

Introduction to Speech Recognition

YUAN-FU LIAO

NATIONAL TAIPEI UNIVERSITY OF TECHNOLOGY

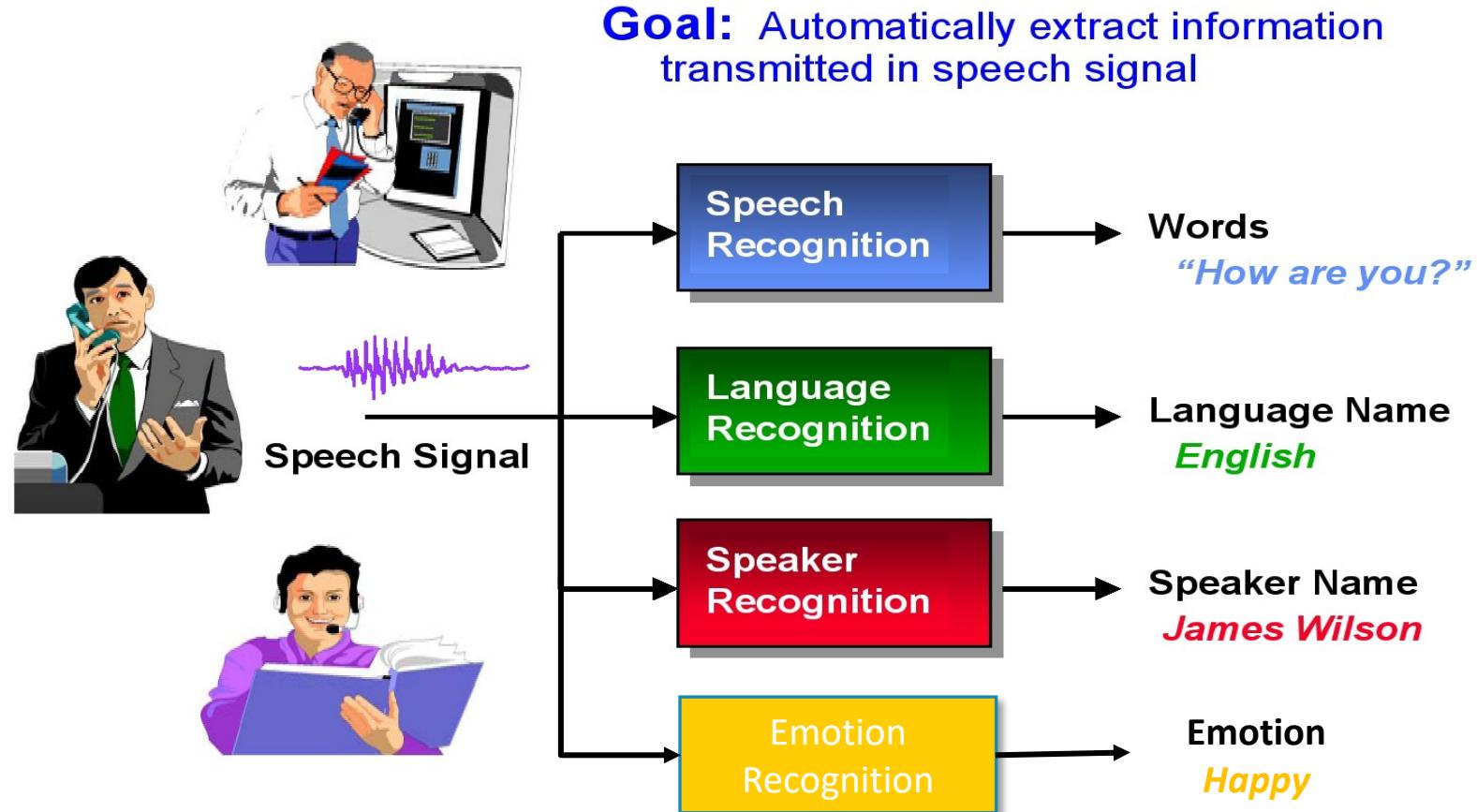
Outline

Overview

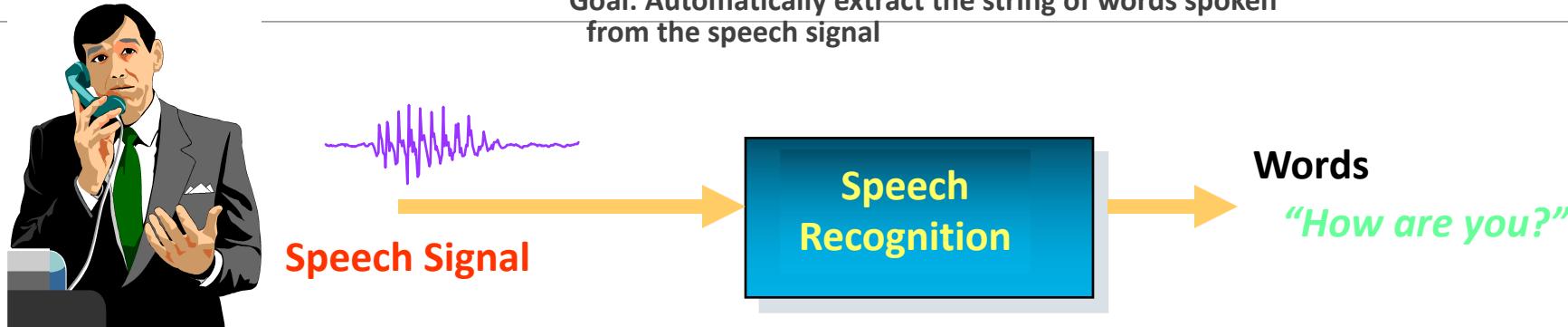
Statistic Speech Recognition

Deep-Learning

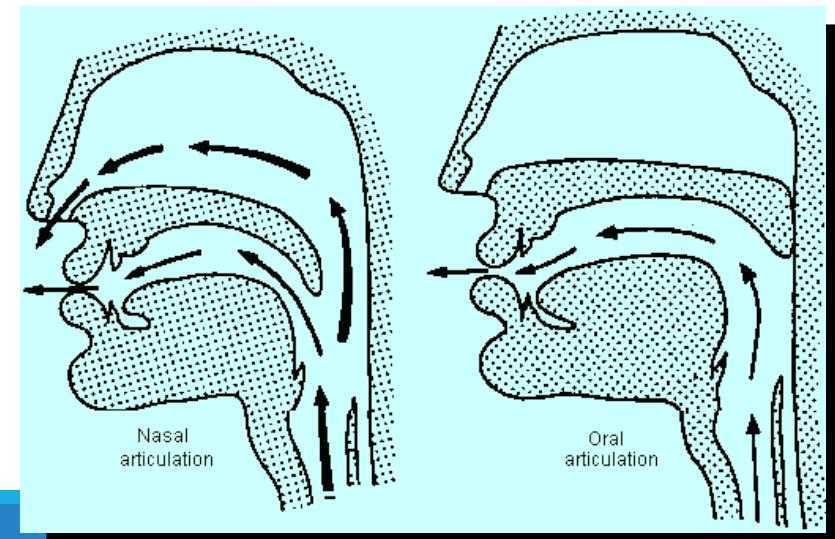
Broad Objectives of Speech Recognition for Machines



