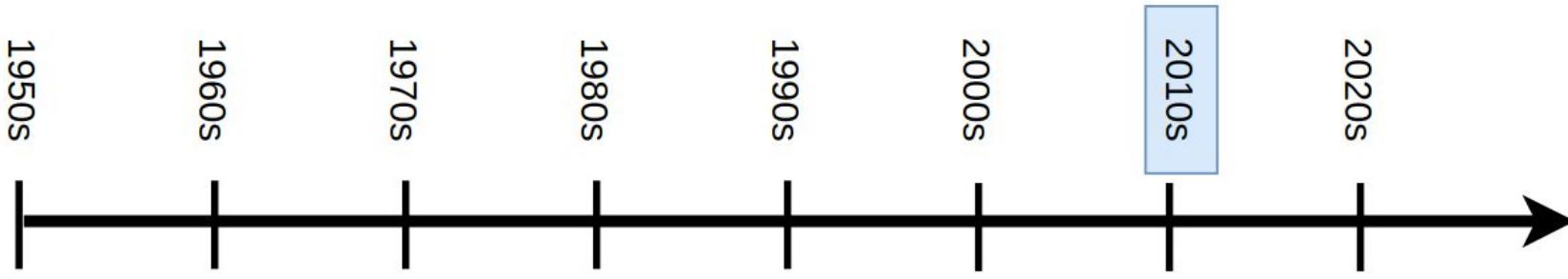


History



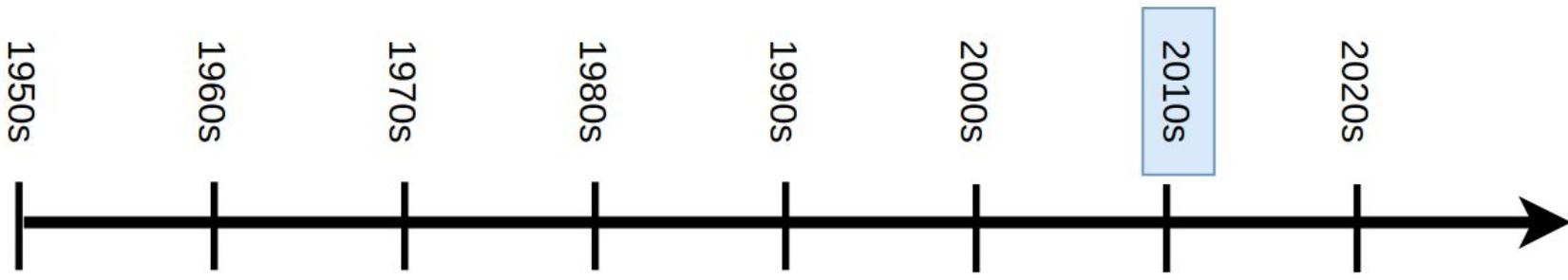
Kaldi (2009)

History



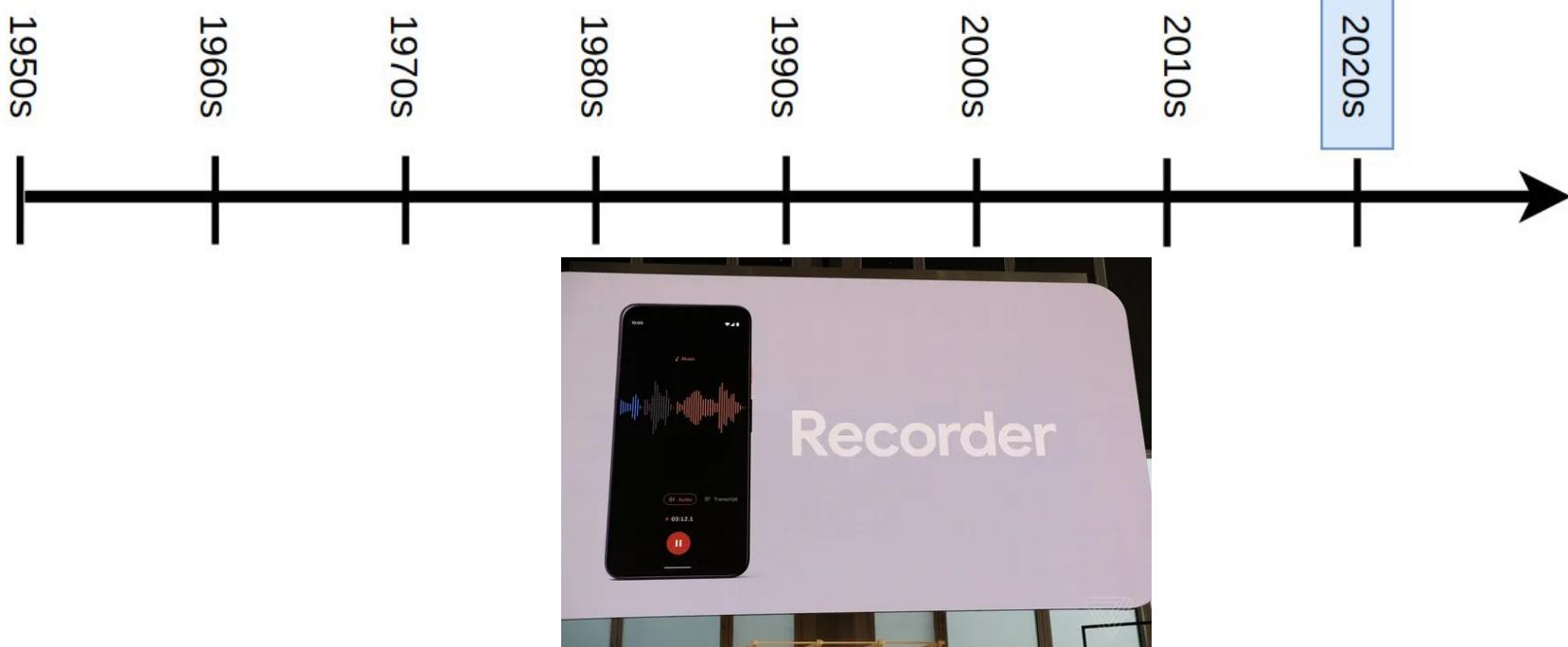
Voice enabled applications (PAs)

History



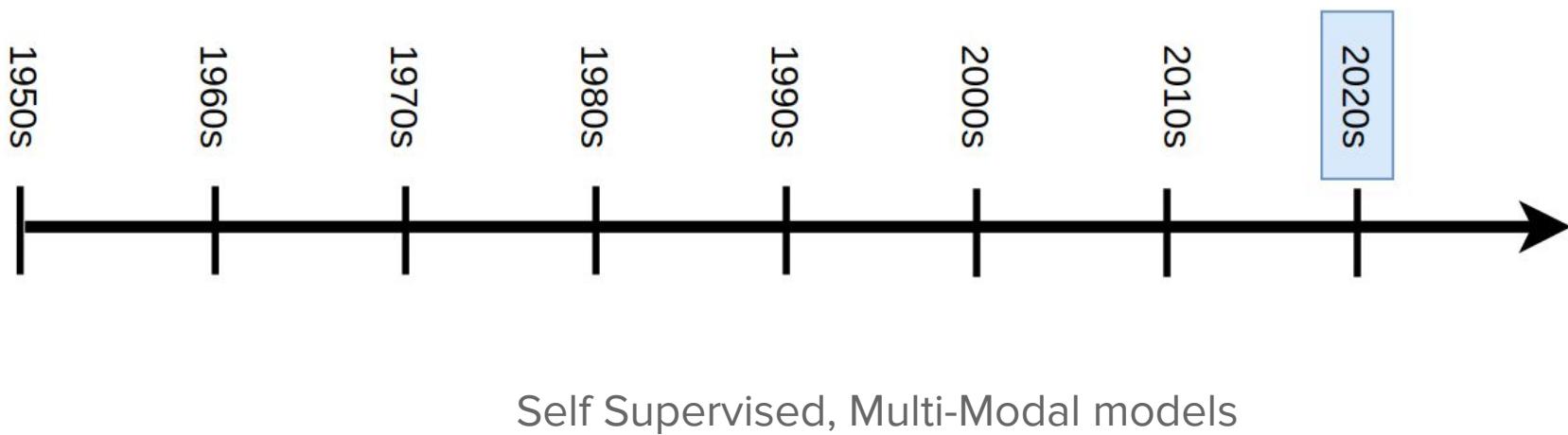
DNNs reached human level (2017)

History



On device embedded ASR

History



Outline

- History
- **ASR as a system**
- Datasets & Evaluation metrics
- ASR levels of difficulty
- Dynamic time warping (DTW)
- ASR problem formulation
- Inference pipeline
 - Decoding
 - Language models
 - Pronunciation models
 - Feature extraction
 - Acoustic models
 - HMM

ASR As a System

So we have the best ASR engine, but now what?

“you know there's something interesting about you i had to come meet you hmm huh what no there's something interesting about you i just wanted what to come meet you i don't know i just got that vibe that's weird is that okay no it's not oh oops”

ASR As a System

Making The Transcript More Intelligible

A: You now there's something interesting about, you I had to come meet you.

B: Hmm?

A: Huh?

B: What?

A: No, there's something interesting about you I just wanted

B: What?

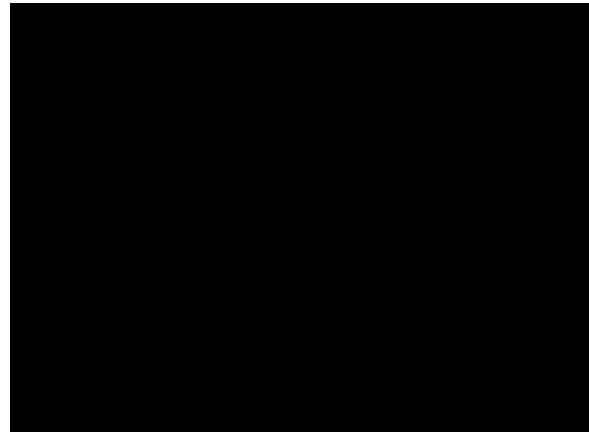
A: to come meet you. I don't know, I just got that vibe

B: That's weird.

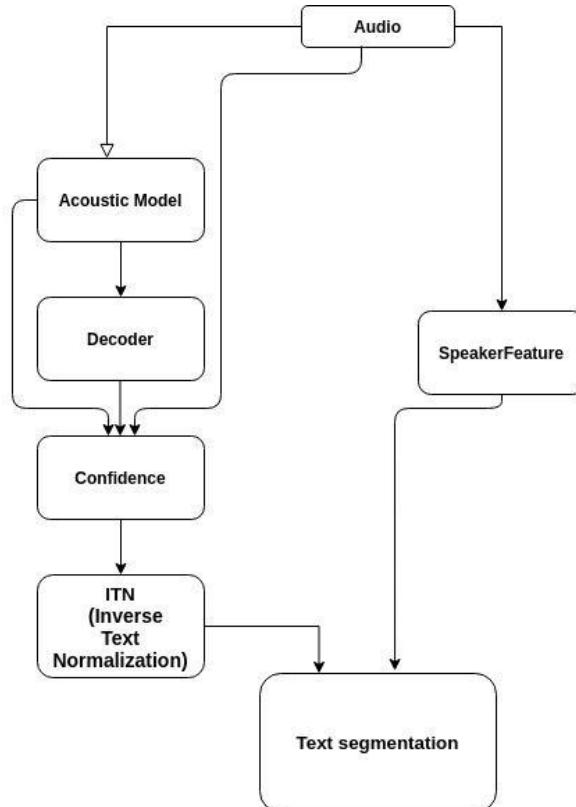
A: Is that okay?

B: No.

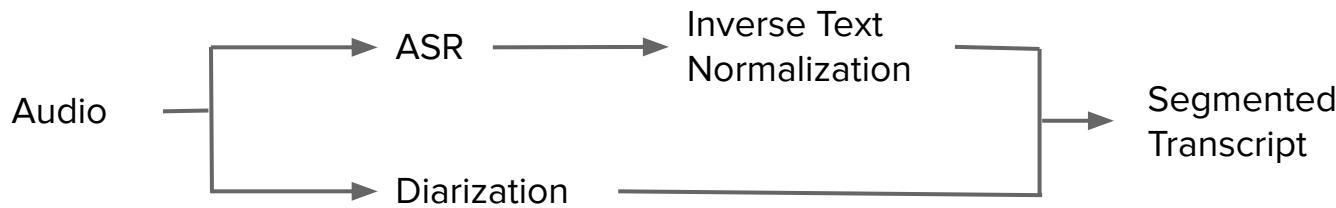
A: it's not. Oh, oops.



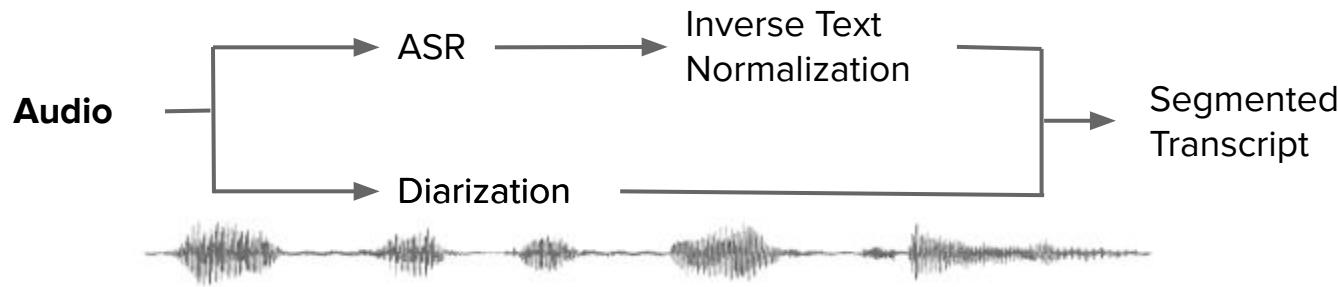
ASR As a System



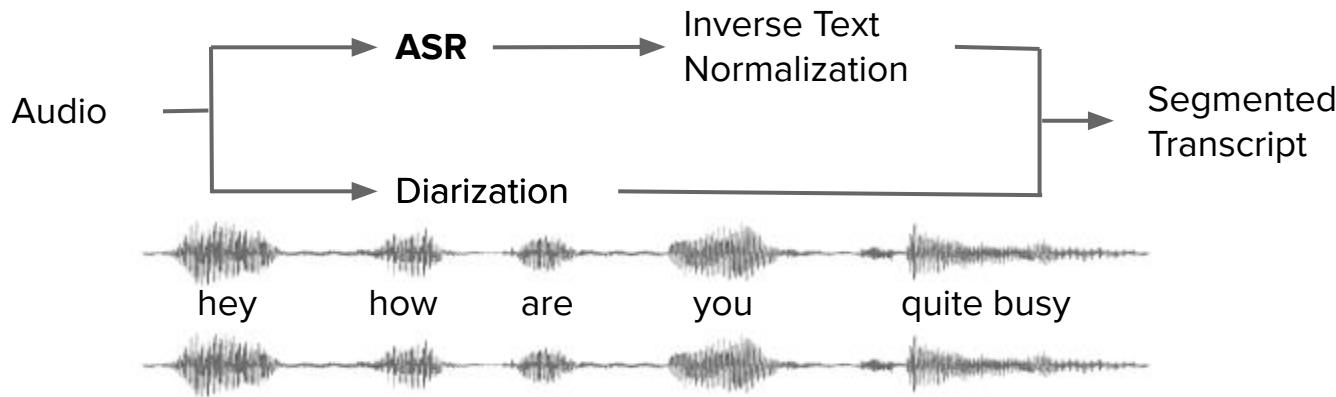
ASR As A System



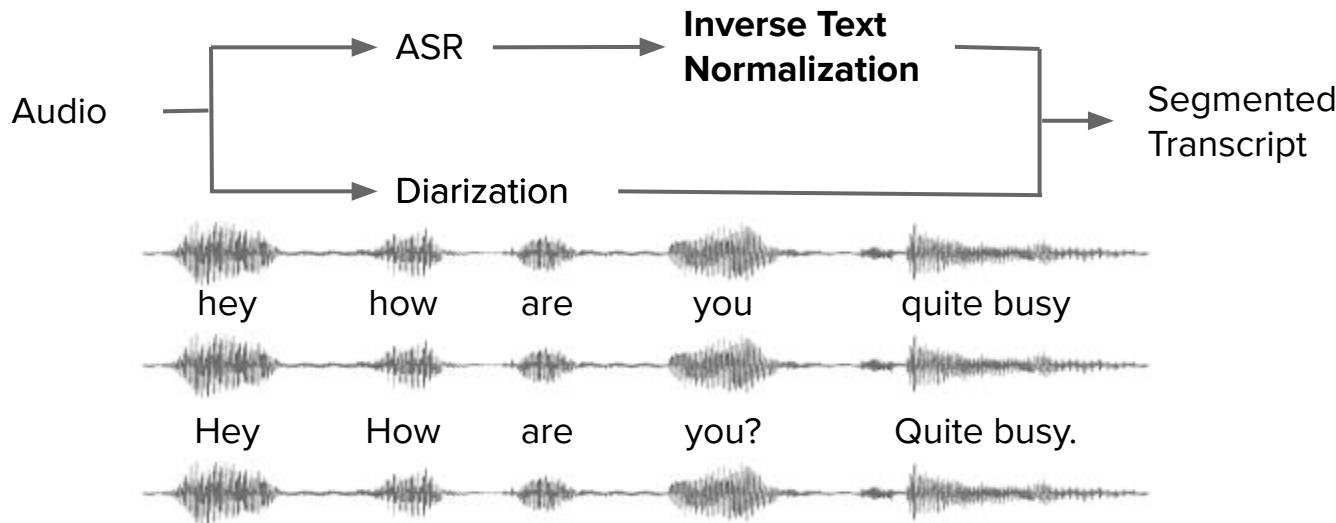
ASR As A System



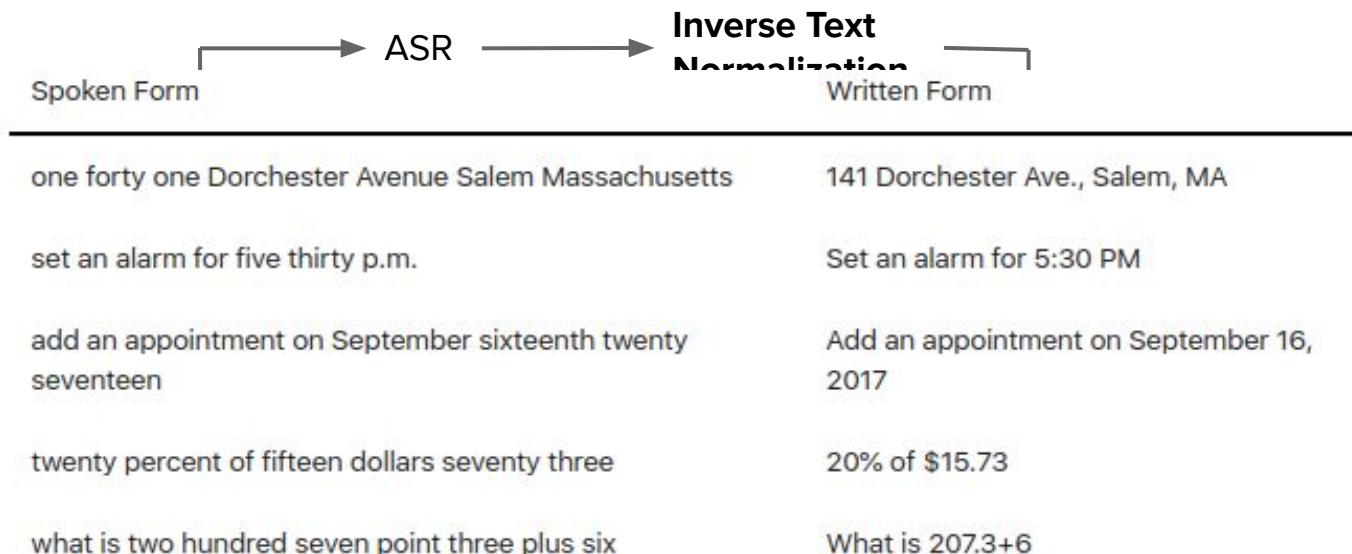
ASR As A System - Transcript



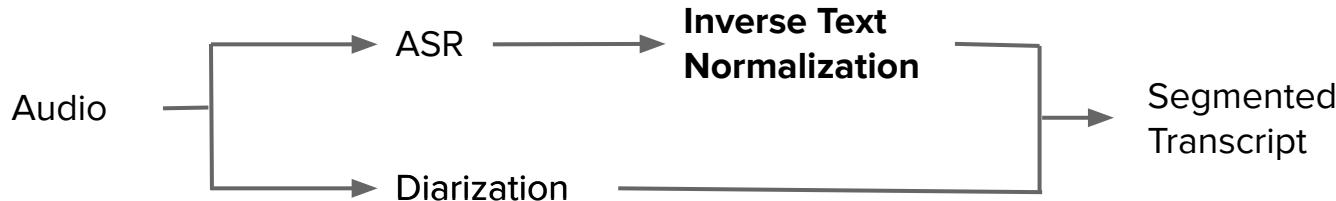
ASR As A System - ITN



ASR As A System - ITN



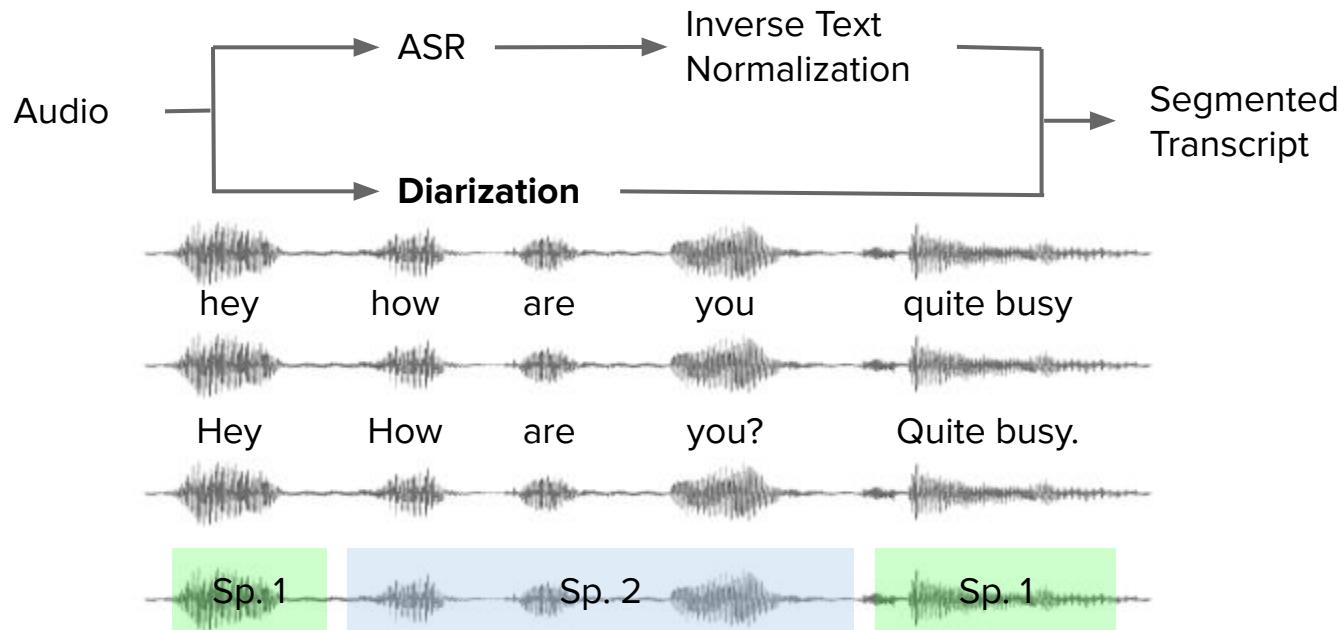
ASR As A System - ITN



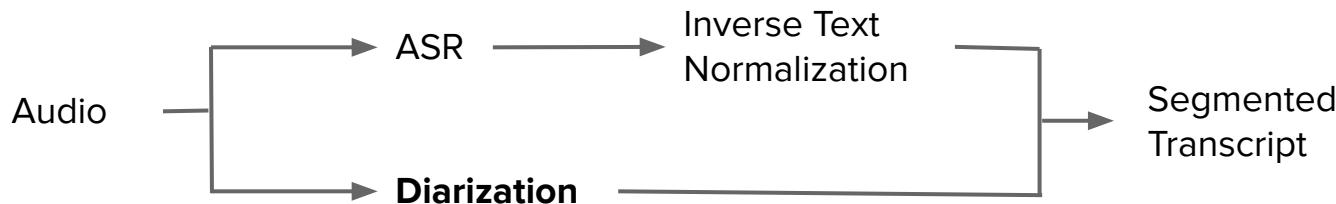
Why Inverse text normalization?

Because usually during the acoustic model training we **normalize** the spoken text into written form.

ASR As A System - Diarization



ASR As A System - Diarization



Diarization = Who spoke when?

