

Introduction to Speech Recognition

YUAN-FU LIAO

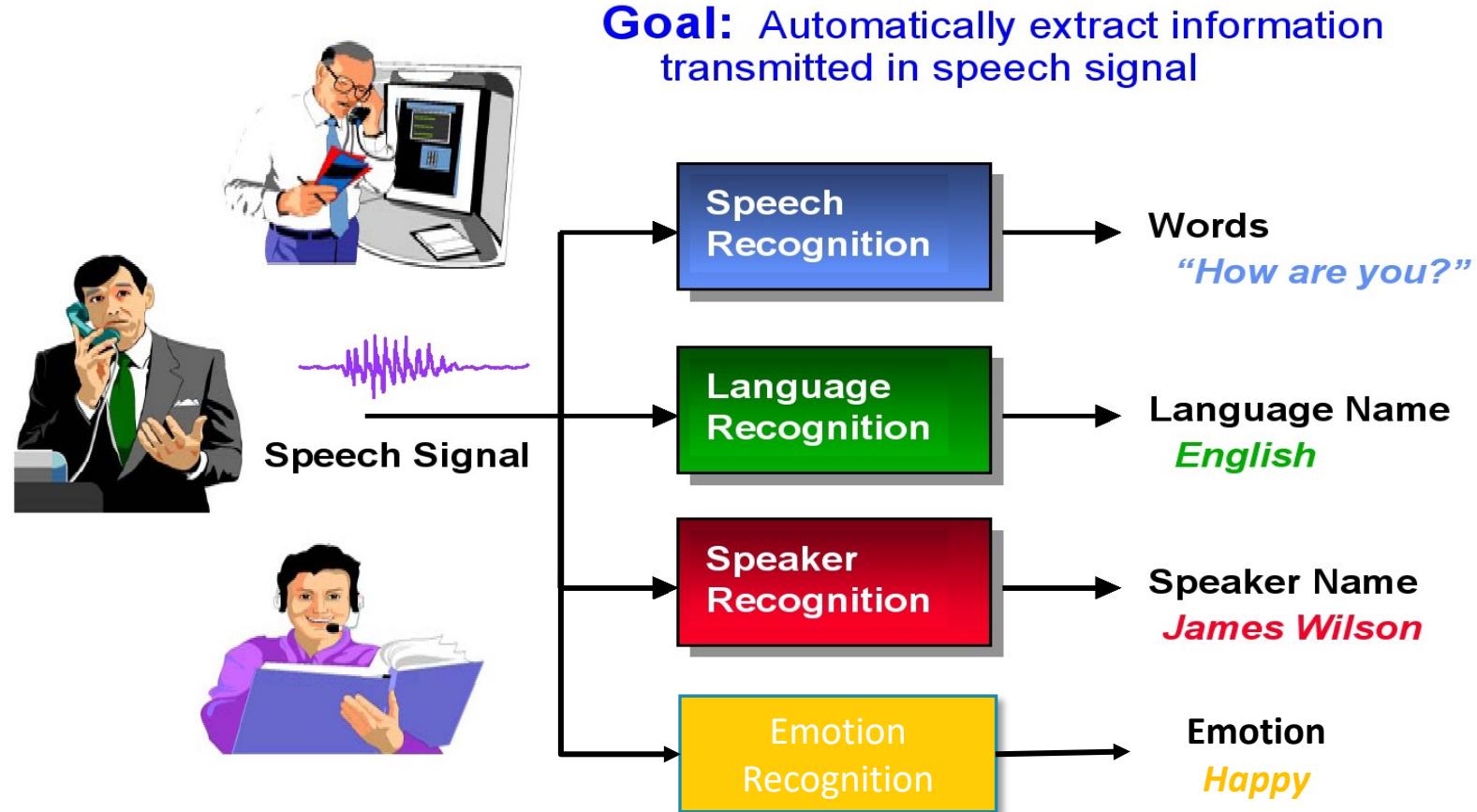
NATIONAL TAIPEI UNIVERSITY OF TECHNOLOGY

Outline

Automatic Speech Recognition (ASR)

- How does ASR work?
- General Architecture of Speech Recognition

Broad Objectives of Speech Recognition for Machines



Communication: A Thing Always Happens in Our Daily Life



What is Behind the Communication?

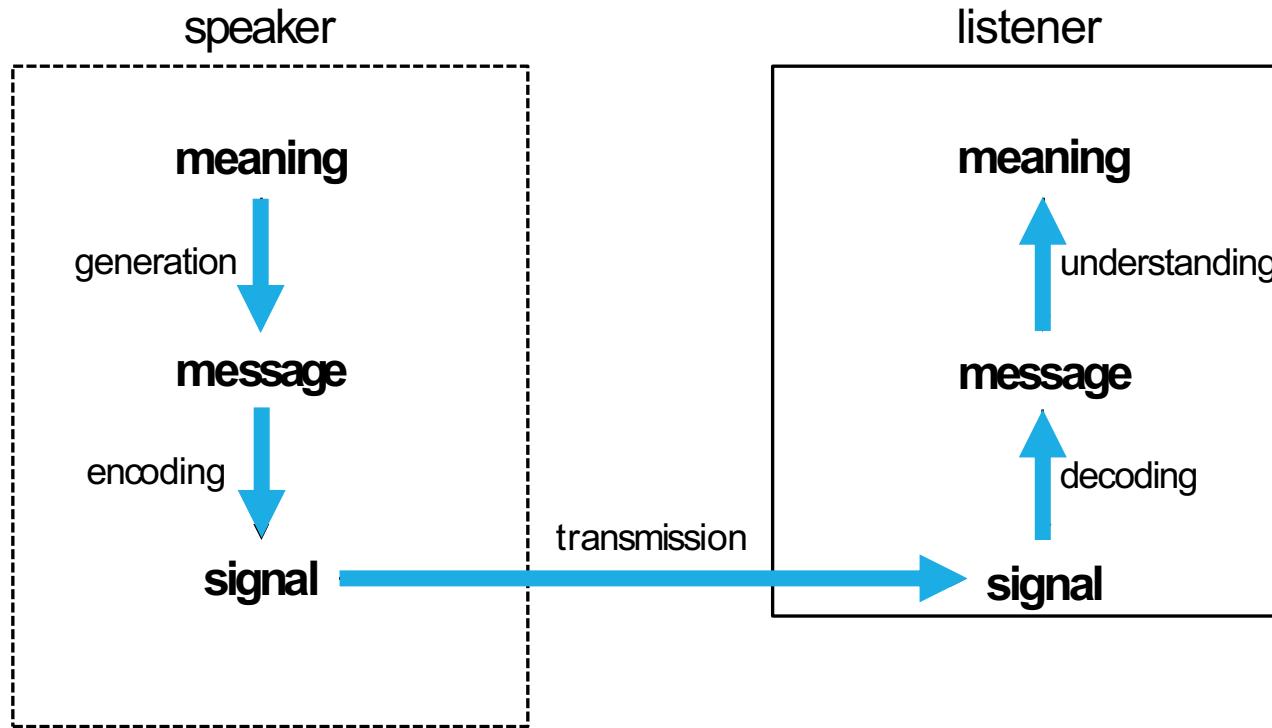
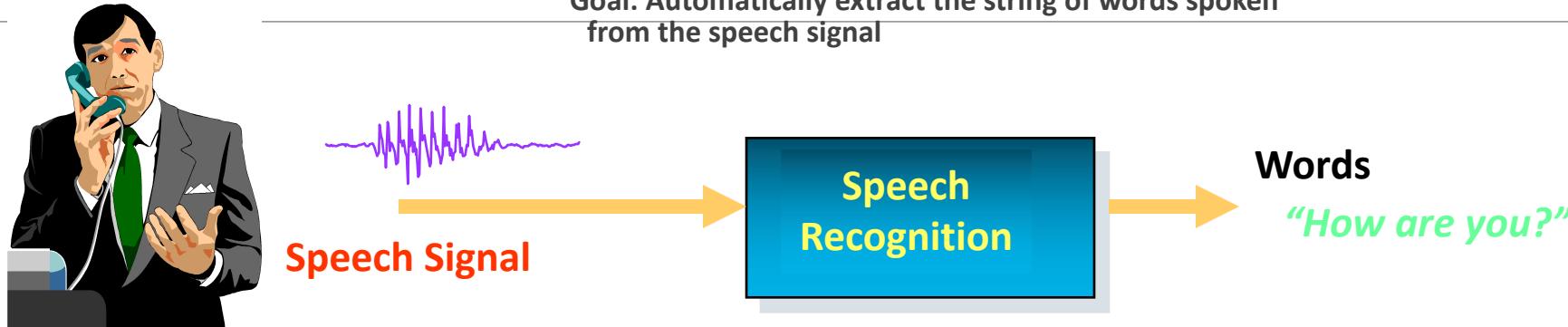
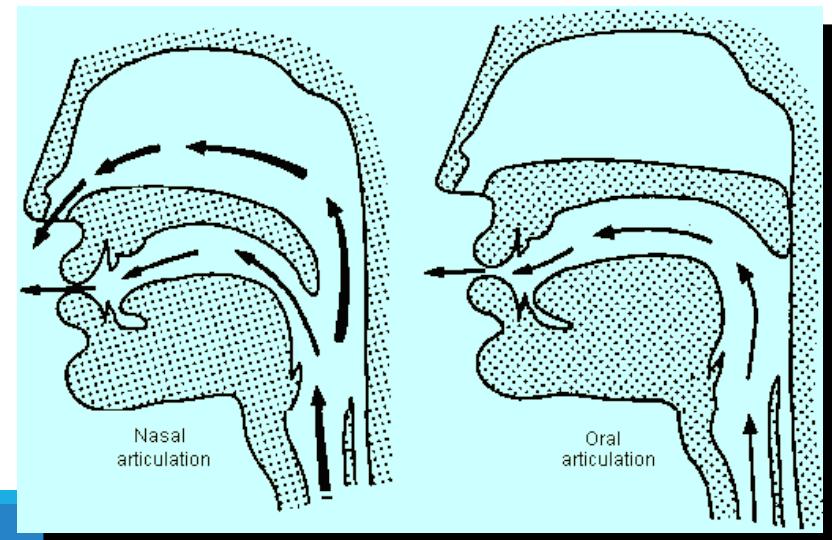


