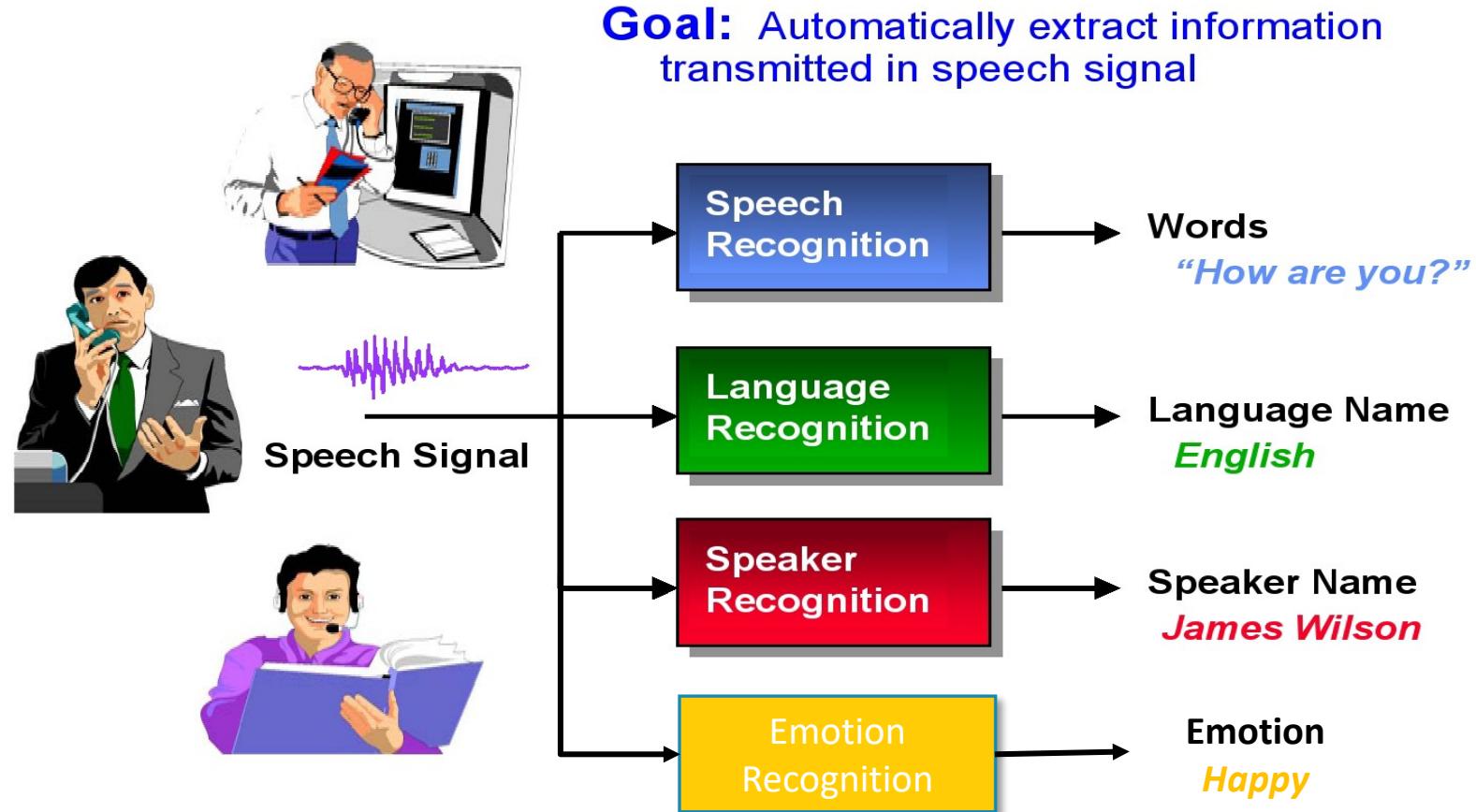


# Introduction to Speech Recognition

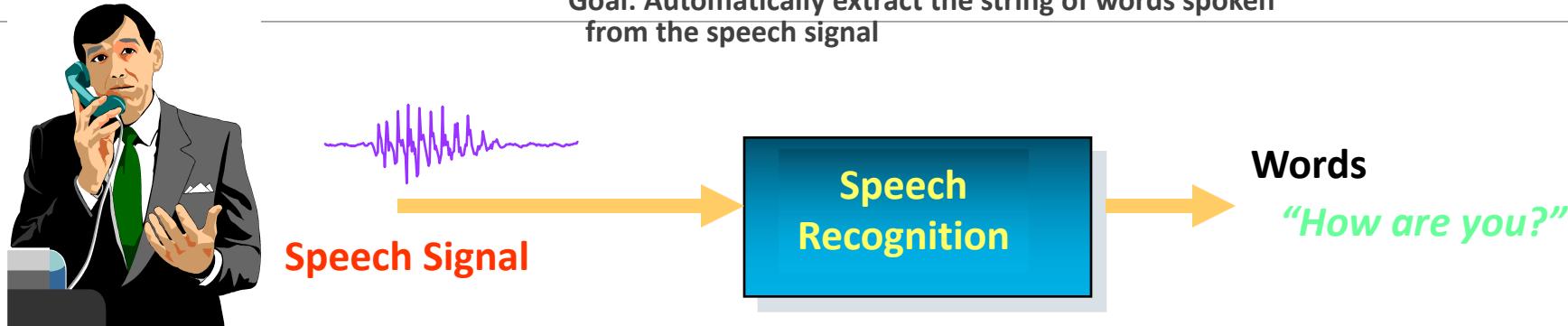
Yuan-Fu Liao

National Taipei University of Technology

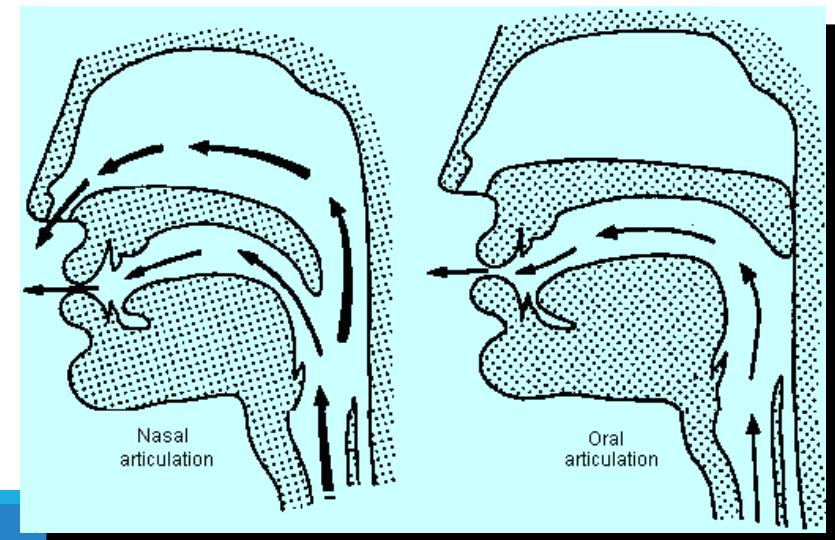
# Broad Objectives of Speech Recognition for Machines