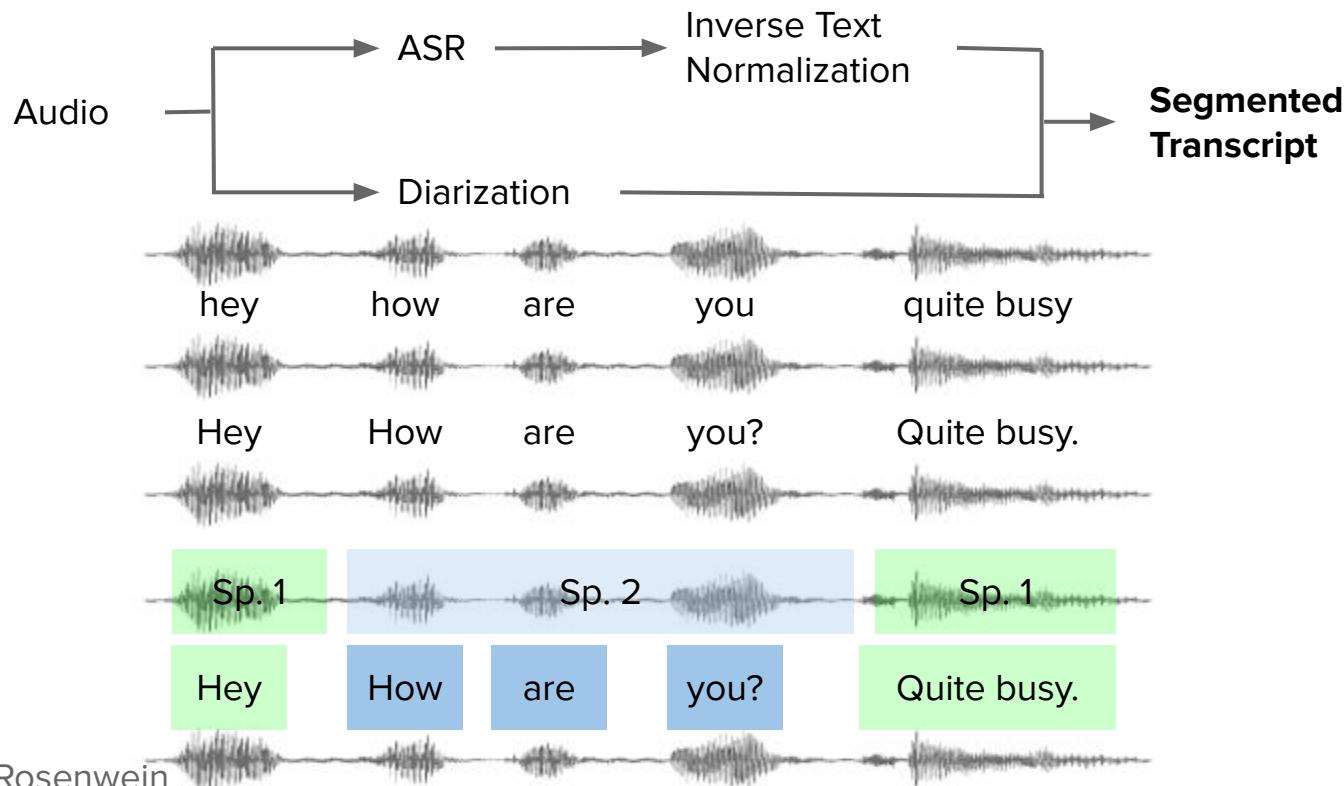
Hey How are you? Quite busy.

Sp. 1

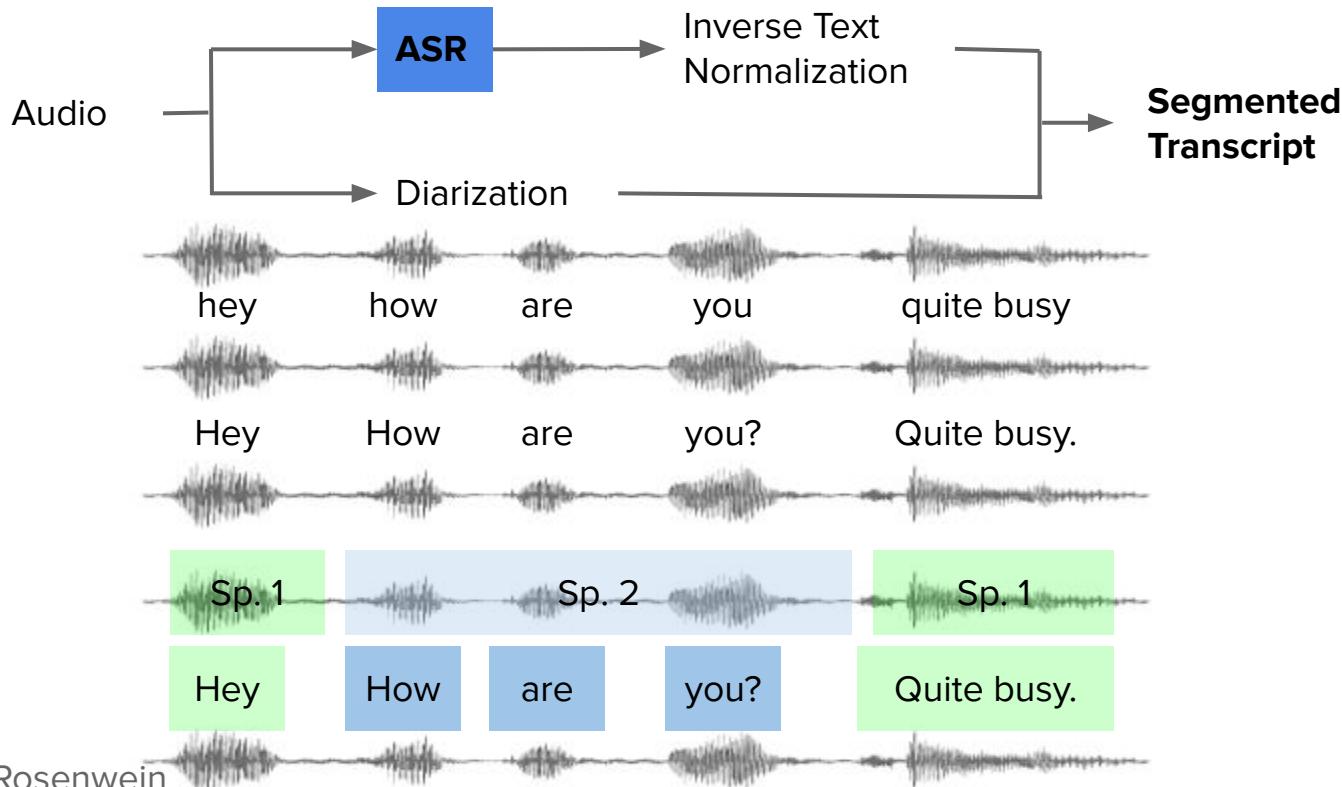
Sp. 2

Sp. 1

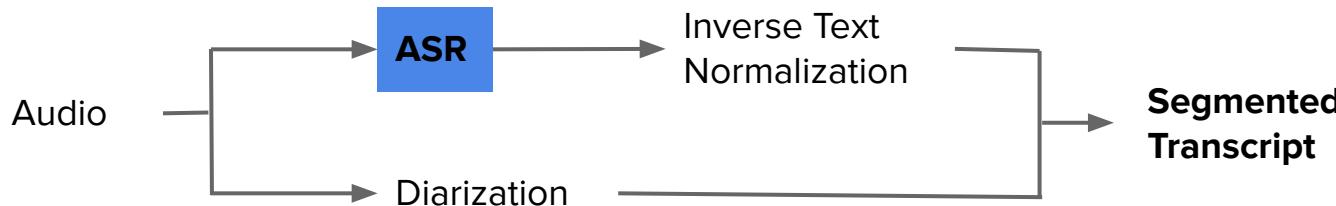
ASR As A System - Final Output



ASR As A System - Our Main Focus



ASR As A System - E2E Training of System



Joint Speech Recognition and Speaker Diarization via Sequence Transduction

Laurent El Shafey, Hagen Soltau, Izhak Shafran

Google

shafey@google.com, soltau@google.com, izhak@google.com

Abstract

Speech applications dealing with conversations require not only recognizing the spoken words, but also determining who spoke when. The task of assigning words to speakers is typically addressed by merging the outputs of two separate systems, namely, an automatic speech recognition (ASR) system and a speaker diarization (SD) system. The two systems are trained independently with different objective functions. Often the SD systems operate directly on the acoustics and are not constrained to respect word boundaries and this deficiency is overcome in an *ad hoc* manner. Motivated by recent advances in sequence to sequence learning, we propose a novel approach to tackle the two tasks by a joint ASR and SD system using a recurrent neural network transducer. Our approach utilizes both

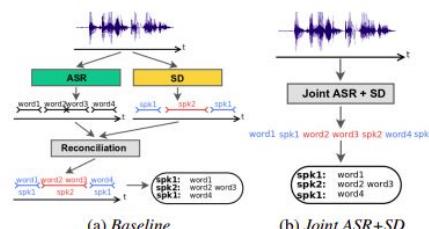


Figure 1: Comparison of the conventional speech recognition and speaker diarization system (Figure 1a) with the proposed approach (Figure 1b), where the task consists of generating a speaker-decorated transcript from raw audio.

hello dr jekyll <spk : pt> hello mr hyde what
brings you here today <spk : dr> I am struggling
again with my bipolar disorder <spk : pt>

Figure 2: Example of an output sequence for our joint ASR and SD RNN-T system. The corresponding input would be the raw audio signal. Speaker turns are displayed in different colors.

Outline

- History
- ASR as a system
- **Datasets & Evaluation metrics**
- ASR levels of difficulty
- Dynamic time warping (DTW)
- ASR problem formulation
- Inference pipeline
 - Decoding
 - Language models
 - Pronunciation models
 - Feature extraction
 - Acoustic models
 - HMM

Datasets

Datasets

8,964 machine learning datasets

Share your dataset with the ML community!

162 dataset results for **Speech**

Search for datasets

Best match

Filter by Modality (clear)

- Speech
- Images 2538
- Texts 2414
- Videos 820
- Audio 975
- Medical 318
- 3D 281

Filter by Task

- Speech Recognition 52
- Speech Synthesis 10
- Automatic Speech Recognition (ASR) 9
- Speech Enhancement 9
- Speech Separation 8
- Text-To-Speech Synthesis 8

Filter by Language

- English 75

LibriSpeech
The LibriSpeech corpus is a collection of approximately 1,000 hours of audiobooks that are a part of the LibriVox project. Most of the audiobooks come from the Project Gutenberg. The...
1,838 PAPERS • 11 BENCHMARKS

Speech Commands
Speech Commands is an audio dataset of spoken words designed to help train and evaluate keyword spotting systems.
317 PAPERS • 4 BENCHMARKS

LibriTTS
LibriTTS is a multi-speaker English corpus of approximately 585 hours of read English speech at 24kHz sampling rate, prepared by Heiga Zen with the assistance of Google Speech and...
168 PAPERS • 1 BENCHMARK

MUSAN
MUSAN is a corpus of music, speech and noise. This dataset is suitable for training models for voice activity detection (VAD) and music/speech discrimination. The dataset consists of mu...
159 PAPERS • NO BENCHMARKS YET

AISHELL-1
AISHELL-1 is a corpus for speech recognition research and building speech recognition systems for Mandarin.
155 PAPERS • 1 BENCHMARK

WSJ0-2mix
WSJ0-2mix is a speech recognition corpus of speech mixtures using utterances from the Wall Street Journal (WSJ0) corpus.
138 PAPERS • 2 BENCHMARKS

VoxPopuli
VoxPopuli is a large-scale multilingual corpus providing 100K hours of unlabelled speech data in 23 languages. It is the largest open data to date for unsupervised representation learnin...
74 PAPERS • 5 BENCHMARKS

LibriMix
LibriMix is an open-source alternative to wsj0-2mix. Based on LibriSpeech, LibriMix consists of

Leaderboard



The 🎉 Open ASR Leaderboard ranks and evaluates speech recognition models on the Hugging Face Hub.
We report the Average [WER](#) (⬆) and [RTE](#) (⬇) - lower the better. Models are ranked based on their Average WER, from lowest to highest. Check the 📈 Metrics tab to understand how the models are evaluated.
If you want results for a model that is not listed here, you can submit a request for it to be included 📝💡.
The leaderboard currently focuses on English speech recognition, and will be expanded to multilingual evaluation in later versions.

model	Average WER	RTF (1e-3)	AMI	Earnings22	Gigaspeech	LS Clean	LS Other	SPGISpeech	Tedli
openai/whisper-large-v3	7.7	10.3	16.01	11.3	10.02	2.03	3.91	2.95	3.9
nvidia/stt_en_fastconformer_transducer_xlarge	8.06	12.3	18.28	16.37	11.58	1.5	2.88	4.4	4.49
openai/whisper-large-v2	8.06	10.5	16.82	12.02	10.57	2.56	5.16	3.77	4.01
nvidia/stt_en_fastconformer_transducer_xxlarge	8.07	14.4	18.81	16.66	11.95	1.38	2.52	4.98	4.74
distil-whisper/distil-large-v2	8.31	4.93	14.65	12.12	10.31	2.95	6.39	3.28	4.3
nvidia/stt_en_fastconformer_ctc_xxlarge	8.34	5	17.62	16.44	11.61	1.69	3.4	4.91	4.64
nvidia/stt_en_conformer_ctc_large	8.39	7.5	15.97	15.83	11.59	2.06	4.16	5.6	4.41
openai/whisper-medium.en	8.5	10.7	16.43	12.59	11.13	3.02	5.84	3.41	4.15
nvidia/stt_en_fastconformer_ctc_xlarge	8.52	2.9	18.41	17.89	11.84	1.73	3.47	5.04	4.71
nvidia/stt_en_fastconformer_ctc_large	8.9	1.8	18.59	18.67	12.15	1.95	4.04	5.03	4.74
nvidia/stt_en_fastconformer_transducer_large	8.94	10.4	20.2	19.36	12.16	1.67	3.64	4.43	4.44
openai/whisper-large	9.2	10.5	17.9	15.77	11.84	3.04	6.01	3.99	4.69
nvidia/stt_en_conformer_transducer_large	9.27	21.8	22.27	19.91	12.5	1.64	3.51	4.96	5.15
distil-whisper/distil-medium.en	9.32	3.95	16.02	12.89	11.26	3.6	7.67	3.77	4.8

Evaluation Metrics - Edit Distance

Edit Distance (ED): sums the Insertion, Deletion & Substitution errors

E X P O N E N T I A L

P O L Y N O M I A L

Deleted Substituted Added

$$ED = \frac{Ins. + Subs. + Del.}{|GT|}$$

Character Error Rate (CER): ED Between Characters

Word Error Rate (WER): ED Between Words

CER < WER

Evaluation Metrics - Edit Distance

WER sometimes correlates poorly with user perception of transcription quality:

- Does not consider semantic correctness
- Weights all errors equally (stop words, Nouns...)- The words that are rare count more.
- Ignores multi-hypothesis (don't vs do not, going to vs gonna).

Evaluation Metrics- Semantic Distance

Table 2: Examples of Ref/Hyp that has the biggest gap between WER and SemDist.

Gap	WER	SemDist	Ref/Hyp
12011	16.67	0.00	Ref: hey portal play mister blue sky Hyp: hey portal play mr blue sky.
8742	50.00	0.00	Ref: I smell hot dogs Hyp: I smell hotdogs.
6390	6.67	0.01	Ref: keep away it is a nightmare thank God we are separated about four thousand kilometres Hyp: keep away, it is a nightmare. thank god we are separated about four thousand kilometers.
3807	20.00	0.02	Ref: keep time zones in mind for the next zoom call Hyp: keep timezones in mind for the next Zoom call.
2725	66.67	3.14	Ref: I'm eagerly waiting Hyp: I am eagerly waiting.

(a) top-5 (Rank_{WER} - Rank_{SemDist}) gap

Gap	WER	SemDist	Ref/Hyp
2486	6.25	302.73	Ref: of course in the first time but one by one I able to handle those complaints Hyp: of course, in the first time, birth one by one, I able to handle those complaints.
2388	10.00	473.47	Ref: okay then it's kind of a great weekly for him Hyp: okay, then it's kind of a great victory for him.
2345	10.00	426.67	Ref: do you know the most whom he used to admire Hyp: do you know the most home he used to admire?
2314	7.69	286.03	Ref: sure but uh I think you will have to be patient because usually it's something like once a year or a little bit more but not much Hyp: sure but uh I think you will have to be patient because you read something like once a year or a little bit more but not much.
2300	10.00	375.39	Ref: yeah I think I've seen that in the news before Hyp: yeah, I think I've seen that in the morning before

Outline

- History
- ASR as a system
- Datasets & Evaluation metrics
- **ASR levels of difficulty**
- Dynamic time warping (DTW)
- ASR problem formulation
- Inference pipeline
 - Decoding
 - Language models
 - Pronunciation models
 - Feature extraction
 - Acoustic models
 - HMM

ASR Levels Of Difficulty

Keywords
spotting

Closed Set
Voice
Commands

Open
Domain
Questions

'Sterile'
LVCSR

Real Life
LVCSR



Hey Siri



Say "Hey Siri, send a
message."



Image source [1](#) [2](#) [3](#) [4](#) [5](#)

ASR Levels Of Difficulty

Keywords spotting Closed Set Voice Open Domain ‘Sterile’ LVCSR Real Life LVCSR

Models that solve downstream tasks other than transcription, are usually trained/deployed because of system constraints

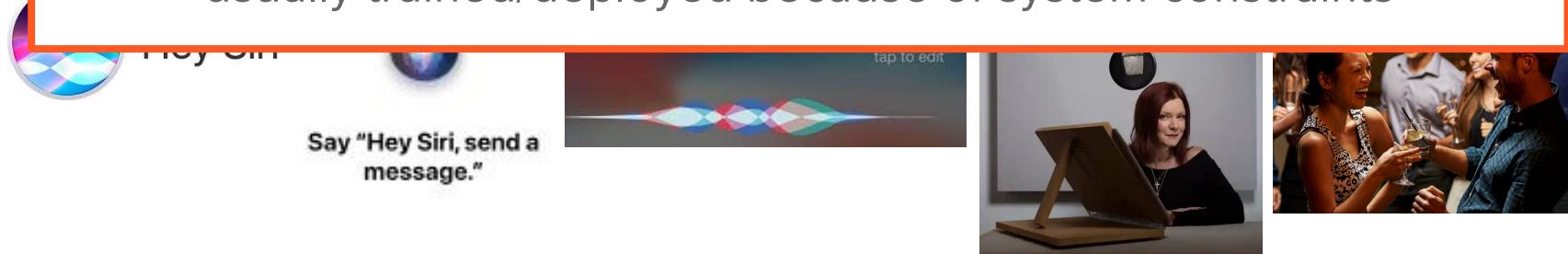


Image source [1](#) [2](#) [3](#) [4](#) [5](#)

Levels Of Difficulty: Why ASR Is Difficult?

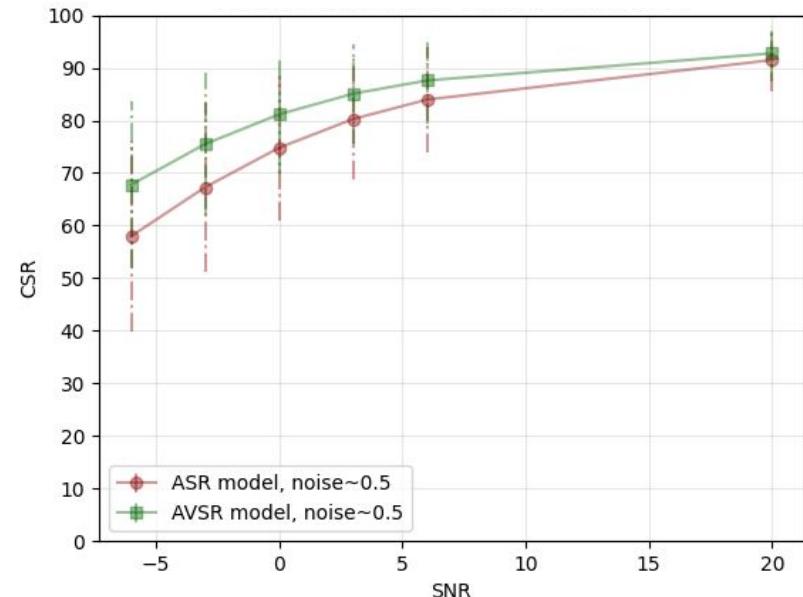
- Variety of Speech Attributes:
 - Accent
 - Rate
 - Gain
 - Gender

Levels Of Difficulty: Why ASR Is Difficult?

- **Variety of Speech Attributes:**
 - Accent
 - Rate
 - Gain
 - Gender
- **Real-Life Scenarios:**
 - Background Noise / Cocktail Party
 - Structured Vs Unstructured Speech

Levels Of Difficulty: Why ASR Is Difficult?

- Variety of Speech Attributes:
 - Accent
 - Rate
 - Gain
 - Gender
- Real-Life Scenarios:
 - Background Noise / Cocktail Party
 - Structured Vs Unstructured Speech



Levels Of Difficulty: Why ASR Is Difficult?

- **Variety of Speech Attributes:**
 - Accent
 - Rate
 - Gain
 - Gender
- **Real-Life Scenarios:**
 - Background Noise / Cocktail Party
 - Structured Vs Unstructured Speech
- **Audio Is Not Sufficient To Determine The Transcription:**
 - Large Context is Needed

Levels Of Difficulty: Why ASR Is Difficult?

- **Variety of Speech Attributes:**
 - Accent
 - Rate
 - Gain
 - Gender
- **Real-Life Scenarios:**
 - Background Noise / Cocktail Party
 - Structured Vs Unstructured Speech
- **Audio Is Not Sufficient To Determine The Transcription:**
 - Large Context is Needed
- **Non-Lexical Conversational Sounds (NLCS) are hard to recognize for ASR**

Levels Of Difficulty: Why ASR Is Difficult?

- **Variety**

JOURNAL ARTICLE

- “Mm-hm,” “Uh-uh”: are non-lexical conversational sounds deal breakers for the ambient clinical documentation technology?

- **Real-time**

[Get access >](#)

- Brian D Tran, Kareem Latif, Tera L Reynolds, Jihyun Park, Jennifer Elston Lafata, Ming Tai-Seale, Kai Zheng 

- **Audiobooks**

Journal of the American Medical Informatics Association, ocad001,
<https://doi.org/10.1093/jamia/ocad001>

- **Non-native speakers**

Published: 23 January 2023 [Article history ▾](#)

for ASR

Levels Of

- **Variety of**
 - Accer
 - Rate
 - Gain
 - Gende
- **Real-Life S**
 - Backg
 - Struct
- **Audio Is N**
 - Large
- **Non-Lexic**

Materials and Methods

In this study, we evaluated 2 contemporary ASR engines, Google Speech-to-Text Clinical Conversation (“Google ASR”), and Amazon Transcribe Medical (“Amazon ASR”), both of which have their language models specifically tailored to clinical conversations. The empirical data used were from 36 primary care encounters. We conducted a series of quantitative and qualitative analyses to examine the word error rate (WER) and the potential impact of misrecognized NLCS on the quality of clinical documentation.

Results

Out of a total of 135 647 spoken words contained in the evaluation data, 3284 (2.4%) were NLCS. Among these NLCS, 76 (0.06% of total words, 2.3% of all NLCS) were used to convey clinically relevant information. The overall WER, of all spoken words, was 11.8% for Google ASR and 12.8% for Amazon ASR. However, both ASR engines demonstrated poor

performance in recognizing NLCS: the WERs across frequently used NLCS were 40.8% (Google) and 57.2% (Amazon), respectively; and among the NLCS that conveyed clinically relevant information, 94.7% and 98.7%, respectively.

?

ize for ASR

Outline

- History
- ASR as a system
- Datasets & Evaluation metrics
- ASR levels of difficulty
- **Dynamic time warping (DTW)**
- ASR problem formulation
- Inference pipeline
 - Decoding
 - Language models
 - Pronunciation models
 - Feature extraction
 - Acoustic models
 - HMM

Dynamic Time Warping (DTW)

- Used for aligning two sequences
 - Two instances of the same / different speaker are not exactly the same.