Figure. Processes involved in communication between two speakers

Speech Recognition



How is SPEECH produced?
⇒ Characteristics of
Acoustic Signal



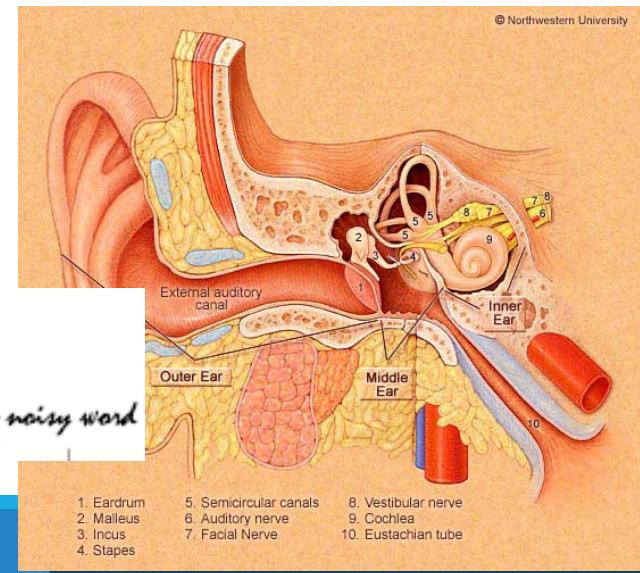
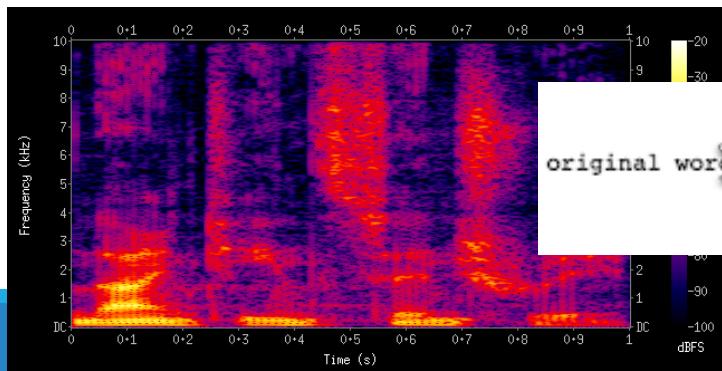
Speech Recognition



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



How is SPEECH perceived?
=> Important Features



Speech Recognition



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



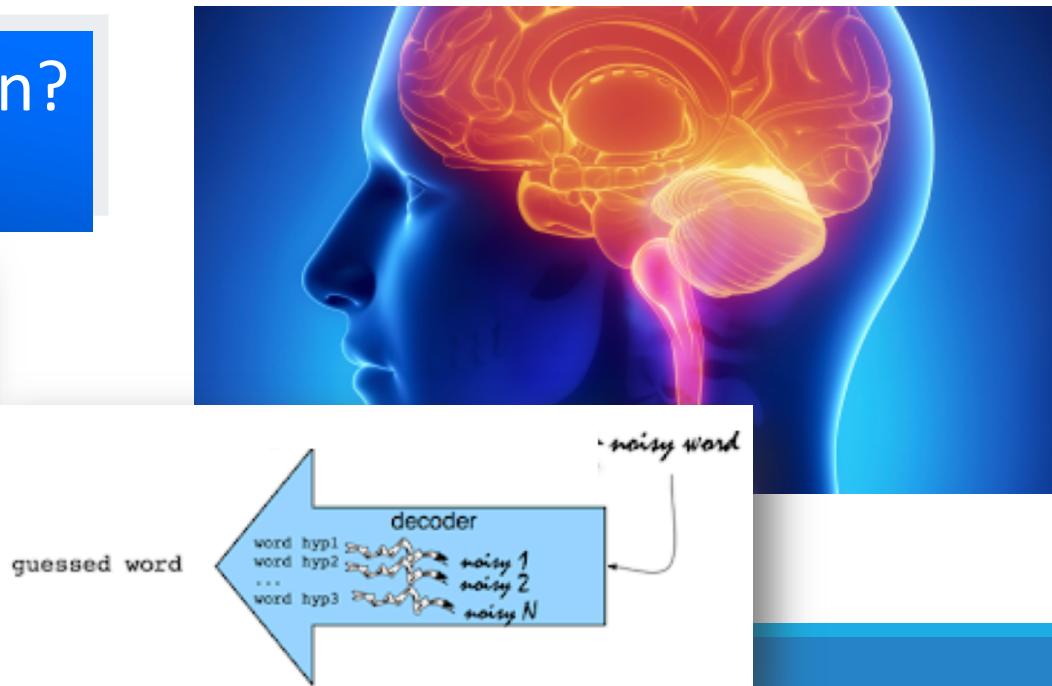
What Sentence is Spoken?
=> Language Model

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?'