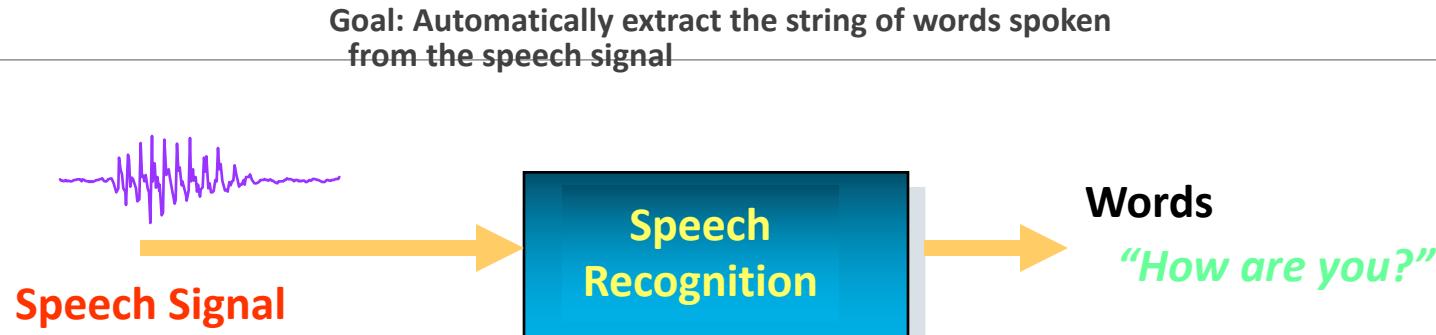
# Speech Recognition



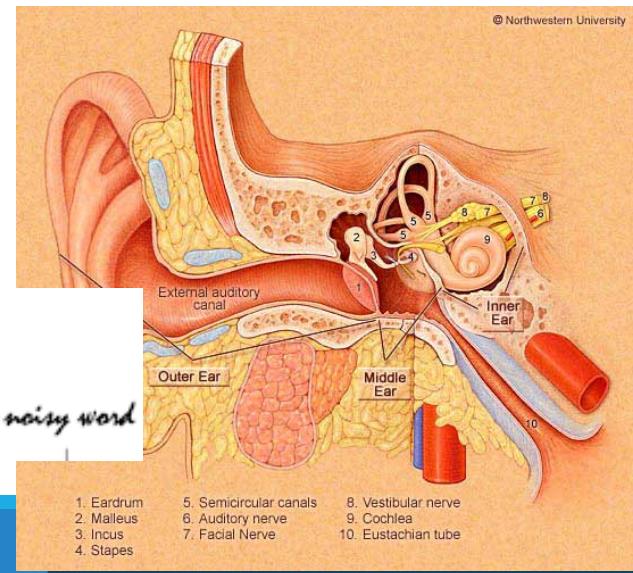
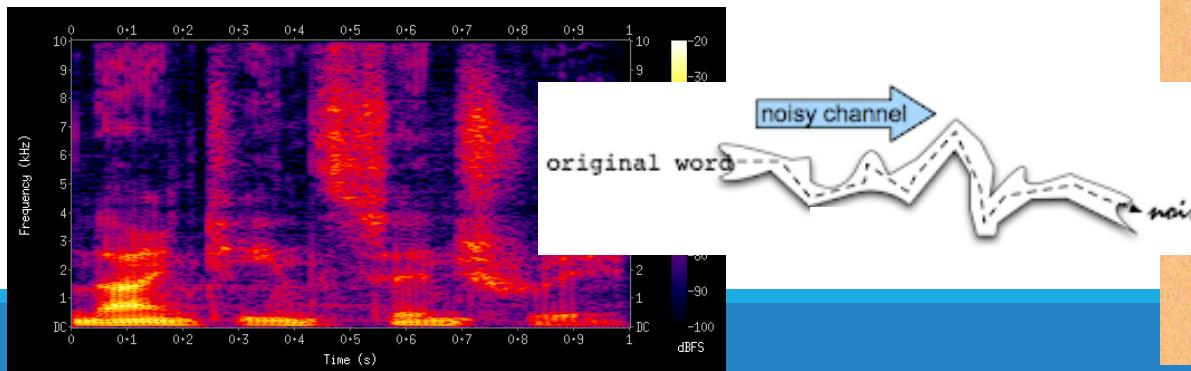
How is SPEECH produced?  
⇒ Characteristics of  
Acoustic Signal



# Speech Recognition



How is SPEECH perceived?  
=> Important Features



# Speech Recognition



Goal: Automatically extract the string of words spoken from the speech signal



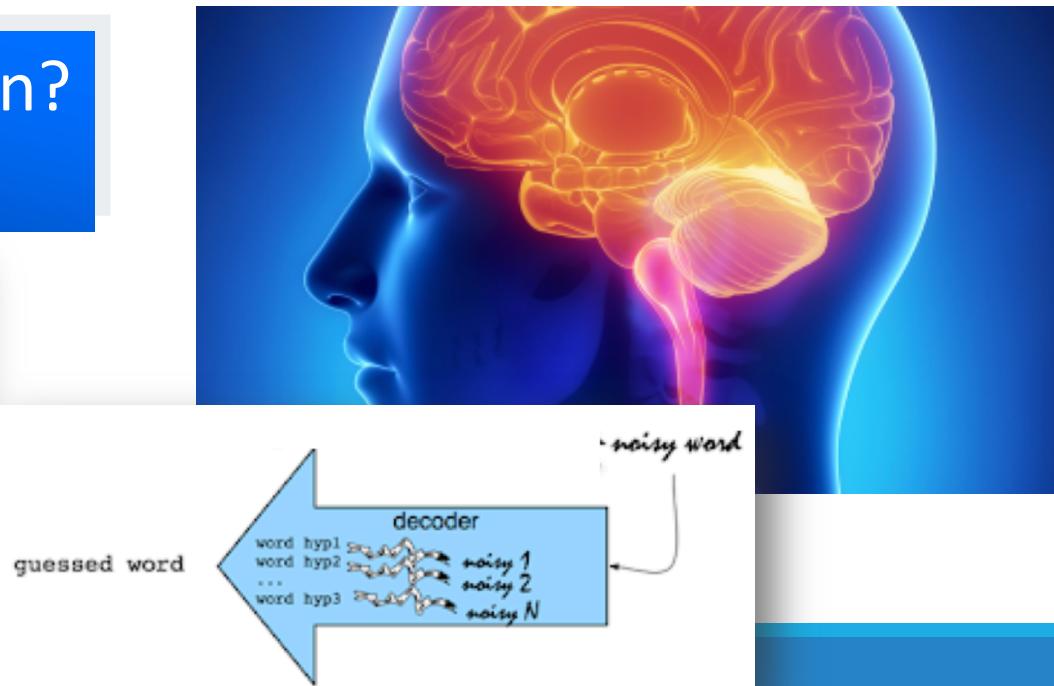
What Sentence is Spoken?  
=> Language Model

Alice was beginning to get very tired of sitting by her sister on the bank, and of having nothing to do: once or twice she had peeped into the book her sister was reading, but it had no pictures or conversations in it, 'and what is the use of a book,' thought Alice 'without pictures or conversation?'

$P(w_{i+1} = \text{of} \mid w_i = \text{tired}) = 1$   
 $P(w_{i+1} = \text{of} \mid w_i = \text{use}) = 1$   
 $P(w_{i+1} = \text{sister} \mid w_i = \text{her}) = 1$   
 $P(w_{i+1} = \text{beginning} \mid w_i = \text{was}) = 1/2$   
 $P(w_{i+1} = \text{reading} \mid w_i = \text{was}) = 1/2$

$P(w_{i+1} = \text{bank} \mid w_i = \text{the}) = 1/3$   
 $P(w_{i+1} = \text{book} \mid w_i = \text{the}) = 1/3$   
 $P(w_{i+1} = \text{use} \mid w_i = \text{the}) = 1/3$

$$P(X|Y) = \frac{P(X,Y)}{P(Y)}$$



# Aim of Automatic Speech Recognition

---

Find the most likely sentence (word sequence)  $\mathbf{W}$ , which transcribes the speech audio  $\mathbf{A}$ :

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

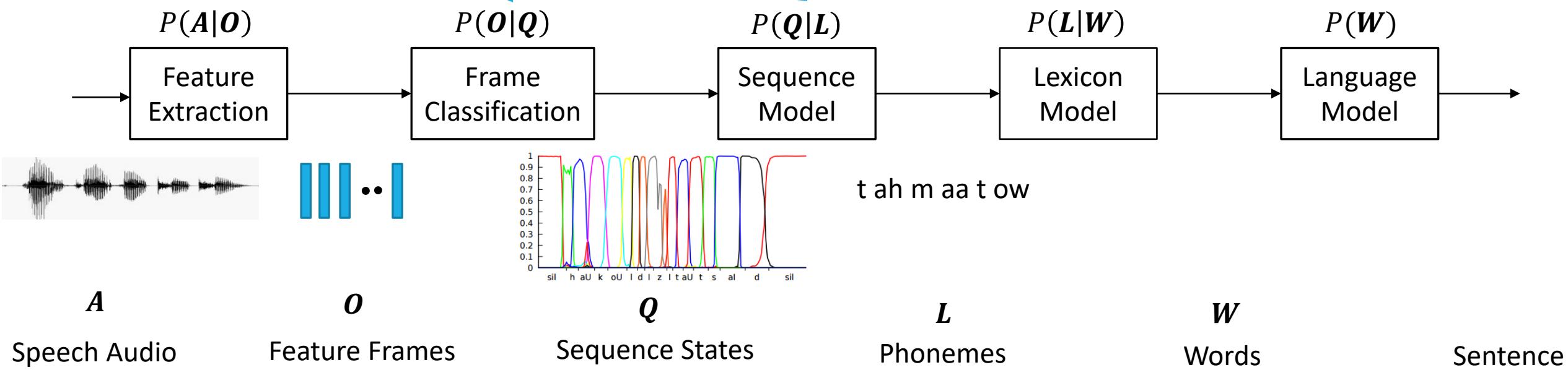
- Acoustic model  $P(\mathbf{A}|\mathbf{W})$
- Language model  $P(\mathbf{W})$