Dynamic Time Warping (DTW)

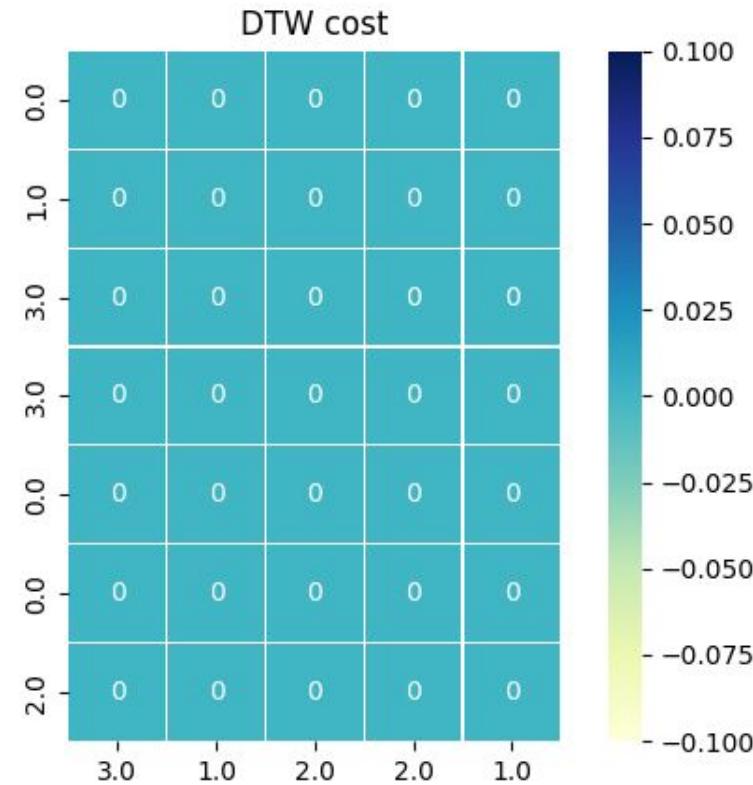
- Used for aligning two sequences
 - Two instances of the same / different speaker are not exactly the same.
- How do we compare 2 files/sequences with different lengths
 - Align (warp) the sequences
 - Calculate the distance

Dynamic Time Warping (DTW)

Given two sequences:

$$x = [3, 1, 2, 2, 1]$$

$$y = [2, 0, 0, 3, 3, 1, 0]$$

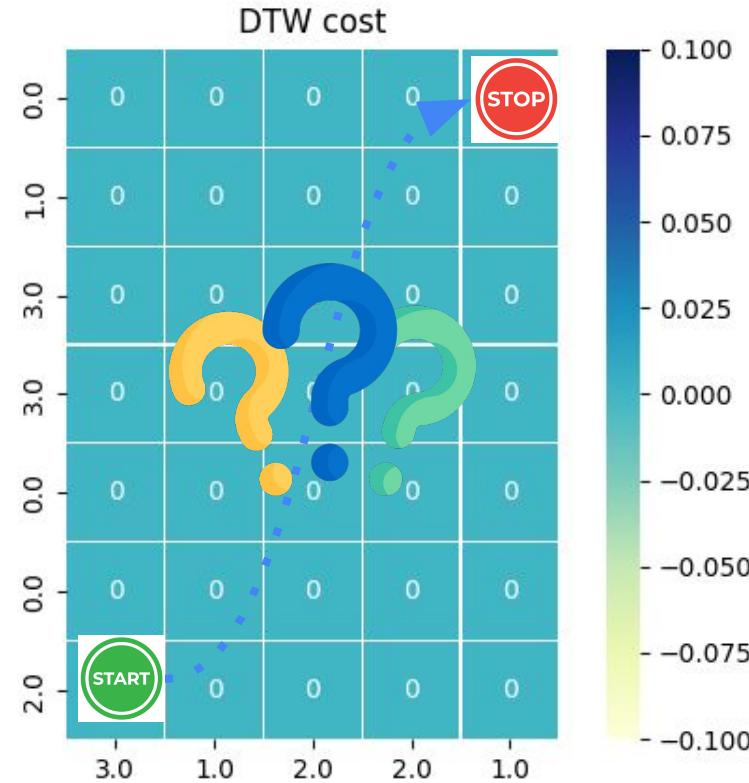


Dynamic Time Warping (DTW)

Given two sequences:

$$x = [3, 1, 2, 2, 1]$$

$$y = [2, 0, 0, 3, 3, 1, 0]$$

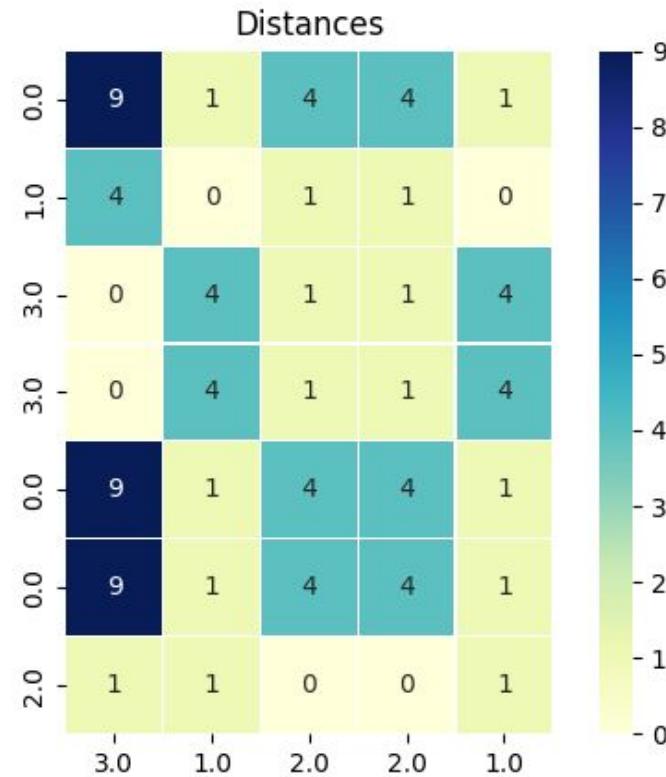


Dynamic Time Warping (DTW)

Given two sequences:

$$x = [3, 1, 2, 2, 1]$$

$$y = [2, 0, 0, 3, 3, 1, 0]$$



Dynamic Time Warping (DTW)

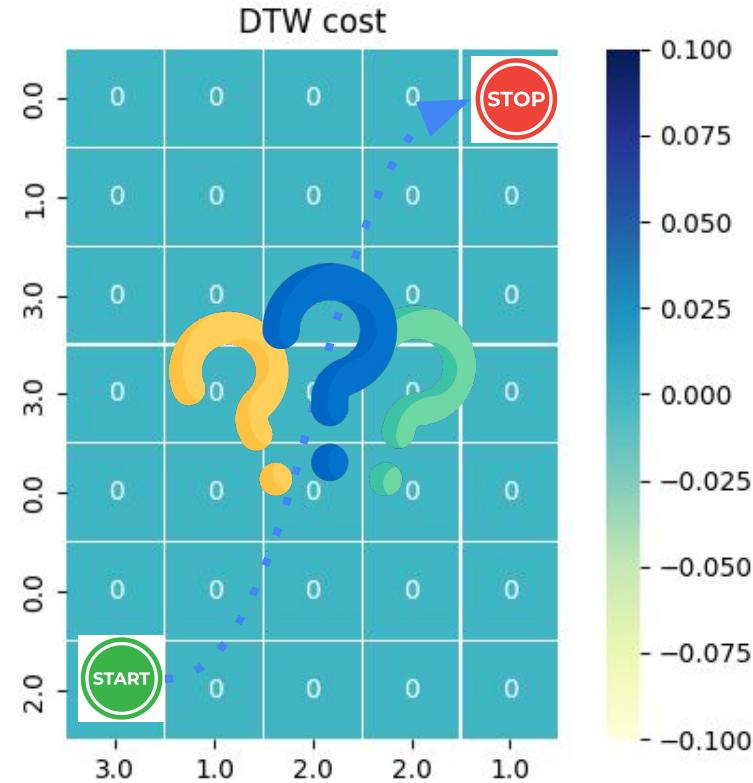
Given two sequences:

$$x = [3, 1, 2, 2, 1]$$

$$y = [2, 0, 0, 3, 3, 1, 0]$$

Constraints: Monotonic alignment

Using DP for efficiency



Dynamic Time Warping (DTW)

Given two sequences:

$$x = [3, 1, 2, 2, 1]$$

$$y = [2, 0, 0, 3, 3, 1, 0]$$

$$DTW[0, 0] = d(0, 0)$$

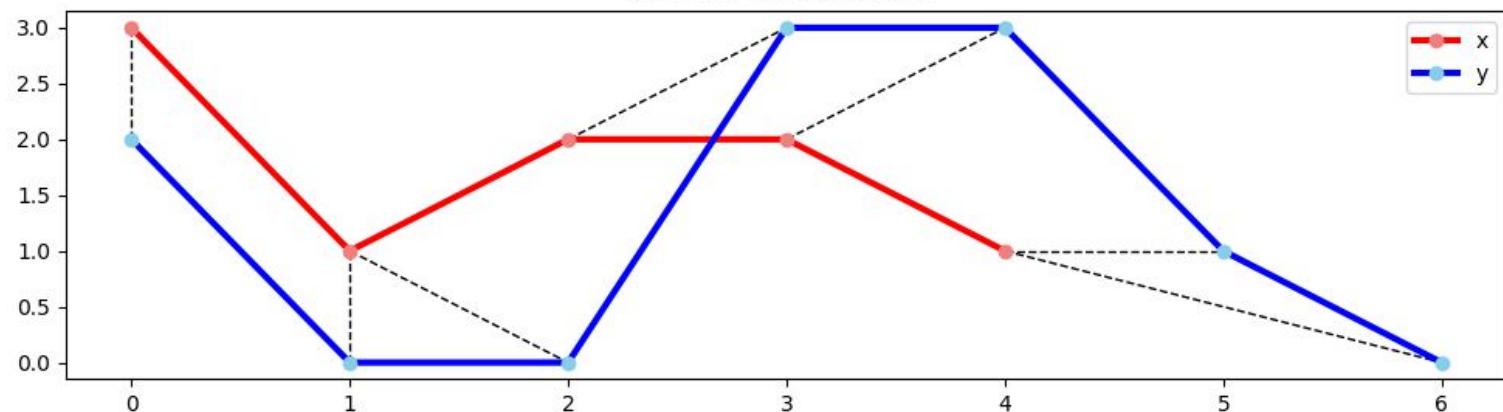
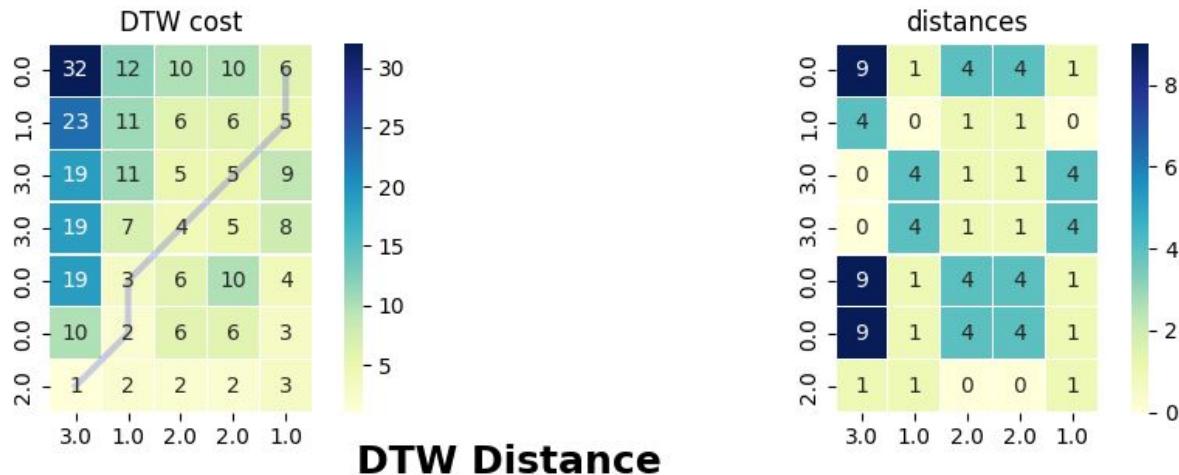
$$DTW[i, j] = d(i, j) + \min \begin{bmatrix} DTW[i - 1, j - 1] \\ DTW[i, j - 1] \\ DTW[i - 1, j] \end{bmatrix}$$

Dy

Gi

X :

y :

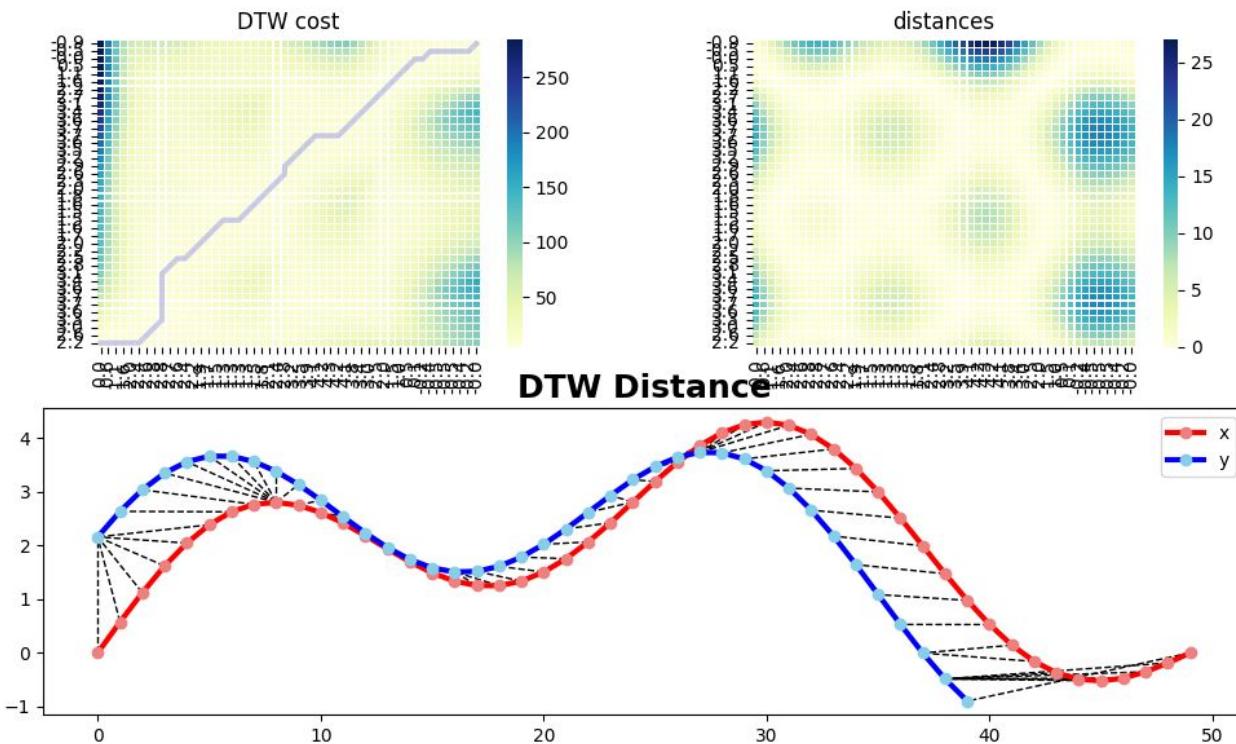


Dynamic Time Warping (DTW)

Given tw

$$x = [3, 1, 2]$$

$$y = [2, 0, 1]$$

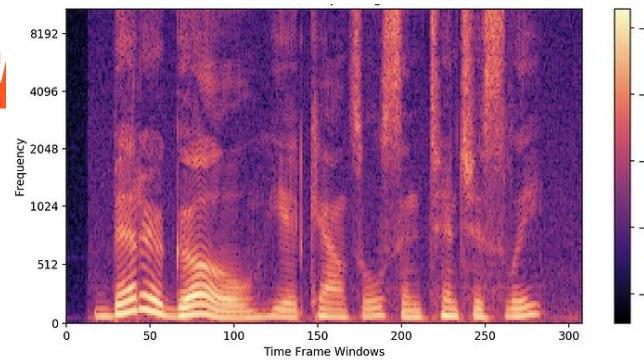


Dynamic Time Warping (DTW) - Inference

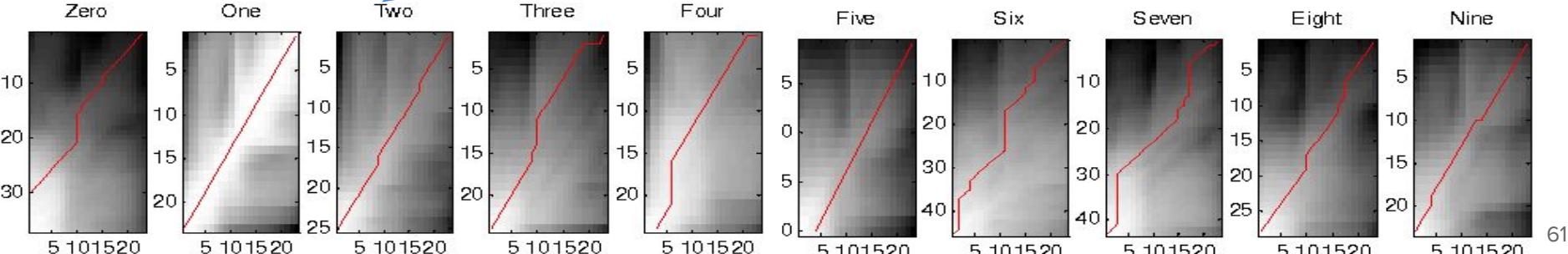


Dynamic Time V

[inferer] ?



$$class = \arg \min_{c \in Lexicon} DTW(Query, L_C)$$



Dynamic Time Warping (DTW) - Limitations

DTW has several limitations:

- How to represent a class?
- Classifying a whole word → Limited dictionary size
- Pre processing is required
- Exhaustive search
- Non parametric / non-probabilistic

Outline

- History
- ASR as a system
- Datasets & Evaluation metrics
- ASR levels of difficulty
- Dynamic time warping (DTW)
- **ASR problem formulation**
- Inference pipeline
 - Decoding
 - Language models
 - Pronunciation models
 - Feature extraction
 - Acoustic models
 - HMM

ASR Pipeline

DTW's limitations:

- Large multi-class
- Non probabilistic / learnable approach.

ASR pipeline consists of 3 main modules:

- **Acoustic model:** maps from audio signal to phonemes.
- **Language model:** computes to probability of a given word sequence.
- **Pronunciation model:** different ways of pronunciation a given word.

ASR Pipeline - Problem Formulation

What is the most likely word sequence \hat{W} out of all sentences in the language L given some acoustic input O ?

$$\widehat{\hat{W}} = \underset{W}{\operatorname{argmax}} P(W|O)$$

ASR Pipeline - Problem Formulation

What is the most likely word sequence \hat{W} out of all sentences in the language L given some acoustic input O ?

$$\widehat{W} = \underset{W}{\operatorname{argmax}} P(W|O)$$

Define:

1. Acoustic input O - a sequence of individual observations: $O = o_1, o_2, \dots, o_t$
2. A sentence as a sequence of words: $W = w_1, w_2, \dots, w_n$
3. Phoneme sequence $P = p_1, p_2, \dots, p_n$, for example:
 $W: w_1 = \text{banana}, P: p_1 = \text{buh}, p_2 = \text{na}, p_3 = \text{nuh},$
4. Aligned phoneme sequence $Q = q_1, q_2, \dots, q_T$, for example:
 $W: w_1 = \text{banana}, P: p_1 = \text{buh}, p_2 = \text{na}, p_3 = \text{nuh}, Q: q_1 = \text{buh}, q_2 = \text{buh}, q_3 = \text{na}, q_4 = \text{na}, q_5 = \text{nah}$

ASR Pipeline - Problem Formulation

What is the most likely word sequence \hat{W} out of all sentences in the language L given some acoustic input O ?

$$\widehat{W} = \underset{W}{\operatorname{argmax}} P(W|O)$$

$$\widehat{W} = \underset{W \in L}{\operatorname{argmax}} P(W|O) = \underset{W \in L}{\operatorname{argmax}} \frac{P(O|W)P(W)}{P(O)}$$

ASR Pipeline - Problem Formulation

What is the most likely word sequence \hat{W} out of all sentences in the language L given some acoustic input O ?

$$\widehat{W} = \underset{W}{\operatorname{argmax}} P(W|O)$$

$$\widehat{W} = \underset{W \in L}{\operatorname{argmax}} P(W|O) = \underset{W \in L}{\operatorname{argmax}} \frac{P(O|W)P(W)}{P(O)}$$

$$\widehat{W} = \underset{W \in L}{\operatorname{argmax}} P(O|W)P(W)$$

ASR Pipeline - Problem Formulation

What is the most likely word sequence \hat{W} out of all sentences in the language L given some acoustic input O ?

$$\widehat{W} = \underset{W \in L}{\operatorname{argmax}} P(O|W)P(W)$$

$$P(O|W) = \sum_P P(O|P, W)P(P|W)$$

ASR Pipeline - Problem Formulation

What is the most likely word sequence \hat{W} out of all sentences in the language L given some acoustic input O ?

$$\widehat{W} = \underset{W \in L}{\operatorname{argmax}} P(O|W)P(W)$$

$$P(O|W) = \sum_P P(O|P, W)P(P|W)$$

For example, The word ‘tomato’ can be pronounced in several ways ([Google search](#)):

- Tuh May tow (American English)
- Tuh Maa tow (British English)

ASR Pipeline - Problem Formulation

What is the most likely word sequence \hat{W} out of all sentences in the language L given some acoustic input O ?

$$\widehat{W} = \underset{W \in L}{\operatorname{argmax}} P(O|W)P(W)$$

$$P(O|W) = \sum_P P(O|P, W)P(P|W)$$
$$\approx \max_P P(O|P, W)P(P|W)$$

ASR Pipeline - Problem Formulation

What is the most likely word sequence \hat{W} out of all sentences in the language L given some acoustic input O ?

$$\widehat{W} = \underset{W \in L}{\operatorname{argmax}} P(O|W)P(W)$$

$$P(O|W) = \sum_P P(O|P, W)P(P|W)$$

$$\approx \underset{P}{\max} P(O|P, W)P(P|W)$$

$$\approx \underset{P}{\max} P(O|P)P(P|W)$$

ASR Pipeline - Problem Formulation

What is the most likely word sequence \hat{W} out of all sentences in the language L given some acoustic input O ?

$$\widehat{W} = \underset{W}{\operatorname{argmax}} P(W|O)$$

$$\widehat{W} = \underset{W \in L}{\operatorname{argmax}} P(O|W)P(W)$$

$$\widehat{W} = \arg \max_{W} \max_{P} P(O|P)P(P|W)P(W)$$

ASR Pipeline - Problem Formulation

What is the most likely word sequence \hat{W} out of all sentences in the language L given some acoustic input O ?

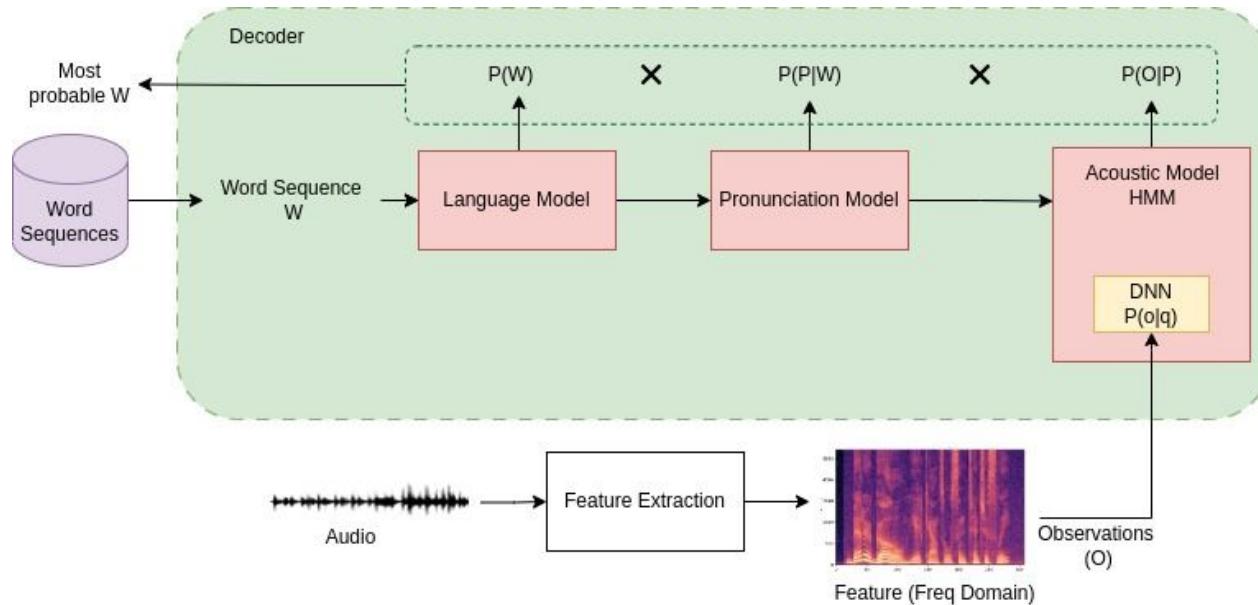
- **P(O|P)**: Acoustic model- the probability that this audio was generated by this phoneme
- **P(P|W)**: Pronunciation model- the probability that this phoneme sequence generates the word sequence.
- **P(W)**: Language model- the probability that this word sequence to be generated.

$$\widehat{W} = \arg \max_W \max_P P(O|P)P(P|W)P(W)$$

Outline

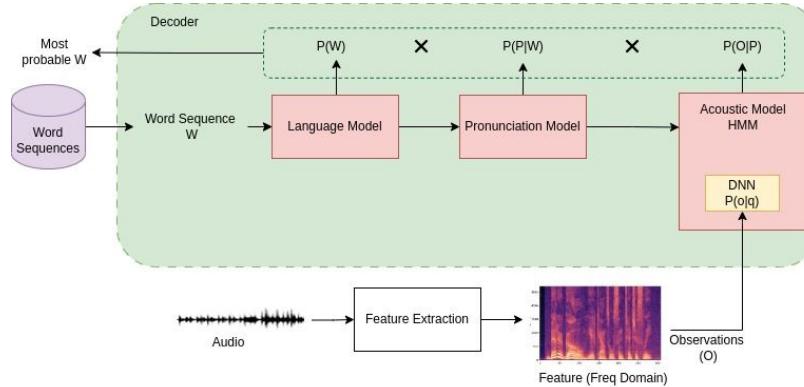
- History
- ASR as a system
- Datasets & Evaluation metrics
- ASR levels of difficulty
- Dynamic time warping (DTW)
- ASR problem formulation
- **Inference pipeline**
 - Decoding
 - Language models
 - Pronunciation models
 - Feature extraction
 - Acoustic models
 - HMM