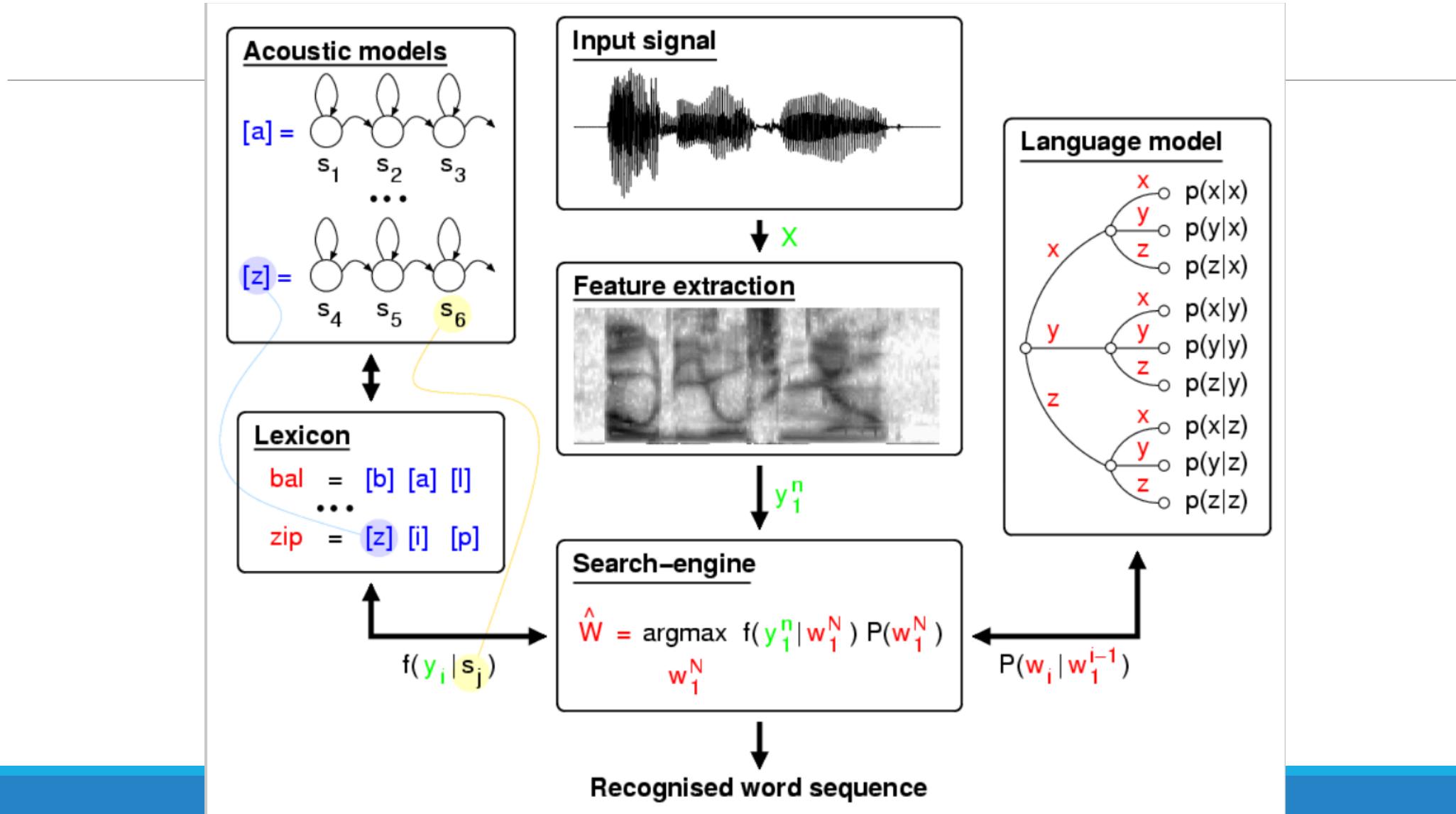
$P(w_{i+1} = \text{of} \mid w_i = \text{tired}) = 1$
 $P(w_{i+1} = \text{of} \mid w_i = \text{use}) = 1$
 $P(w_{i+1} = \text{sister} \mid w_i = \text{her}) = 1$
 $P(w_{i+1} = \text{beginning} \mid w_i = \text{was}) = 1/2$
 $P(w_{i+1} = \text{reading} \mid w_i = \text{was}) = 1/2$