Training: find parameters for acoustic and language model separately

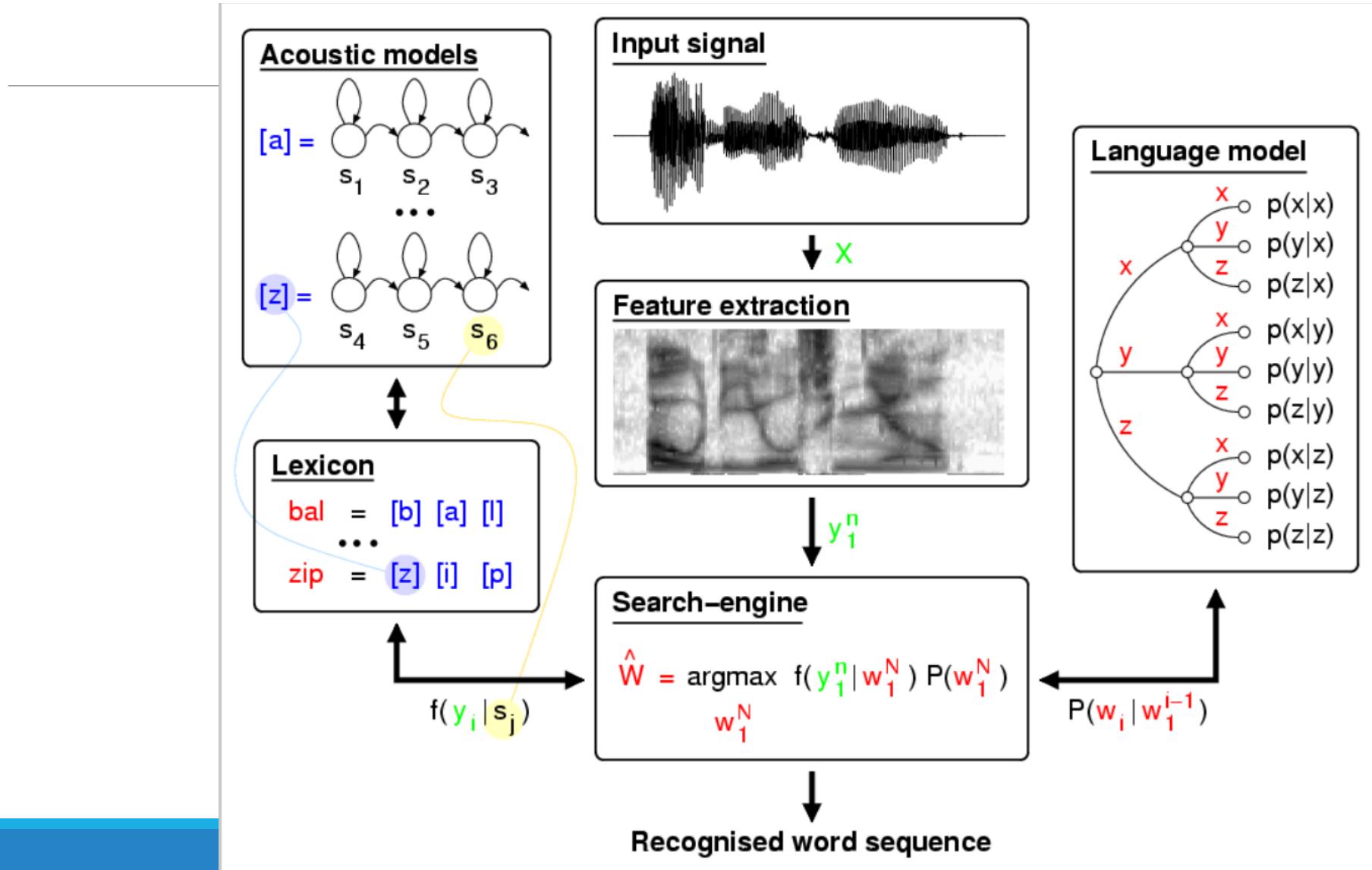
- Speech Corpus: speech waveform and human-annotated transcriptions
- Language model: with extra data (prefer daily expressions corpus for spontaneous speech)

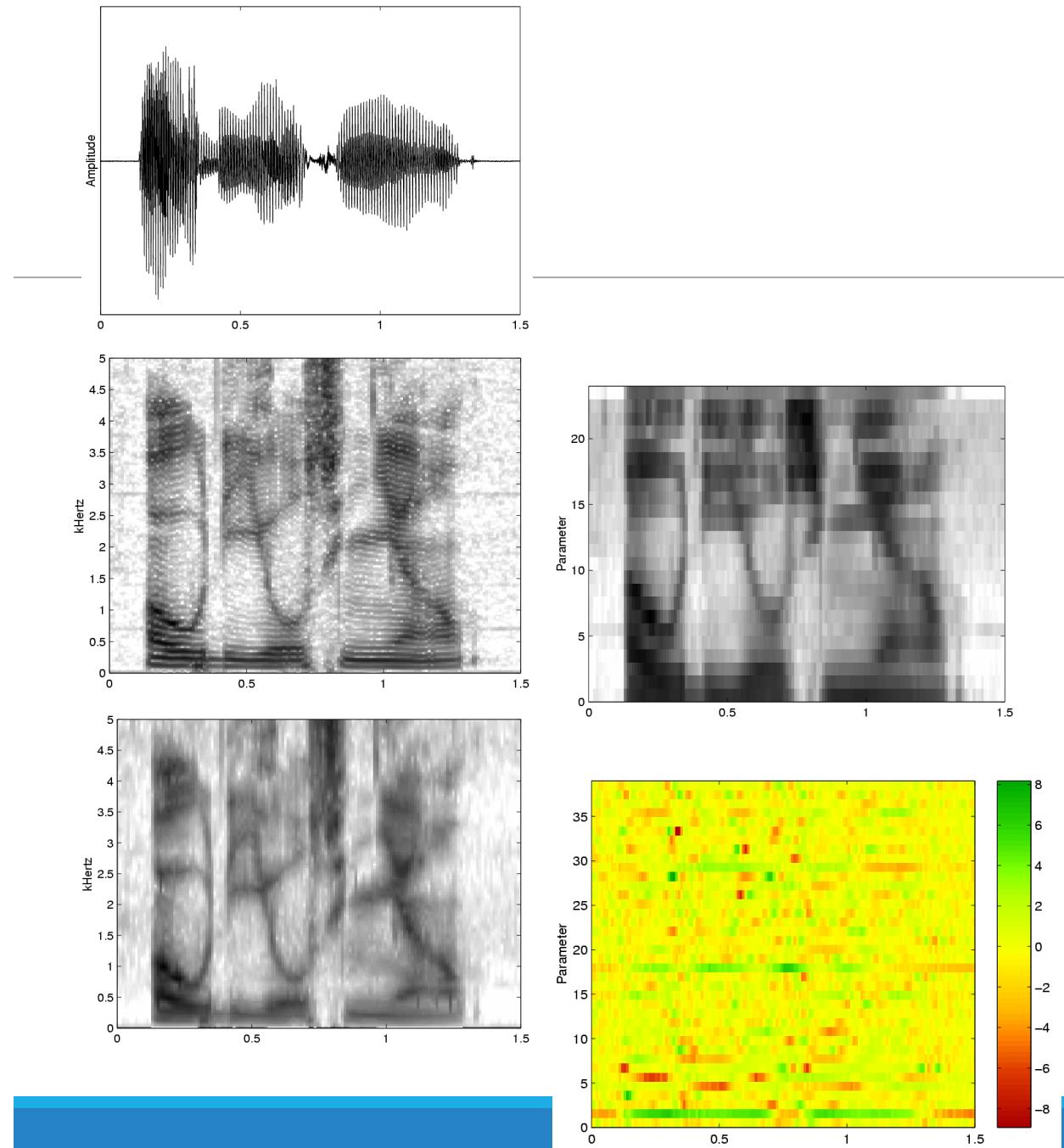
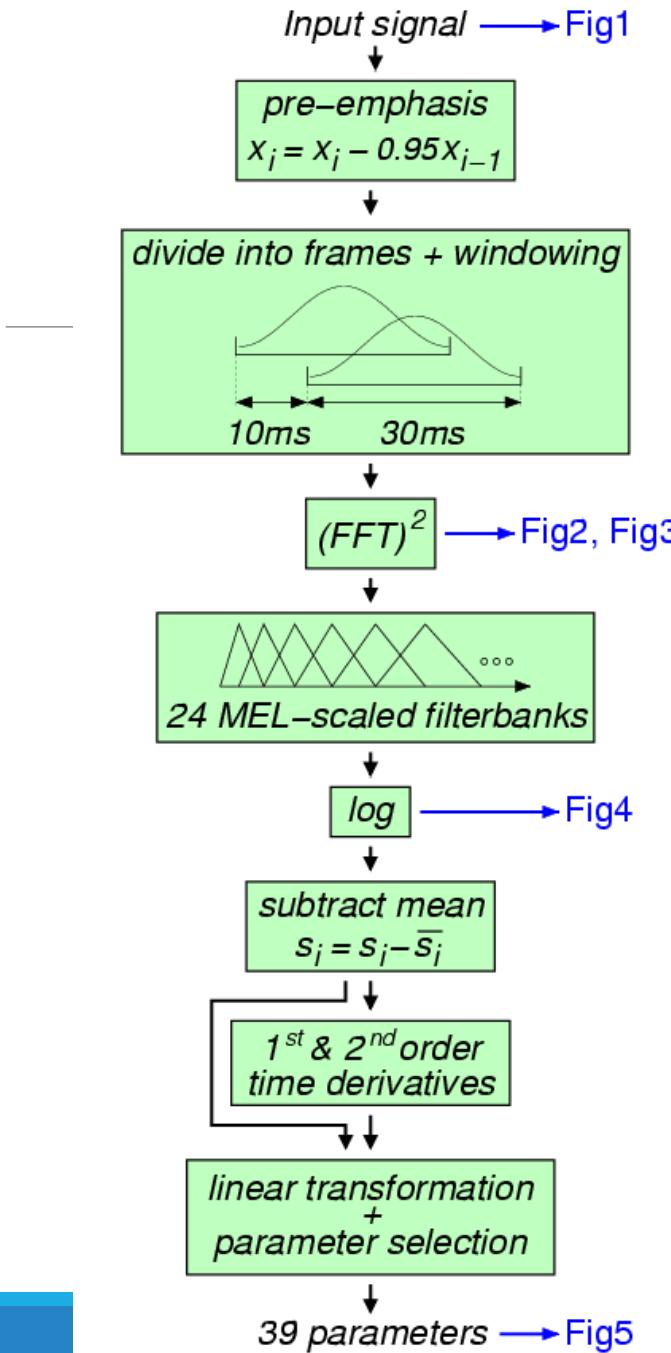
# Automatic Speech Recognition

$$\hat{W} = \underset{W}{\operatorname{argmax}} P(W|O) = \underset{W}{\operatorname{argmax}} P(A|O)P(O|Q)P(Q|L)P(L|W)P(W)$$