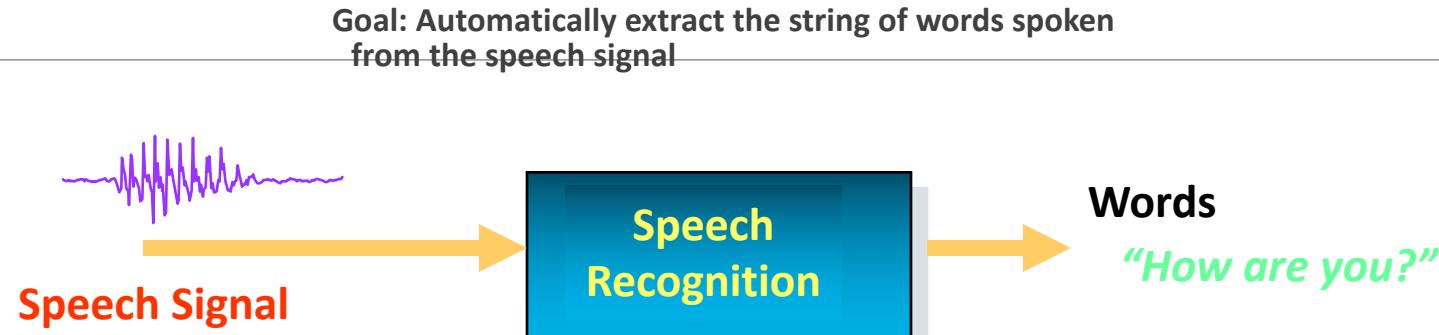
Speech Recognition



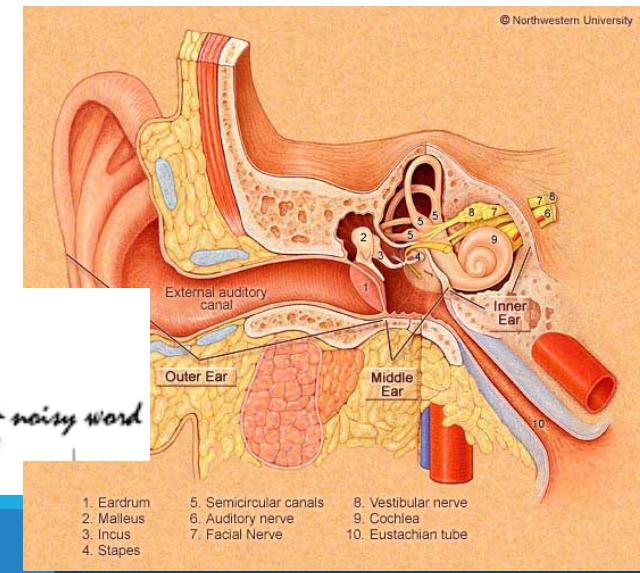
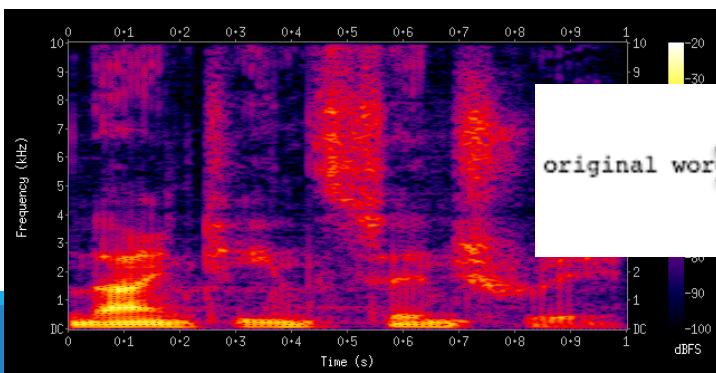
How is SPEECH produced?
⇒ Characteristics of
Acoustic Signal



Speech Recognition



How is SPEECH perceived?
=> Important Features



Speech Recognition



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



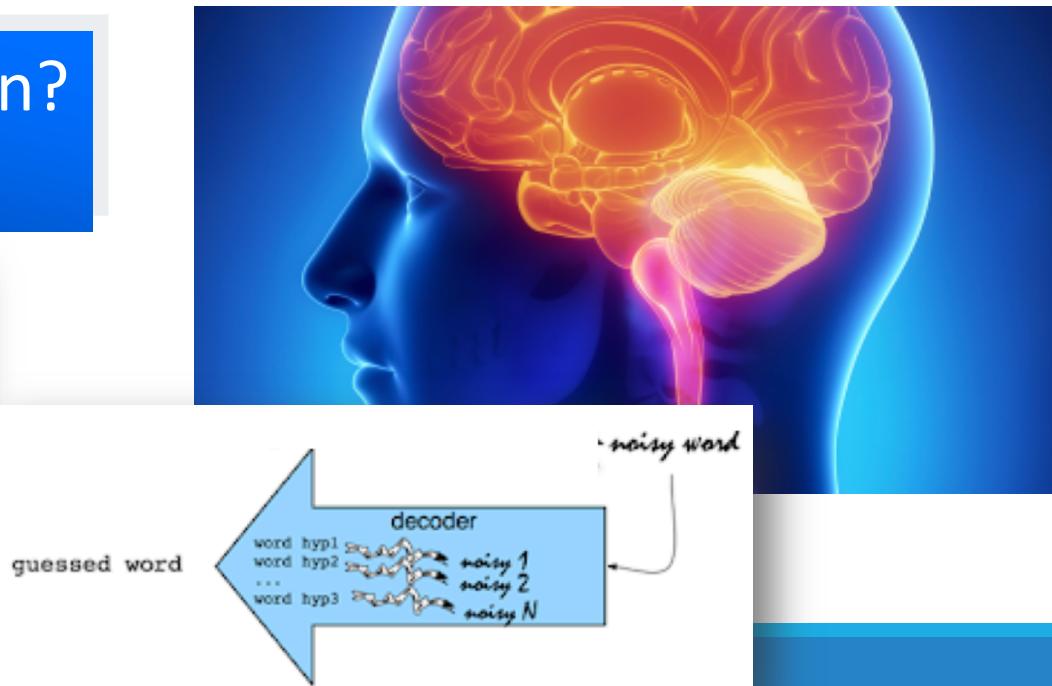
What Sentence is Spoken?
=> Language Model

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?'

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

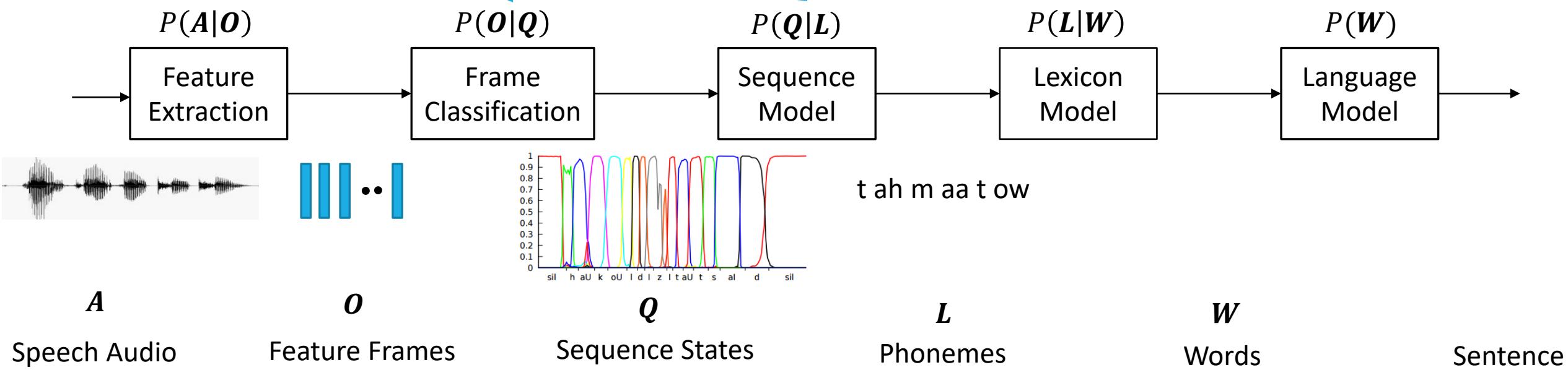
$P(w_{i+1} = \text{bank} | w_i = \text{the}) = 1/3$
 $P(w_{i+1} = \text{book} | w_i = \text{the}) = 1/3$
 $P(w_{i+1} = \text{use} | w_i = \text{the}) = 1/3$