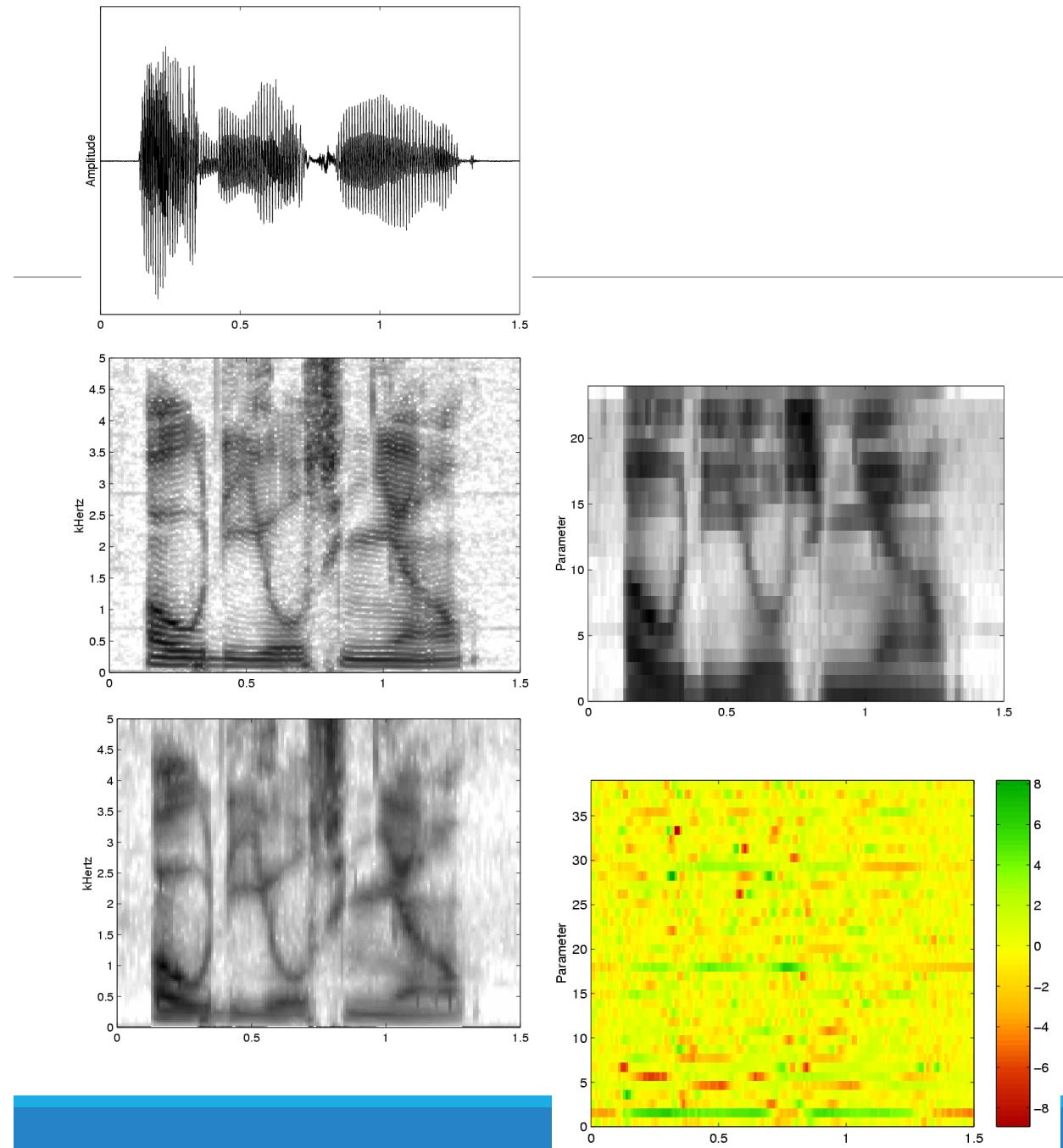
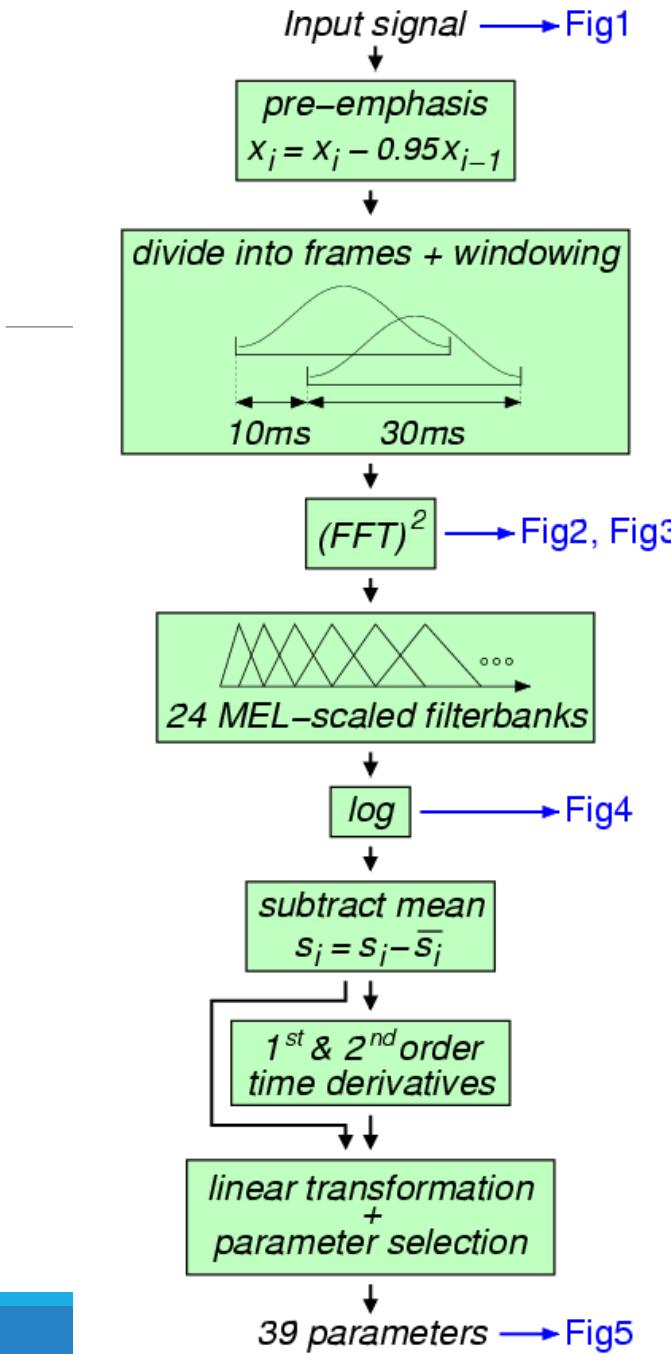
$P(w_{i+1} = \text{bank} \mid w_i = \text{the}) = 1/3$
 $P(w_{i+1} = \text{book} \mid w_i = \text{the}) = 1/3$
 $P(w_{i+1} = \text{use} \mid w_i = \text{the}) = 1/3$

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

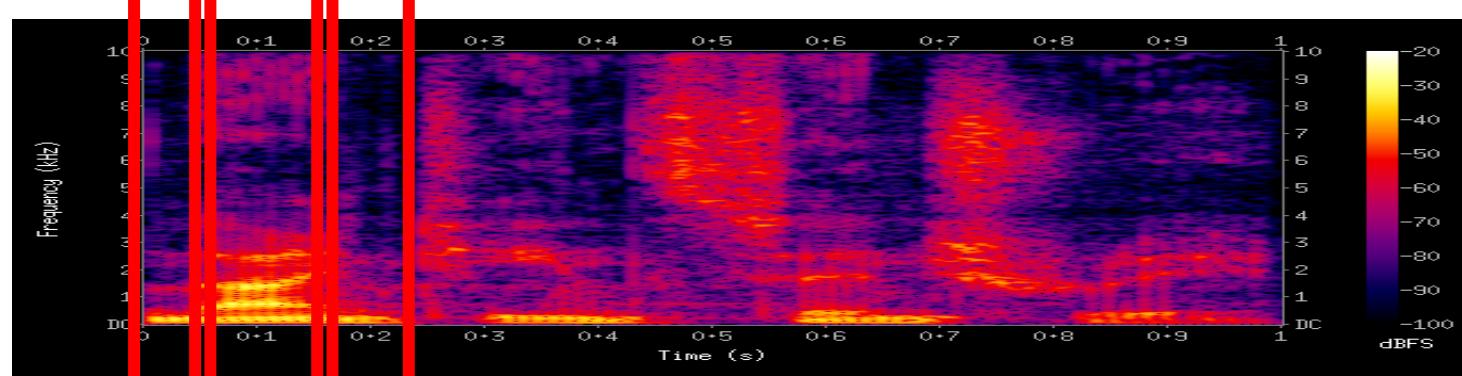
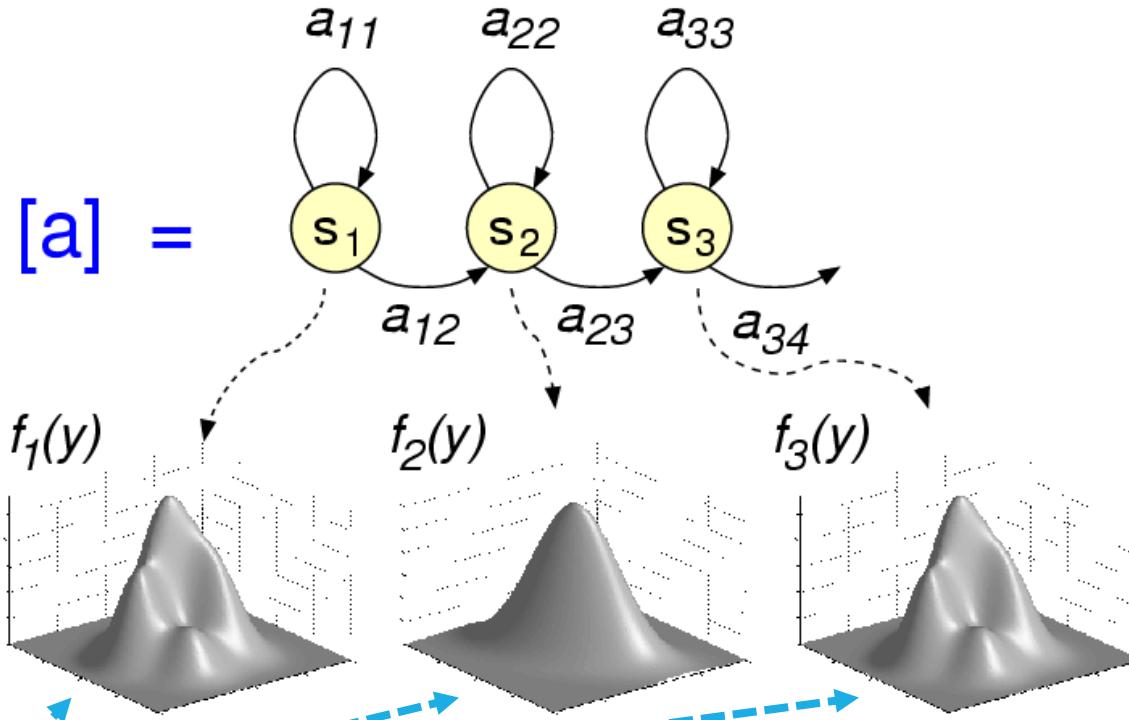