ASR Pipeline - Inference Pipeline



$$\widehat{W} = \arg \max_W \max_P P(O|P)P(P|W)P(W)$$

ASR Pipeline - Inference Pipeline



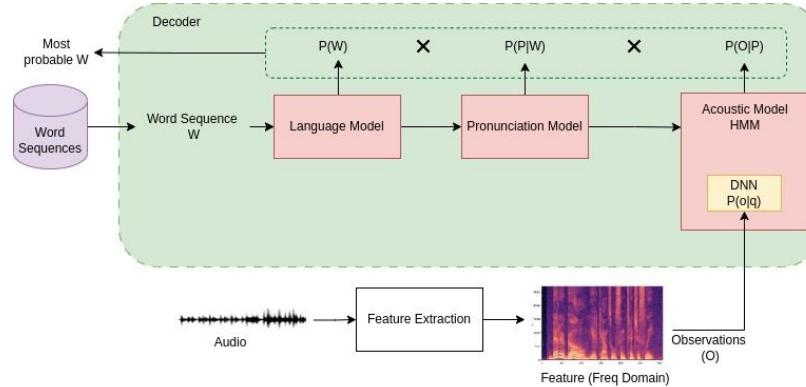
HMM-GMM has 5 components:

1. **Feature extraction:** from audio to spectral features
2. **Acoustic models:**
 - a. GMM for computing $p(olq)$, where q is the HMM hidden state.
 - b. HMM that wraps the GMM and outputs the $p(olp)$
3. **Lexicon/Pronunciation Model:** prob. for phoneme seq given the word seq $p(plw)$.
4. **Language Model:** captures semantics. Usually using N-grams
5. **Decoder:** Find the most probable word sequence using DP / search algorithm

Outline

- History
- ASR as a system
- Datasets & Evaluation metrics
- ASR levels of difficulty
- Dynamic time warping (DTW)
- ASR problem formulation
- Inference pipeline
 - **Decoding**
 - Language models
 - Pronunciation models
 - Feature extraction
 - Acoustic models
 - HMM

ASR Pipeline - Decoding

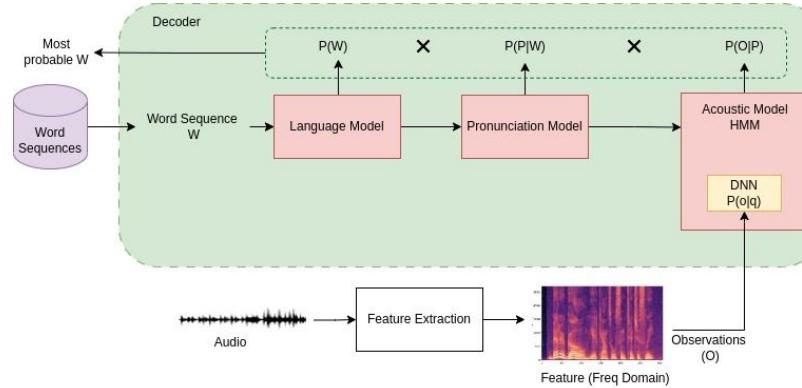


Searching over all possible word sequences is too exhaustive
Solution: Search over paths with signal from the acoustics

Outline

- History
- ASR as a system
- Datasets & Evaluation metrics
- ASR levels of difficulty
- Dynamic time warping (DTW)
- ASR problem formulation
- Inference pipeline
 - Decoding
 - **Language models**
 - Pronunciation models
 - Feature extraction
 - Acoustic models
 - HMM

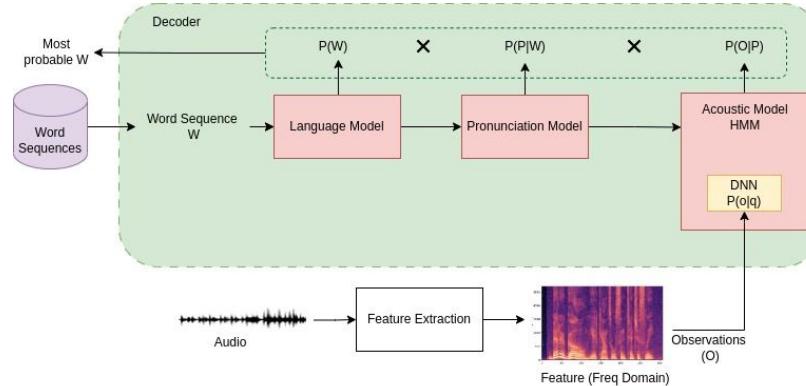
ASR Pipeline - Language Model



Language models are extremely important during decoding as they:

- Enlarge the acoustic model's context and add semantics into decoding.
- Enforce transcript: Speech Commands- “Hey Siri Call X”
- Recover from acoustic model errors.

ASR Pipeline - Language Model



Language Model **injects A-Prior information** into the decoding process, where it **enforces (by some extent)** the **output sequence to make sense**

$$P(w_1, \dots, w_n) = \prod_{i=1}^m P(w_i | w_1, \dots, w_{i-1})$$

Causal

ASR Pipeline - Language Model

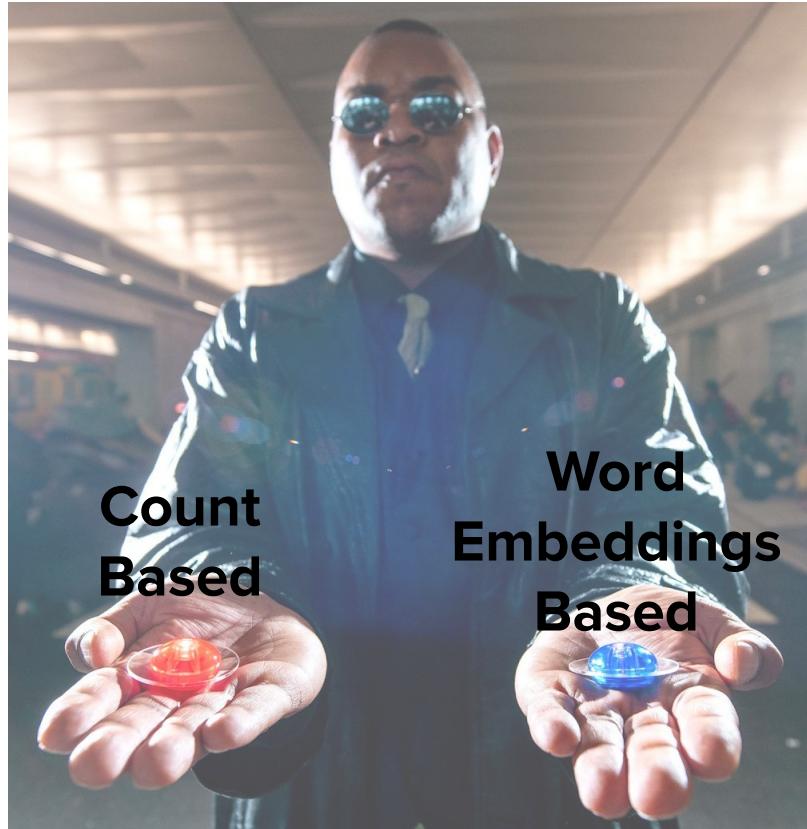


Image [source](#)

83

ASR Pipeline - Language Model - Count Based

“What You See Is What You Get”

All Out Of Vocabulary (OOV) Have The Same Prob.

$$P(\text{"bla bla fds"}) = P(\text{"gotten into his"}) = P(\text{"got iphone X"})$$

Lookup Tables (GBs)

Main tools:

- KenLM: <https://github.com/kpu/kenlm>
- SRILM: <http://www.speech.sri.com/projects/srilm/>

-1.5299898	looking into his
-1.6576592	get into his
-1.5772834	way into his
-1.5478122	them into his
-1.4195554	got into his
-0.5696159	year into his
-0.49294835	reached into his
-1.7273724	went into his
-1.0968437	well into his
-1.4932973	back into his

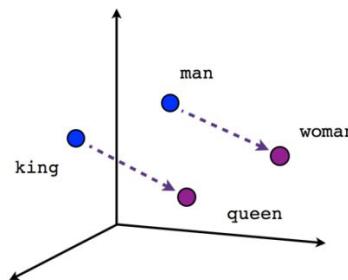
ASR Pipeline - Language Model - Embed. Based



ASR Pipeline - Language Model - Embed. Based

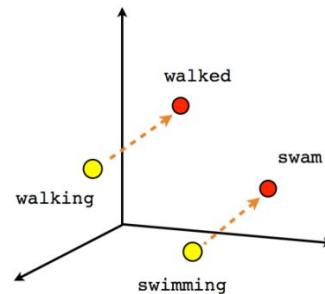
Word2Vec: Discrete words into continuous (embedding) space.

Can improves performances on OOV words (generalization).

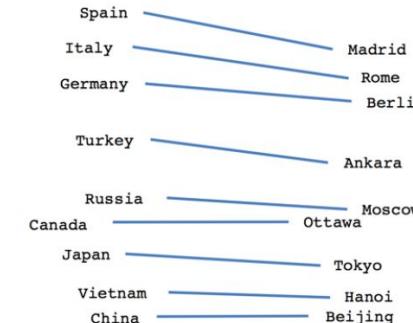


Male-Female

$$\text{king} - \text{man} + \text{woman} \approx \text{queen}$$



Verb tense



Country-Capital

Image source [1.2](#)

ASR Pipeline - Language Model - Embed. Based

Contextualized Word Embedding: Words in different contexts / Homonyms (words that have multiple meanings) should have different representations based on their context:

- I eat a **date** every day
- Let's go on a **date** sometime

Pioneering works:

- **ELMo:** 2018 Deep contextualized word representations.
- **BERT:** 2019 BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding

ASR Pipeline - Language Model - Embed. Based

Attention is the mechanism that takes into account the context of a given word, but doesn't learn a specific weight for the exact sequence.

It is affected by ‘the animal’.

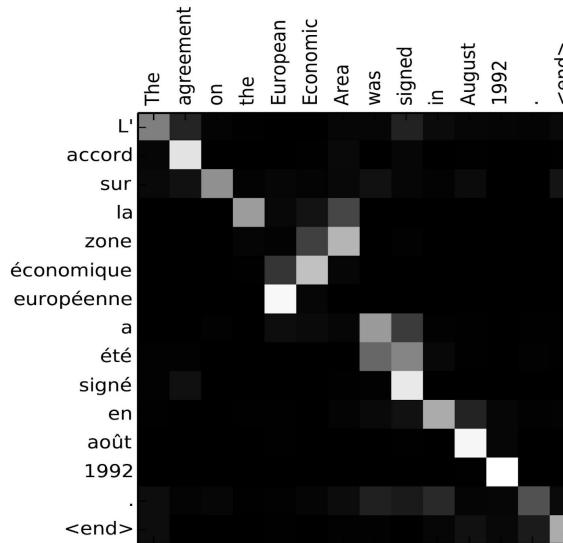
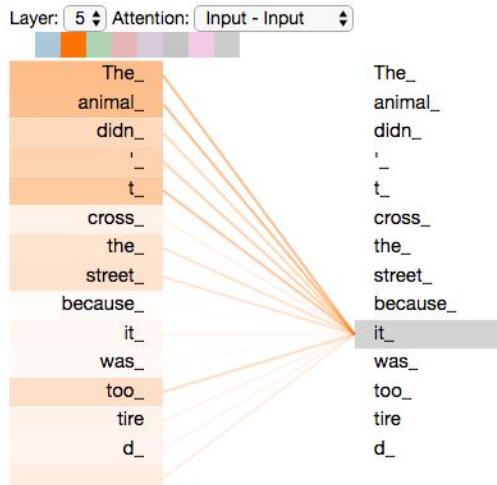


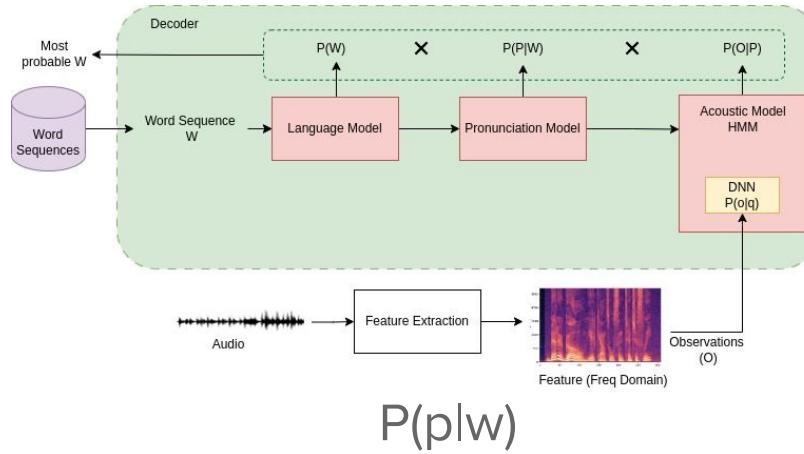
Image [source](#)

Outline

- History
- ASR as a system
- Datasets & Evaluation metrics
- ASR levels of difficulty
- Dynamic time warping (DTW)
- ASR problem formulation
- Inference pipeline
 - Decoding
 - Language models
 - **Pronunciation models**
 - Feature extraction
 - Acoustic models
 - HMM

ASR Pipeline - Lexicon/Pronunciation Model

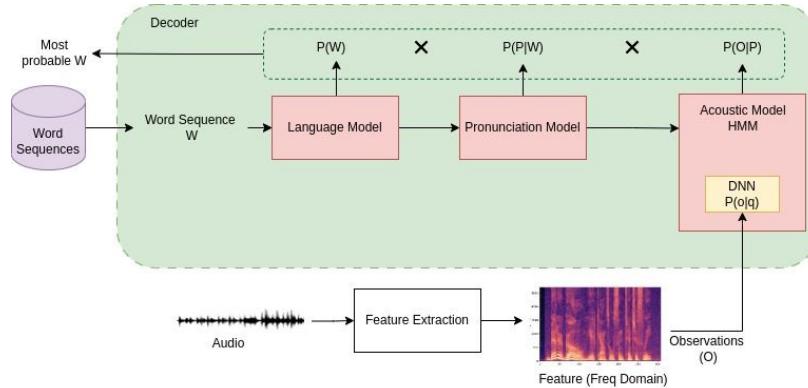
ADVANCING	AH D V AE N S IH NG
ADVANI	AE D V AA N IY
ADVANTA	AE D V AE N T AH
ADVANTA'S	AE D V AE N T AH Z
ADVANTA'S(2)	AH D V AE N T AH Z
ADVANTA(2)	AH D V AE N T AH
ADVANTAGE	AE D V AE N T IH JH
ADVANTAGE(2)	AH D V AE N T IH JH
ADVANTAGE(3)	AE D V AE N IH JH
ADVANTAGE(4)	AH D V AE N AH JH
ADVANTAGED	AE D V AE N T IH JH D
ADVANTAGED(2)	AH D V AE N T IH JH D
ADVANTAGED(3)	AE D V AE N IH JH D
ADVANTAGED(4)	AH D V AE N IH JH D
ADVANTAGEOUS	AE D V AH N T EY JH AH S
ADVANTAGES	AE D V AE N T IH JH IH Z
ADVANTAGES(2)	AH D V AE N T IH JH IH Z
ADVANTAGES(3)	AE D V AE N IH JH IH Z
ADVANTAGES(4)	AH D V AE N IH JH IH Z
ADVANTEST	AE D V AE N T AH S T
ADVANTEST(2)	AH D V AE N T AH S T
ADVECTION	AE D V EH K SH AH N



Outline

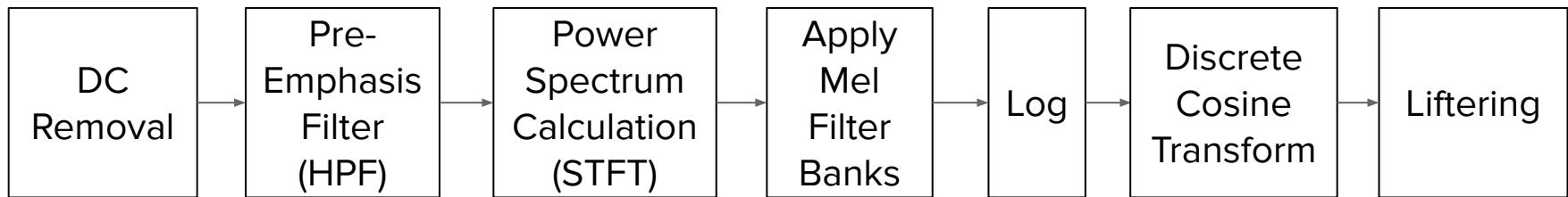
- History
- ASR as a system
- Datasets & Evaluation metrics
- ASR levels of difficulty
- Dynamic time warping (DTW)
- ASR problem formulation
- Inference pipeline
 - Decoding
 - Language models
 - Pronunciation models
 - **Feature extraction**
 - Acoustic models
 - HMM

ASR Pipeline - Feature Extraction



MFCCs / Mel Filter Banks, Raw audio

MFCCs - Mel Frequency Cepstral Coefficients

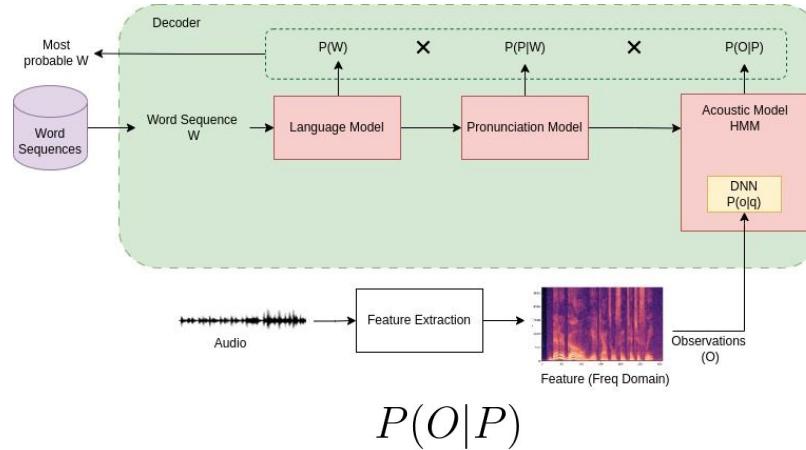


Source [1](#), [2](#)

Outline

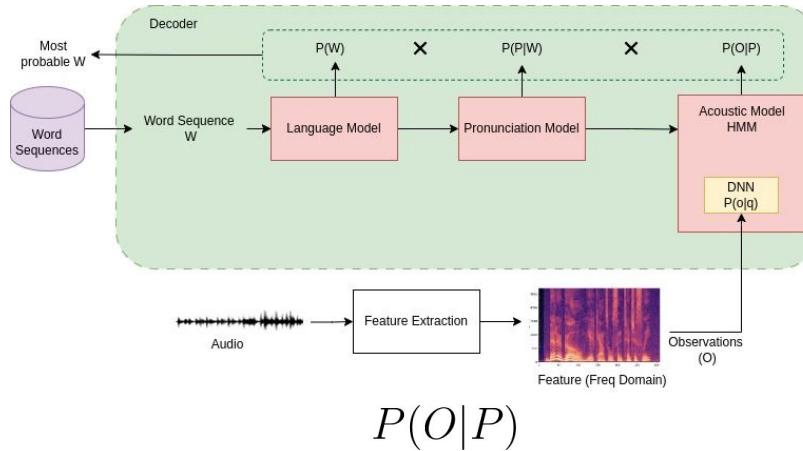
- History
- ASR as a system
- Datasets & Evaluation metrics
- ASR levels of difficulty
- Dynamic time warping (DTW)
- ASR problem formulation
- Inference pipeline
 - Decoding
 - Language models
 - Pronunciation models
 - Feature extraction
 - **Acoustic models**
 - HMM

ASR Pipeline - Acoustic Model

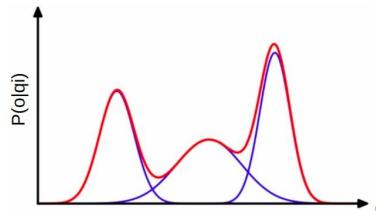


The probability of observing the feature vector O given the phoneme sequence P

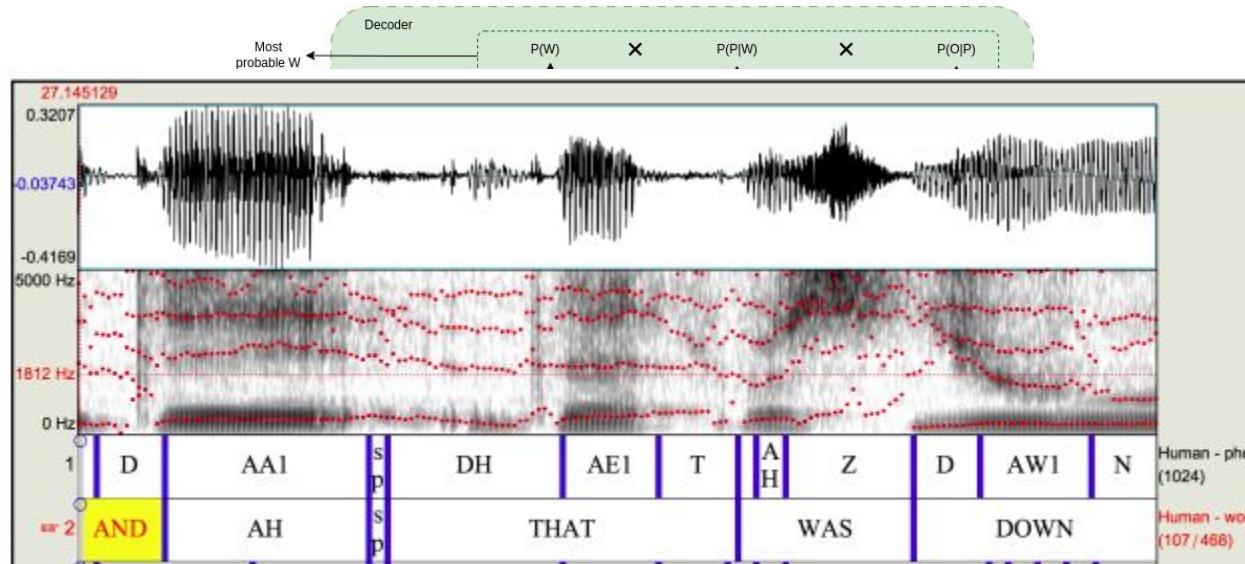
ASR Pipeline - Acoustic Model - Inference



- A single model for all states (multiclass over all states Q)- $p(q_i|o)$.
- A model per each state- $p(o|q_i)$.

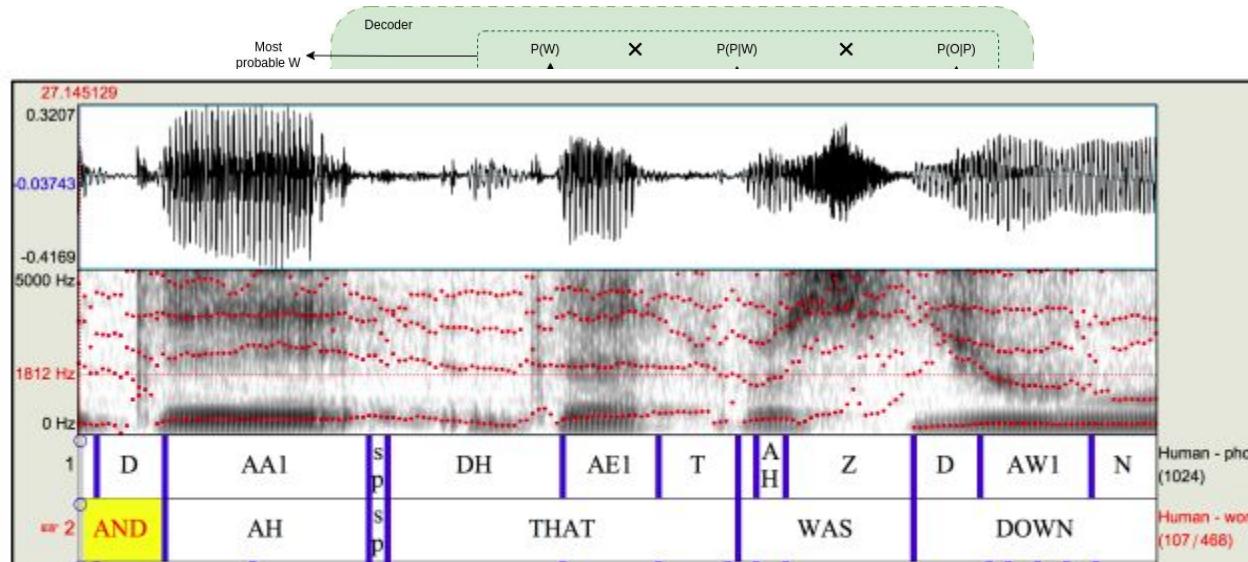


ASR Pipeline - Acoustic Model - Inference



Let's define an alignment Q

ASR Pipeline - Acoustic Model - Inference

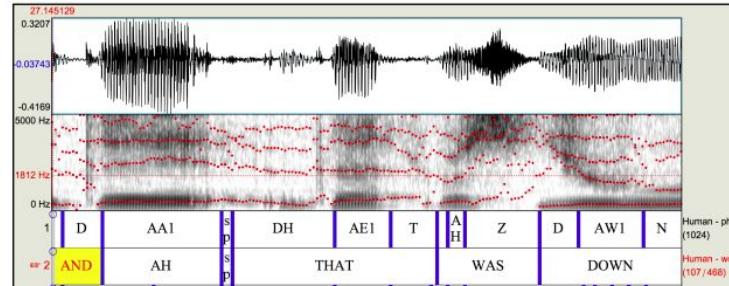


Let's define an alignment Q

P is the phoneme sequence and Q is the **aligned** phoneme sequence.