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

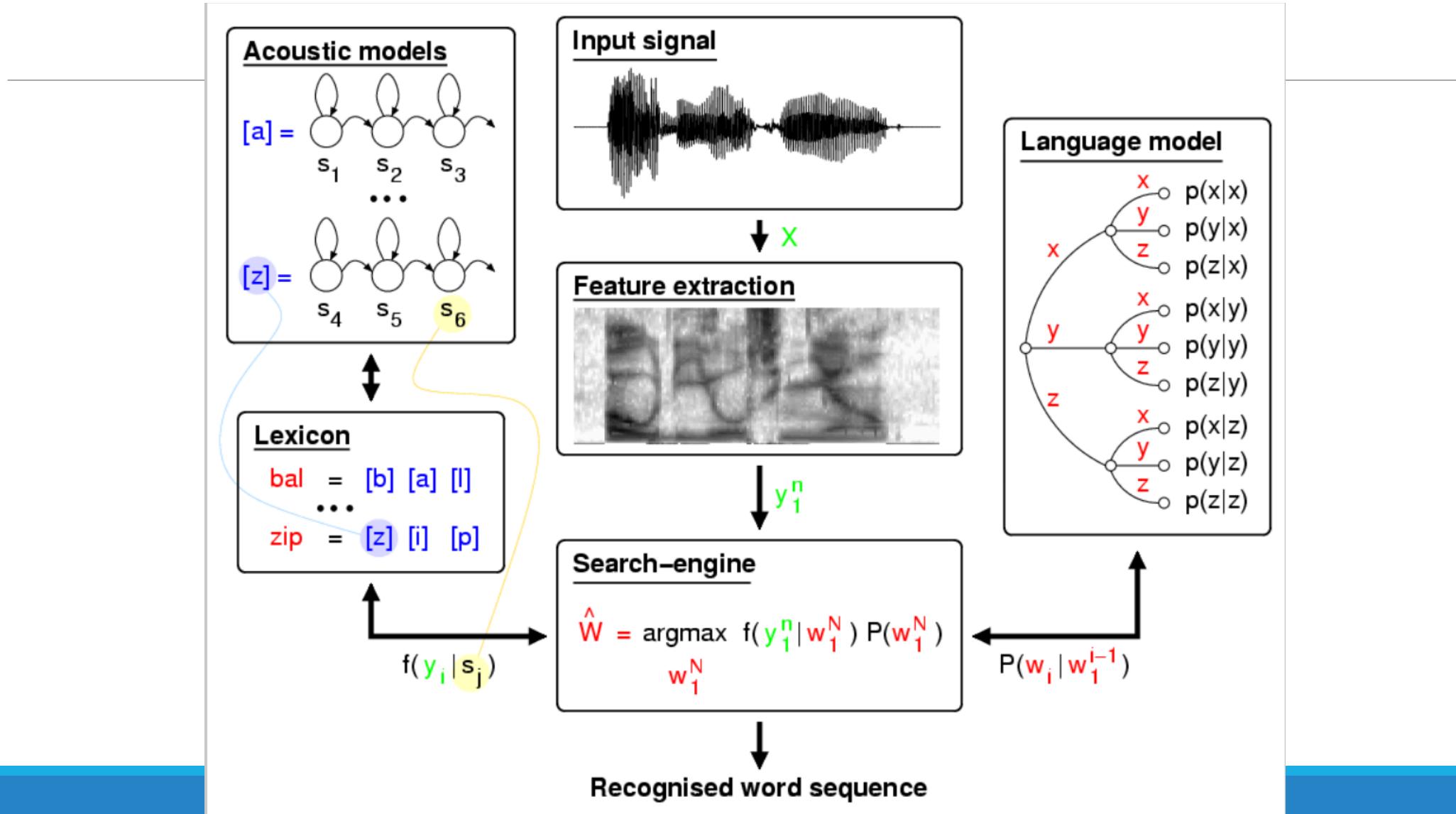


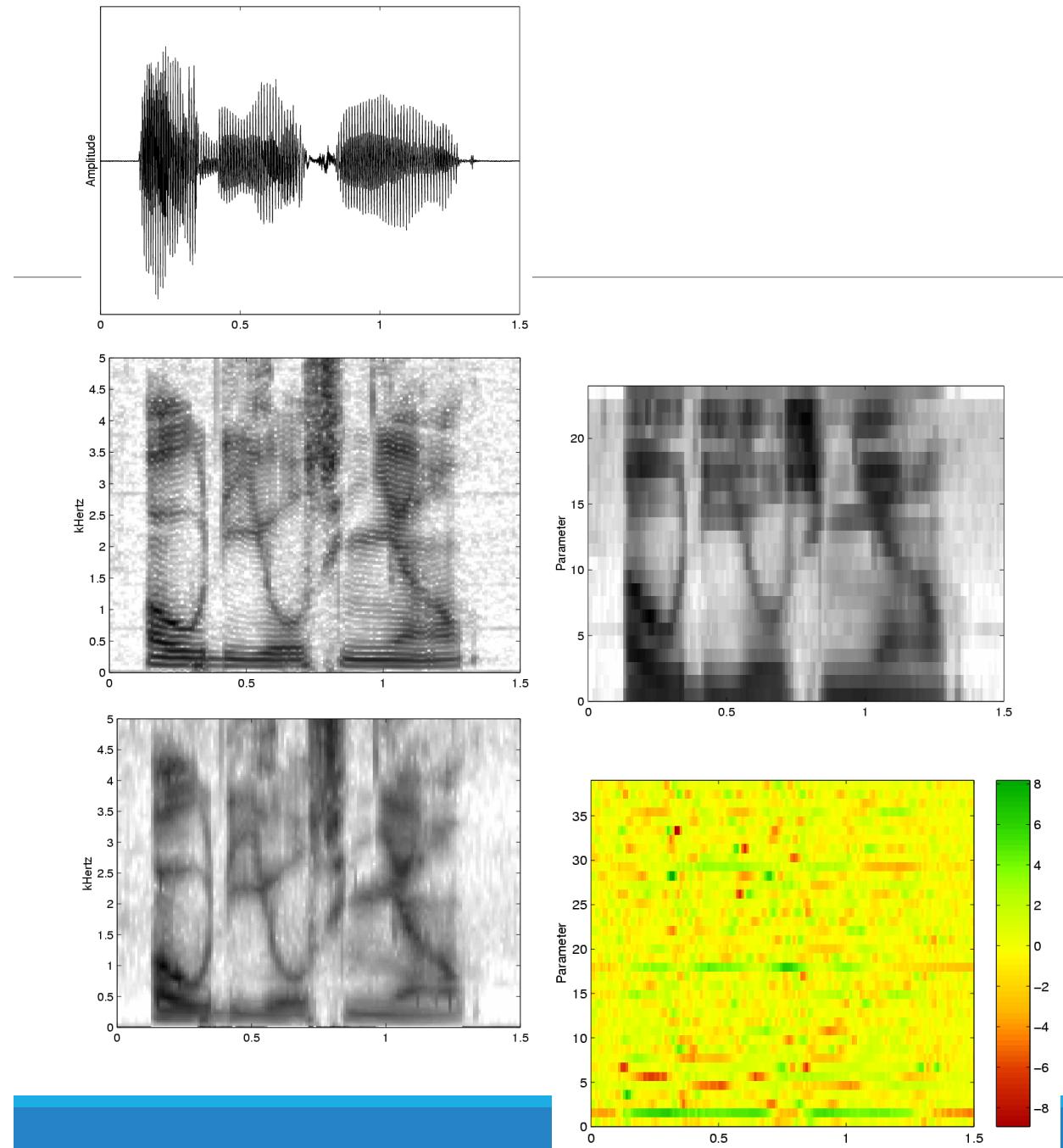
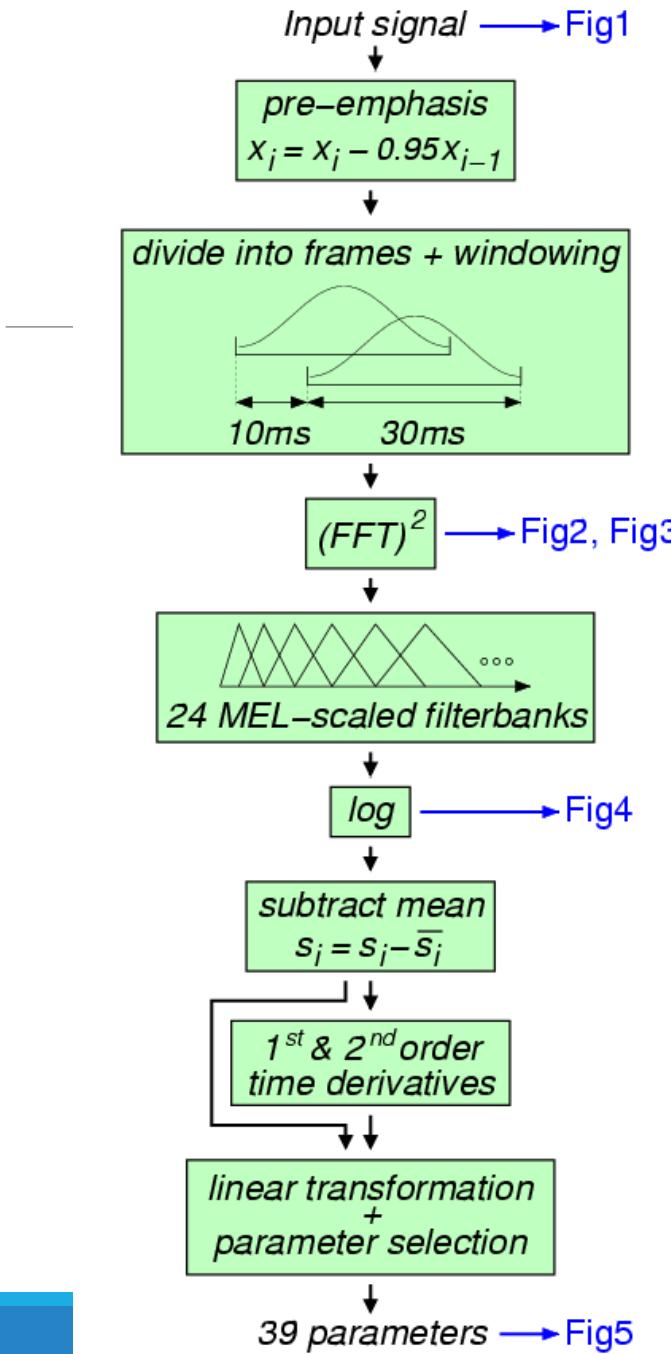
Automatic Speech Recognition