Automatic Speech Recognition: General Architecture

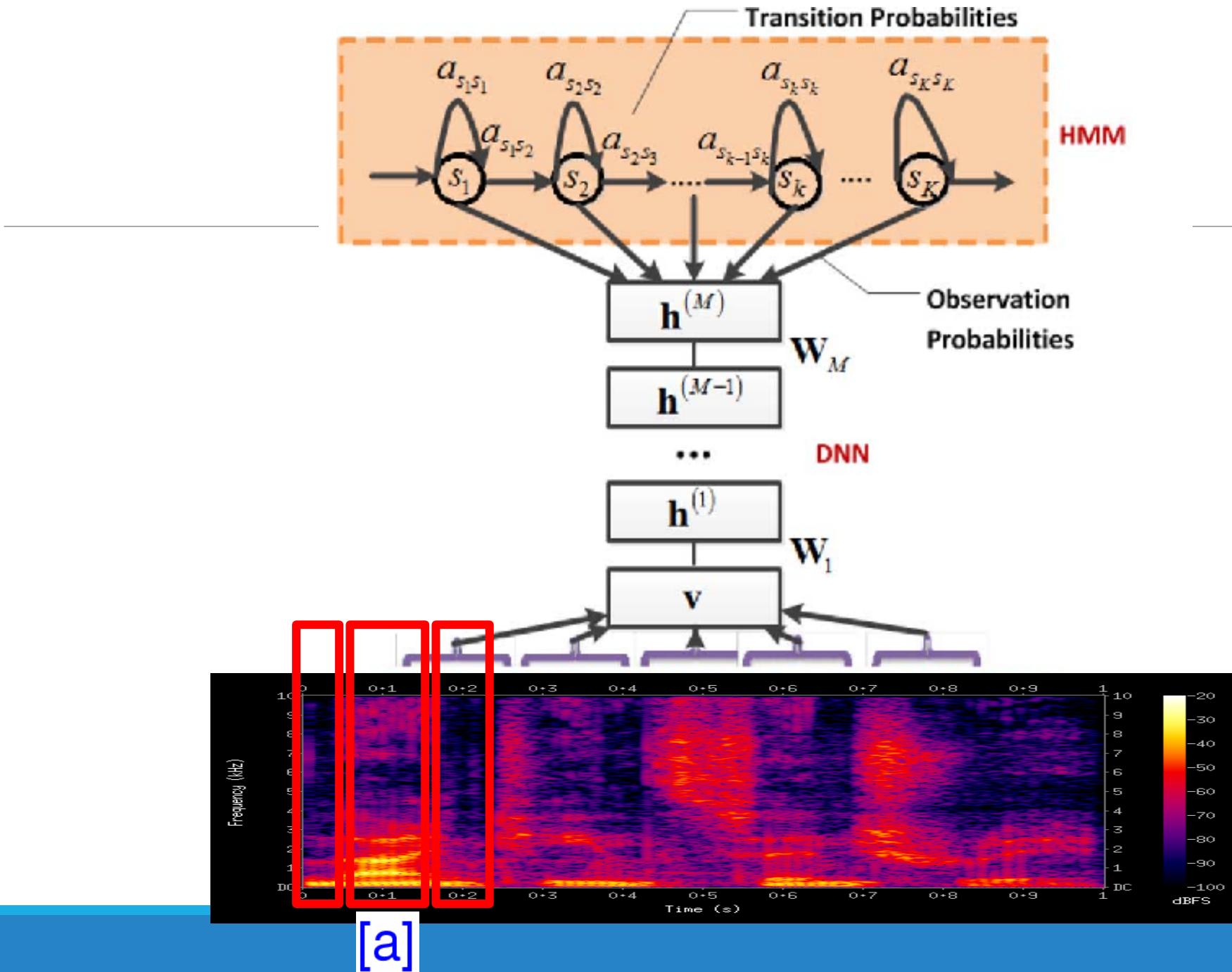




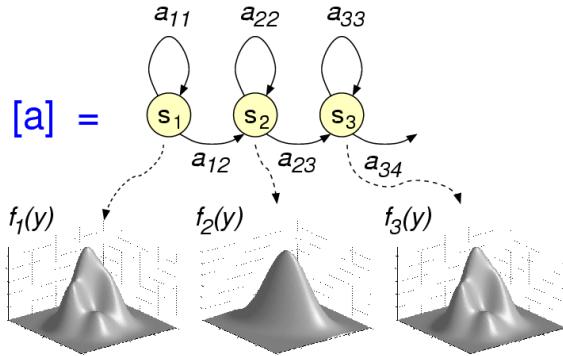
Hidden Markov Models



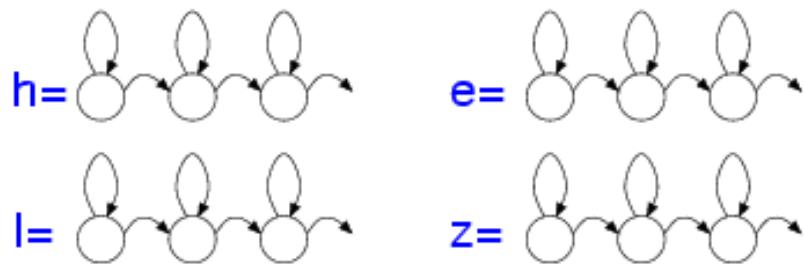
[a]