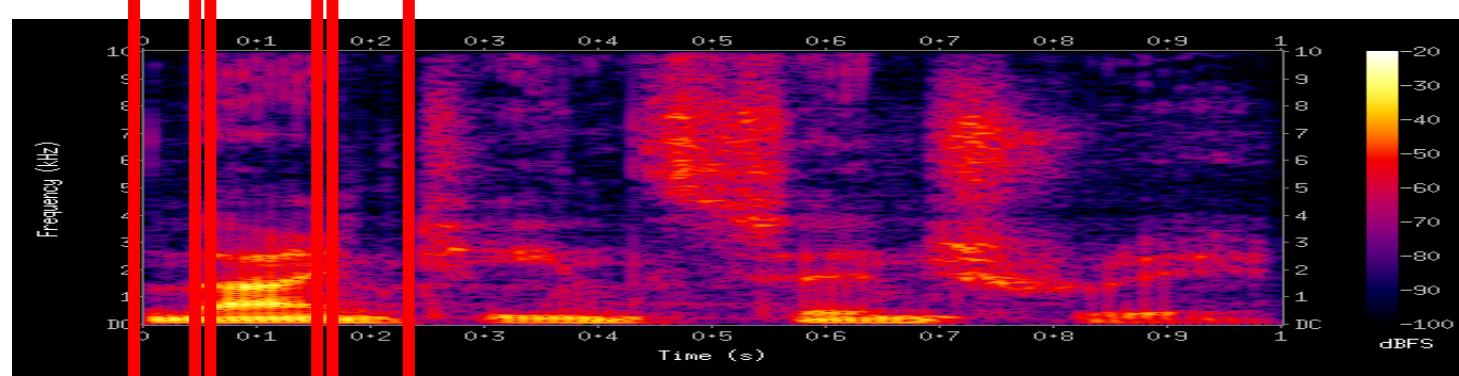
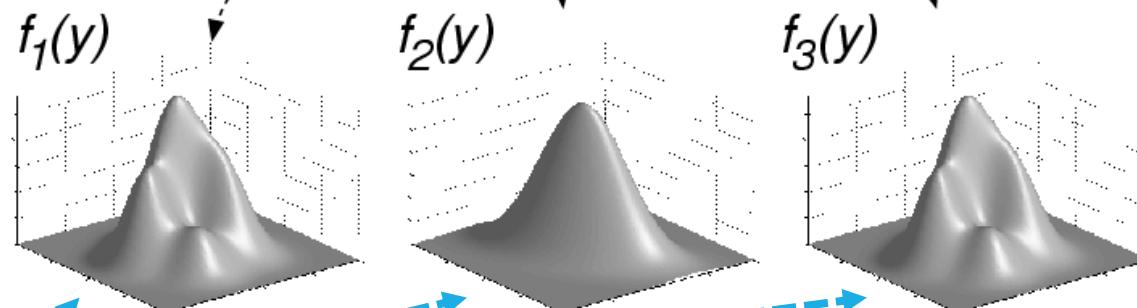
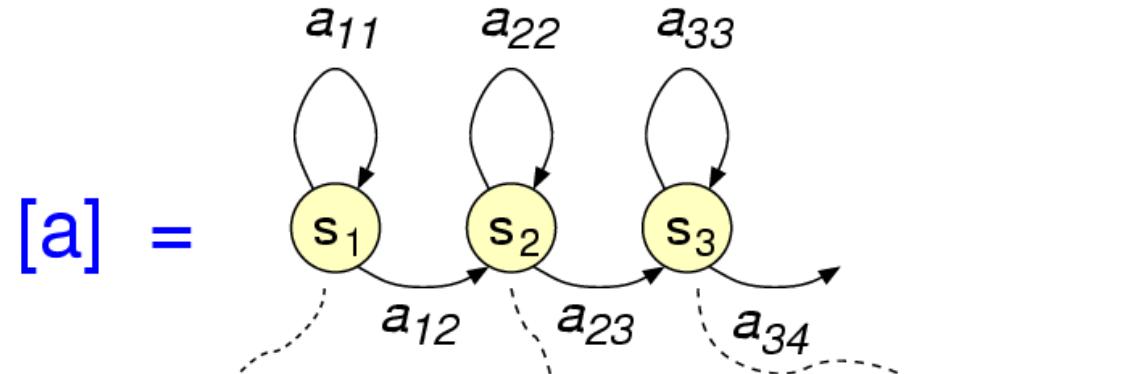


# Automatic Speech Recognition: Short Introduction

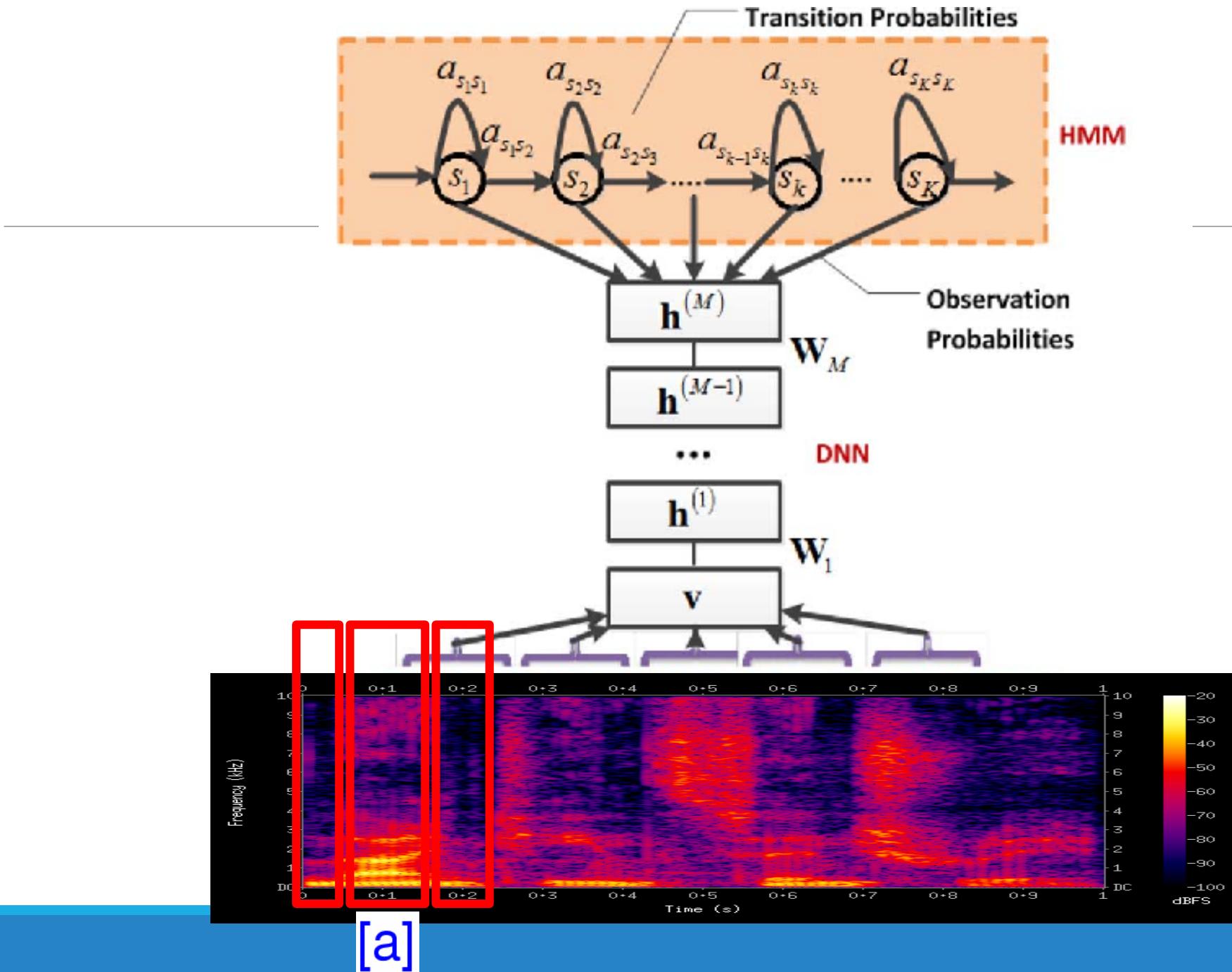




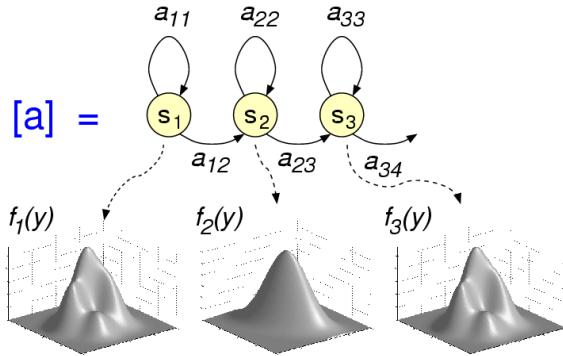
# Hidden Markov Models



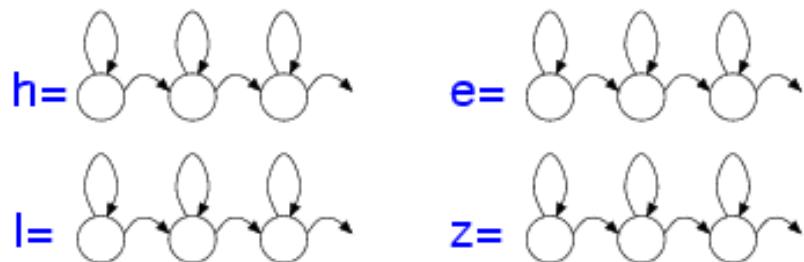
[a]