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

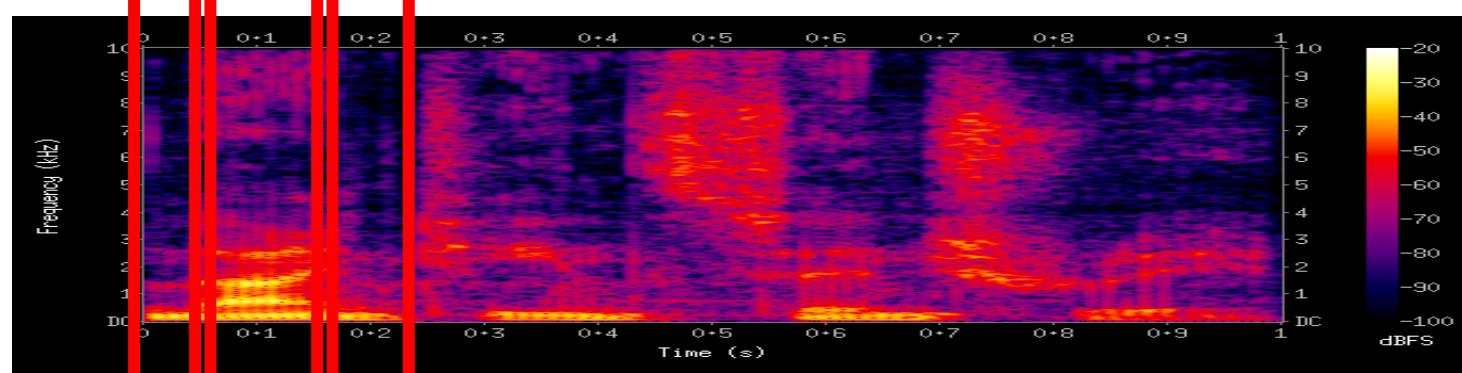
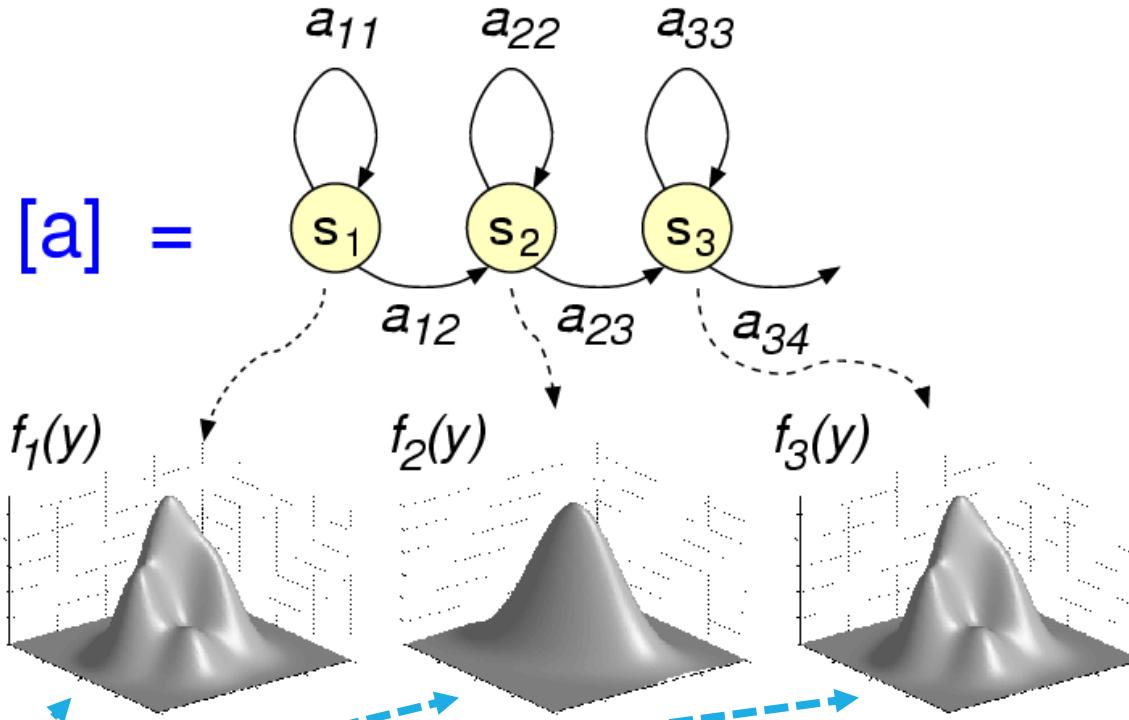