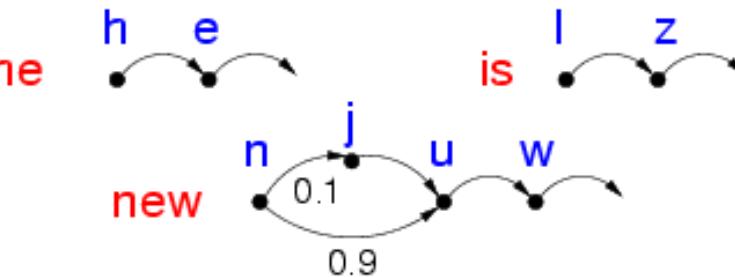
Hidden Markov Models



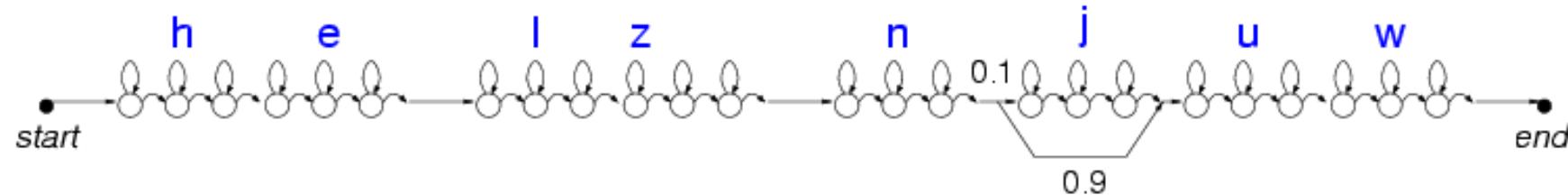
HMM phone models



Lexicon



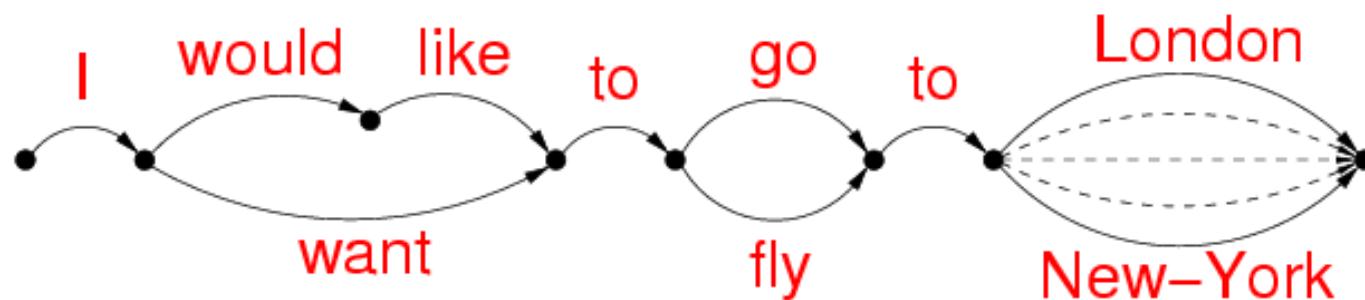
Sentence model: 'he is new'



Grammar:

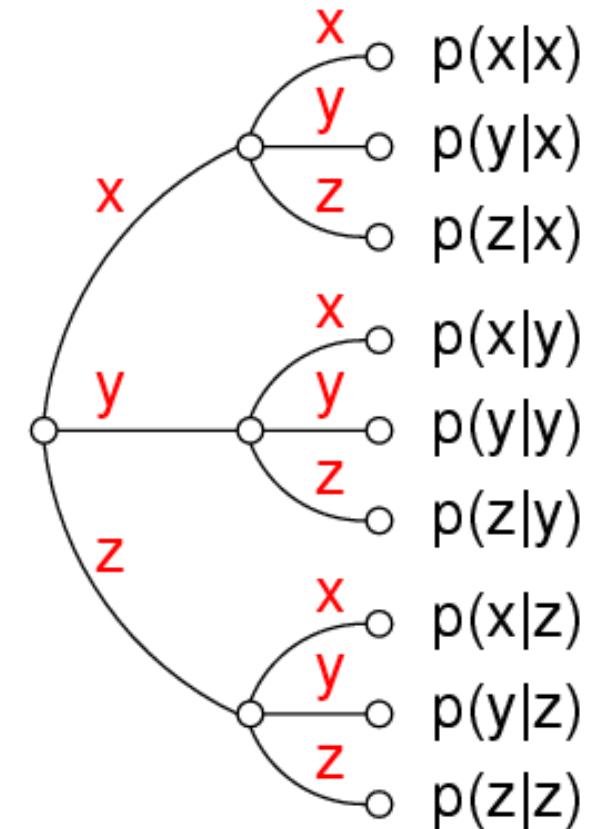
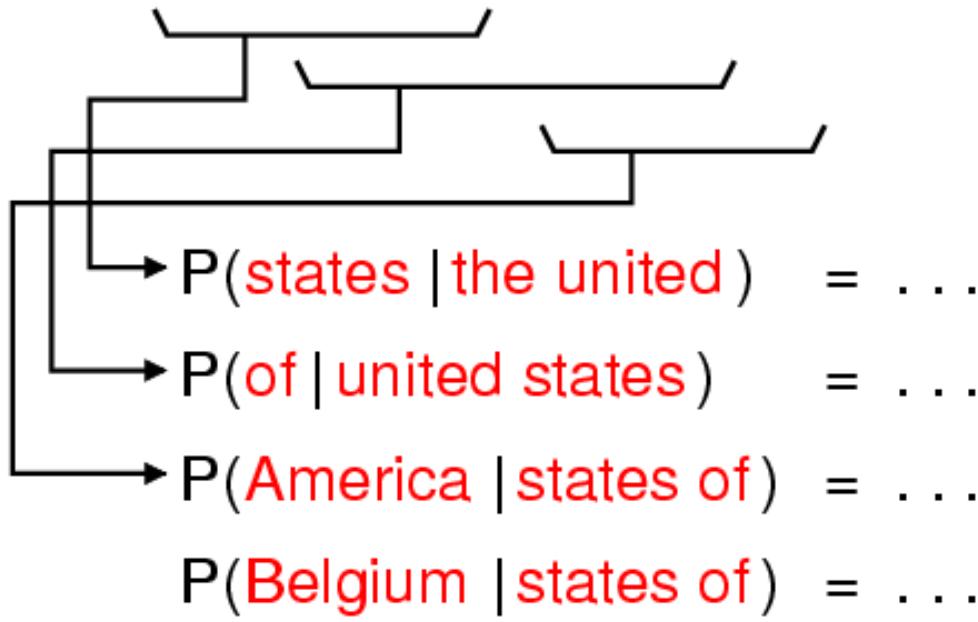
$$\langle \text{sentence}_1 \rangle = I \left\{ \begin{array}{c} \text{would like} \\ \text{want} \end{array} \right\} \text{to} \left\{ \begin{array}{c} \text{go} \\ \text{fly} \end{array} \right\} \text{to} \langle \text{airport} \rangle$$
$$\langle \text{sentence}_2 \rangle = \dots$$
$$\langle \text{airport}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

Finite-state representation:



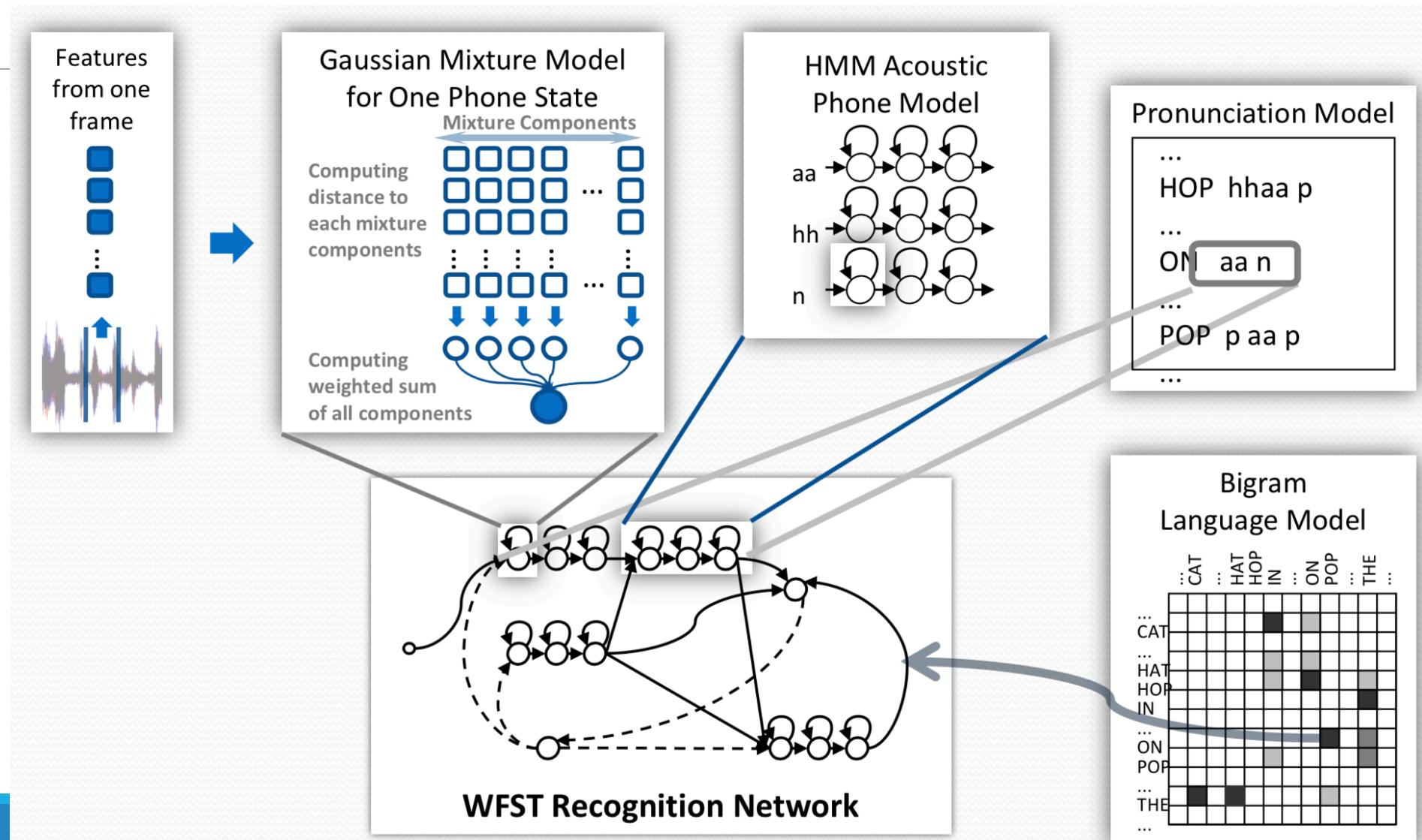
N -gram language models

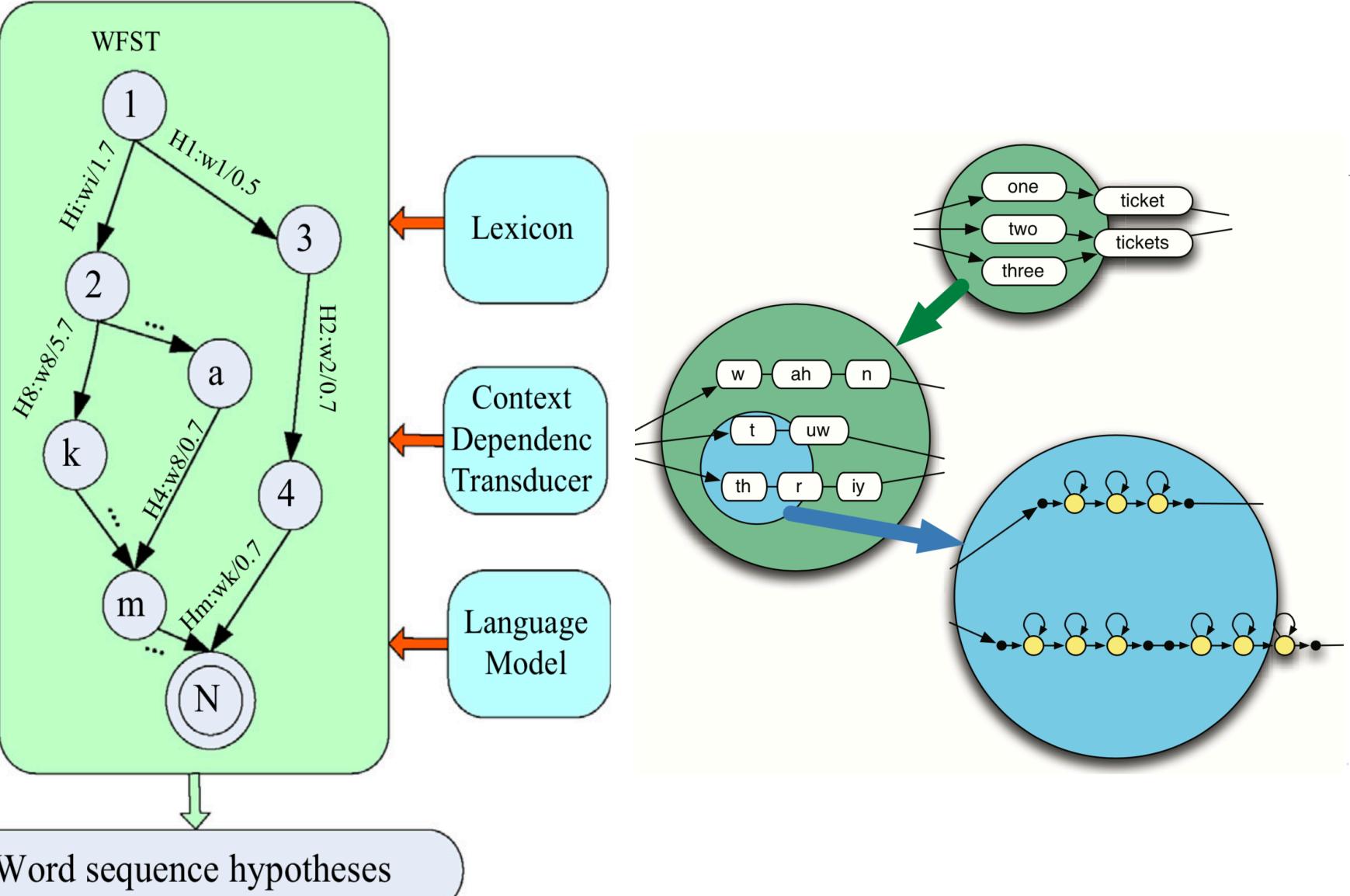
... the united states of ???