- For example: $p=[\text{cat}]$. $\text{len}(p) = 3$.
- $q=[\text{ccccaaaaaa.....ttttt}]$. $\text{len}(q) = T$.

ASR Pipeline - Acoustic Model - Inference



$$P(O|P)$$

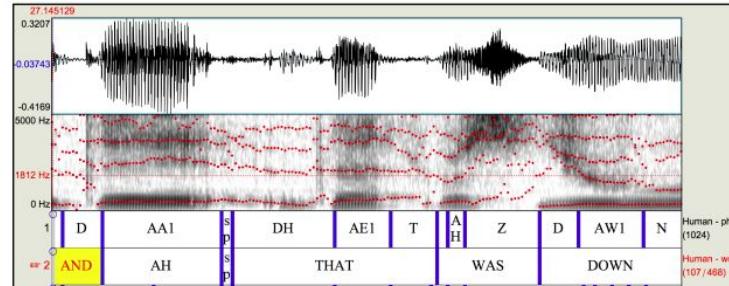
P is the phoneme sequence and Q is the **aligned** phoneme sequence.

- For example: p=[cat]. len(p) = 3.
- q=[ccccaaaaaa.....ttttt]. len(q) = T.

Had we known the alignment Q, the probability of the observation given the alignment can be calculated as follows:

$$P(O|Q) = \prod_{t=1}^T P(o_t|q_t)$$

ASR Pipeline - Acoustic Model - Inference



$$P(O|P)$$

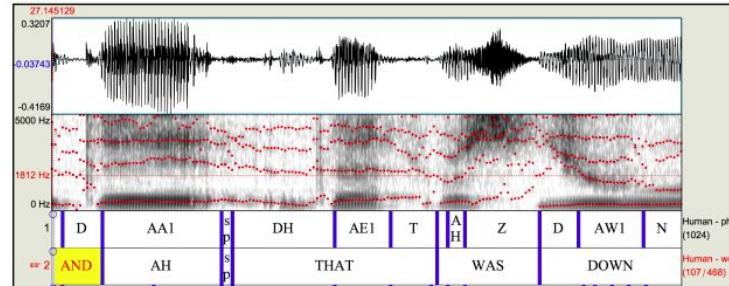
P is the phoneme sequence and Q is the **aligned** phoneme sequence.

- For example: p=[cat]. len(p) = 3.
- q=[ccccaaaaaa.....ttttt]. len(q) = T.

In reality we do not know the alignment, hence:

$$P(O|P) = \sum_{Q \in A(P)} P(O|Q)P(Q|P)$$

ASR Pipeline - Acoustic Model - Inference



$$P(O|P)$$

P is the phoneme sequence and Q is the **aligned** phoneme sequence.

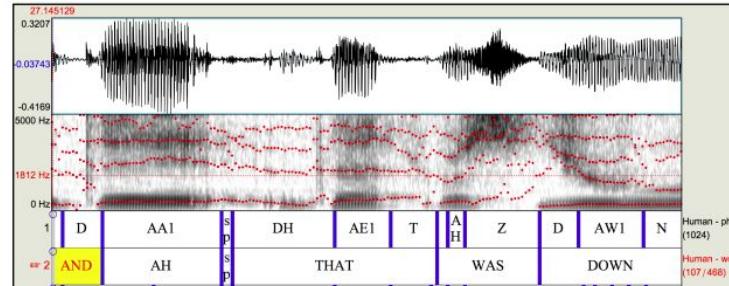
- For example: p=[cat]. len(p) = 3.
- q=[ccccaaaaaa.....ttttt]. len(q) = T.

In reality we do not know the alignment, hence: $P(O|P) = \sum_{Q \in A(P)} P(O|Q)P(Q|P)$

Applying First order markov assumption:

$$P(q_t = S_j | q_{t-1} = S_i, q_{t-2} = S_k, \dots, q_1 = S_l) = P(q_t = S_j | q_{t-1} = S_i)$$

ASR Pipeline - Acoustic Model - Inference



$$P(O|P)$$

P is the phoneme sequence and Q is the **aligned** phoneme sequence.

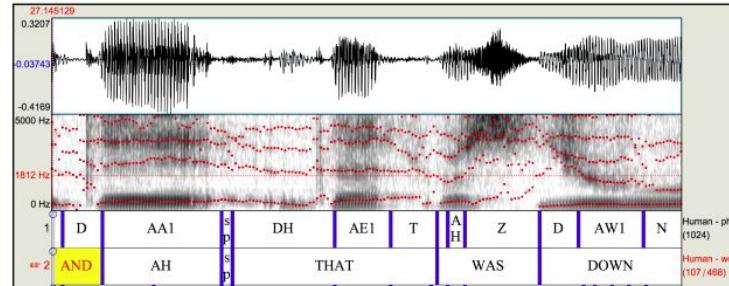
- For example: $p=[\text{cat}]$. $\text{len}(p) = 3$.
- $q=[\text{ccccaaaaaa.....ttttt}]$. $\text{len}(q) = T$.

In reality we do not know the alignment, hence:

$$P(O|P) = \sum_{Q \in A(P)} P(O|Q)P(Q|P)$$

Applying first order markov assumption: $P(q_t = S_j | q_{t-1} = S_i, q_{t-2} = S_k, \dots, q_1 = S_l) = P(q_t = S_j | q_{t-1} = S_i)$

ASR Pipeline - Acoustic Model - Inference



$$P(O|P)$$

P is the phoneme sequence and Q is the **aligned** phoneme sequence.

- For example: p=[cat]. len(p) = 3.
- q=[ccccaaaaaa.....ttttt]. len(q) = T.

$$P(O|P) = \sum_{Q \in A(P)} P(O|Q)P(Q|P)$$

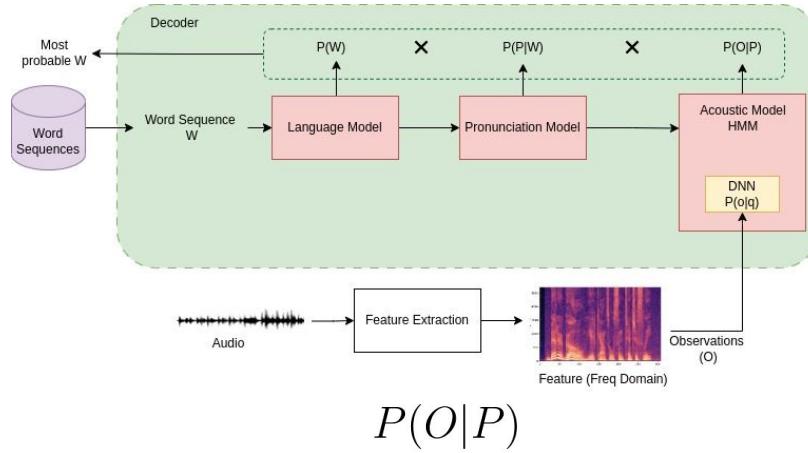
In reality we do not know the alignment, hence:

Applying first order markov assumption: $P(q_t = S_j | q_{t-1} = S_i, q_{t-2} = S_k, \dots, q_1 = S_l) = P(q_t = S_j | q_{t-1} = S_i)$

Results in:

$$P(O|P) = \sum_{Q \in A(P)} P(O|Q)P(Q|P) = \sum_{Q \in A(P)} \prod_{t=1}^T P(o_t | q_t)P(q_t | q_{t-1}, P)$$

ASR Pipeline - Acoustic Model - Training



Components to be trained:

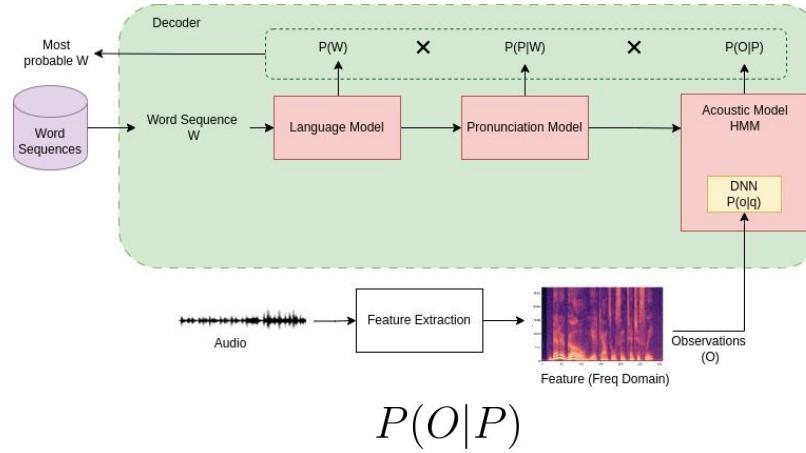
- Acoustic model
- Transition probability (HMM)

$$P(O|P) = \sum_{Q \in A(P)} P(O|Q)P(Q|P) = \sum_{Q \in A(P)} \prod_{t=1}^T P(o_t|q_t)P(q_t|q_{t-1}, P)$$

[1](#) [2](#)

103

ASR Pipeline - Acoustic Model - Training



Components to be trained:

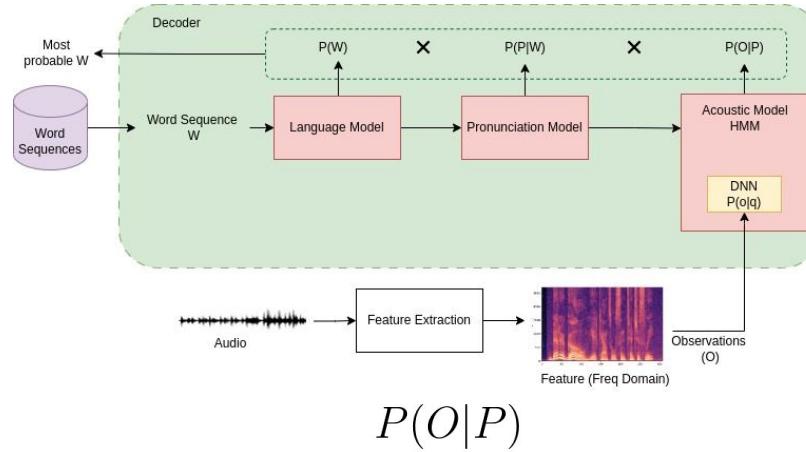
- **Acoustic model**
- Transition probability (HMM)

$$P(O|P) = \sum_{Q \in A(P)} P(O|Q)P(Q|P) = \sum_{Q \in A(P)} \prod_{t=1}^T P(o_t|q_t)P(q_t|q_{t-1}, P)$$

[1](#), [2](#)

104

ASR Pipeline - Acoustic Model - Training

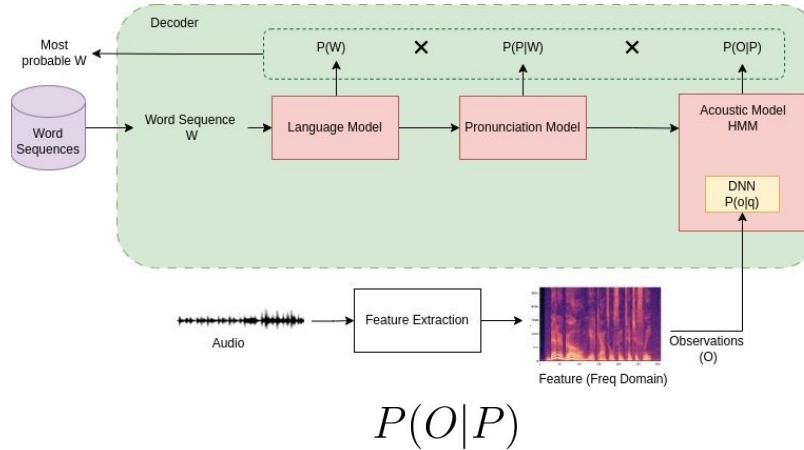


Acoustic model

$P(o|q)$ is trained in a supervised manner

- o and q are known/labeled for each frame.
- Collected manually and/or using Force alignment

ASR Pipeline - Acoustic Model - Training



Acoustic model

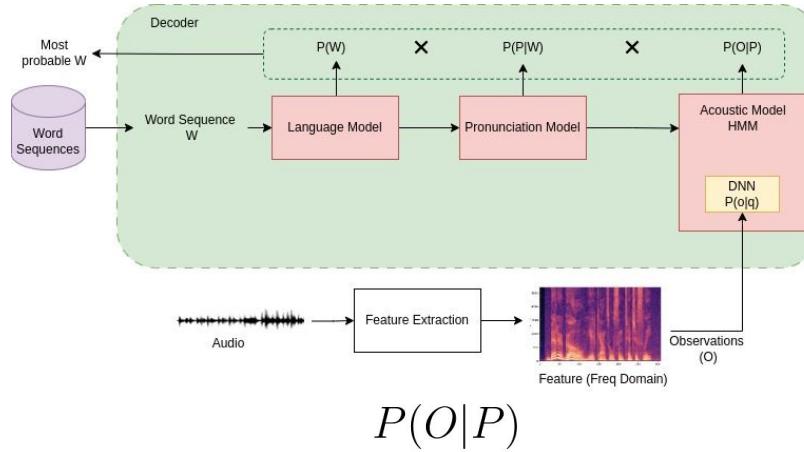
$P(o|q)$ is trained in a supervised manner

- o and q are known/labeled for each frame.
- Collected **manually** and/or using Force alignment

“She had your dark suit in greasy wash water all year”

0	3050	h#	22920	23271	g
3050	4559	sh	23271	24229	r
4559	5723	ix	24229	25566	ix
5723	6642	hv	25566	27156	s
6642	8772	eh	27156	28064	ix
8772	9190	dcl	28064	29660	w
9190	10337	jh	29660	31719	ao
10337	11517	ih	31719	33360	sh
11517	12500	dcl	33360	33754	epi
12500	12640	d	33754	34715	w
12640	14714	ah	34715	36080	ao
14714	15870	kcl	36080	36326	dx
15870	16334	k	36326	37556	axr
16334	18088	s	37556	39561	ao
18088	20417	ux	39561	40313	l
20417	21199	q	40313	42059	y
21199	22560	en	42059	43479	th
22560	22920	gcl	43479	44586	axr
			44586	46720	h#

ASR Pipeline - Acoustic Model - Training



$$P(O|P)$$

Acoustic model

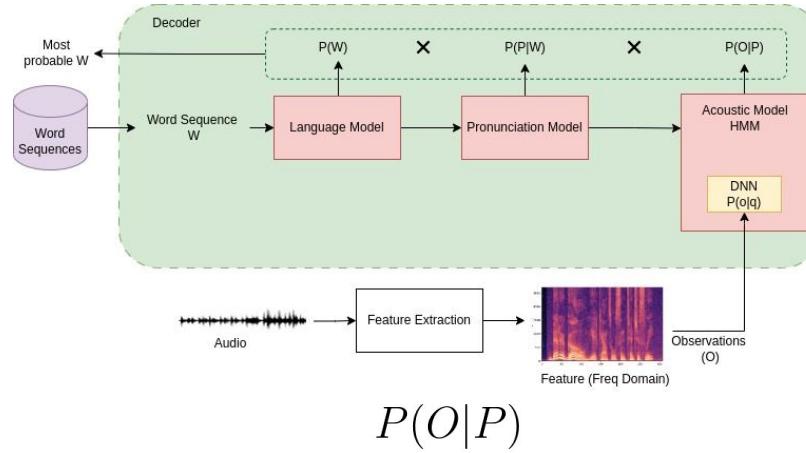
$P(o|q)$ is trained in a supervised manner

- o and q are known/labeled for each frame.
- Collected **manually** and/or using Force alignment

[Link](#)

“MARGUERITE TO BE UNABLE TO LIVE APART FROM ME IT WAS THE DAY AFTER THE EVENING WHEN SHE
CAME TO SEE ME THAT I SENT HER MANON LESCAUT FROM THAT TIME SEEING THAT I COULD NOT CHANGE
MY MISTRESS'S LIFE I CHANGED MY OWN” - [Librispeech](#)

ASR Pipeline - Acoustic Model - Training

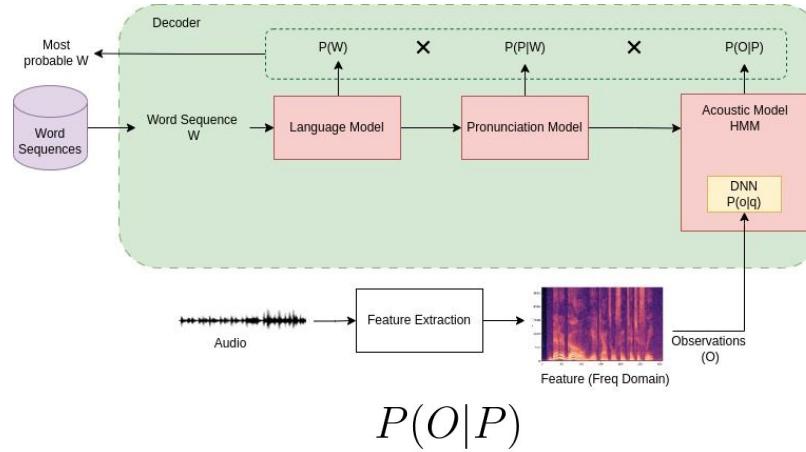


Acoustic model

$P(o|q)$ is trained in a supervised manner

- o and q are known/labeled for each frame.
- Collected manually and/or using **Force alignment**

ASR Pipeline - Acoustic Model - Training

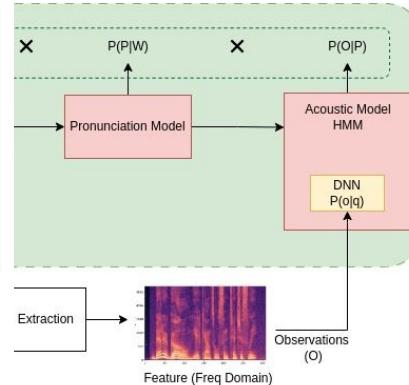
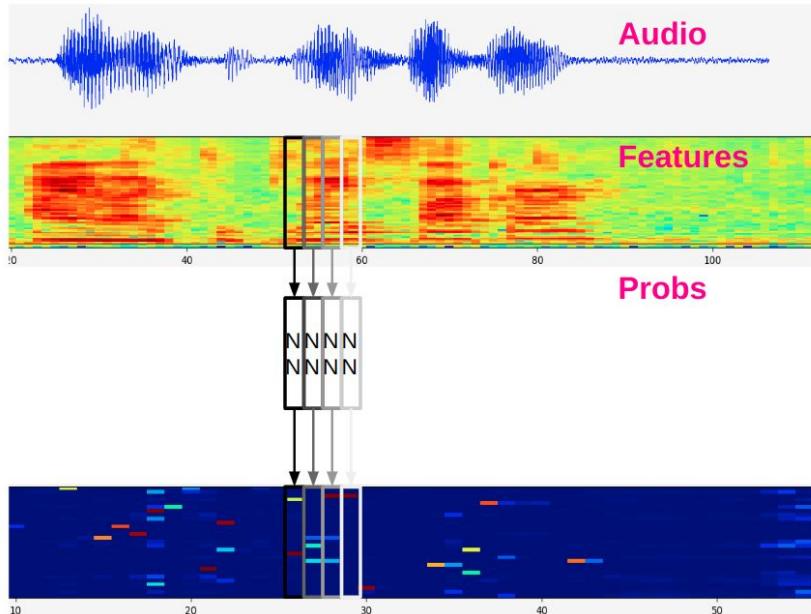


Acoustic model

$P(o|q)$ is trained in a supervised manner

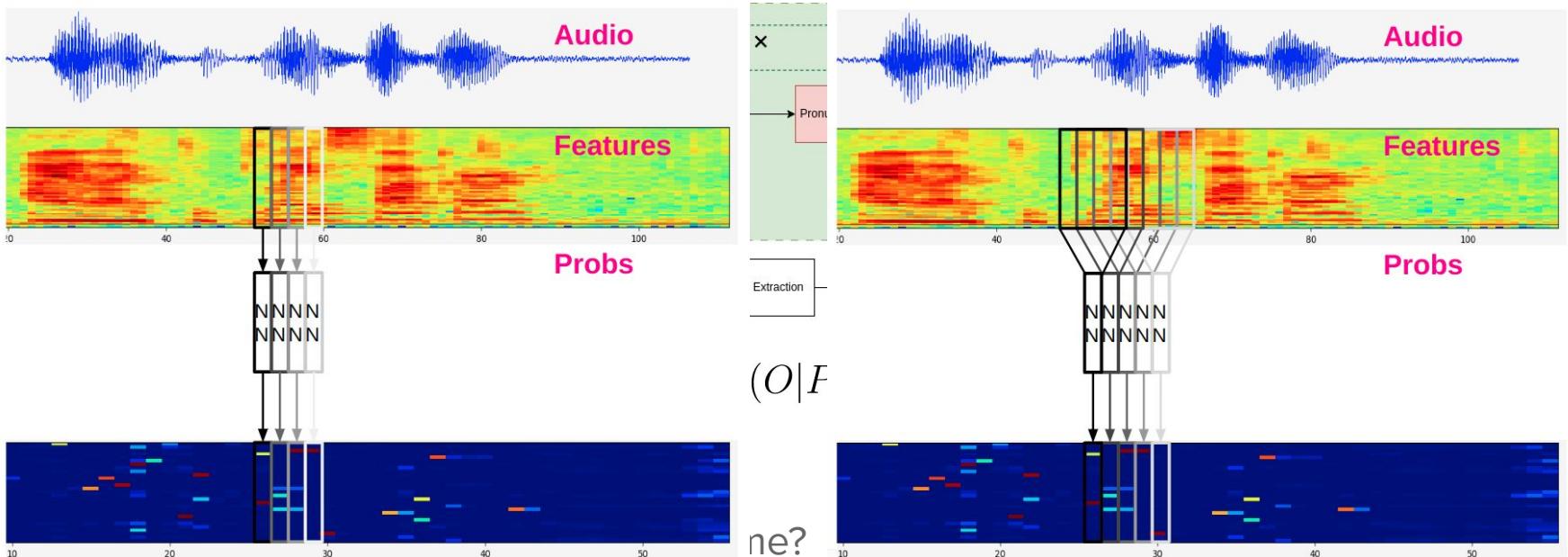
Is one frame enough for classifying a phoneme?

ASR Pipeline - Acoustic Model - Training

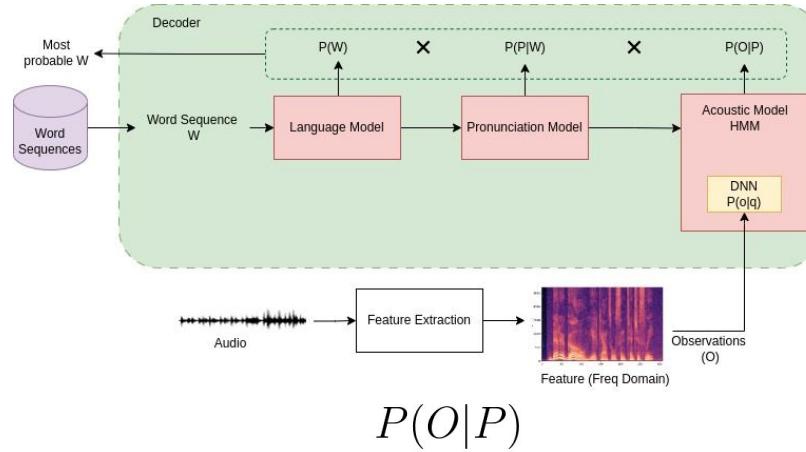


$$(O|P)$$

ASR Pipeline - Acoustic Model - Training



ASR Pipeline - Acoustic Model - Training



Acoustic model

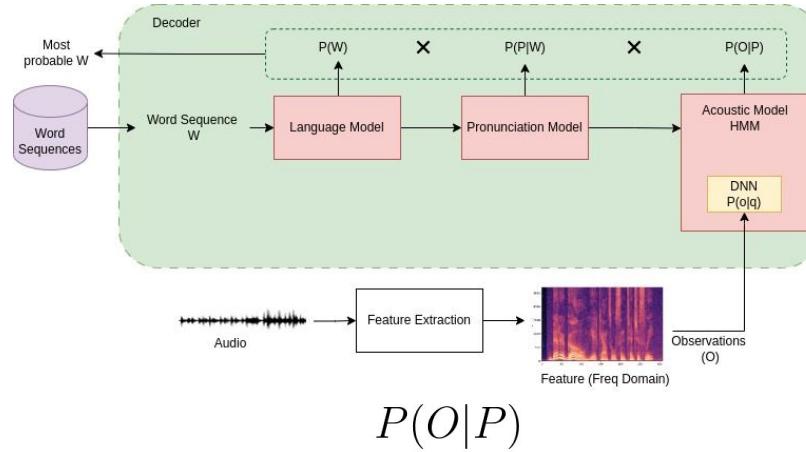
$P(o|q)$ is trained in a supervised manner

Is one frame enough for classifying a phoneme?

