

4 Speech Recognition

Chapter 4 Speech Recognition

- 4.1 Overview of Speech Recognition
- 4.2 Acoustic Model
- 4.3 Language Model
- 4.4 Decoding Algorithm for Speech Recognition
- 4.5 Prospect of Speech Recognition

4.1 Overview of Speech Recognition

- 4.1.1 Basic Concepts of Speech Recognition
- 4.1.2 History and Main Contents of Speech Recognition

What is speech recognition?

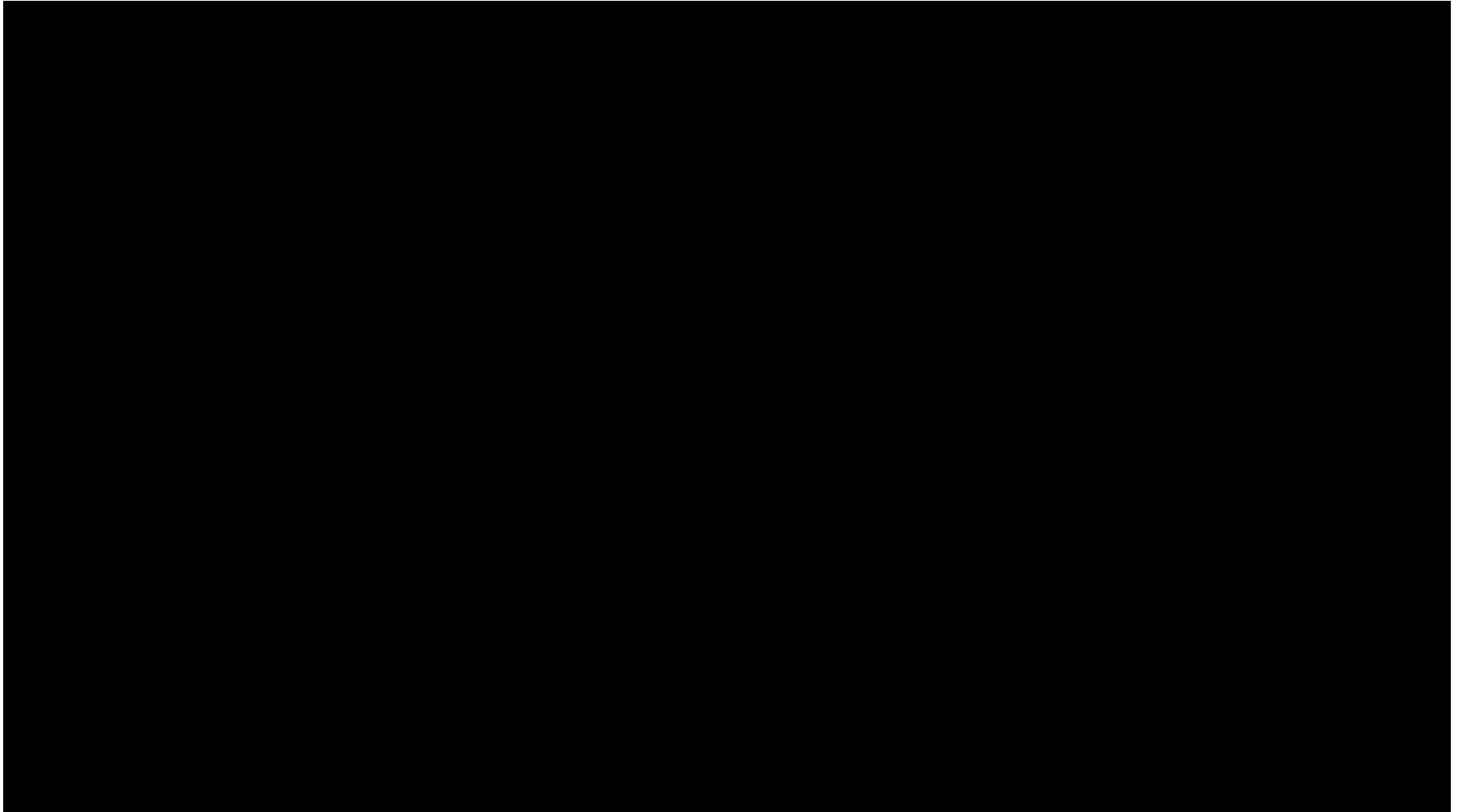
Speech Recognition

Conversion of spoken words into text (words, syllables, phones, etc.), also known as “Automatic Speech Recognition” (ASR).