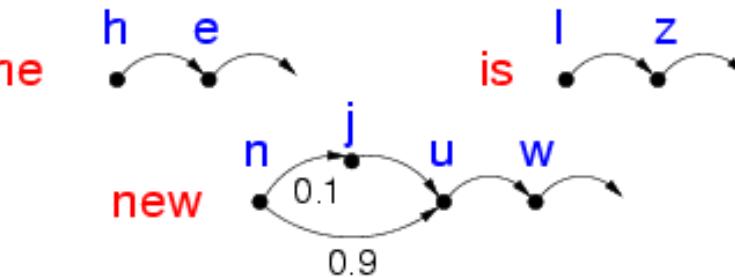
## Hidden Markov Models



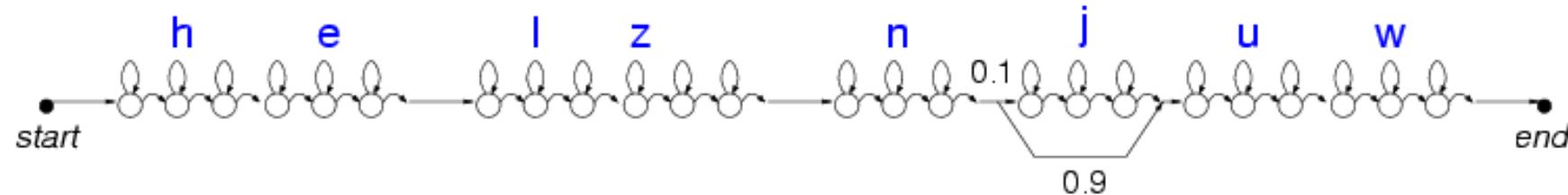
### HMM phone models



### Lexicon



### Sentence model: 'he is new'

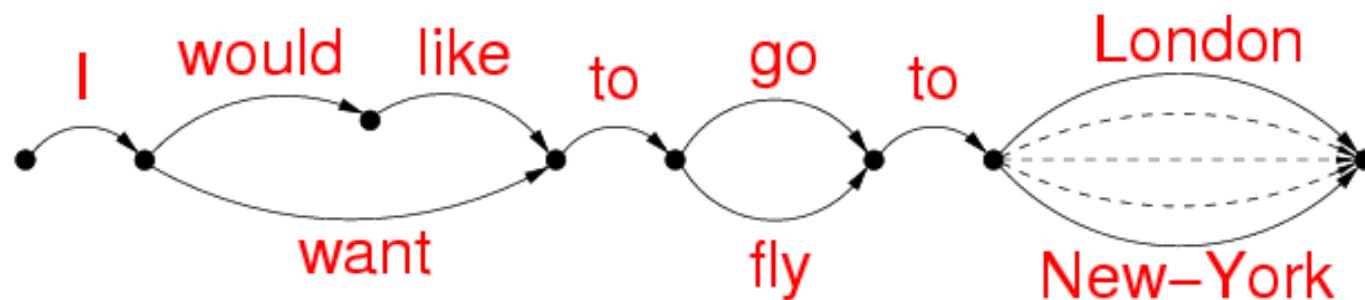


## Grammar:

---

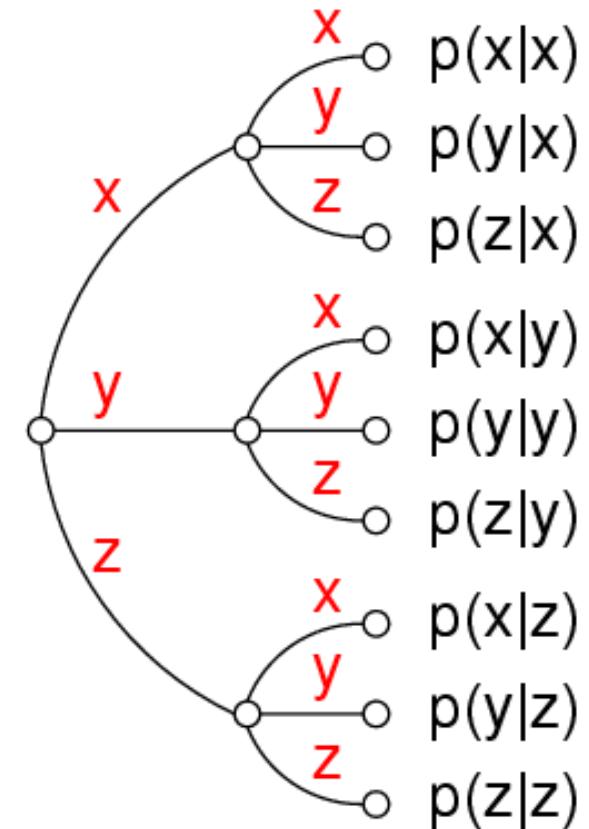
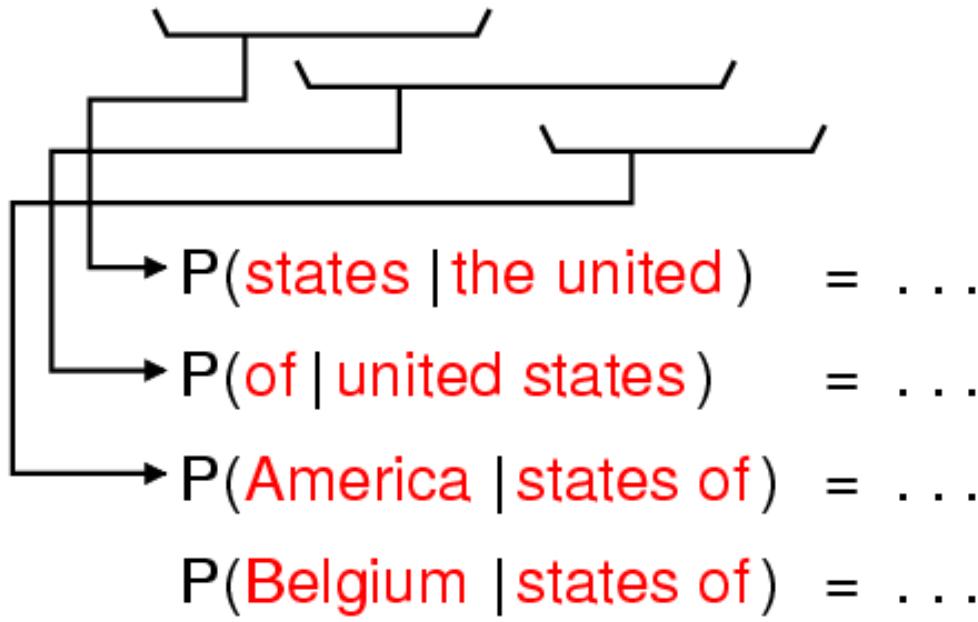
$$\langle \text{sentence}_1 \rangle = I \left\{ \begin{array}{c} \text{would like} \\ \text{want} \end{array} \right\} \text{to} \left\{ \begin{array}{c} \text{go} \\ \text{fly} \end{array} \right\} \text{to} \langle \text{airport} \rangle$$
$$\langle \text{sentence}_2 \rangle = \dots$$
$$\langle \text{airport} \rangle = \{ \text{London}, \text{New-York}, \dots \}$$

## Finite-state representation:



## $N$ -gram language models

... the united states of ???



Predict the next word, give a set of predecessor words

# Language model

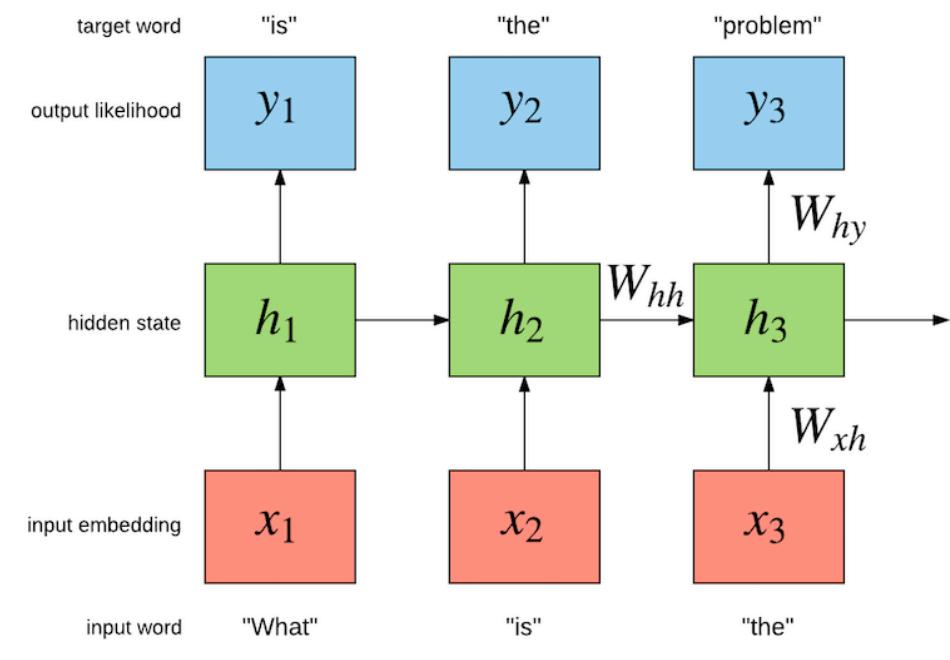
Language model is a probabilistic model used to