A single classifier for all phonemes (states)- discriminative model

$$L_{\Theta}(Q|O) = \sum_t -\log(P(q = q_t|O_t))$$

ASR Pipeline - Acoustic Model - Training



Acoustic model

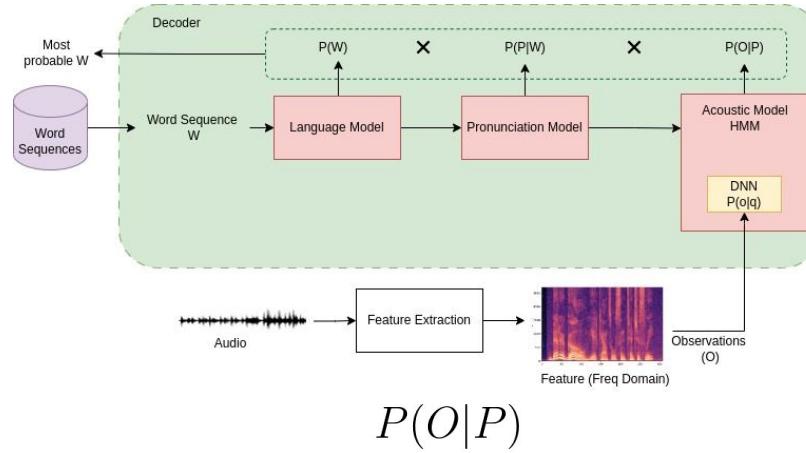
$P(o|q)$ is trained in a supervised manner

Is one frame enough for classifying a phoneme?

A single classifier for all phonemes (states)- discriminative model

$$L_{\Theta}(Q|O) = \sum_t -\log(P(q = q_t|O_t)) \quad P(O|Q) = \frac{P(Q|O)P(O)}{P(Q)}$$

ASR Pipeline - Acoustic Model - Training



Components to be trained:

- **Acoustic model**
- Transition probability (HMM)

$$P(O|P) = \sum_{Q \in A(P)} P(O|Q)P(Q|P) = \sum_{Q \in A(P)} \prod_{t=1}^T P(o_t|q_t)P(q_t|q_{t-1}, P)$$

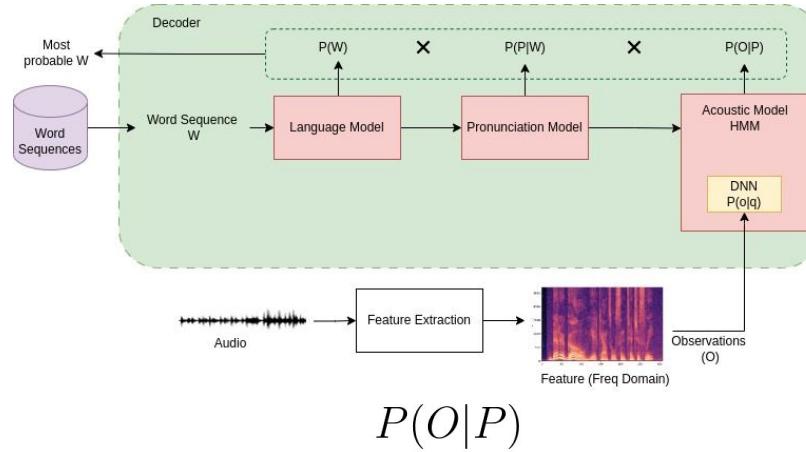
[1](#), [2](#)

114

Outline

- History
- ASR as a system
- Datasets & Evaluation metrics
- ASR levels of difficulty
- Dynamic time warping (DTW)
- ASR problem formulation
- Inference pipeline
 - Decoding
 - Language models
 - Pronunciation models
 - Feature extraction
 - Acoustic models
 - **HMM**

ASR Pipeline - Acoustic Model - Training

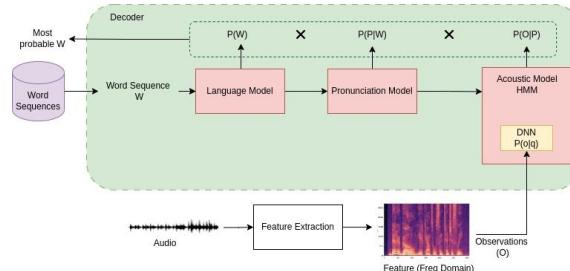


Components to be trained:

- Acoustic model
- **Transition probability (HMM)**

$$P(O|P) = \sum_{Q \in A(P)} P(O|Q)P(Q|P) = \sum_{Q \in A(P)} \prod_{t=1}^T P(o_t|q_t)P(q_t|q_{t-1}, P)$$

ASR Pipeline - Acoustic Model - Training



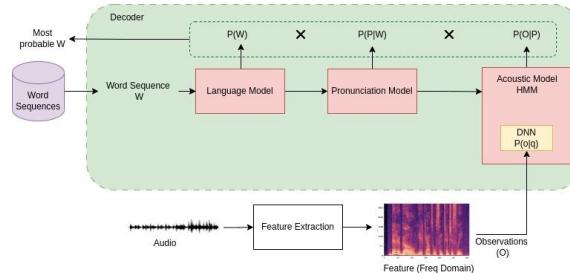
$$P(O|P)$$

First-Order **Observable** Markov Chain

State weather: $Q = q_1, q_2, \dots, q_N$; the state at time t is q_t



ASR Pipeline - Acoustic Model - Training



$$P(O|P)$$

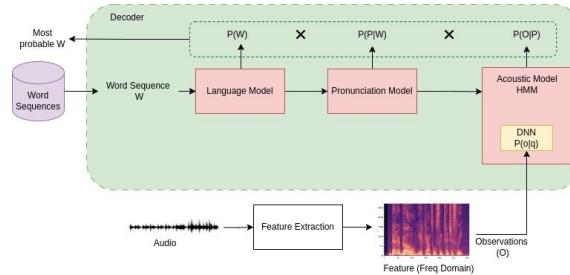
First-Order **Observable** Markov Chain

Current state depends on **previous state only** (not on the entire state history).

$$P(q_t | q_{t-1}, q_{t-2} \dots, q_1) = P(q_t | q_{t-1})$$



ASR Pipeline - Acoustic Model - Training



$$P(O|P)$$

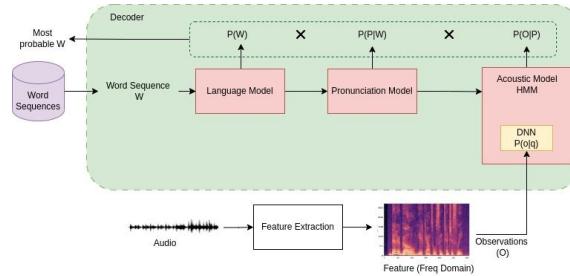
First-Order **Observable** Markov Chain

Current state depends on **previous state only** (not on the entire state history).

Add state transitions



ASR Pipeline - Acoustic Model - Training

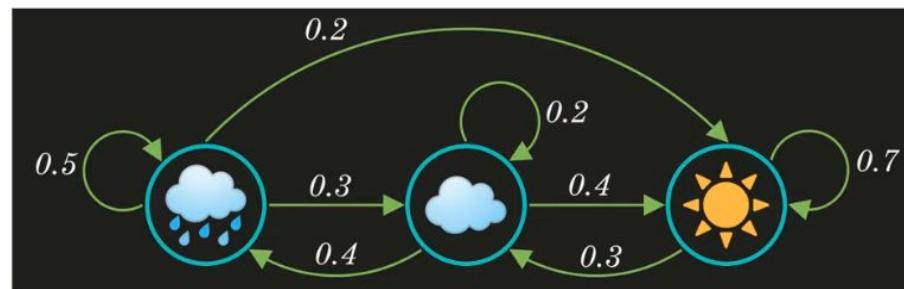


$$P(O|P)$$

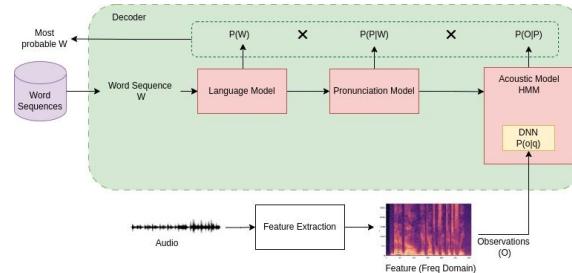
First-Order **Observable** Markov Chain

Current state depends on **previous state only** (not on the entire state history).

Add state transitions & transitions probabilities



ASR Pipeline - Acoustic Model - Training



$$P(O|P)$$

First-Order **Observable** Markov Chain

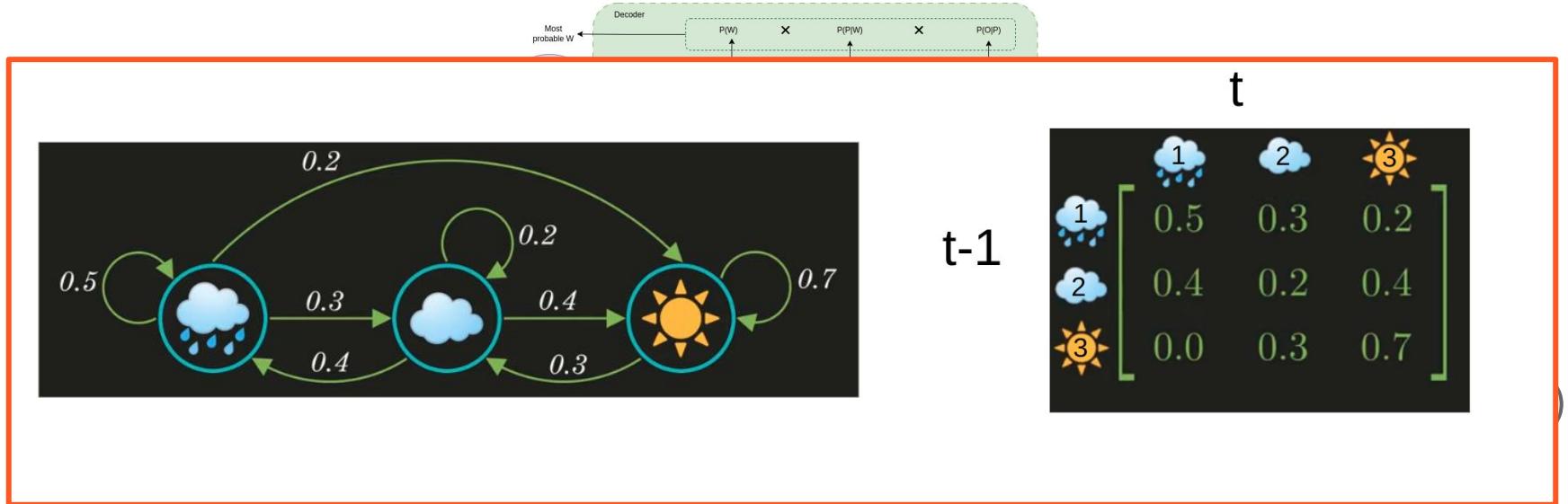
Current state depends on **previous state only** (not on the entire state history).

Add state transitions & transitions probabilities

The set of transition probabilities is annotated as $A = a_{01}, a_{02}, \dots, a_{n1}, \dots, a_{nn}$, where a_{ij} is defined as the probability of transitioning from state i to state j:

$$a_{ij} = P(q_t = j | q_{t-1} = i)$$

ASR Pipeline - Acoustic Model - Training



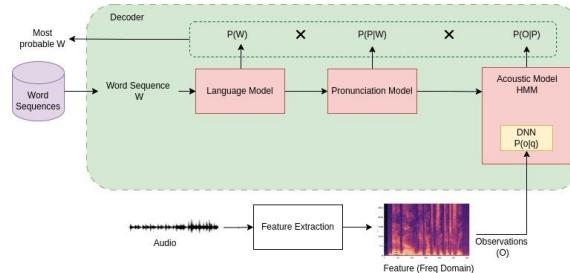
The set of transition probabilities is annotated as $A = a_{01}, a_{02}, \dots, a_{n1}, \dots, a_{nn}$, where a_{ij} is defined as the probability of transitioning from state i to state j :

$$a_{ij} = P(q_t = j | q_{t-1} = i) \quad \sum_{j=1}^N a_{ij} = 1; \quad 1 \leq i \leq N$$

Source [1, 2](#)

122

ASR Pipeline - Acoustic Model - Training



$$P(O|P)$$

First-Order **Observable** Markov Chain

Current state depends on **previous state only** (not on the entire state history).

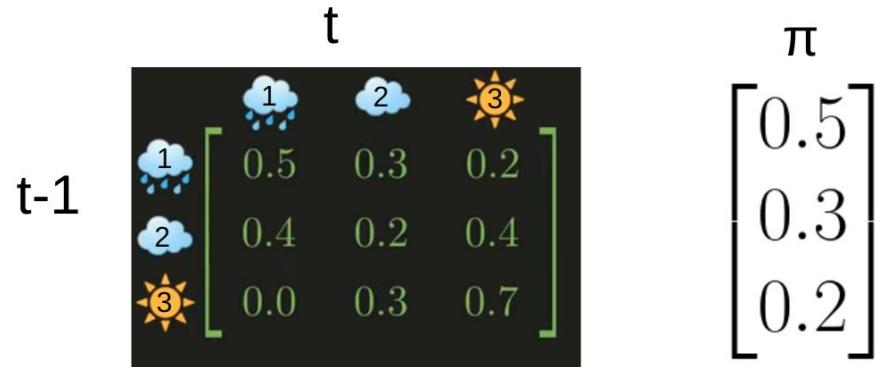
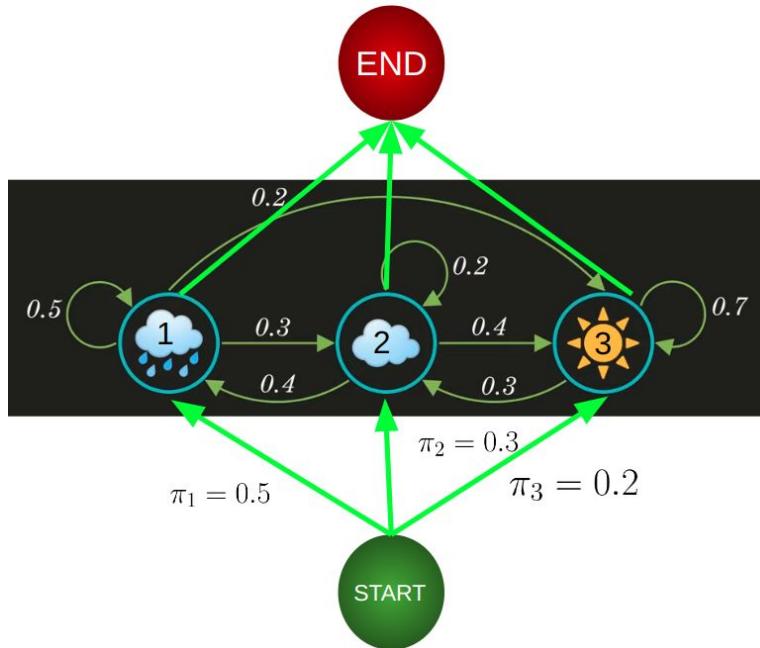
Add state transitions & transitions probabilities

First-Order Observable Markov Chain has distinguished **start** and **end** states

$$\pi_i = P(q_1 = i) \quad 1 \leq i \leq N$$

$$\sum_{j=1}^N \pi_j = 1$$

ASR Pipeline - Acoustic Model - Training



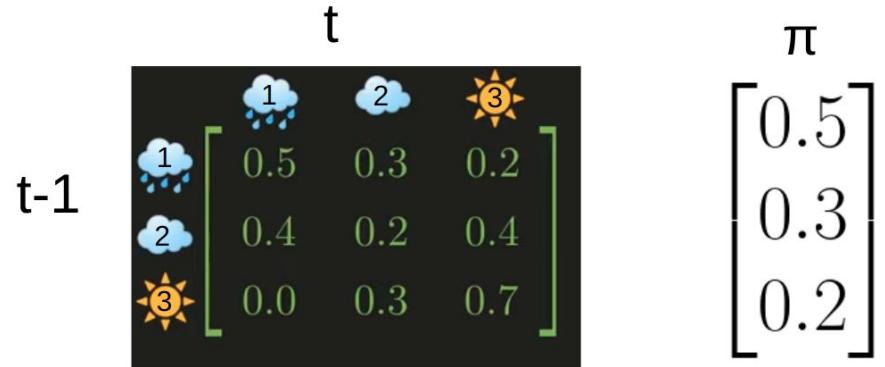
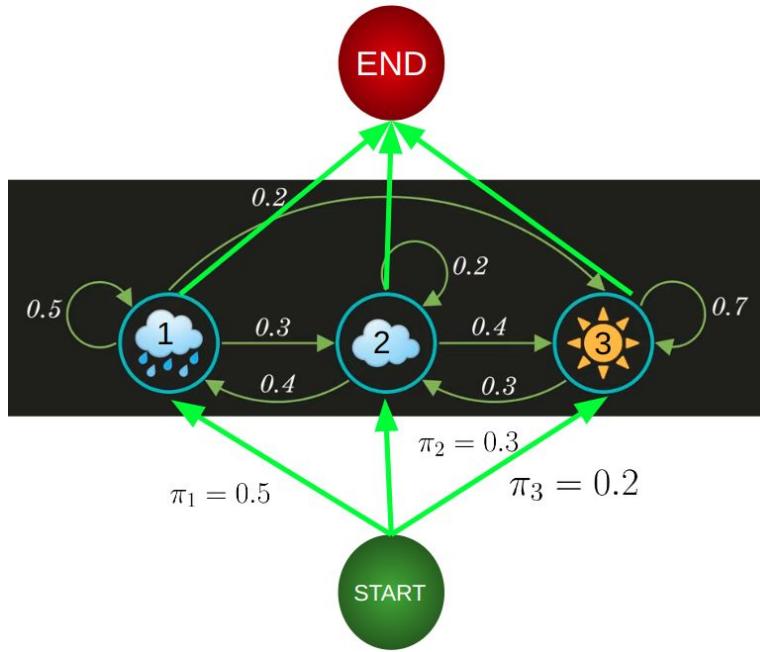
$$\pi_i = P(q_1 = i) \quad 1 \leq i \leq N$$

$$\sum_{j=1}^N \pi_j = 1$$

Source [1](#), [2](#)

124

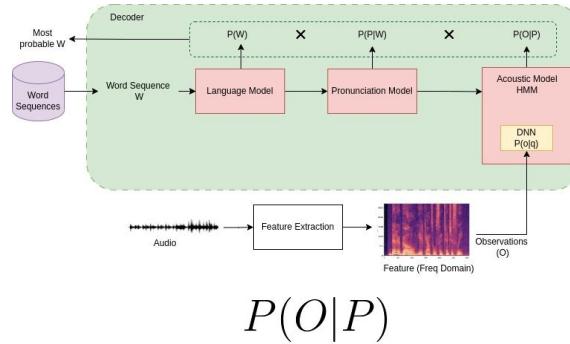
ASR Pipeline - Acoustic Model - Training



The probability of the trajectory 1, 2, 3 ,1

$$P(1, 2, 3, 1) = \pi_1 * a_{21} * a_{32} * a_{13} = 0.5 * 0.3 * 0.4 * 0.0 = 0.0$$

ASR Pipeline - Acoustic Model - Training

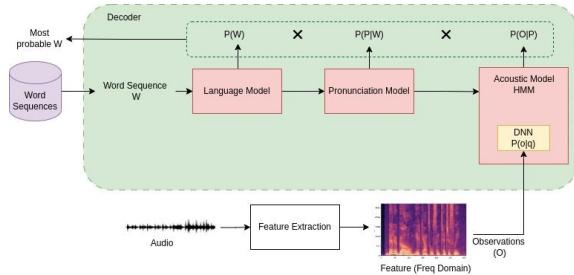


Problem:

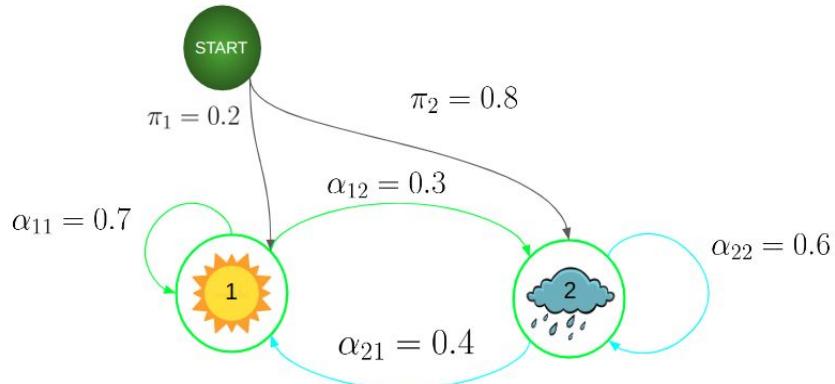
- States are observable in Markov chains → output symbols == state symbols.
- This is not the case in speech recognition. We can observe (hear) the audio, but the transcript (phoneme states) is ‘hidden’ from us.

Solution: extend the markov chain to the **Hidden Markov Model (HMM)** where the states are hidden, meaning we can observe something related to the state, but we cannot access the state itself.

ASR Pipeline - Acoustic Model - Training



$$P(O|P)$$



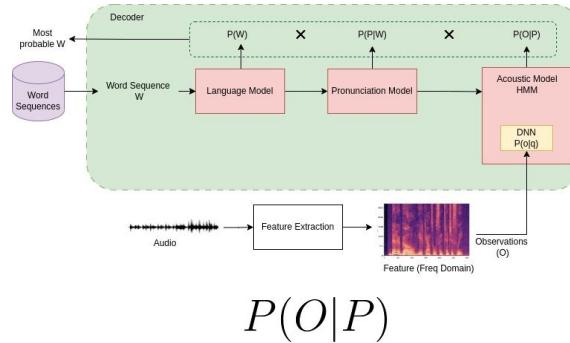
$$a_{ij} = P(q_t = j | q_{t-1} = i)$$

t π

Below the HMM diagram, the state probabilities are shown for time steps $t-1$ and t . At time $t-1$, the sun state (1) has a probability of $[0.7 \ 0.3]$ and the cloud state (2) has a probability of $[0.4 \ 0.6]$. At time t , the sun state (1) has a probability of $[0.2 \ 0.8]$ and the cloud state (2) has a probability of $[0.8 \ 0.2]$.

$$\begin{matrix} t-1 & \begin{bmatrix} 0.7 & 0.3 \\ 0.4 & 0.6 \end{bmatrix} & \begin{bmatrix} 0.2 & 0.8 \\ 0.8 & 0.2 \end{bmatrix} \\ \text{1} & [0.7 \ 0.3] & [0.2 \ 0.8] \\ \text{2} & [0.4 \ 0.6] & [0.8 \ 0.2] \end{matrix}$$

ASR Pipeline - Acoustic Model - Training

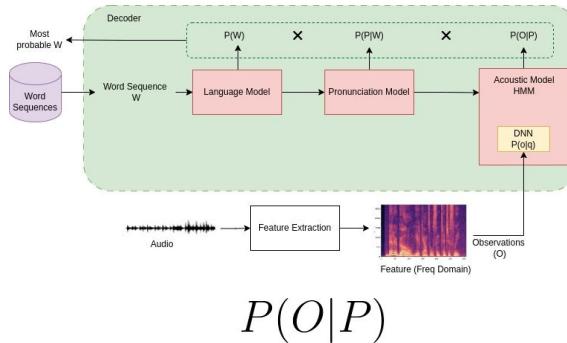


We are **not exposed** to the weather, but we **do know** how many ice-creams Jack ate every day

- Hidden States: weather condition: hot / cold
- Observations: number of ice-creams he ate.



ASR Pipeline - Acoustic Model - Training



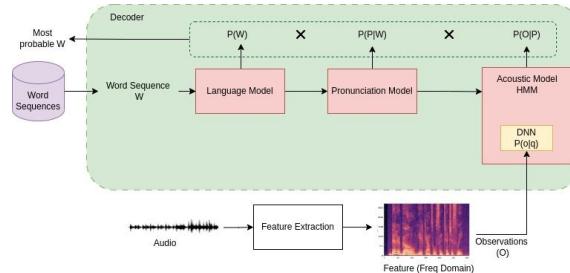
We are **not exposed** to the weather, but we **do know** how many ice-creams Jack ate every day

- Hidden States: weather condition: hot / cold
- Observations: number of ice-creams he ate.



Our job: figure out how is the weather (hidden states)- cold / hot, based on the number of ice-creams (observations) that he ate.