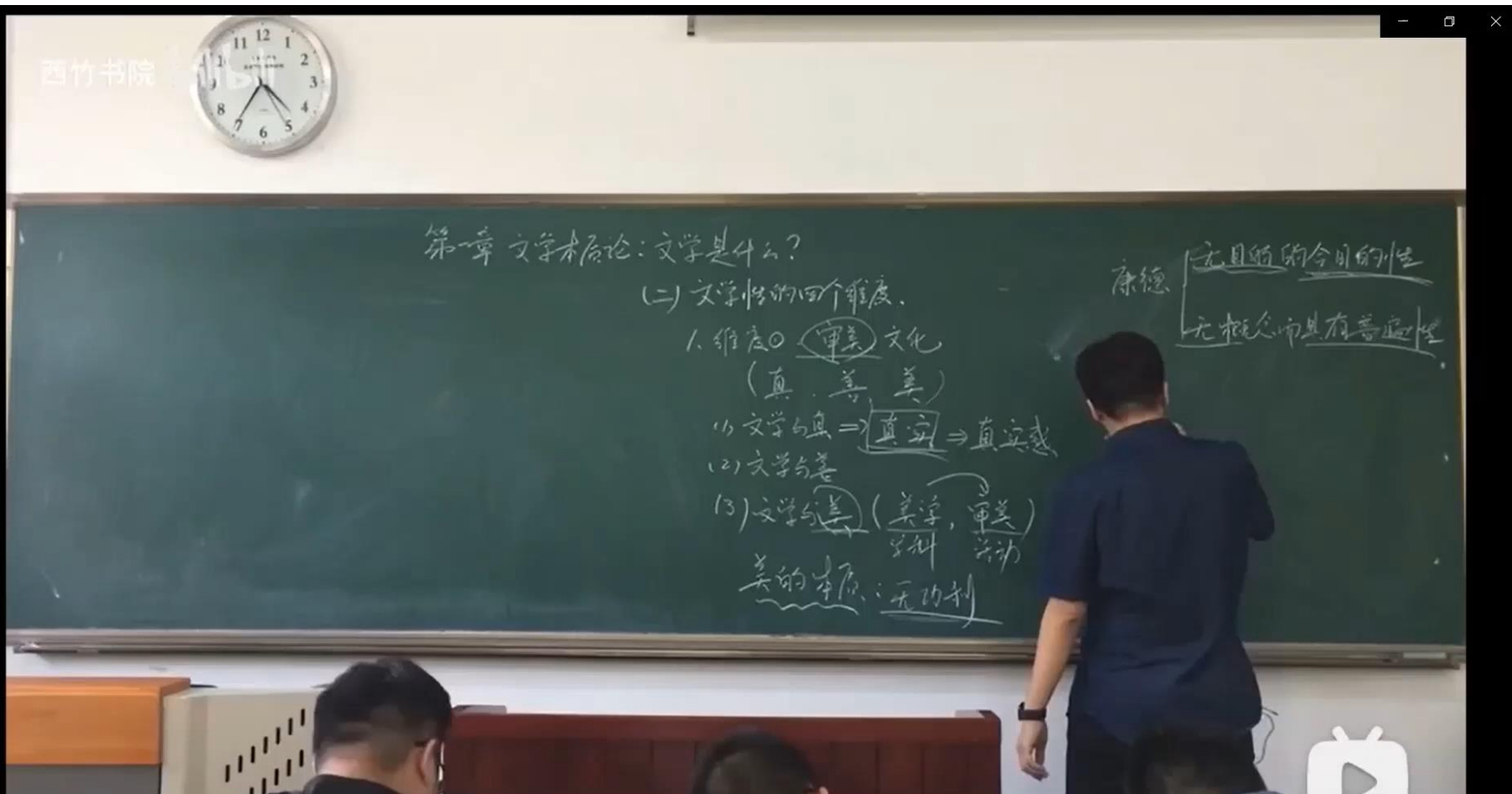
What is speech recognition?

The recognized text can be

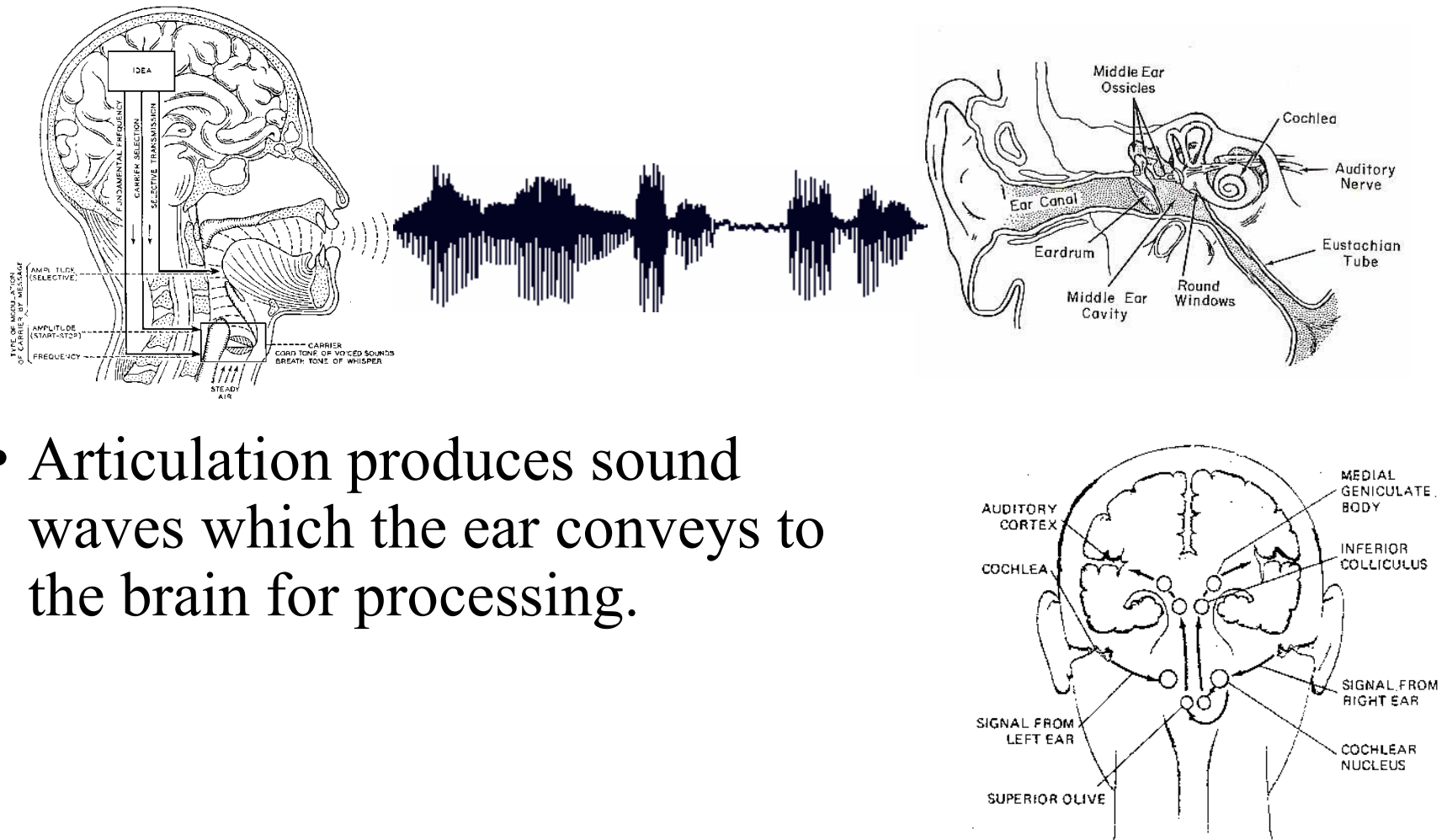


What is speech recognition?

The recognized text can be



How do humans recognize speech?