Automatic Speech Recognition: Short Introduction

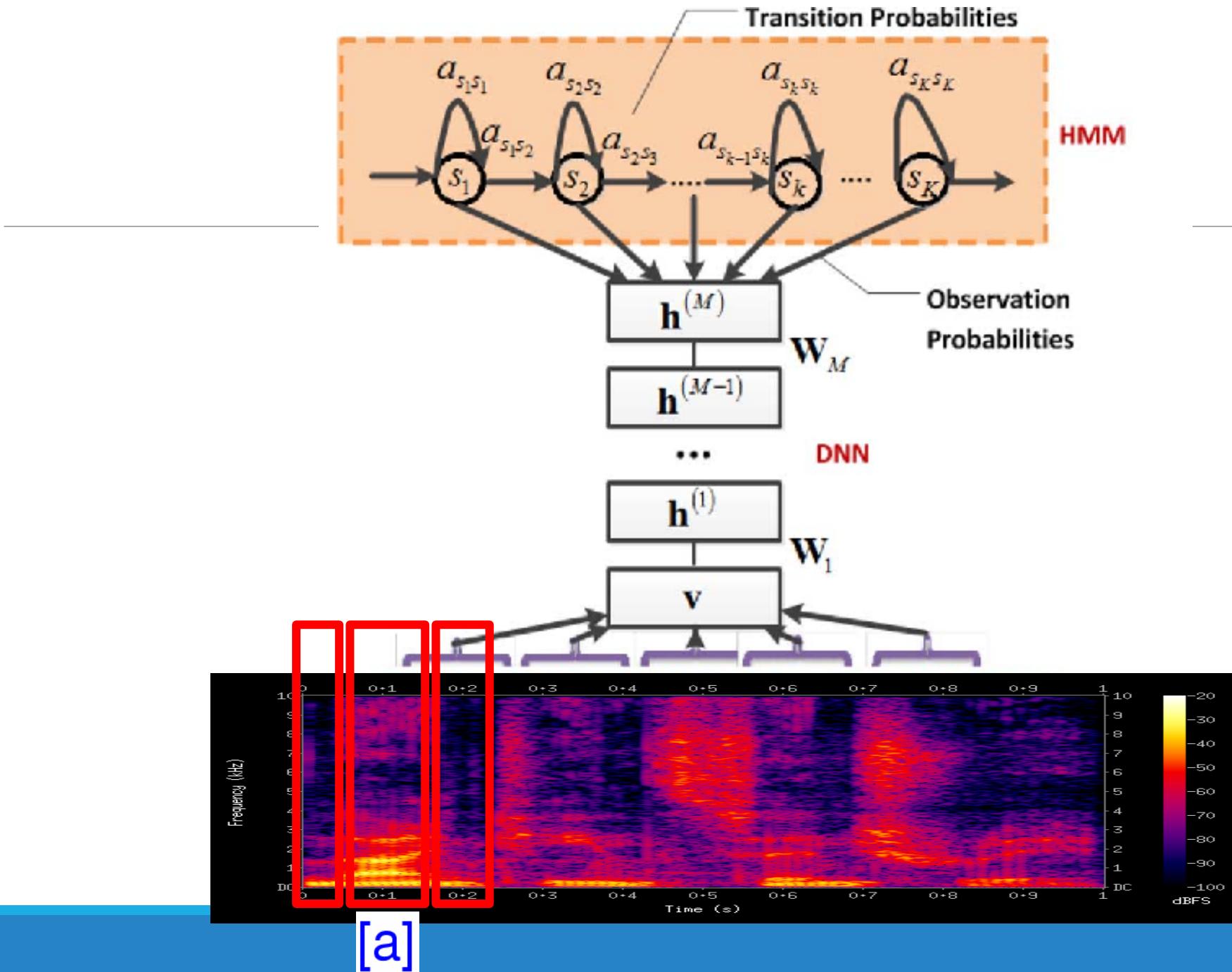




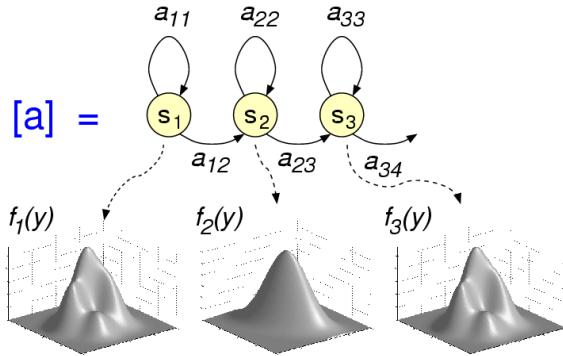
Hidden Markov Models



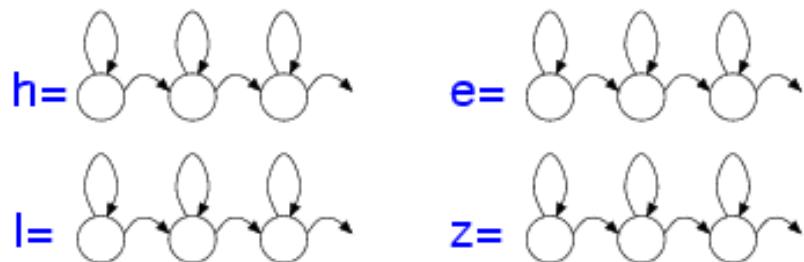
[a]



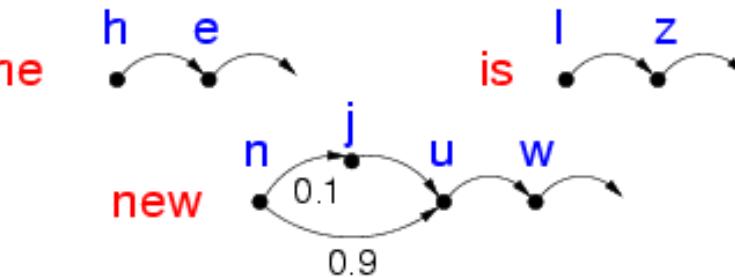
Hidden Markov Models



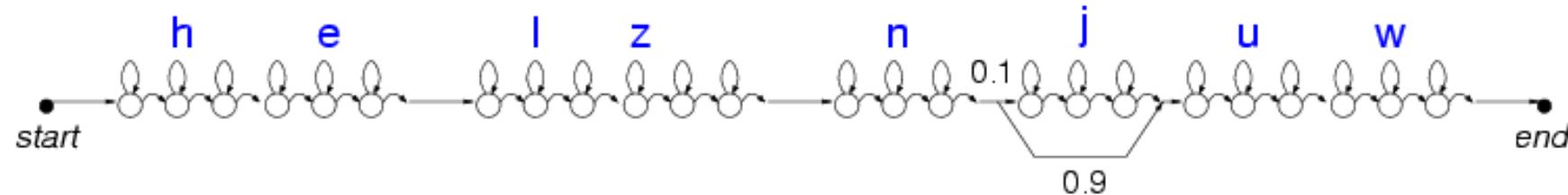
HMM phone models



Lexicon



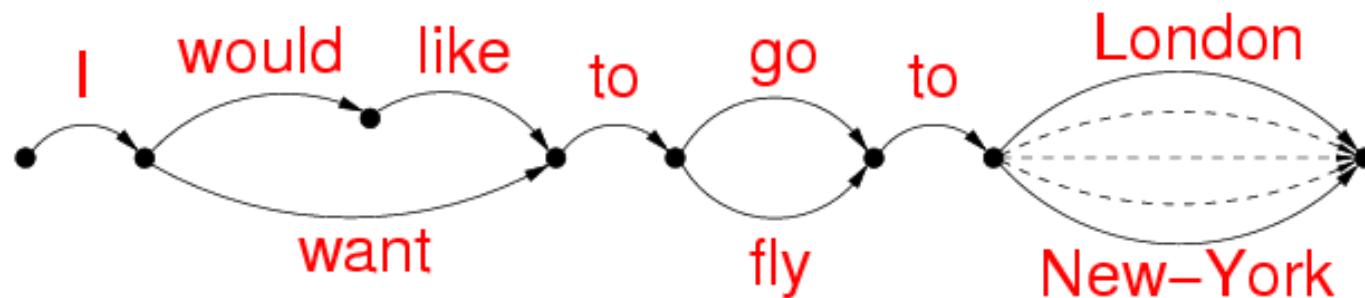
Sentence model: 'he is new'



Grammar:

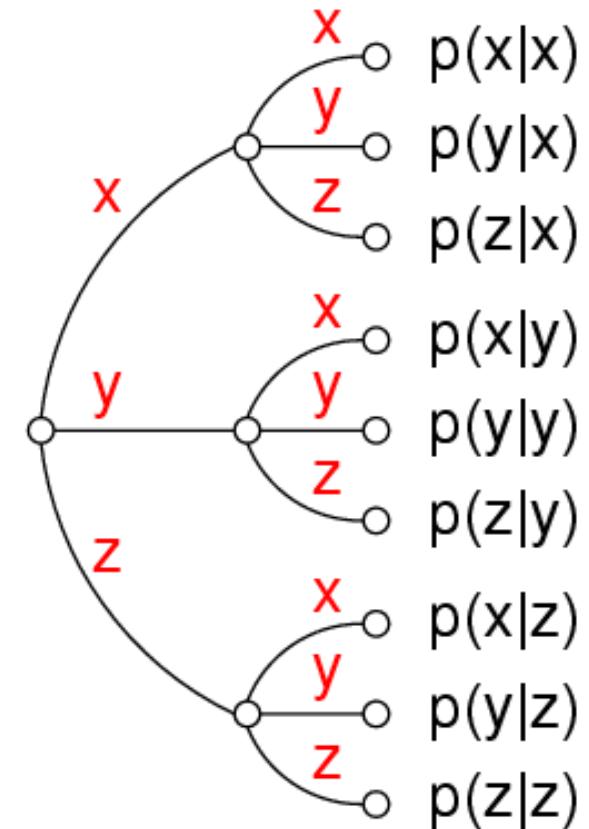
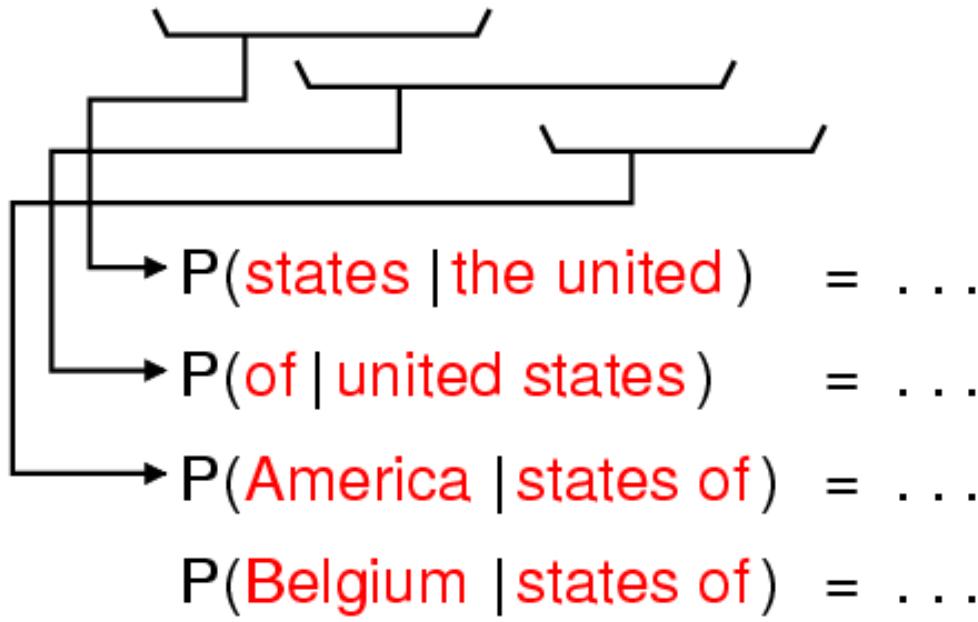
$$\begin{aligned}\langle \text{sentence}_1 \rangle &= I \left\{ \begin{array}{c} \text{would like} \\ \text{want} \end{array} \right\} \text{to} \left\{ \begin{array}{c} \text{go} \\ \text{fly} \end{array} \right\} \text{to} \langle \text{airport} \rangle \\ \langle \text{sentence}_2 \rangle &= \dots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