- Guide the search algorithm (predict next word given history)
- Disambiguate between phrases which are acoustically similar
  - Great wine vs Grey twine

It assigns probability to a sequence of tokens to be finally recognized

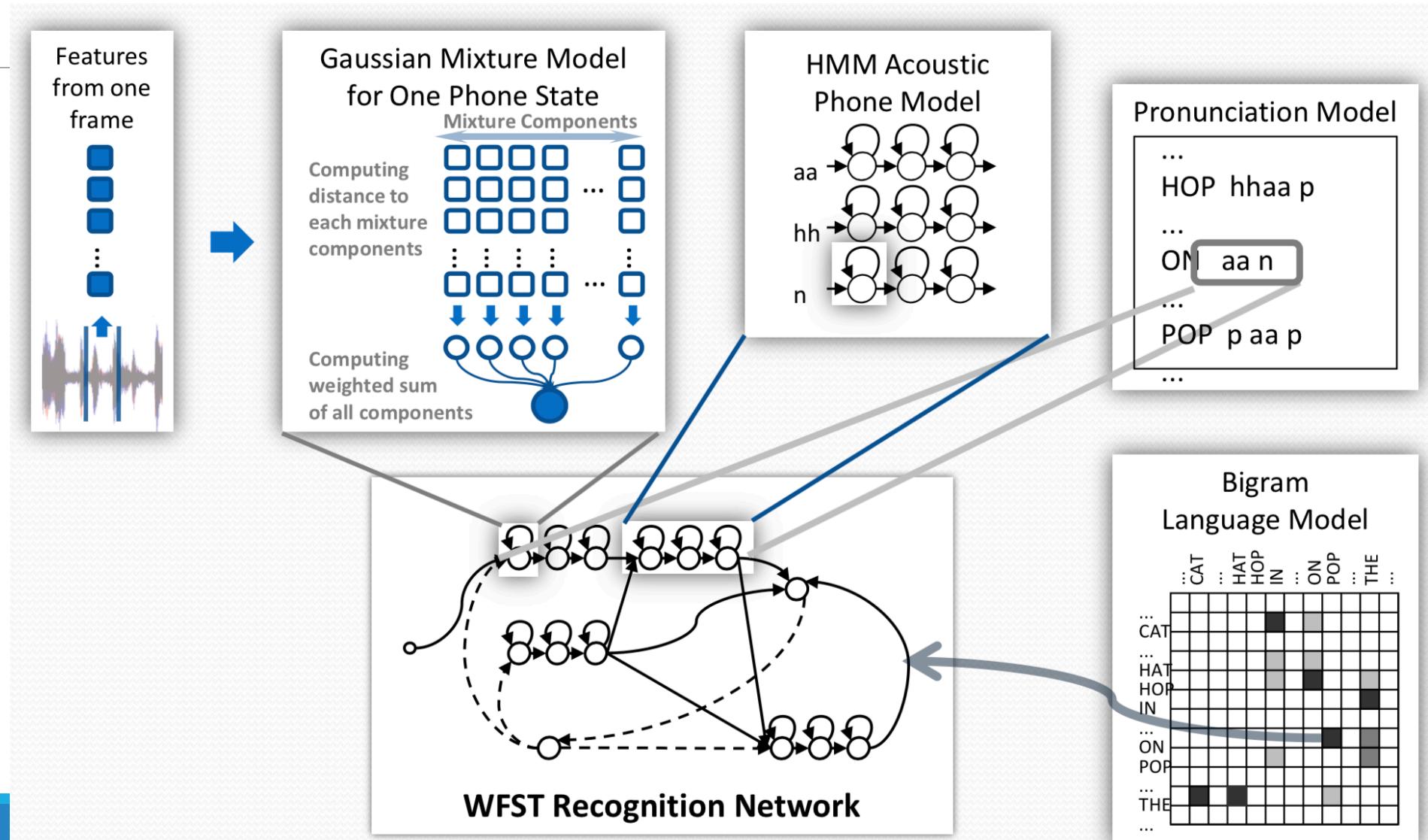
N-gram model  $P(w_N | w_1, w_2, \dots, w_{N-1})$