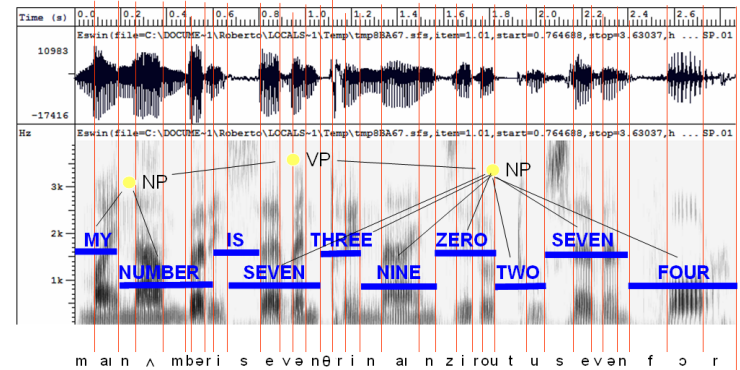
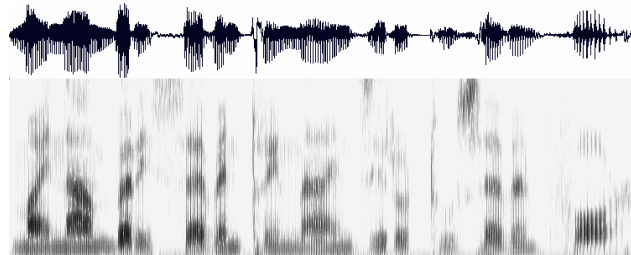


How might computers do it?



Acoustic waveform

Acoustic signal



- Digitization
- Acoustic analysis of the speech signal
- Linguistic interpretation

Speech recognition

Types of Speech Recognition

- Different classes based on types of utterances they are able to recognize
 - 1. Isolated Words
 - “Listen/Not-Listen” states “打开空调”
 - 2. Connected Words
 - Connected digit recognition
“139-1234-5678” (Phone Number)
 - “run-together”
“现在 我 简单 介绍 语音 识别 的 基本 概念。”
 - 3. Continuous Speech Large Vocabulary Continuous Speech Recognition (LVCSR)
 - Natural speech
“现在我简单介绍语音识别的基本概念。”

Types of Speech Recognition

- Different classes based on different vocabulary size
 - 1. Small vocabulary
 - About 10-1000 words
 - 2. Medium vocabulary
 - About 1000-10000 words
 - 3. Large vocabulary
 - More than 10000 words

Large Vocabulary Continuous Speech Recognition (LVCSR)

Types of Speech Recognition

- Different classes based on speaker dependency
 - 1. Speaker dependent
 - Training/Adapting model using user's speech
 - Time consuming for individual speaker/user.
 - 2. Speaker independent
 - Using same model for all speakers/users
 - Need various speakers and styles for training data

ASR Metric: Word Error Rate (WER)

- How to evaluate the word string output by a speech recognizer?

$$\text{Word Error Rate} = \frac{\text{Insertions} + \text{Substitutions} + \text{Deletions}}{\text{Total Words in Correct Transcript}} * 100\%$$

Alignment example:

REF: portable **** PHONE UPSTAIRS last night so

HYP: portable FORM OF STORES last night so

Eval I S S

$$\text{WER} = (1+2+0)/6 * 100\% = 50\%$$

4.1 Overview of Speech Recognition

- 4.1.1 Basic Concepts of Speech Recognition
- 4.1.2 History and Main Contents of Speech Recognition

Trends of ASR

Modeling approaches:

Before mid 70's	Mid 70's – mid 80's	After mid 80's
Heuristic	Template matching	Mathematical
Rule-based and declarative	Deterministic and data-driven	Probabilistic and data-driven

Pattern Recognition Approach

Isolated word

Isolated word

Connected speech

LVCSR

HMM-based

DL-based

Rule-based approach

- Use knowledge of phonetics (语音学) and linguistics (语言学) to guide search process
- Rules express everything (anything) that might help to decode:
 - Phonetics, phonology (音韵学), phonotactics (音位结构学)
 - Syntax (语法)
 - Pragmatics (语用学)

Rule-based approach

- Typical approach is based on “blackboard” architecture:
 - At each decision point, lay out the possibilities
 - Apply rules to determine which sequences are permitted
- Poor performance due to
 - Difficulty to express rules
 - Difficulty to make rules interact
 - Difficulty to know how to improve the system