Finite-state representation:



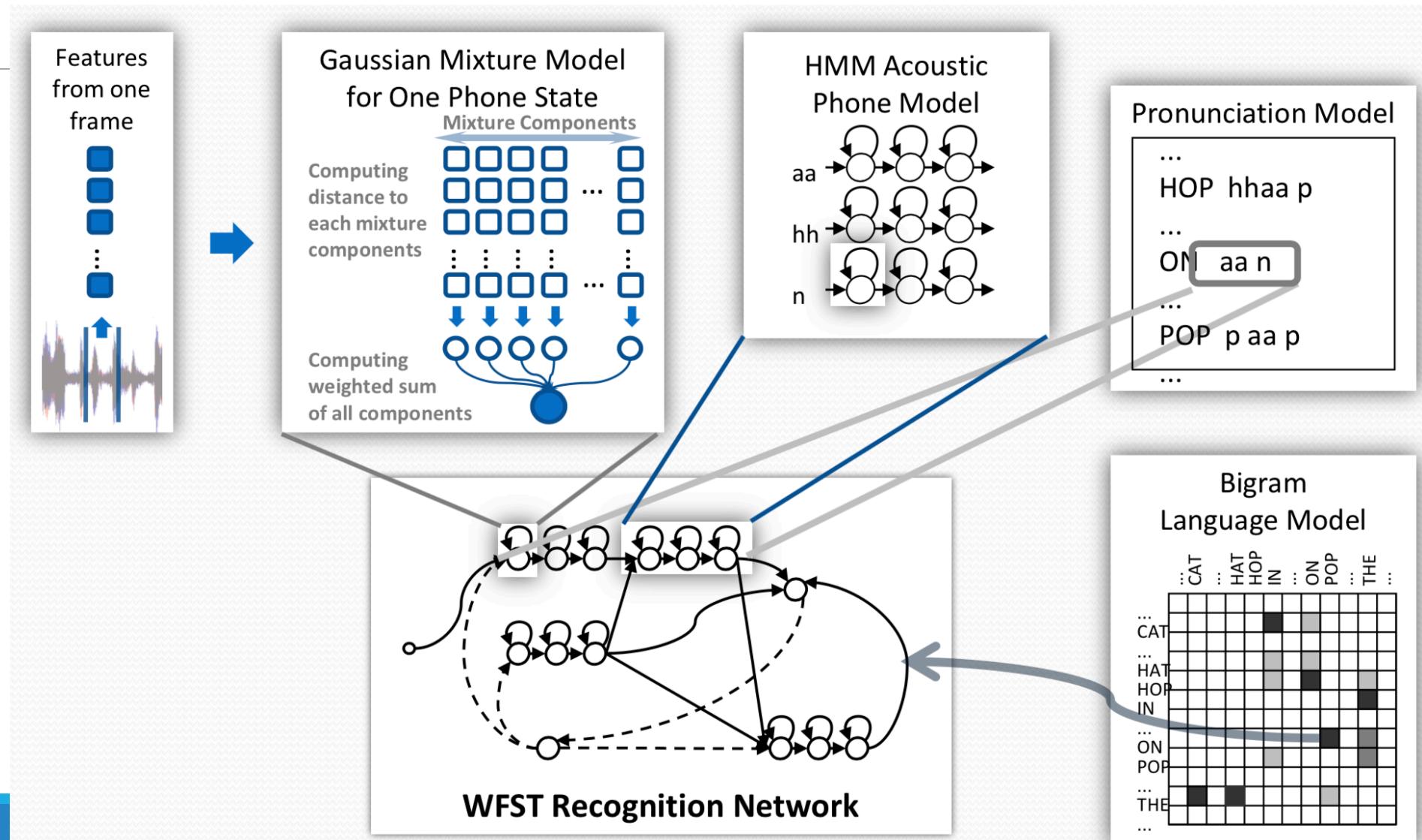
N -gram language models

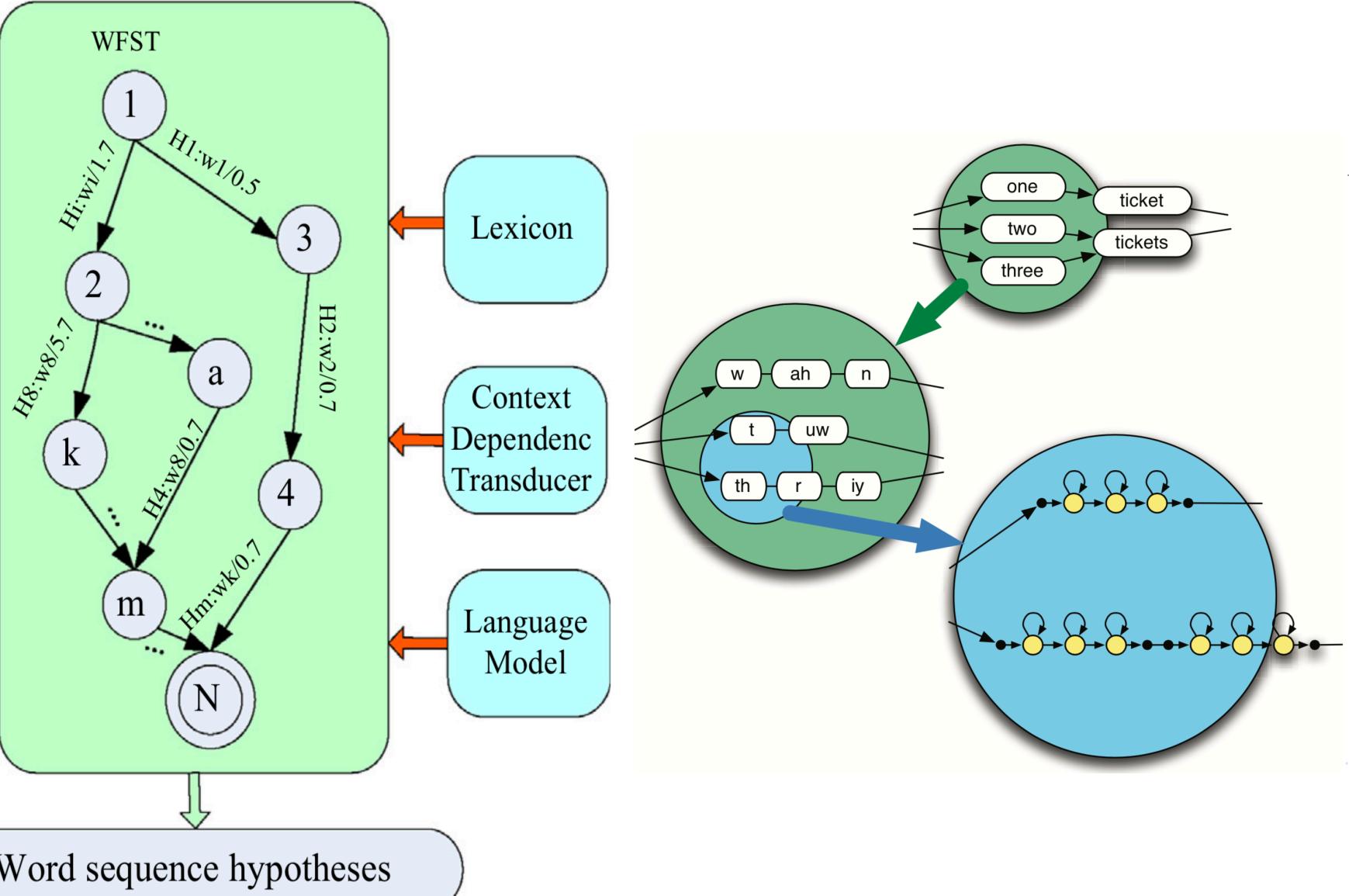
... the united states of ???