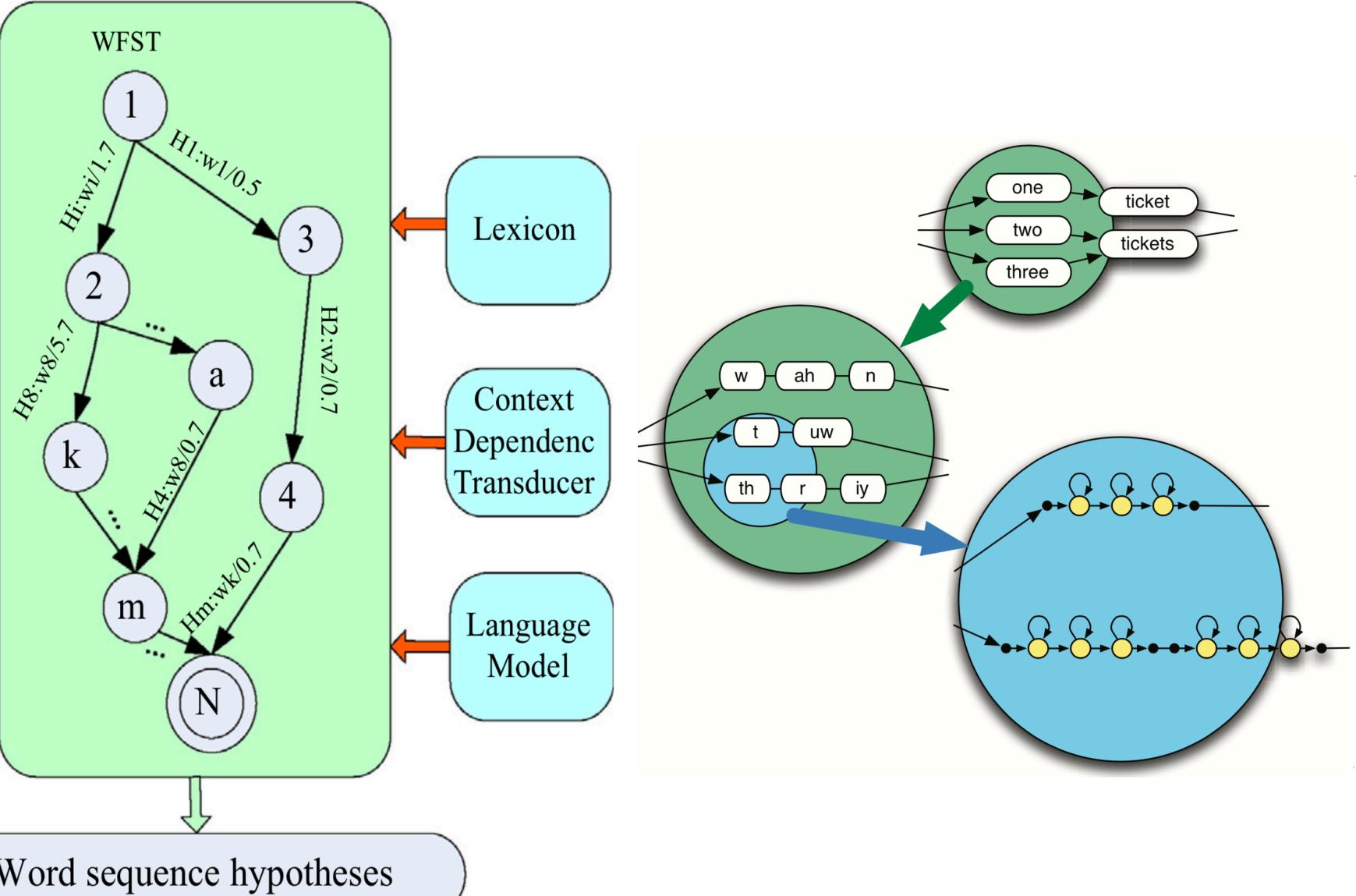
Recurrent neural network



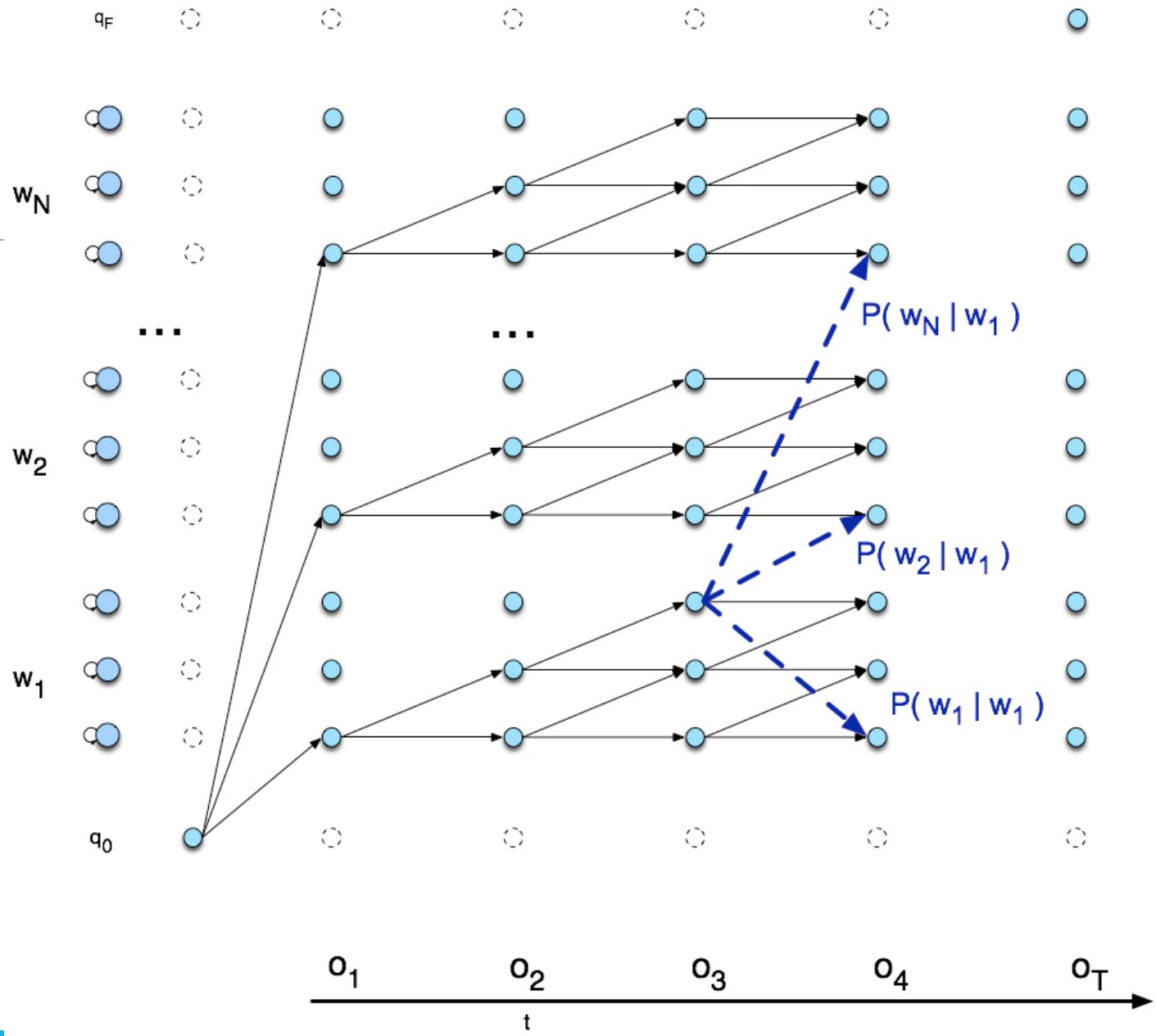
Recurrent neural network based Language model