s k
h
p
t i:
iə
I
s
tʃ
h
s

Pattern Recognition Approach

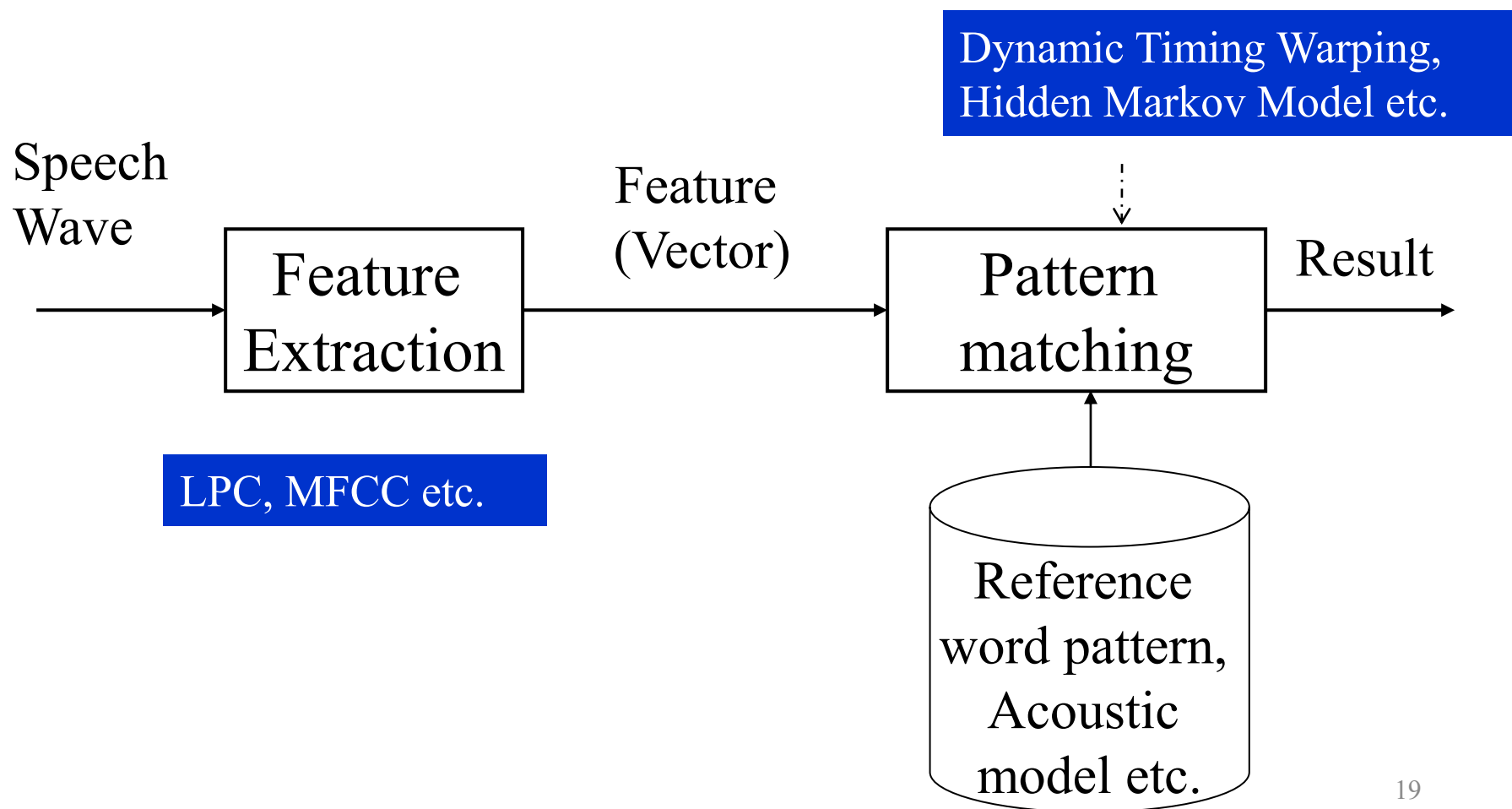
- 2 steps:
 - Pattern Training
 - Pattern Comparison
- Forms:
 - Speech Template **Deterministic**
 - Statistical Model (HMM, DNN) **Probabilistic**
- Goal to determine identity of unknown speech according to how well patterns match

Methods in Pattern Comparison Approach

- Template Based Approach
 - Patterns stored as dictionary of words
 - Match unknown utterance with reference templates
 - Select best matching pattern
- Statistical Approach (HMM、 DNN)
 - Probabilistic Models
 - Uncertainty and Incompleteness

Diagram of data-driven ASR

Example of isolated word recognition



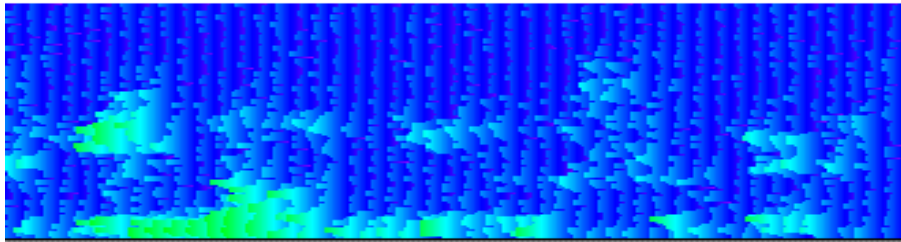
Template-based approach

Template-based approach

- Store examples of units (words, phonemes), then find the example that most closely fits the input
- Extract features from speech signal, then it's “just” a complex similarity matching problem, using solutions developed for all sorts of applications
- OK for discrete utterances, and a single user