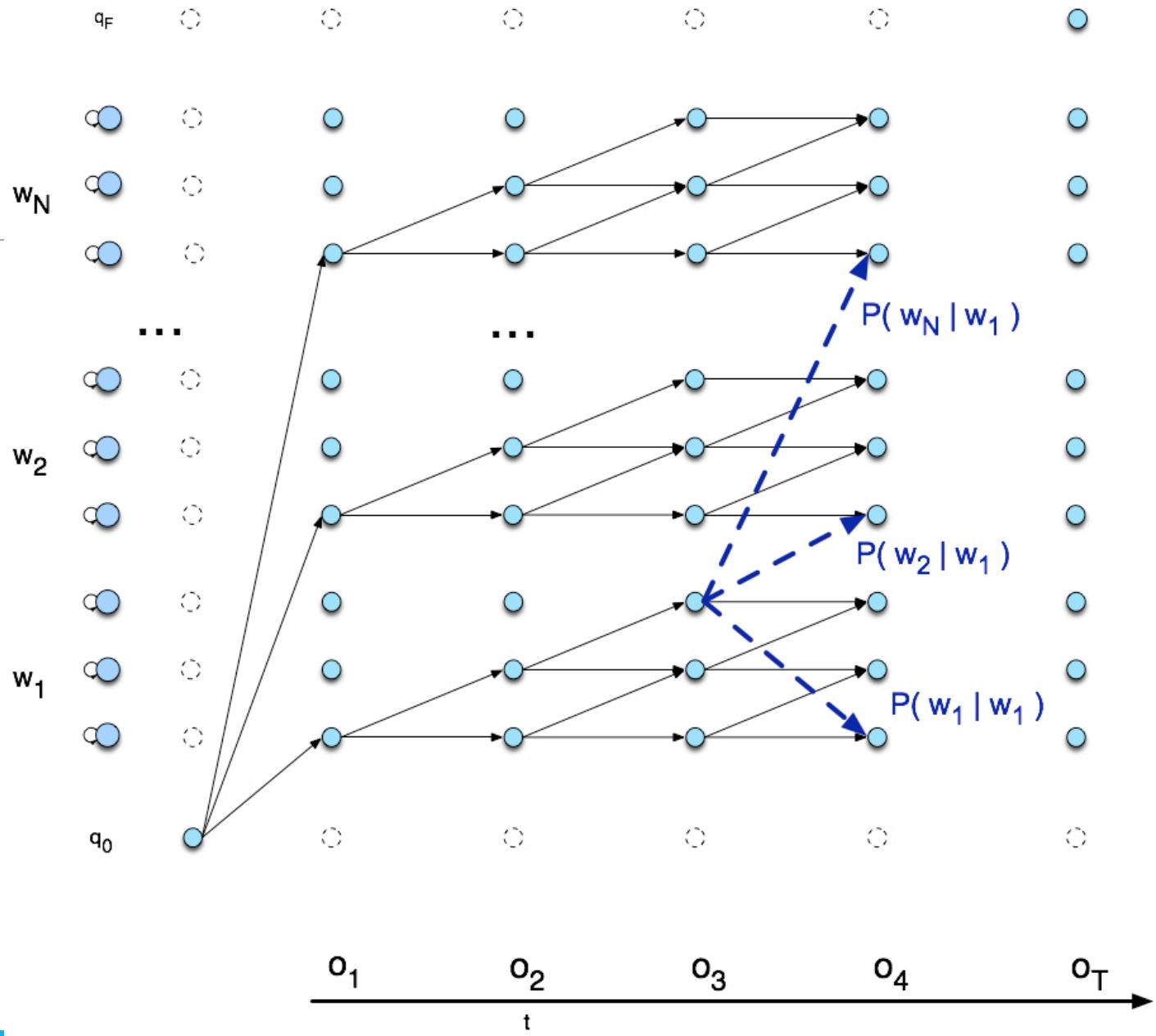
Predict the next word, give a set of predecessor words