# Automatic Speech Recognizer

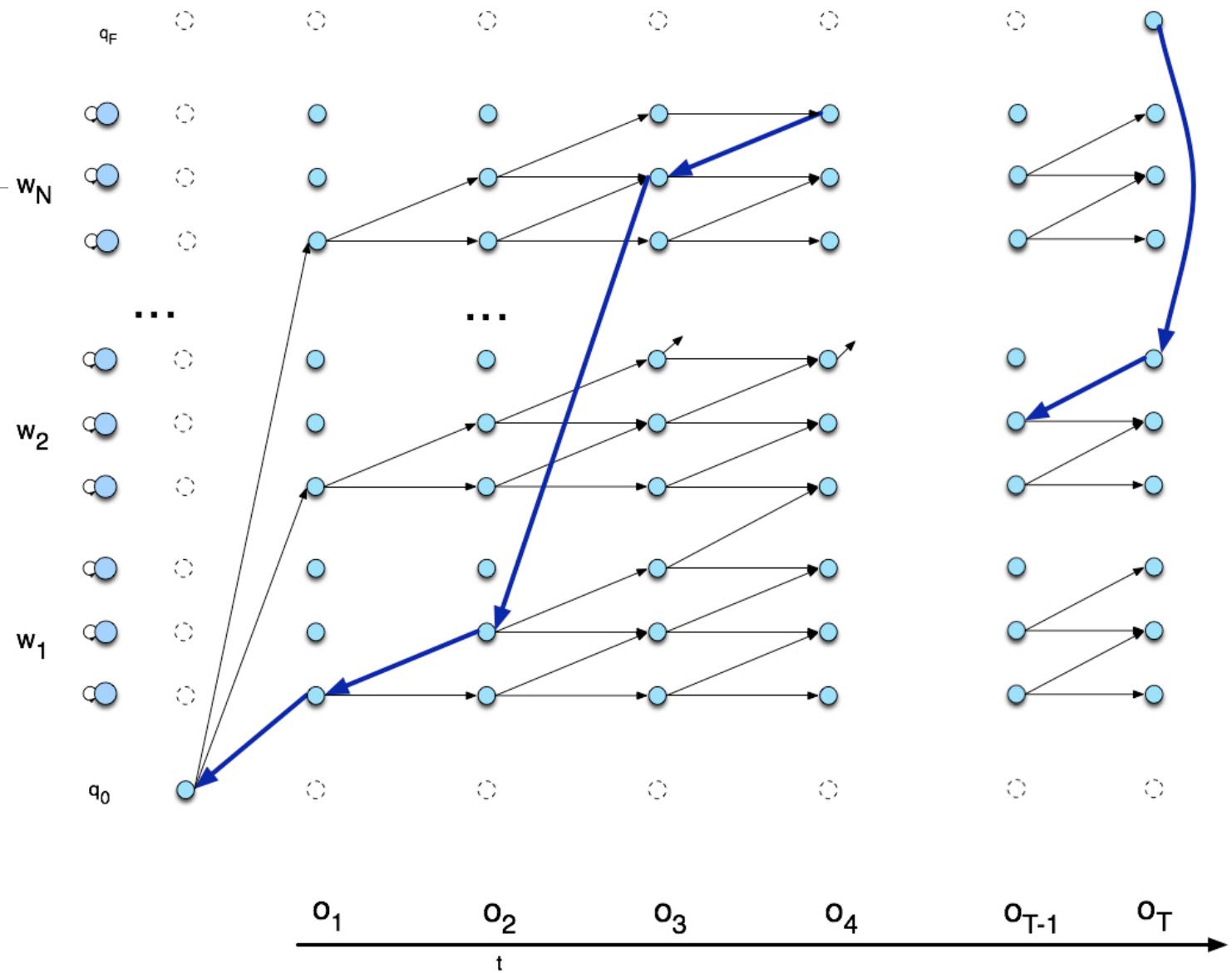




# Viterbi Trellis

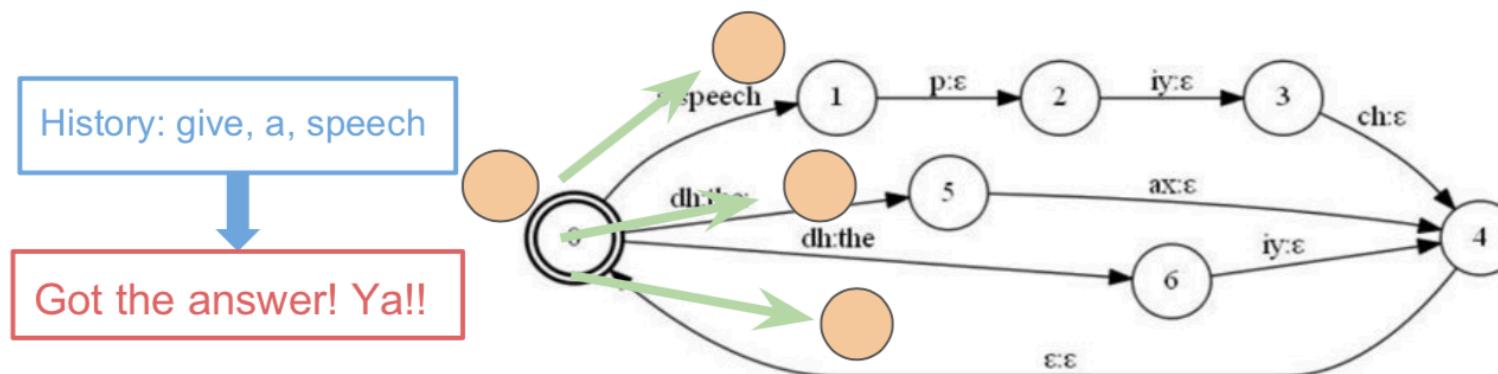


# Viterbi Backtrace

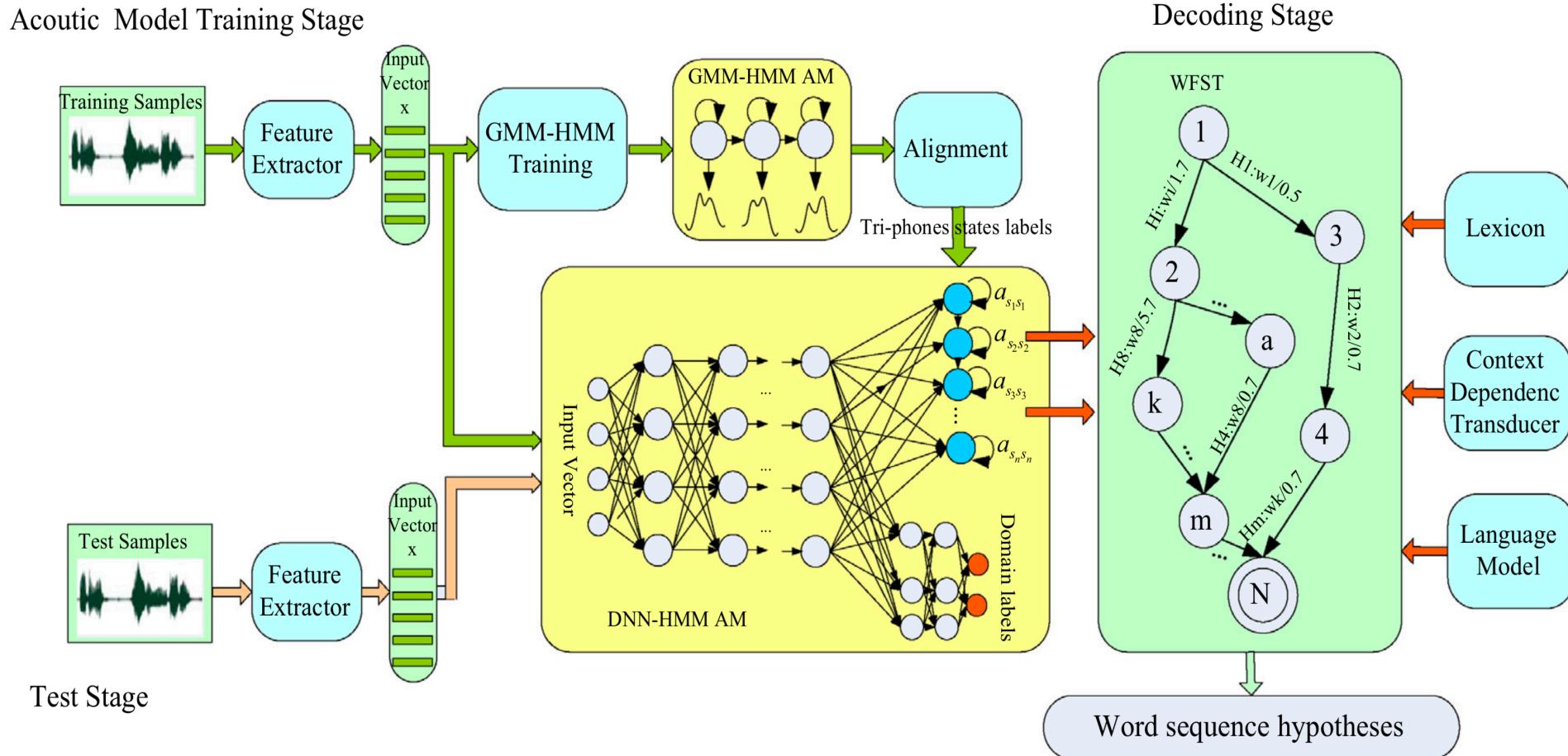


# Token & Beam Pruning

- The loop of Decode:
  - Token Copy
  - Update AM/LM scores
  - Record the **highest\_score** over all tokens
  - If **Token.score < (highest\_score - beam)**, kill it. → beam pruning
- Finally, we choose the token with the highest score.
- Output its history words as answer.



# Automatic Speech Recognizer

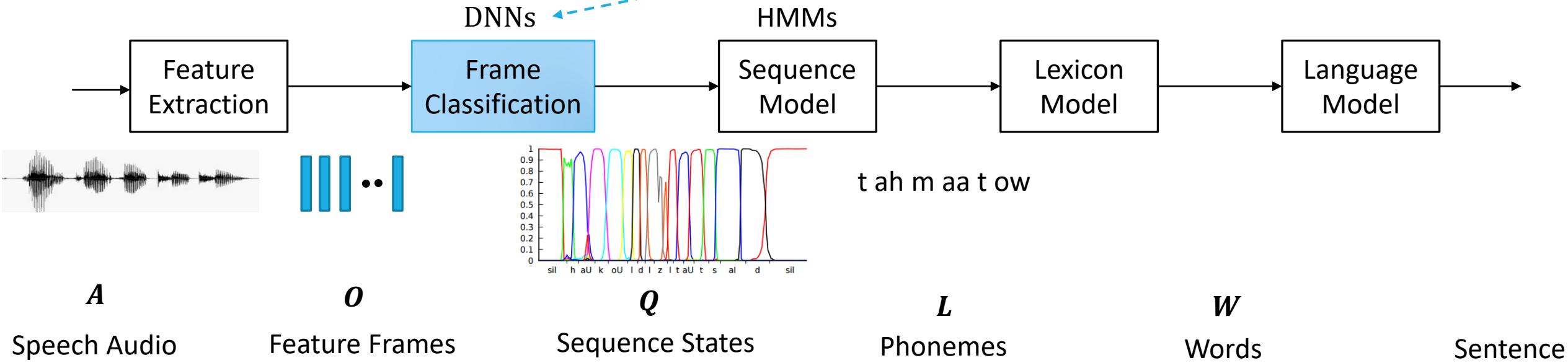


# DNN-HMM

**Kaldi**

DNN: Deep Neural Networks

HMMs: Hidden Markov Models

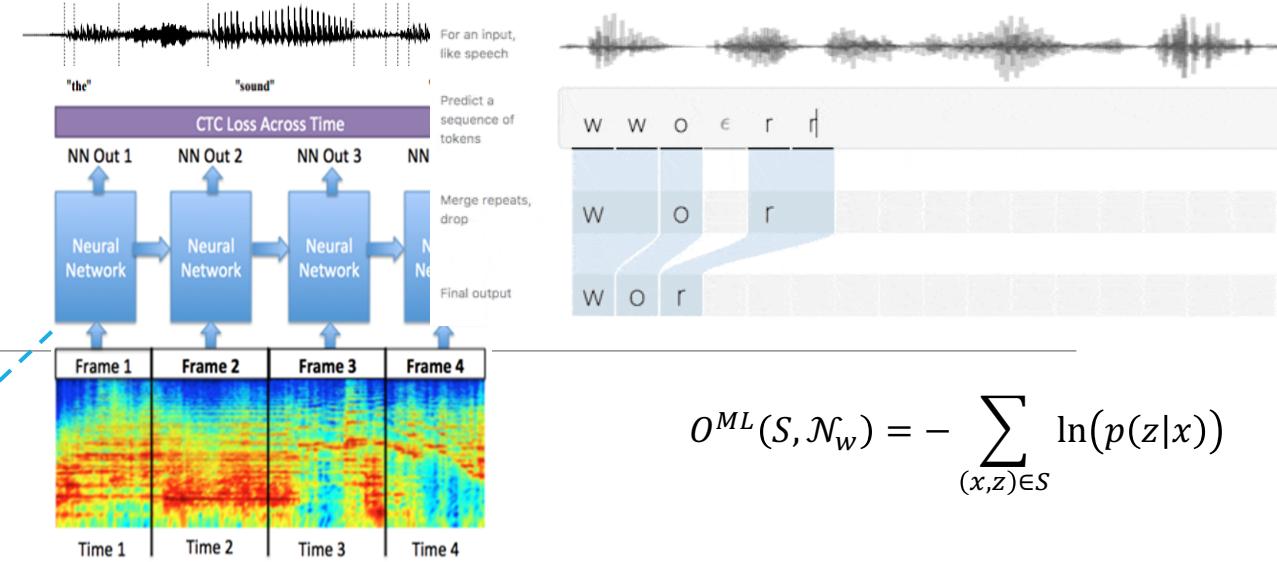
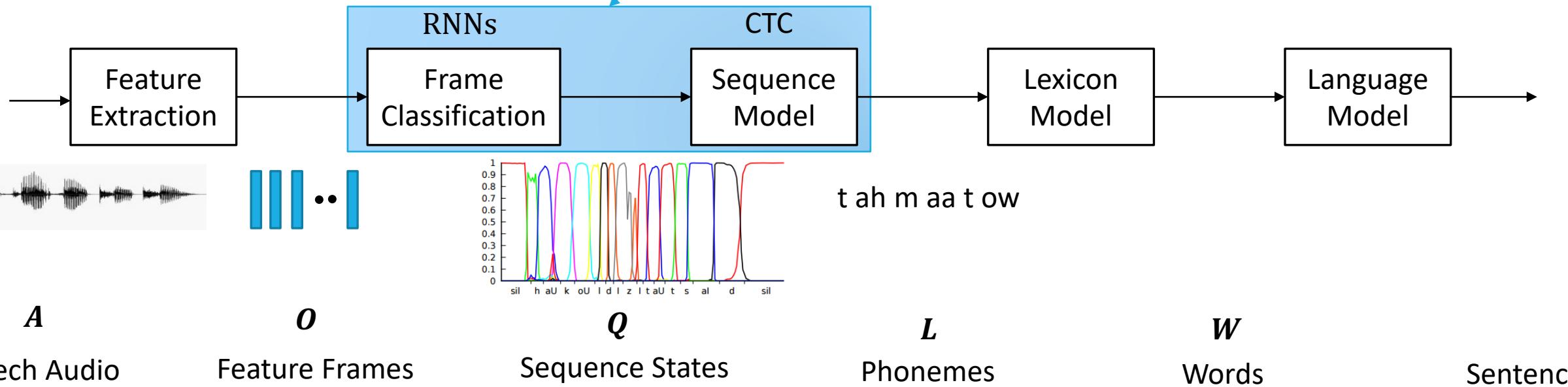


# End-to-End: CTC

## DeepSpeech2

RNN: Recurrent Neural Networks

CTC: Connectionist Temporal Classification



# End-to-End: Attention

## Listen, Attend and Spell

Predict character sequence directly

Sequence-to-Sequence with Attention