Pattern matching with different length

Input
pattern



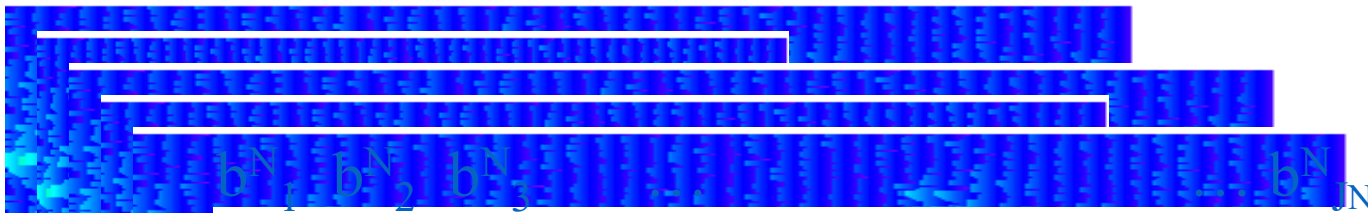
a_1 a_2 a_3

\dots a_I

b^1_1 b^1_2 b^1_3

\dots $b^1_{J_1}$

Reference
Pattern 1



Reference
Pattern N

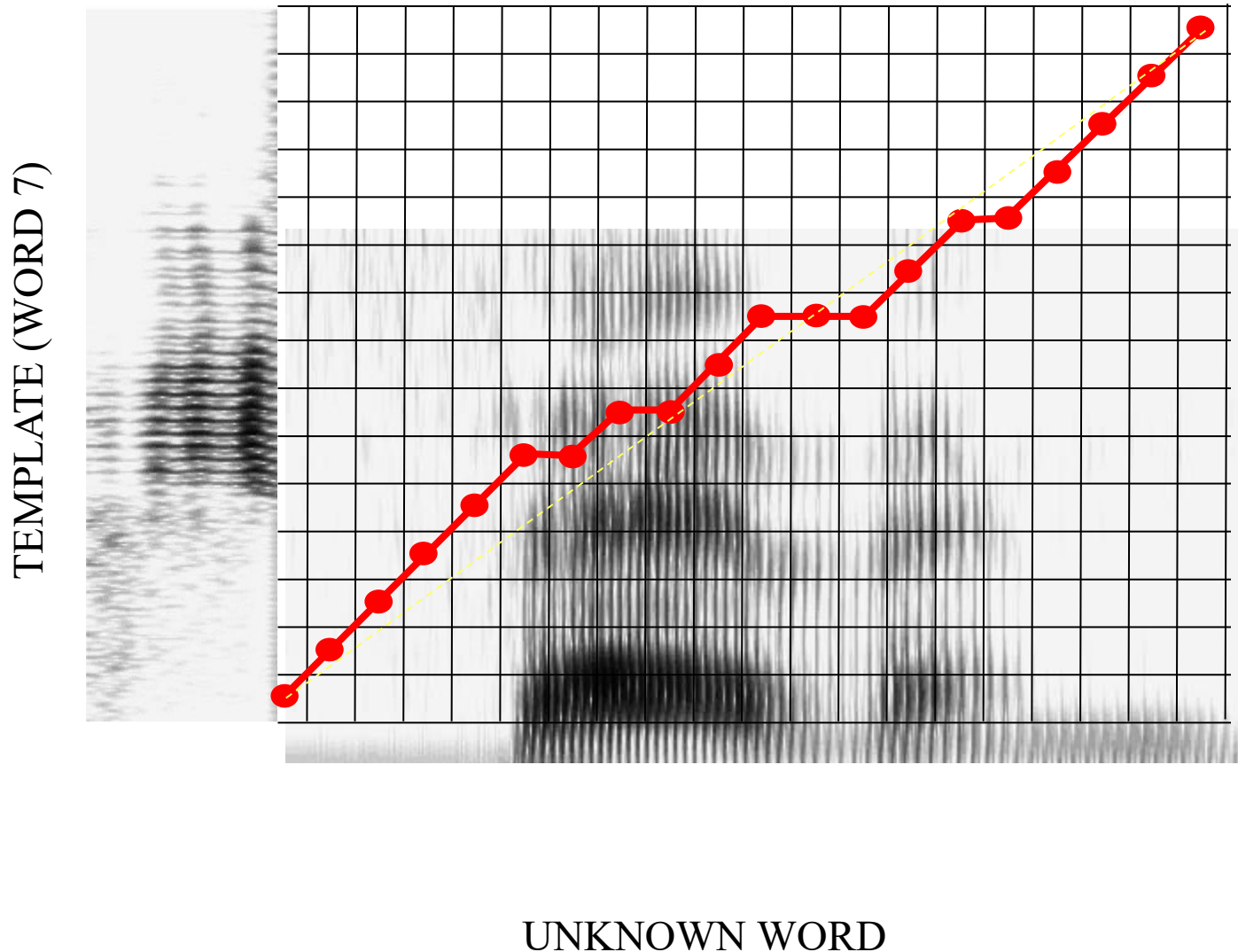
Different length



How to compare?

How about linear
alignment?

Dynamic Time Warping (DTW) /Dynamic Programming (DP) Matching

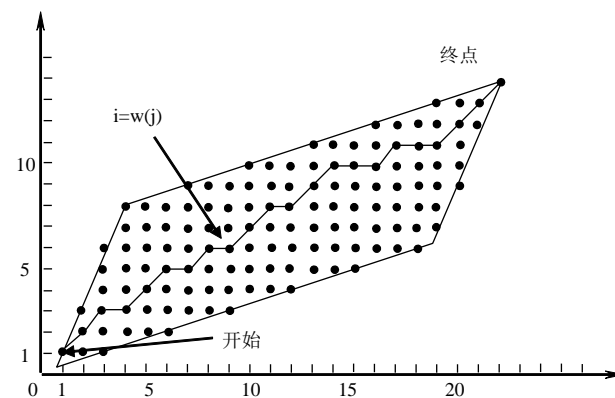


DTW Algorithm

(1) Initialization:

$$i = j = 1, \quad g(1, 1) = 2d(x_1, y_1)$$

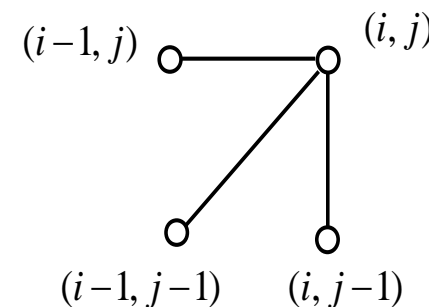
$$g(i, j) = \begin{cases} 0 & (i, j) \in \text{Reg} \\ \text{huge} & (i, j) \notin \text{Reg} \end{cases}$$



Constrained Region

$g(i, j)$: minimum partial accumulated distance

DTW Algorithm



(2) Recursion:

Minimum partial

accumulated distance

Local distance

Path constraints

$$g(i, j) = \min \left\{ \begin{array}{l} g(i-1, j) + d(\mathbf{x}_i, \mathbf{y}_j)W(1) \\ g(i-1, j-1) + d(\mathbf{x}_i, \mathbf{y}_j)W(2) \\ g(i, j-1) + d(\mathbf{x}_i, \mathbf{y}_j)W(3) \end{array} \right\}$$

$$i = 2, 3, \dots, I; j = 2, 3, \dots, J; (i, j) \in \text{Reg}$$

$$W(1) = W(3) = 1, \quad W(2) = 2$$

(3) Termination:

$$\frac{g(I, J)}{\sum W}$$

Problem of template-based approach

- Hard to distinguish very similar templates
- And quickly degrades when acoustic conditions between input and templates are different
- Therefore needs techniques to mitigate these degradations:
 - More subtle matching techniques
 - Multiple templates which are aggregated
- Taken together, these suggested ...

Statistic-based approach

Statistic-based approach

- Can be seen as extension of template-based approach, using more powerful mathematical and statistical models
- Sometimes seen as “anti-linguistic” approach

Statistic-based approach

- Collect a large corpus of transcribed speech recordings
- Train the computer to learn the correspondences (“machine learning”)
- At run time, apply statistical processes to search through the space of all possible solutions, and pick the statistically most likely one

Noisy Channel in a Picture



- Search through space of all possible sentences
- Pick the one that is most probable given the waveform

Noisy Channel Model

- In speech recognition, you observe an acoustic signal ($A=a_1, \dots, a_n$) and you want to determine the most likely sequence of words ($W=w_1, \dots, w_n$): $P(W | A)$
- Assume that the acoustic signal (A) is already segmented
- $P(W | A)$ could be computed as

$$P(W | A) = \prod_{a_i} \max_{w_i} P(w_i | a_i)$$

- Problem: Finding the most likely word corresponding to an acoustic representation depends on the context