Automatic Speech Recognizer

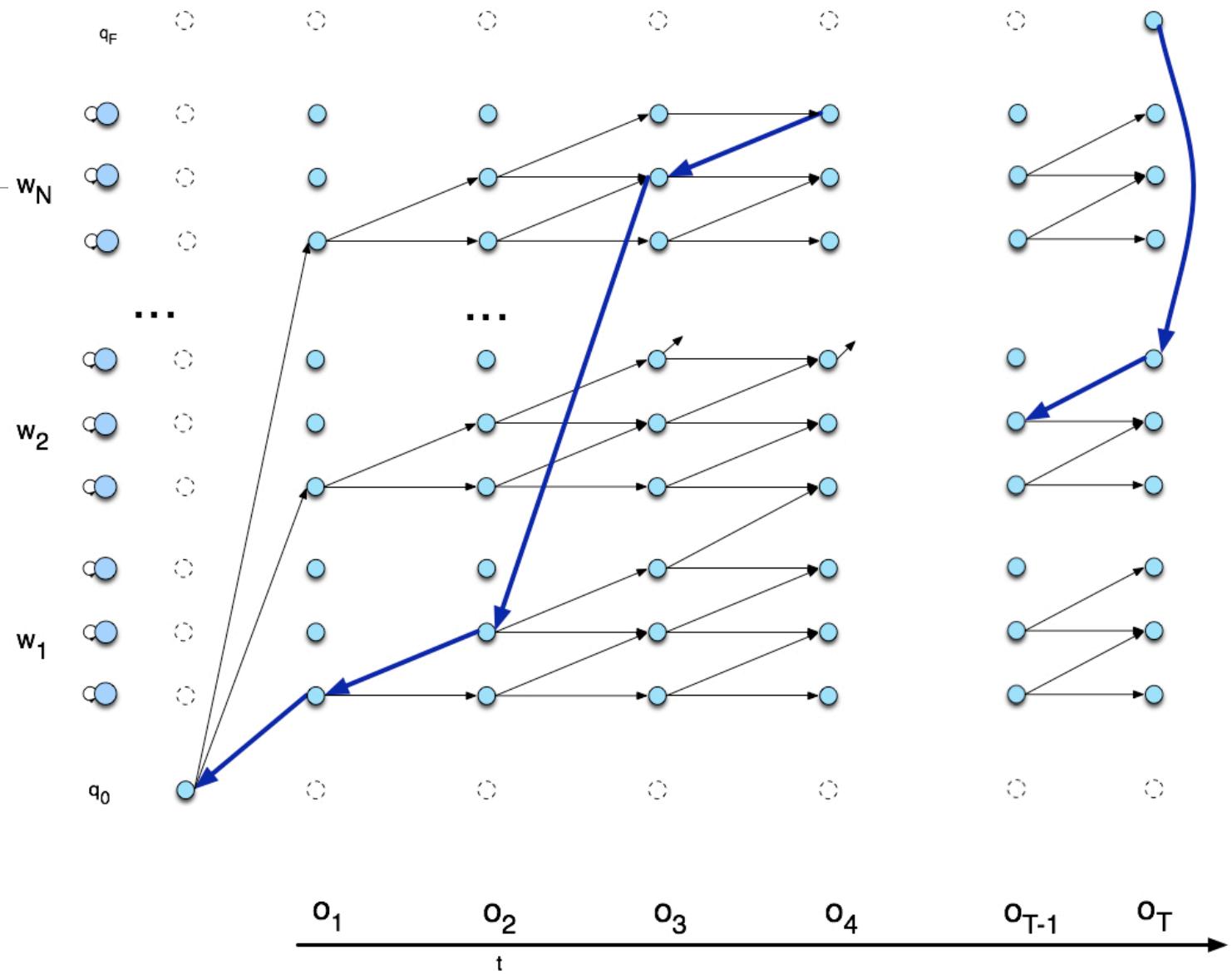




Viterbi Trellis

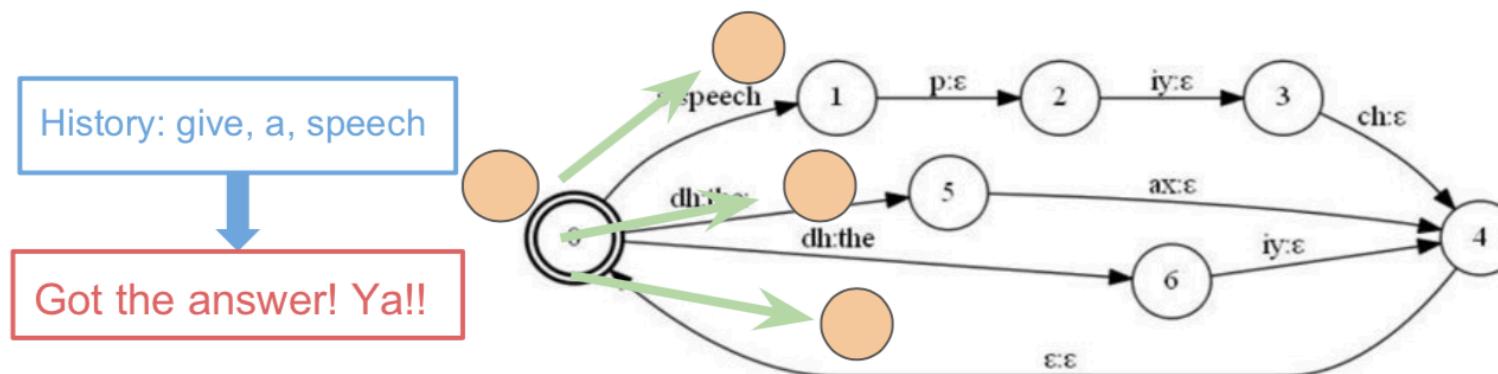


Viterbi Backtrace

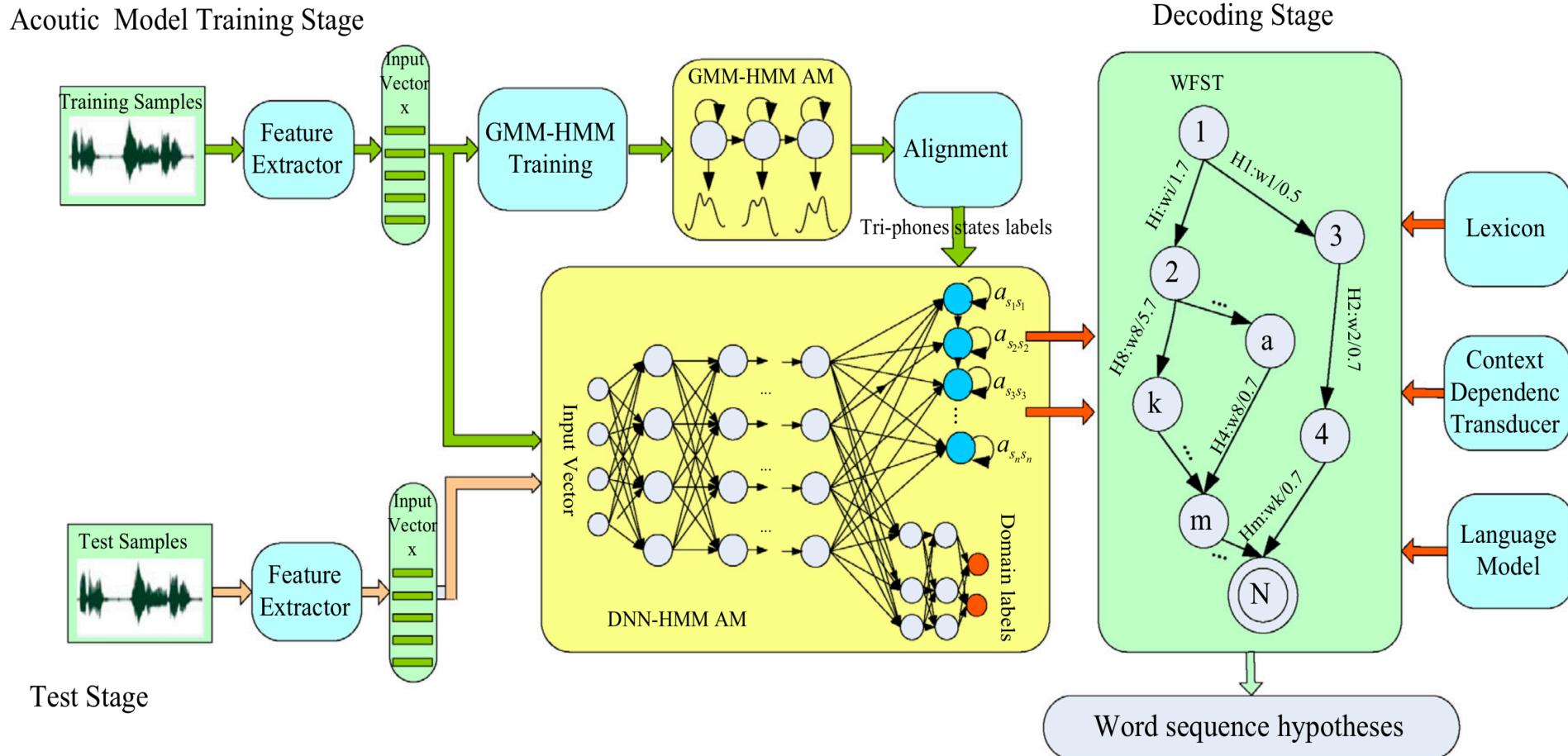


Token & Beam Pruning

- 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

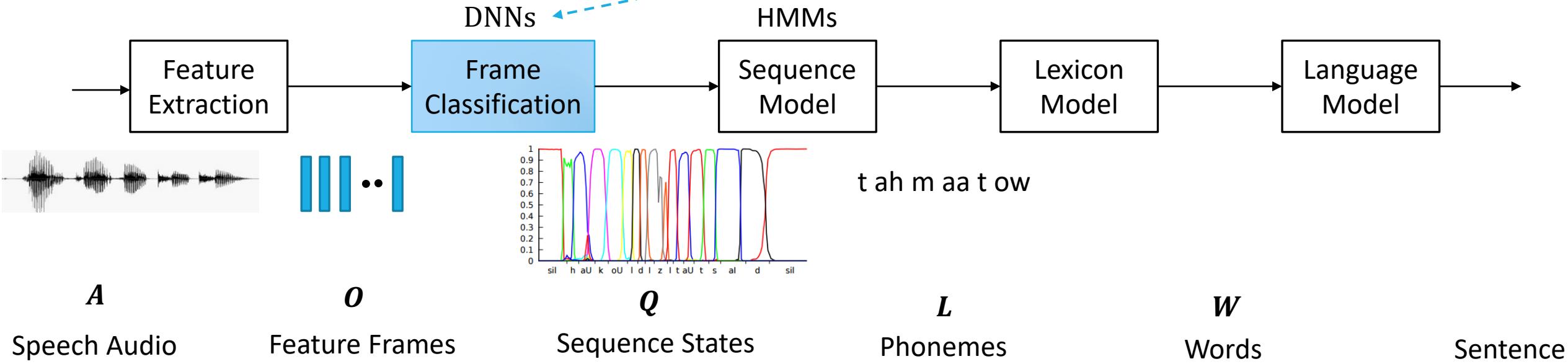


DNN-HMM

Kaldi

DNN: Deep Neural Networks

HMMs: Hidden Markov Models

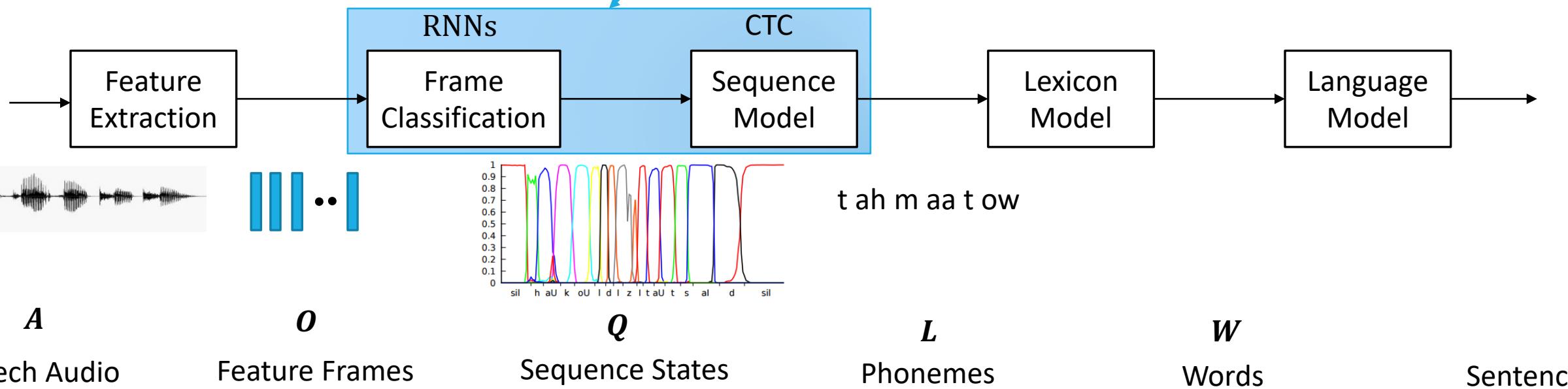


End-to-End: CTC

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