ASR Pipeline - Acoustic Model - Training



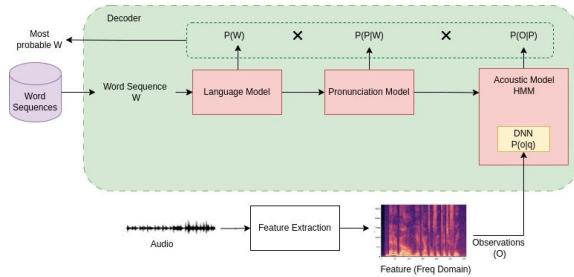
$$P(O|P)$$

Define emission probability:

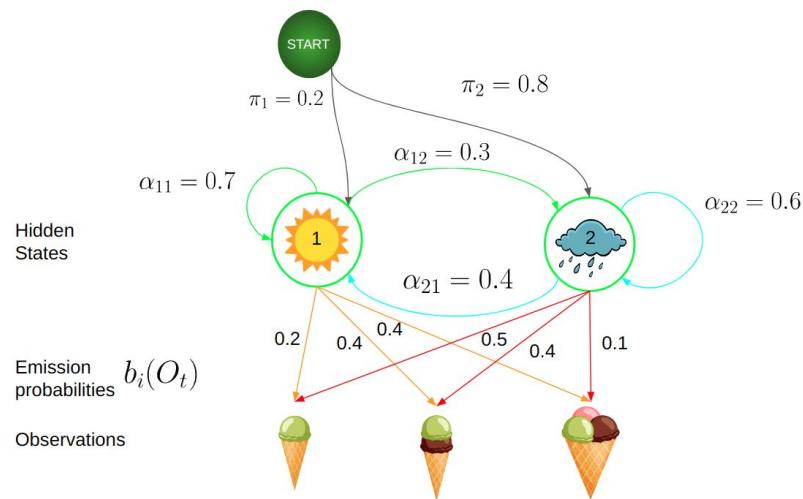
$$b_i(t) = P(o_t | q_t = i)$$

Emission probabilities are **stochastic** and **NOT deterministic**

ASR Pipeline - Acoustic Model - Training



$$P(O|P)$$

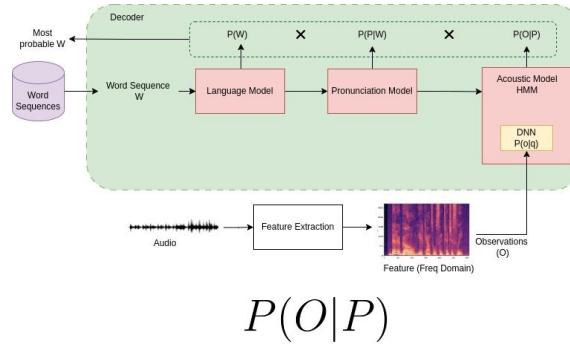


t-1	t	π	B1	B2
		$\begin{bmatrix} 1 \\ 2 \end{bmatrix}$	$\begin{bmatrix} P(\text{Sun} \text{Sun}) & P(\text{Cloud} \text{Sun}) \\ P(\text{Sun} \text{Cloud}) & P(\text{Cloud} \text{Cloud}) \end{bmatrix} = \begin{bmatrix} 0.7 & 0.3 \\ 0.4 & 0.6 \end{bmatrix}$	$\begin{bmatrix} P(\text{Sun} \text{Sun}) \\ P(\text{Cloud} \text{Sun}) \end{bmatrix} = \begin{bmatrix} 0.2 \\ 0.8 \end{bmatrix}$
		$\begin{bmatrix} 0.2 \\ 0.8 \end{bmatrix}$	$\begin{bmatrix} P(\text{Ice Cream} \text{Sun}) \\ P(\text{Ice Cream} \text{Cloud}) \end{bmatrix} = \begin{bmatrix} 0.4 \\ 0.6 \end{bmatrix}$	$\begin{bmatrix} P(\text{Ice Cream} \text{Sun}) \\ P(\text{Ice Cream} \text{Cloud}) \end{bmatrix} = \begin{bmatrix} 0.5 \\ 0.4 \end{bmatrix}$

$$a_{ij}(t) = P(q_t = j | q_{t-1} = i)$$

$$b_i(t) = P(o_t | q_t = i)$$

ASR Pipeline - Acoustic Model - Training



Output-independence assumption

$$P(o_t | O_1^{t-1}, q_1^t) = P(o_t | q_t)$$

The number of ice-creams that he ate (observation) depends only on the weather of that particular day (hidden state).

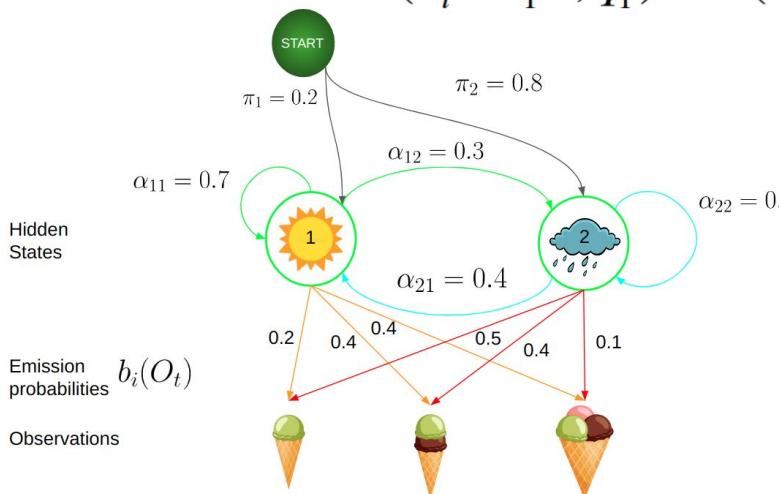
ASR Pipeline - Acoustic Model - Training

First order Markov assumption:

$$P(q_t | q_{t-1}, q_{t-2} \dots, q_1) = P(q_t | q_{t-1})$$

Output-independence assumption

$$P(o_t | O_1^{t-1}, q_t) = P(o_t | q_t)$$

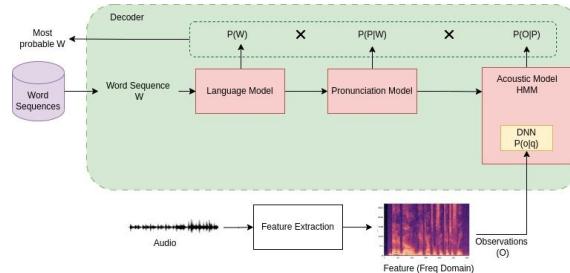


$t-1$	t	π	$B1$	$B2$
$\begin{matrix} 1 \\ 2 \end{matrix}$	$\begin{matrix} 1 & 2 \\ 1 & 2 \end{matrix}$	$\begin{bmatrix} 0.2 \\ 0.8 \end{bmatrix}$	$\begin{bmatrix} P(\text{Ice Cream 1} 1) \\ P(\text{Ice Cream 2} 1) \\ P(\text{Ice Cream 3} 1) \end{bmatrix} = \begin{bmatrix} 0.2 \\ 0.4 \\ 0.4 \end{bmatrix}$	$\begin{bmatrix} P(\text{Ice Cream 1} 2) \\ P(\text{Ice Cream 2} 2) \\ P(\text{Ice Cream 3} 2) \end{bmatrix} = \begin{bmatrix} 0.5 \\ 0.4 \\ 0.1 \end{bmatrix}$

$$a_{ij}(t) = P(q_t = j | q_{t-1} = i)$$

$$b_i(t) = P(o_t | q_t = i)$$

ASR Pipeline - Acoustic Model - Training

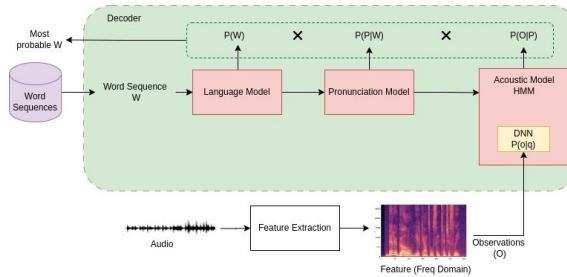


$$P(O|P)$$

HMM solves the following tasks:

- **Evaluation (likelihood):** Given the observation sequence $O = \{o_1, o_2, \dots, o_T\}$, and an HMM model $\lambda = \{A, B\}$, how do we efficiently compute $P(O|\lambda)$, the probability of the observed sequence, given the model?
- **Decoding:** Given the observation sequence $O = \{o_1, o_2, \dots, o_T\}$, and an HMM model $\lambda = \{A, B\}$, how do we choose a corresponding state sequence $Q = \{q_1, q_2, \dots, q_T\}$ that is optimal, i.e., best explains the observations?
- **Learning:** How do we learn the parameters $\lambda = \{A, B\}$ to maximize $P(O|\lambda)$

ASR Pipeline - Acoustic Model - Training

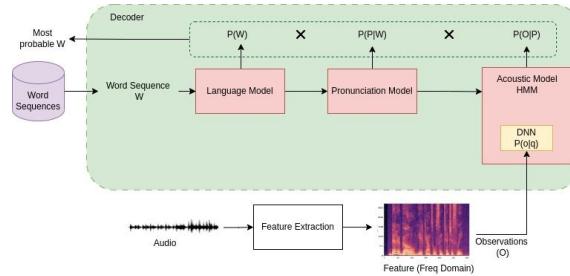


$$P(O|P)$$

HMM solves the following tasks:

- **Evaluation (likelihood):** Given the observation sequence $O = \{o_1, o_2, \dots, o_T\}$, and an HMM model $\lambda = \{A, B\}$, how do we efficiently compute $P(O|\lambda)$, the probability of the observed sequence, given the model?
- **Decoding:** Given the observation sequence $O = \{o_1, o_2, \dots, o_T\}$, and an HMM model $\lambda = \{A, B\}$, how do we choose a corresponding state sequence $Q = \{q_1, q_2, \dots, q_T\}$ that is optimal, i.e., best explains the observations? (This task is highlighted with a red border.)
- **Learning:** How do we learn the parameters $\lambda = \{A, B\}$ to maximize $P(O|\lambda)$

ASR Pipeline - Acoustic Model - Training



$$P(O|P)$$

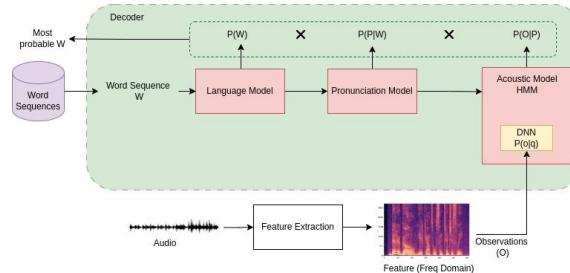
Decoding

Decoding solves the followings:

Given the observation sequence $O = \{o_1, o_2, \dots, o_T\}$, and an HMM model $\lambda = \{A, B\}$, how do we choose a corresponding state sequence $Q = \{q_1, q_2, \dots, q_T\}$ that is optimal, i.e., best explains the observations?

$$\hat{Q} = \arg \max_{Q=\{q_1, \dots, q_N\}} P(O, Q | \lambda)$$

ASR Pipeline - Acoustic Model - Training



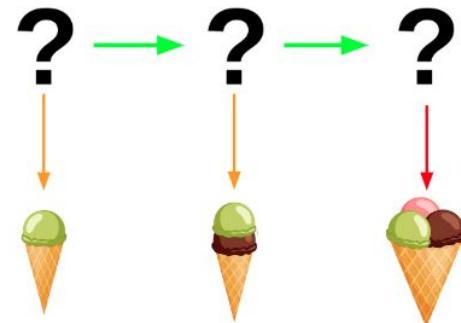
$$P(O|P)$$

Decoding

Decoding solves the followings:

Given the observation sequence

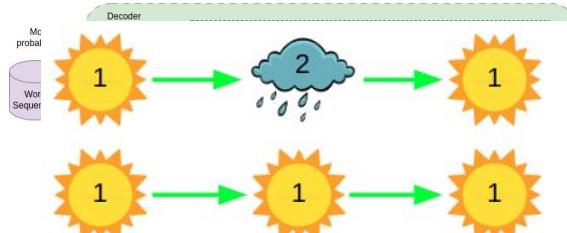
we choose a corresponding



an HMM model $\lambda = \{A, B\}$, how do we choose a corresponding q_1, \dots, q_T that is optimal, i.e., best is?

$$\hat{Q} = \arg \max_{Q=\{q_1, \dots, q_N\}} P(O, Q | \lambda)$$

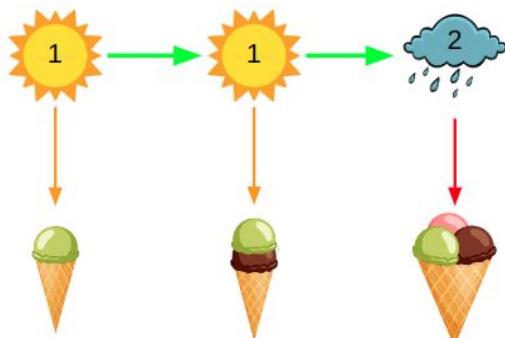
ASR Pipeline - Acoustic Model - Training



Decoding

Decoding solves the followings:

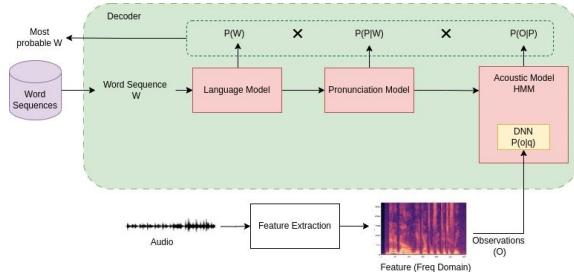
Given the observation sequence
we choose a corresponding



an HMM model $\lambda = \{A, B\}$, how do we choose a corresponding q_1, q_2, \dots, q_T that is optimal, i.e., best fits?

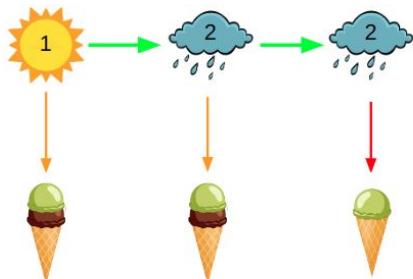
$$\hat{Q} = \arg \max_{Q=\{q_1, \dots, q_N\}} P(O, Q | \lambda)$$

ASR Pipeline - Acoustic Model - Training



$$P(O|P)$$

Decoding



$$P(2, 2, 1, H, C, C | A, B) = ?$$

	t	π	$B1$	$B2$
$t-1$	$\begin{bmatrix} 1 \\ 2 \end{bmatrix}$ $\begin{bmatrix} 0.7 & 0.3 \\ 0.4 & 0.6 \end{bmatrix}$	$\begin{bmatrix} 0.2 \\ 0.8 \end{bmatrix}$	$\begin{bmatrix} P(\text{ice cream} 1) \\ P(\text{ice cream} 1) \\ P(\text{ice cream} 1) \end{bmatrix} = \begin{bmatrix} 0.2 \\ 0.4 \\ 0.4 \end{bmatrix}$	$\begin{bmatrix} P(\text{ice cream} 2) \\ P(\text{ice cream} 2) \\ P(\text{ice cream} 2) \end{bmatrix} = \begin{bmatrix} 0.5 \\ 0.4 \\ 0.1 \end{bmatrix}$
			$a_{ij}(t) = P(q_t = j q_{t-1} = i)$	
			$b_i(t) = P(o_t q_t = i)$	

$$\begin{aligned}
 P(o_1 = 2, o_2 = 2, o_3 = 1, q_1 = H, q_2 = C, q_3 = C) &= \\
 \pi_H * p(2|H) * p(C|H) * p(2|C) * p(C|C) * p(1|C) &= \\
 0.2 * 0.4 * 0.3 * 0.4 * 0.6 * 0.5 &= 0.00288
 \end{aligned}$$

ASR Pipeline - Acoustic Model - Training

Decoding solves the followings:

Given the observation sequence $O = \{o_1, o_2, \dots, o_T\}$, and an HMM model $\lambda = \{A, B\}$, how do we choose a corresponding state sequence $Q = \{q_1, q_2, \dots, q_T\}$ that is optimal, i.e., best explains the observations?

General Solution

$$P(O, Q) = P(O|Q) \times P(Q) = \prod_{i=1}^n P(o_i|q_i) \times \prod_{i=1}^n P(q_i|q_{i-1})$$

Finding the state sequence that maximizes the probability of the observed sequence:

- Using the naive approach: Iterate over all state sequences Q:
 - HHH, HHC, CCC, CHC, etc.
 - For each state sequence calculate $P(O, Q)$
 - Take the argmax
- This exhaustive search has **high computational complexity \rightarrow infeasible**

ASR Pipeline - Acoustic Model - Training

Decoding solves the followings:

Given the observation sequence $O = \{o_1, o_2, \dots, o_T\}$, and an HMM model $\lambda = \{A, B\}$, how do we choose a corresponding state sequence $Q = \{q_1, q_2, \dots, q_T\}$ that is optimal, i.e., best explains the observations?

General Solution

$$P(O, Q) = P(O|Q) \times P(Q) = \prod_{i=1}^n P(o_i|q_i) \times \prod_{i=1}^n P(q_i|q_{i-1})$$

Solution: Viterbi algorithm

Apply **dynamic programming** to efficiently calculate the max probability.

Viterbi finds the state sequence that maximizes the probability of observed sequence.

ASR Pipeline - Acoustic Model - Training

Decoding solves the followings:

Given the observation sequence $O = \{o_1, o_2, \dots, o_T\}$, and an HMM model $\lambda = \{A, B\}$, how do we choose a corresponding state sequence $Q = \{q_1, q_2, \dots, q_T\}$ that is optimal, i.e., best explains the observations?

General Solution

$$P(O, Q) = P(O|Q) \times P(Q) = \prod_{i=1}^n P(o_i|q_i) \times \prod_{i=1}^n P(q_i|q_{i-1})$$

Solution: Viterbi algorithm

$$v_t(j) = \max_{q_0, q_1, \dots, q_{t-1}} P(q_0, q_1 \dots q_{t-1}, o_1, o_2 \dots o_t, q_t = j | \lambda)$$

ASR Pipeline - Acoustic Model - Training

Decoding solves the followings:

Given the observation sequence $O = \{o_1, o_2, \dots, o_T\}$, and an HMM model $\lambda = \{A, B\}$, how do we choose a corresponding state sequence $Q = \{q_1, q_2, \dots, q_T\}$ that is optimal, i.e., best explains the observations?

General Solution

$$P(O, Q) = P(O|Q) \times P(Q) = \prod_{i=1}^n P(o_i|q_i) \times \prod_{i=1}^n P(q_i|q_{i-1})$$

Solution: Viterbi algorithm

$$v_t(j) = \max_{q_0, q_1, \dots, q_{t-1}} P(q_0, q_1 \dots q_{t-1}, o_1, o_2 \dots o_t, q_t = j | \lambda)$$

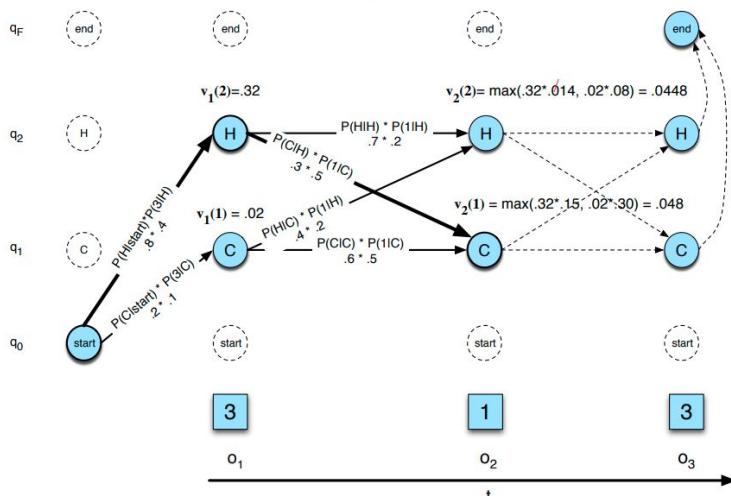
$$v_t(q_j) = \max_{i=1 \dots N} v_{t-1}(q = i) \alpha_{ij} b_j(o_t)$$

ASR Pipeline - Acoustic Model - Training

Decoding solves the followings:

Given the observation sequence $O = \{o_1, o_2, \dots, o_T\}$, and an HMM model $\lambda = \{A, B\}$, how do we choose a corresponding state sequence $Q = \{q_1, q_2, \dots, q_T\}$ that is optimal, i.e., best explains the observations?

Viterbi trellis



t	π	B1	B2
1 2	$\begin{bmatrix} 0.7 & 0.3 \\ 0.4 & 0.6 \end{bmatrix}$	$\begin{bmatrix} P(\text{ice cream} 1) \\ P(\text{ice cream} 2) \\ P(\text{rain} 1) \end{bmatrix} = \begin{bmatrix} 0.2 \\ 0.4 \\ 0.4 \end{bmatrix}$	$\begin{bmatrix} P(\text{ice cream} 2) \\ P(\text{ice cream} 2) \\ P(\text{rain} 2) \end{bmatrix} = \begin{bmatrix} 0.5 \\ 0.4 \\ 0.1 \end{bmatrix}$
$t-1$	π	$B1$	$B2$

$$a_{ij}(t) = P(q_t = j | q_{t-1} = i)$$

$$b_i(t) = P(o_t | q_t = i)$$

ASR Pipeline - Acoustic Model - Training

Decoding solves the followings:

Given the observation sequence $O = \{o_1, o_2, \dots, o_T\}$, and an HMM model $\lambda = \{A, B\}$, how do we choose a corresponding state sequence $Q = \{q_1, q_2, \dots, q_T\}$ that is optimal, i.e., best explains the observations?

Solution: Viterbi algorithm

Apply **dynamic programming** to efficiently calculate the max probability

1. **Initialization:**

$$\begin{aligned} v_1(j) &= a_{0j} b_j(o_1) \quad 1 \leq j \leq N \\ bt_1(j) &= 0 \end{aligned}$$

2. **Recursion** (recall that states 0 and q_F are non-emitting):

$$\begin{aligned} v_t(j) &= \max_{i=1}^N v_{t-1}(i) a_{ij} b_j(o_t); \quad 1 \leq j \leq N, 1 < t \leq T \\ bt_t(j) &= \operatorname{argmax}_{i=1}^N v_{t-1}(i) a_{ij} b_j(o_t); \quad 1 \leq j \leq N, 1 < t \leq T \end{aligned}$$

3. **Termination:**

$$\text{The best score: } P_* = v_t(q_F) = \max_{i=1}^N v_T(i) * a_{i,F}$$

$$\text{The start of backtrace: } q_T^* = bt_T(q_F) = \operatorname{argmax}_{i=1}^N v_T(i) * a_{i,F}$$

ASR Pipeline - Acoustic Model - Training

Decoding solves the followings:

Given the observation sequence $O = \{o_1, o_2, \dots, o_T\}$, and an HMM model $\lambda = \{A, B\}$, how do we choose a corresponding state sequence $Q = \{q_1, q_2, \dots, q_T\}$ that is optimal, i.e., best explains the observations?

Solution: Viterbi algorithm

Apply dynamic programming to efficiently calculate the max probability

1. Initialization:

$$\begin{aligned} v_1(j) &= a_{0j} b_j(o_1) \quad 1 \leq j \leq N \\ b_{t1}(j) &= 0 \end{aligned}$$

Viterbi finds the state sequence that maximizes the probability of observed sequence.

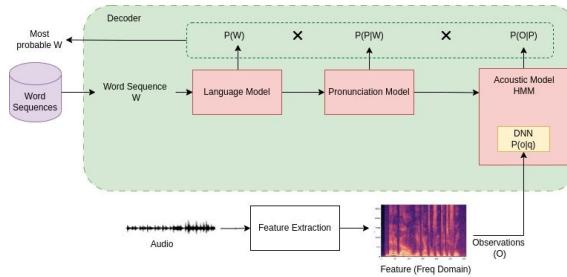
Main Usage: Force alignment

3. Termination:

$$\text{The best score: } P* = v_t(q_F) = \max_{i=1}^N v_T(i) * a_{i,F}$$

$$\text{The start of backtrace: } q_T* = b_{tT}(q_F) = \arg \max_{i=1}^N v_T(i) * a_{i,F}$$

ASR Pipeline - Acoustic Model - Training



$$P(O|P)$$

HMM solves the following tasks:

- **Evaluation (likelihood):** Given the observation sequence $O = \{o_1, o_2, \dots, o_T\}$, and an HMM model $\lambda = \{A, B\}$, how do we efficiently compute $P(O|\lambda)$, the probability of the observed sequence, given the model?
- **Decoding:** Given the observation sequence $O = \{o_1, o_2, \dots, o_T\}$, and an HMM model $\lambda = \{A, B\}$, how do we choose a corresponding state sequence $Q = \{q_1, q_2, \dots, q_T\}$ that is optimal, i.e., best explains the observations?
- **Learning:** How do we learn the parameters $\lambda = \{A, B\}$ to maximize $P(O|\lambda)$

ASR Pipeline - Acoustic Model - Training

Evaluation solves the followings:

Given the observation sequence $O = \{o_1, o_2, \dots, o_T\}$, and an HMM model $\lambda = \{A, B\}$, how do we efficiently compute $P(O|\lambda)$, the probability of the observed sequence, given the model?

The main difference from decoding is that we sum instead of using max.

ASR Pipeline - Acoustic Model - Training

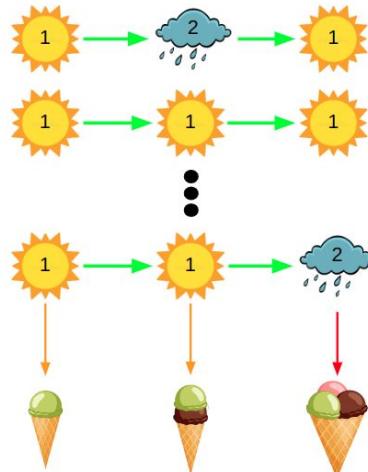
Evaluation solves the followings:

Given the observation sequence $O = \{o_1, o_2, \dots, o_T\}$, and an HMM model $\lambda = \{A, B\}$, how do we efficiently compute $P(O|\lambda)$, the probability of the observed sequence, given the model?

Computing the likelihood of observing 1, 2, 3

$$P(O) = \sum_Q P(O, Q) = \sum_Q P(O|Q)P(Q)$$

$$P(O, Q) = P(O|Q) \times P(Q) = \prod_{i=1}^n P(o_i|q_i) \times \prod_{i=1}^n P(q_i|q_{i-1})$$



ASR Pipeline - Acoustic Model - Training

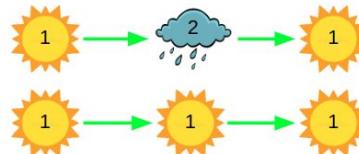
Evaluation solves the followings:

Given the observation sequence $O = \{o_1, o_2, \dots, o_T\}$, and an HMM model $\lambda = \{A, B\}$, how do we efficiently compute $P(O|\lambda)$, the probability of the observed sequence, given the model?