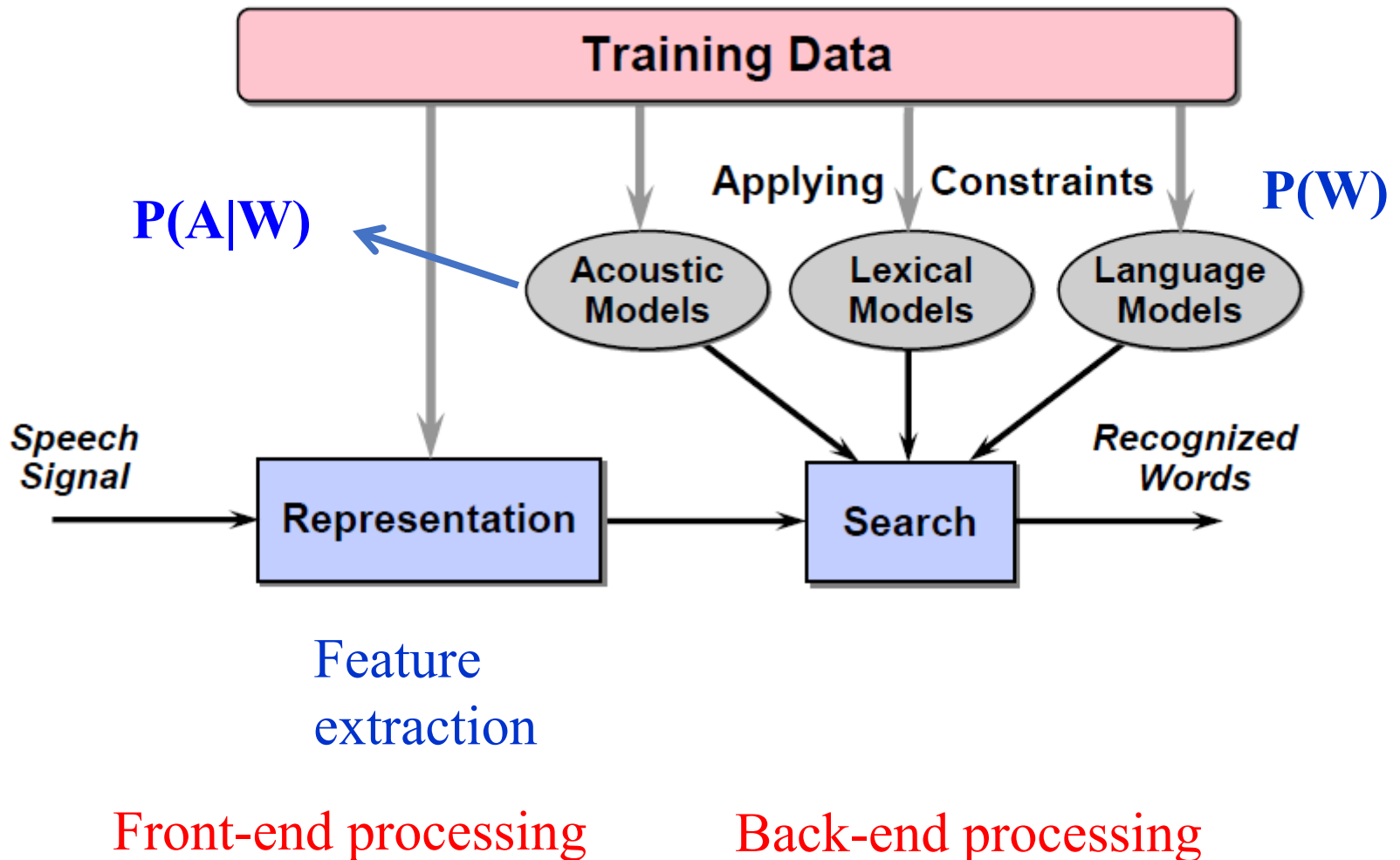
Noisy Channel Model

- Given a candidate sequence W we need to compute $P(W)$ and combine it with $P(W | A)$
- Applying Bayes' rule:

$$\arg \max_W P(W | A) = \arg \max_W \frac{P(A | W)P(W)}{P(A)}$$

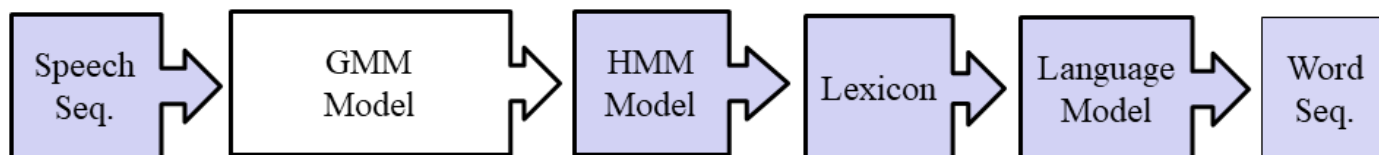
- The denominator $P(A)$ can be dropped, because it is constant for all W

Diagram of ASR system

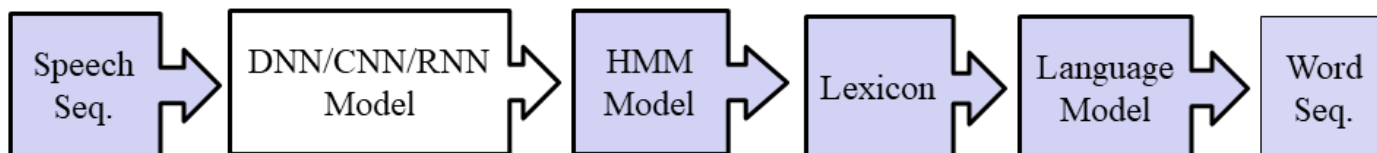


Evolution of speech recognition systems (Statistic model-based)

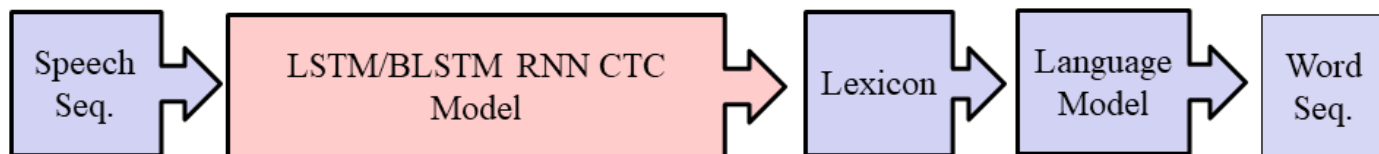
1990s-2009: The GMM-HMM hybrid system. (CU-HTK)



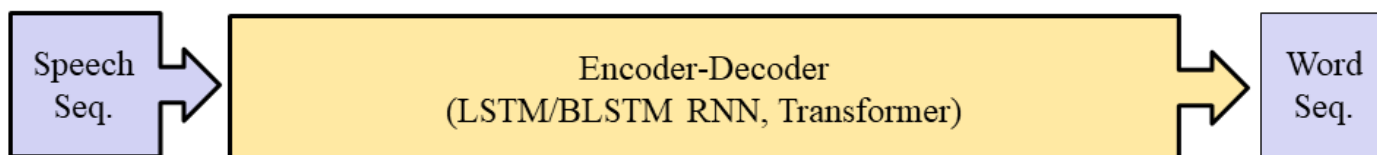
2009 The DNN-HMM hybrid system. (JHU-Kaldi, MS-CNTK)



2014: The CTC End-to-End system. (CMU-EESEN, Baidu-WarpCTC)

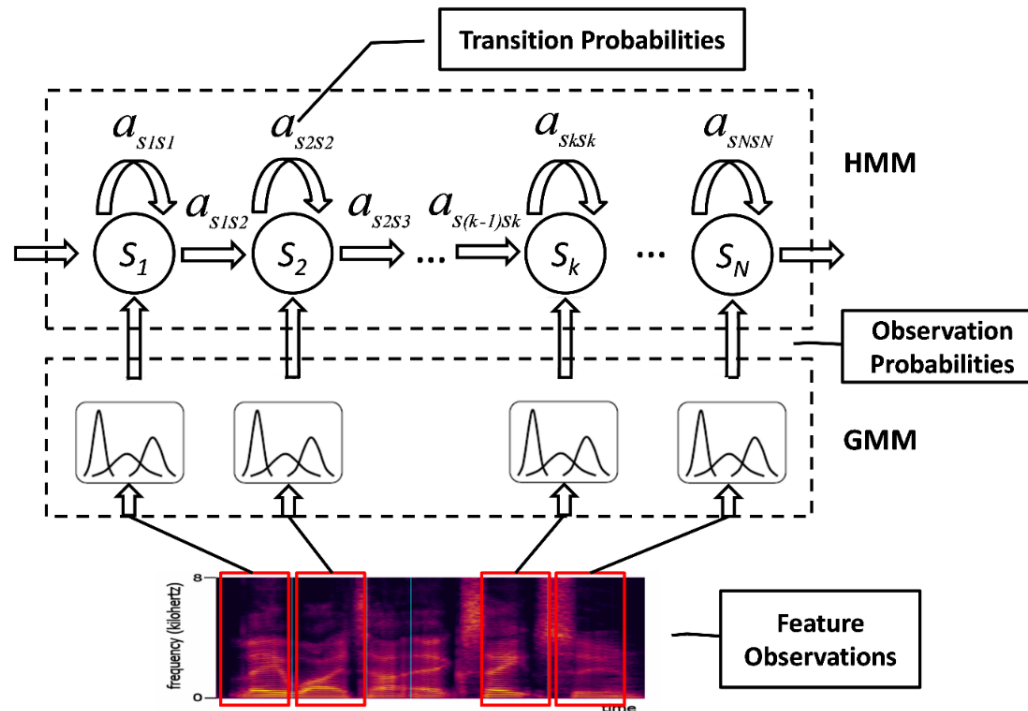


2016: The Encoder-Decoder End-to-End system.
(Google-LAS/Transformer, facebook-wav2letter, JHU-ESPNet)



Acoustic Model

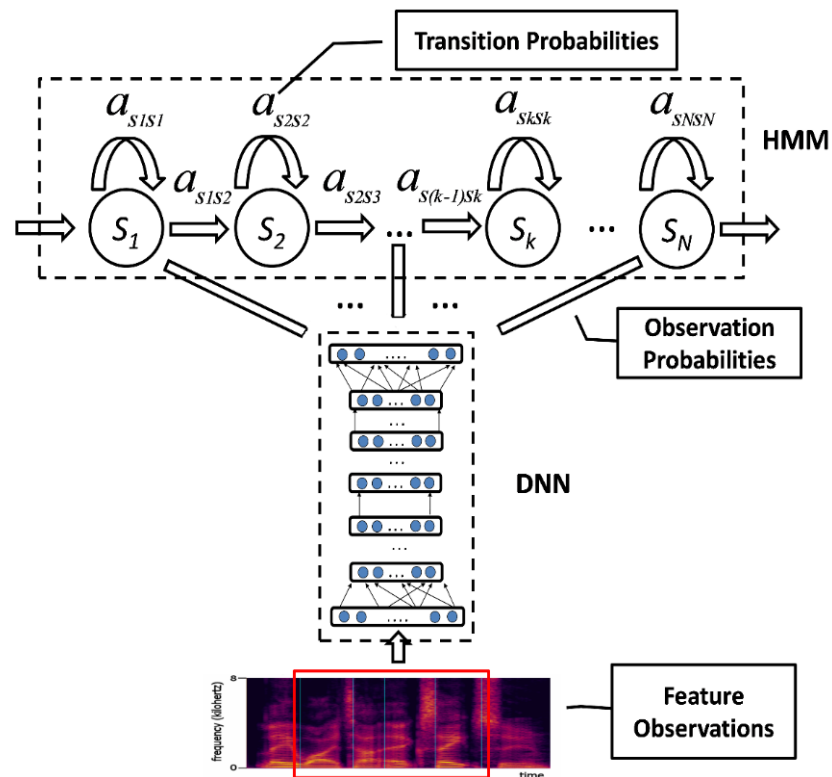
■ GMM-HMM



- L. R. Rabiner. A tutorial on hidden markov models and selected applications in speech recognition. Proc. IEEE, 77(2):257–286, 1989.