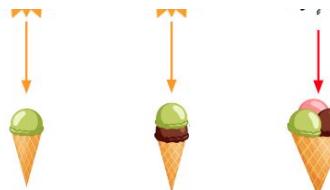
Computing the likelihood of observing 1, 2, 3

$$P(O) = \sum_Q P(O, Q) = \sum_Q P(O|Q)P(Q)$$

$$P(O, Q) = P(O|Q) \times P(Q) = \prod_{i=1}^n P(o_i|q_i) \times \prod_{i=1}^n P(q_i|q_{i-1})$$



Using the naive approach is exhaustive



ASR Pipeline - Acoustic Model - Training

Evaluation solves the followings:

Given the observation sequence $O = \{o_1, o_2, \dots, o_T\}$, and an HMM model $\lambda = \{A, B\}$, how do we efficiently compute $P(O|\lambda)$, the probability of the observed sequence, given the model?

Solution: Forward algorithm

Apply **dynamic programming** to compute **efficiently** the likelihood of the observation sequence by **summing** over all possible hidden state sequences

ASR Pipeline - Acoustic Model - Training

Evaluation solves the followings:

Given the observation sequence $O = \{o_1, o_2, \dots, o_T\}$, and an HMM model $\lambda = \{A, B\}$, how do we efficiently compute $P(O|\lambda)$, the probability of the observed sequence, given the model?

Forward algorithm

$$\alpha_t(j) = P(o_1, o_2 \dots o_t, q_t = j | \lambda)$$

Each cell of the forward algorithm trellis, $\alpha_t(j)$, **Represents the probability of being in state j After seeing the first t observations Given the automaton (regardless the path that it took- assuming all paths)**

ASR Pipeline - Acoustic Model - Training

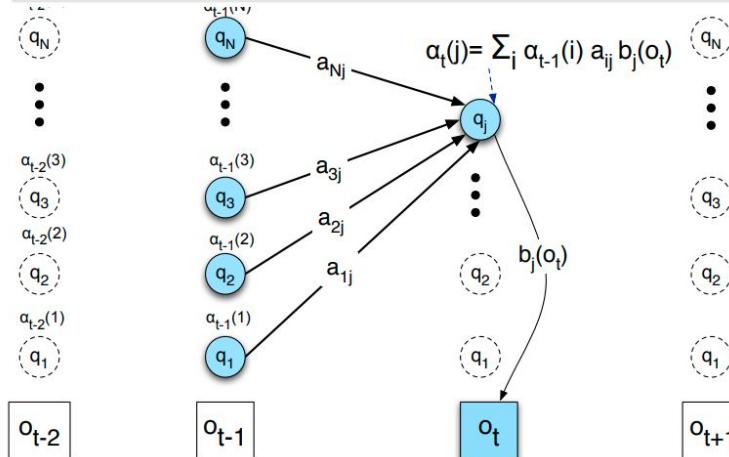
Evaluation solves the followings:

Given the observation sequence $O = \{o_1, o_2, \dots, o_T\}$, and an HMM model $\lambda = \{A, B\}$, how do we efficiently compute the forward path probability?

Forward algorithm

Each cell of the
being in state
(regardless the

$\alpha_{t-1}(i)$	the previous forward path probability from the previous time step
a_{ij}	the transition probability from previous state q_i to current state q_j
$b_j(o_t)$	the state observation likelihood of the observation symbol o_t given the current state j



forward path probability of the sequence, given the

the probability of
even the automaton

ASR Pipeline - Acoustic Model - Training

Evaluation solves the followings:

Given the observation sequence $O = \{o_1, o_2, \dots, o_T\}$, and an HMM model $\lambda = \{A, B\}$, how do we efficiently compute $P(O|\lambda)$, the probability of the observed sequence, given the model?

The Forward Algorithm

1. Initialization

$$\alpha_1(j) = a_{0j} b_j(o_1) \quad 1 \leq j \leq N$$

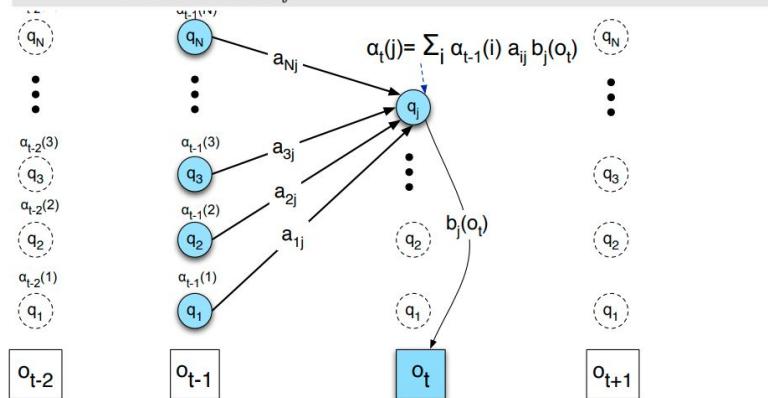
2. Recursion (since state 0 and F are non-emitting):

$$\alpha_t(j) = \sum_{i=1}^N \alpha_{t-1}(i) a_{ij} b_j(o_t); \quad 1 \leq j \leq N, 1 < t \leq T$$

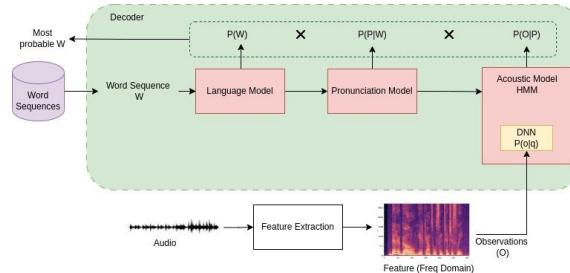
3. Termination:

$$P(O|\lambda) = \alpha_T(q_F) = \sum_{i=1}^N \alpha_T(i) a_{iF}$$

$\alpha_{t-1}(i)$	the previous forward path probability from the previous time step
a_{ij}	the transition probability from previous state q_i to current state q_j
$b_j(o_t)$	the state observation likelihood of the observation symbol o_t given the current state j



ASR Pipeline - Acoustic Model - Training



$$P(O|P)$$

HMM solves the following tasks:

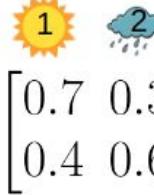
- **Evaluation (likelihood):** Given the observation sequence $O = \{o_1, o_2, \dots, o_T\}$, and an HMM model $\lambda = \{A, B\}$, how do we efficiently compute $P(O|\lambda)$, the probability of the observed sequence, given the model?
- **Decoding:** Given the observation sequence $O = \{o_1, o_2, \dots, o_T\}$, and an HMM model $\lambda = \{A, B\}$, how do we choose a corresponding state sequence $Q = \{q_1, q_2, \dots, q_T\}$ that is optimal, i.e., best explains the observations?
- **Learning:** How do we learn the parameters $\lambda = \{A, B\}$ to maximize $P(O|\lambda)$

ASR Pipeline - Acoustic Model - Training

Learning

How do we learn the parameters $\lambda = \{A, B\}$ to maximize $P(O|\lambda)$

$$O = \{o_1, \dots, o_N\}; \lambda = \{A, B, \pi\}?$$

α_{ij}	π	B1	B2
  $\begin{bmatrix} 0.7 & 0.3 \\ 0.4 & 0.6 \end{bmatrix}$	$\begin{bmatrix} 0.2 \\ 0.8 \end{bmatrix}$	$\begin{bmatrix} P(\text{ice cream} 1) \\ P(\text{ice cream} 2) \\ P(\text{apple} 1) \end{bmatrix} = \begin{bmatrix} 0.2 \\ 0.4 \\ 0.4 \end{bmatrix}$	$\begin{bmatrix} P(\text{ice cream} 2) \\ P(\text{ice cream} 1) \\ P(\text{apple} 2) \end{bmatrix} = \begin{bmatrix} 0.5 \\ 0.4 \\ 0.1 \end{bmatrix}$

ASR Pipeline - Acoustic Model - Training

Learning

How do we learn the parameters $\lambda = \{A, B\}$ to maximize $P(O|\lambda)$

$$O = \{o_1, \dots, o_N\}; \lambda = \{A, B, \pi\}?$$

Define forward and backward variables:

- Forward variable: the probability of seeing the observations o_1, o_2, \dots, o_t and being in state i at time t

$$\alpha_t(i) = P(o_1, o_2, \dots, o_t, q_t = S_i | \lambda)$$

- Backward variable: the probability of the ending partial sequence $q_{t+1}, q_{t+2}, \dots, q_T$ starting state i at time t .

$$\beta_t(i) = P(q_{t+1}, q_{t+2}, \dots, q_T | q_t = S_i, \lambda)$$

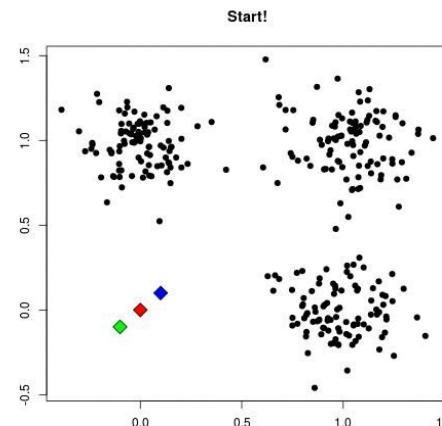
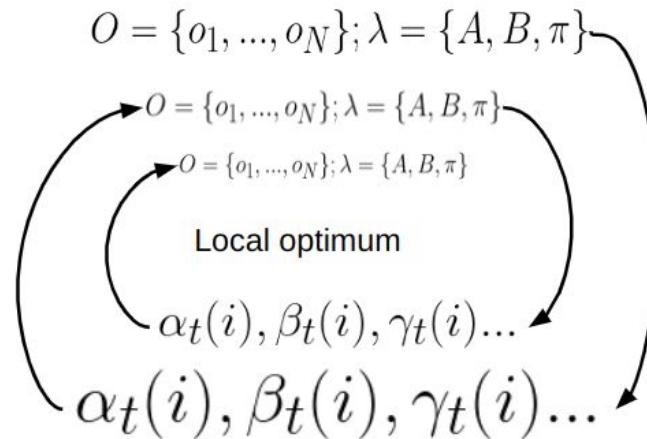
ASR Pipeline - Acoustic Model - Training

Learning

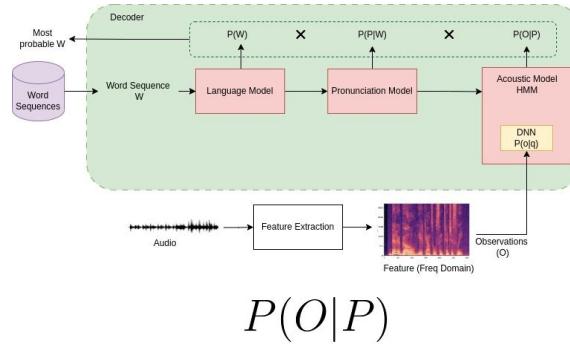
How do we learn the parameters $\lambda = \{A, B\}$ to maximize $P(O|\lambda)$

$$O = \{o_1, \dots, o_N\}; \lambda = \{A, B, \pi\}?$$

Baum Welch Algorithm - EM



ASR Pipeline - Acoustic Model - Training



$$P(O|P)$$

HMM solves the following tasks:

- **Evaluation (likelihood):** Given the observation sequence $O = \{o_1, o_2, \dots, o_T\}$, and an HMM model $\lambda = \{A, B\}$, how do we efficiently compute $P(O|\lambda)$, the probability of observing the sequence O given the model λ ?

HMM for speech recognition

Given a model $\lambda = \{A, B\}$, how do we choose a corresponding state sequence $Q = \{q_1, q_2, \dots, q_T\}$ that is optimal, i.e., best explains the observations?

- **Learning:** How do we learn the parameters $\lambda = \{A, B\}$ to maximize $P(O|\lambda)$

ASR Pipeline - Acoustic Model - Notes

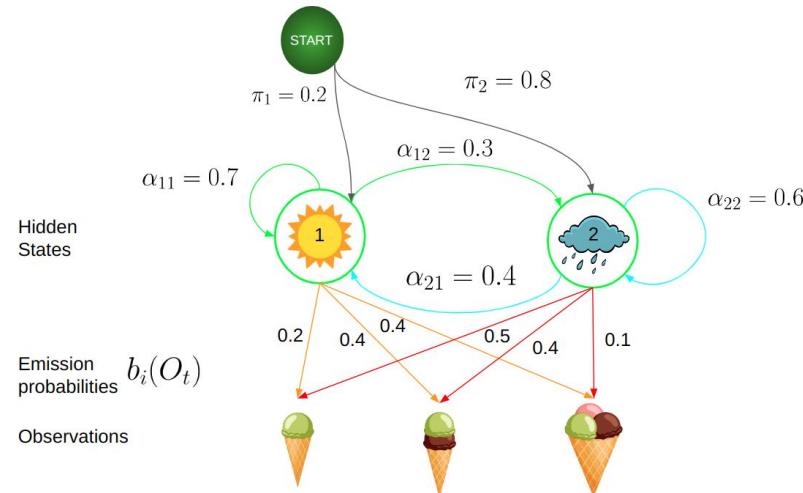
HMM for speech recognition

For a given phone/word string W, we will evaluate $P(O|W)$:

- Observation sequence O- a series of MFCC vectors
- Hidden states W- phones and words

A and B defines:

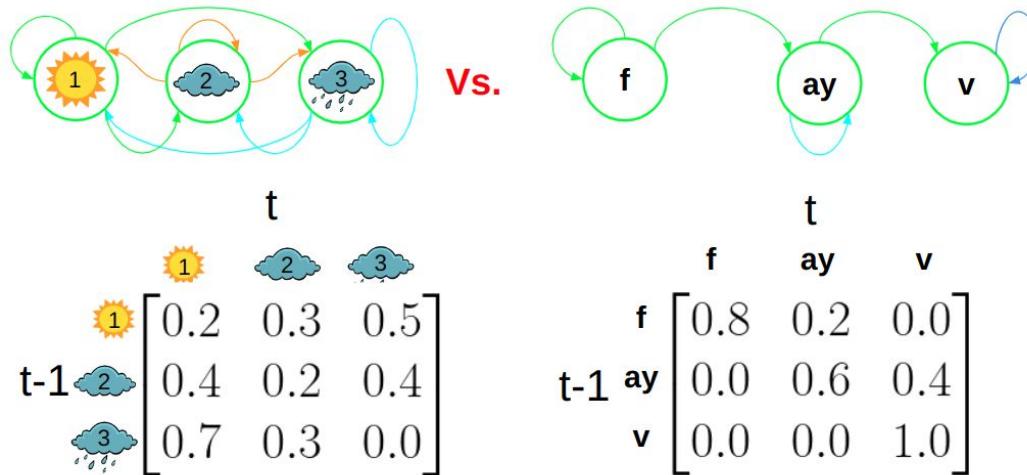
- A- **what** will be **said** (defines the phoneme sequence (hidden states) and the transition probabilities).
- B- **How** it will be **pronounced**



ASR Pipeline - Acoustic Model - Notes

HMM for speech recognition

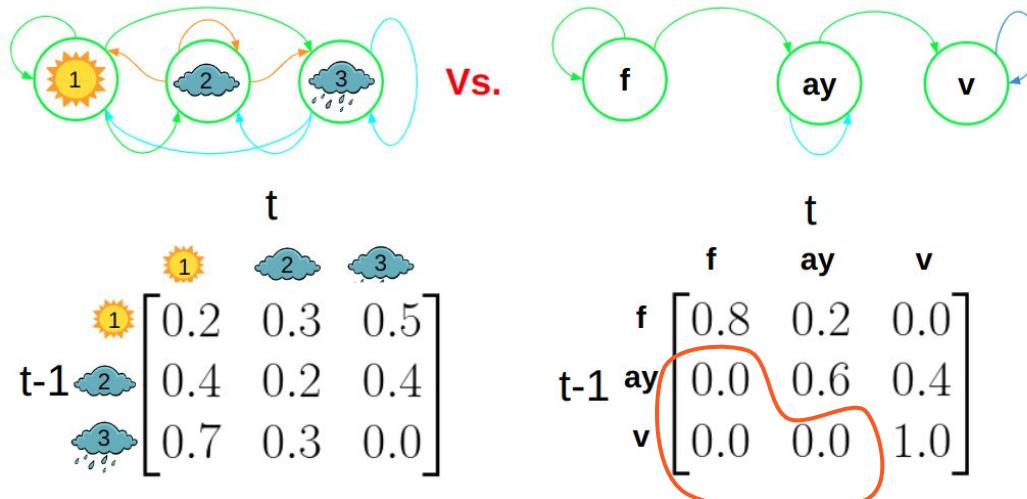
For a given phone/word string W, we will evaluate $P(O|W)$



ASR Pipeline - Acoustic Model - Notes

HMM for speech recognition

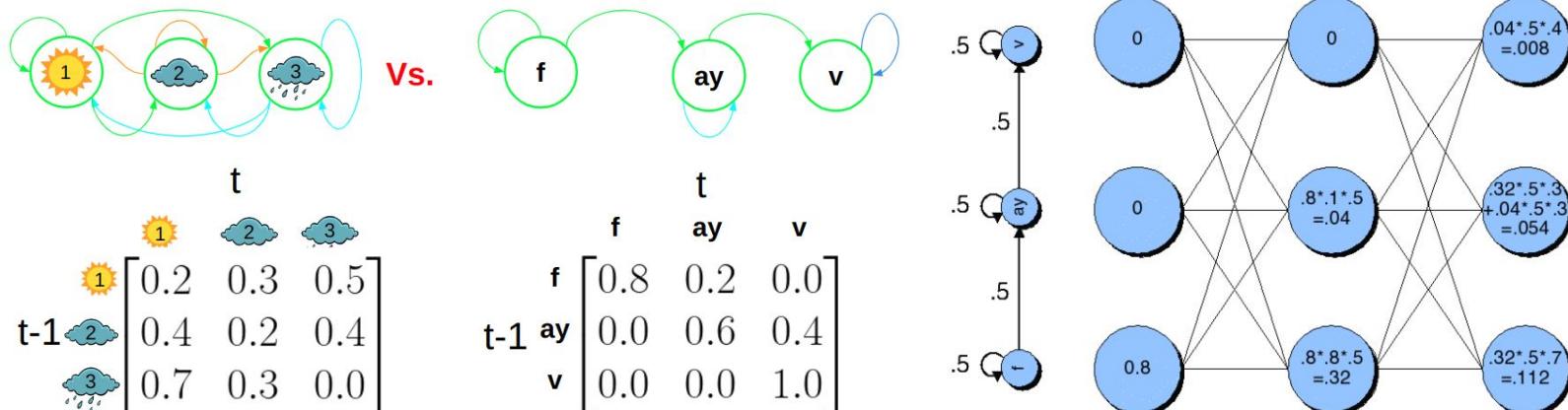
For a given phone/word string W, we will evaluate $P(O|W)$ - evaluation/likelihood.



ASR Pipeline - Acoustic Model - Notes

HMM for speech recognition

For a given phone/word string W, we will evaluate $P(O|W)$:



ASR Pipeline - Acoustic Model - Data

Alignment

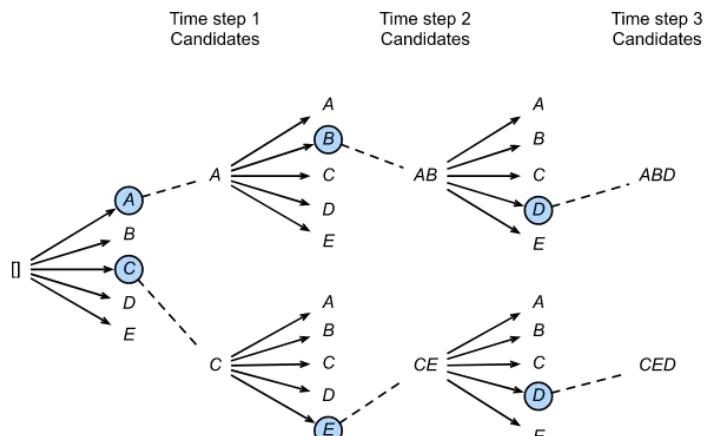
Because aligned data is expensive and hard to achieve (need phoneme level annotations), how can we obtain more training data?

- Training an acoustic model + HMM on a phone-level annotated data
- Take data with utterance level transcription
- Apply Viterbi (force-align) to find the best alignment while enforcing the ground truth text.
- Train an improved acoustic model + HMM using this pseudo labels.

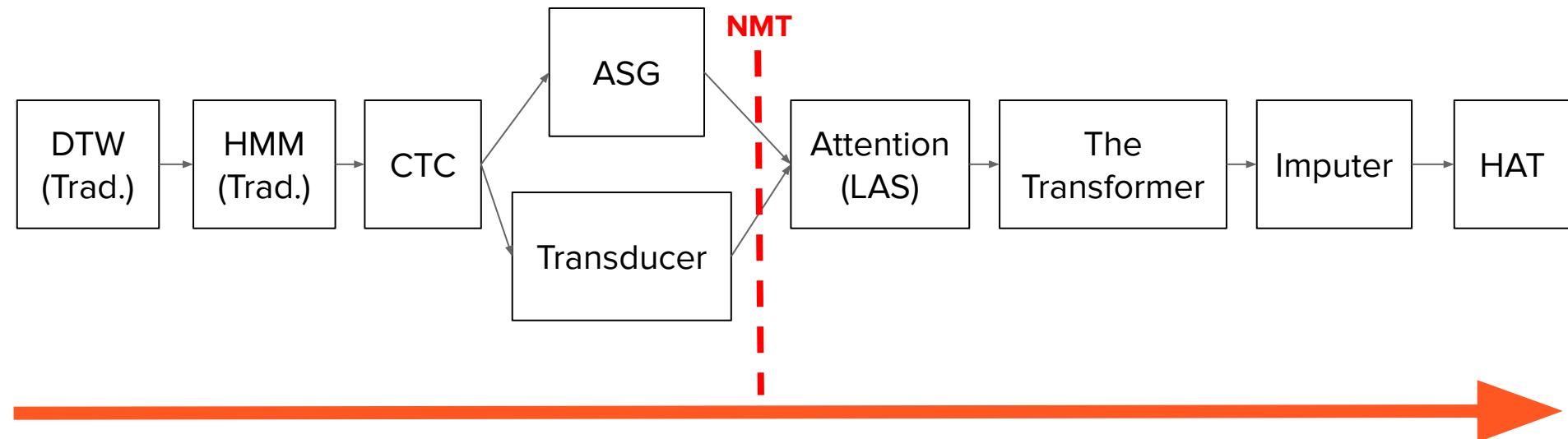
ASR Pipeline - Acoustic Model - Inference

A word on the search

- In practice, it's not practical to search over all word sequences. The search space is too large and exhaustive.
- Solution: beam search



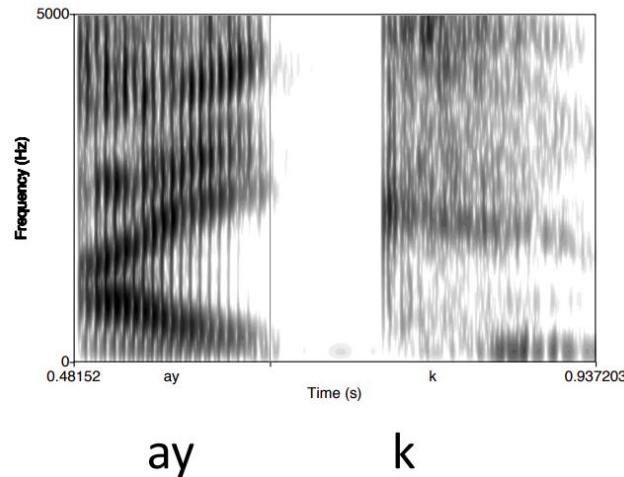
ASR Models - Overview



ASR Models - GMM- HMM

GMM-HMM Drawbacks

- Phonemes are not homogeneous
 - Each phone is divided into 3 substates (begin, middle, final)



ASR Models - GMM- HMM

GMM-HMM Drawbacks

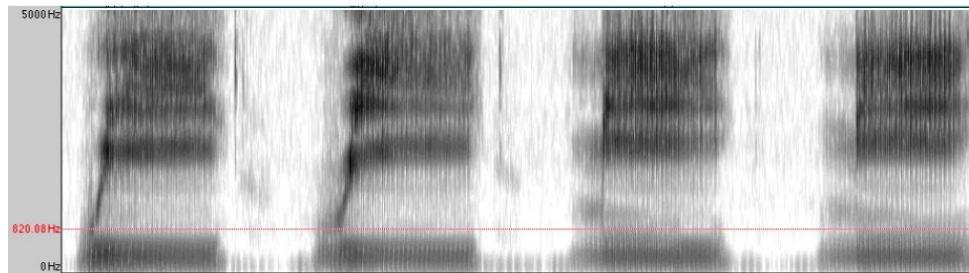
- Phonemes are not homogeneous
 - Each phone is divided into 3 substates (begin, middle, final)
- The strongest factor affecting phonetic variability is the neighboring phone
 - Context-Independent (CI) Vs Context-Dependent (CD) phones: Triphones

w iy

r iy

m iy

n iy



ASR Models - DNN-HMM

- Acoustic Model Is A **Single** DNN For All States
- Improves Performances:
 - From 27% WER to 18% WER on SwitchBoard. Deep Neural Networks for Acoustic Modeling in Speech Recognition, Hinton et al (2012).
 - For comparison: current SOTA using End2End method is ~5% WER.