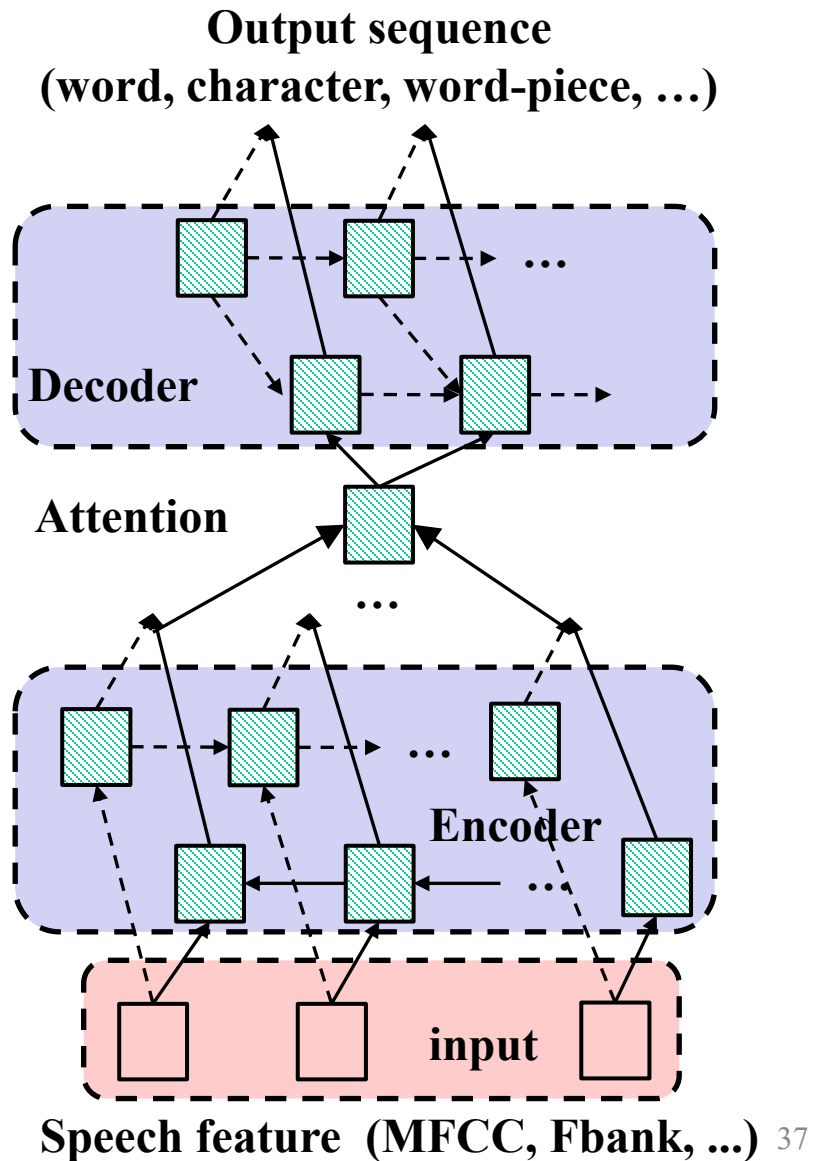
Acoustic Model

■ DNN-HMM



End-to-End Model

- Jointly learn acoustic model and language model



Language model

- Finite state automaton grammar
- **N-gram models:** (Models likelihood of word given previous word(s))
 - Build the model by calculating bigram or trigram probabilities from text training corpus
 - Smoothing issues
- Recurrent Neural Network-based model

Decoding Algorithm

- Calculation of best hypothesis of word sequence

$$P(W | Y) = \arg \max_{\{w_1^N\} \{t_1^N\}} \left\{ \sum_{n=1}^N \log(P_{acoust}(y_{t_{n-1}+1}^{t_n} | w_n)) \right. \\ \left. + \lambda \cdot \sum_{n=1}^N \log(P_{lang}(w_n | w_{n-1}) + \delta) \right\}$$

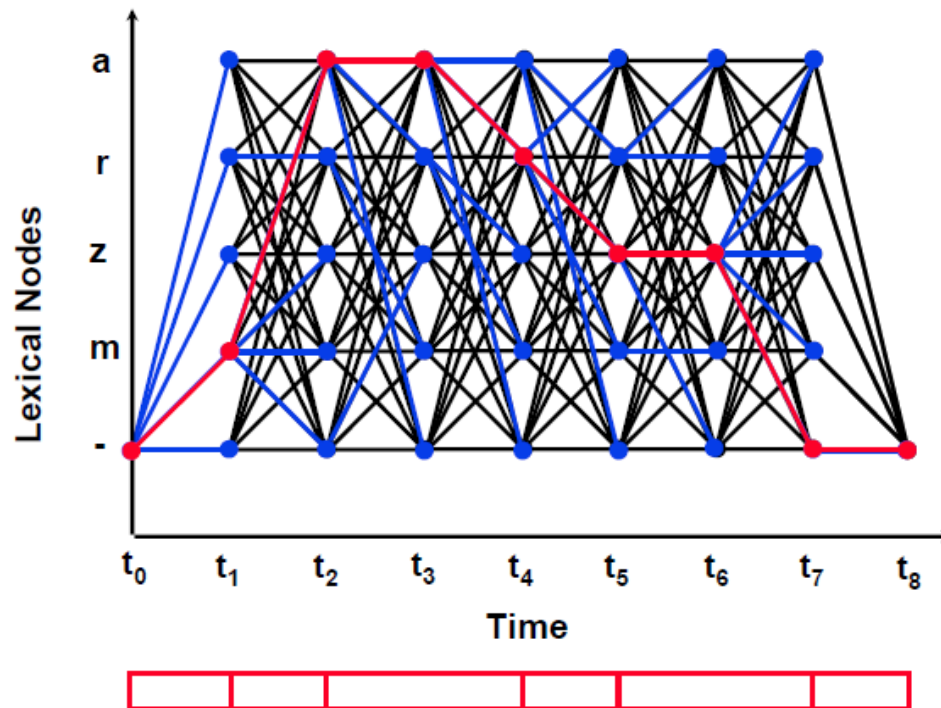
$P_{acoust}(y_{t_{n-1}+1}^{t_n} | w_n)$: likelihood of AM

$P_{lang}(w_n | w_{n-1})$: likelihood of LM

λ : Weighting δ : Penalty of insertion

Decoding Algorithm

- Viterbi Search



- WFST (Weighted finite state transducer)