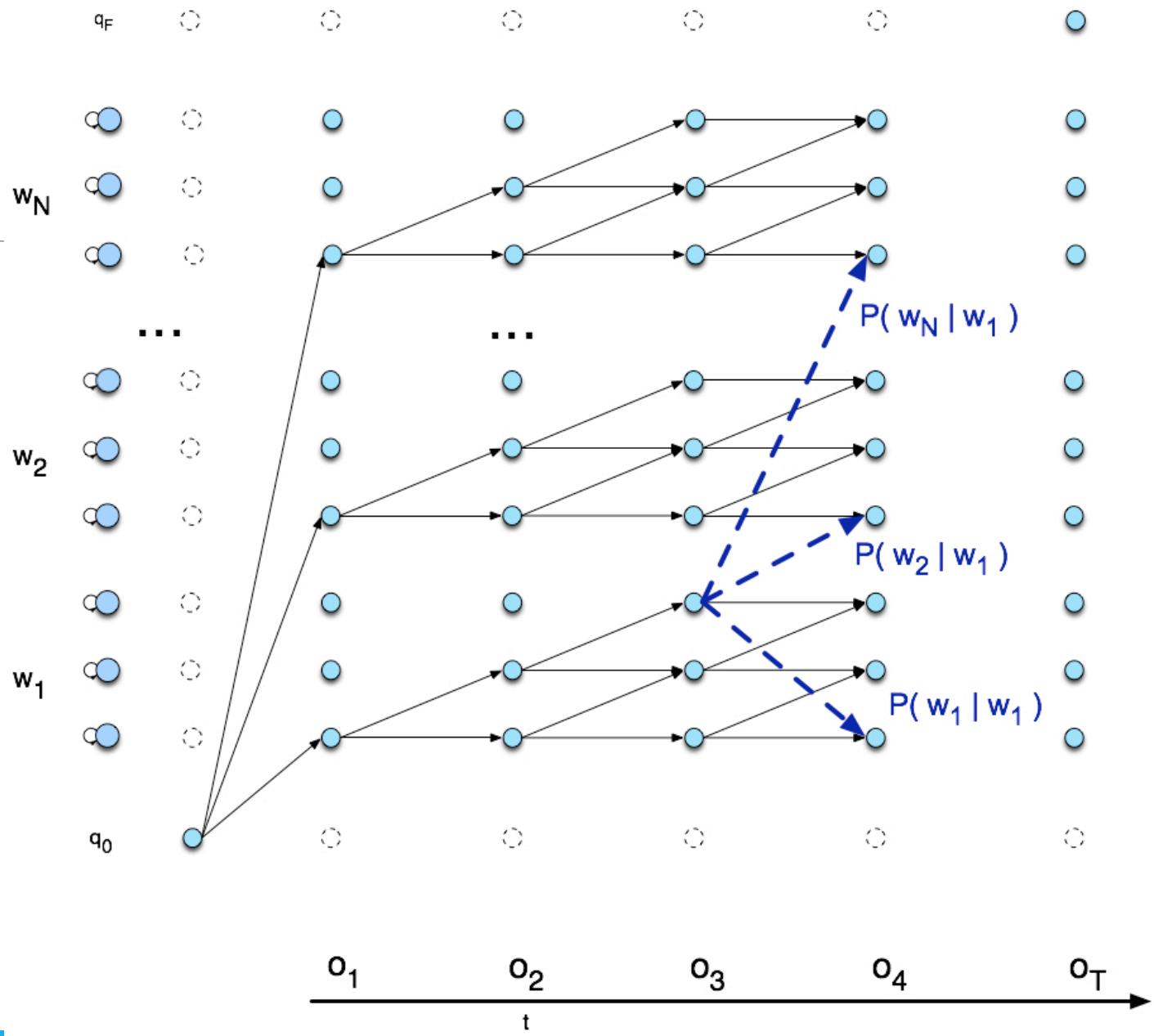
Predict the next word, give a set of predecessor words

Automatic Speech Recognizer

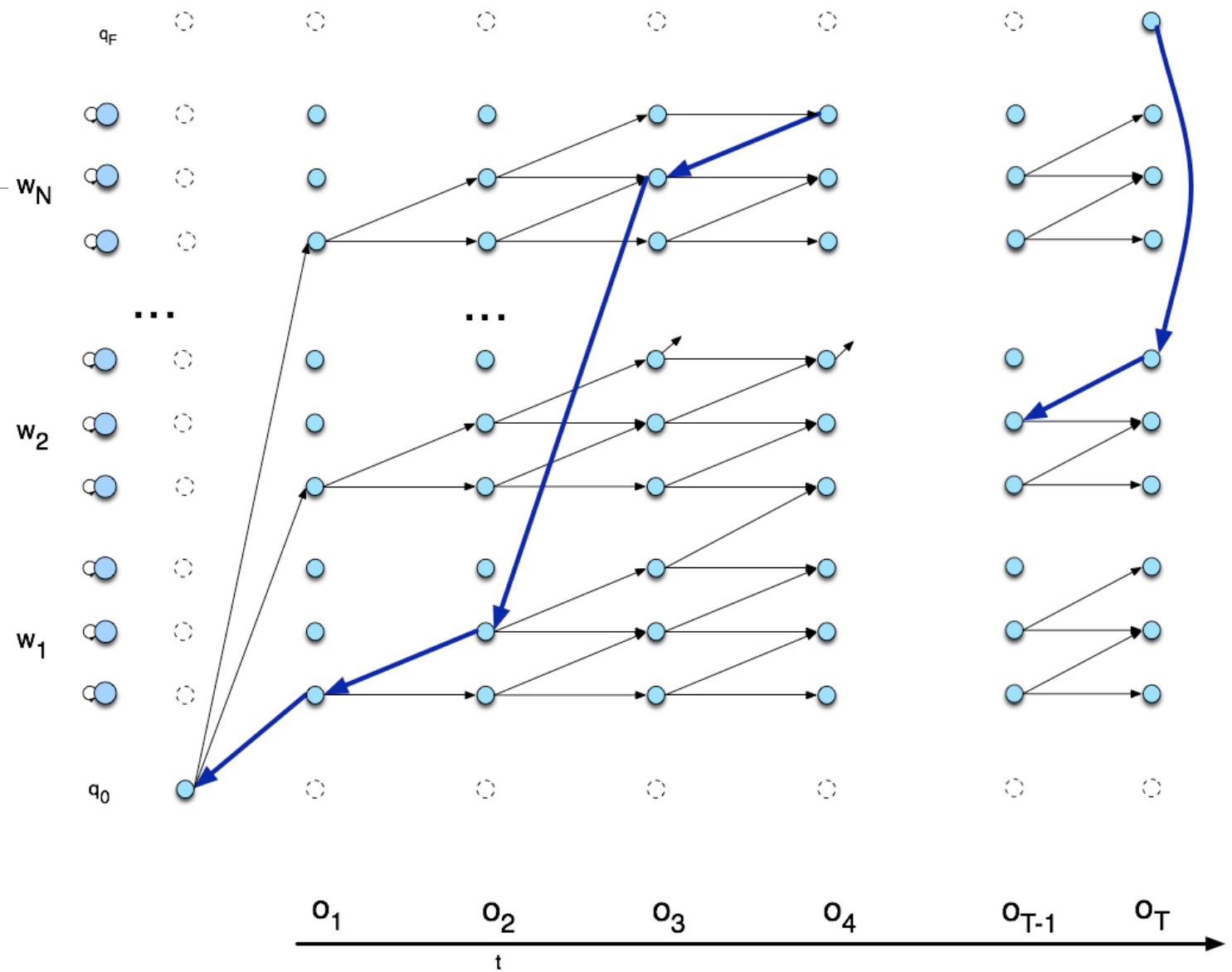




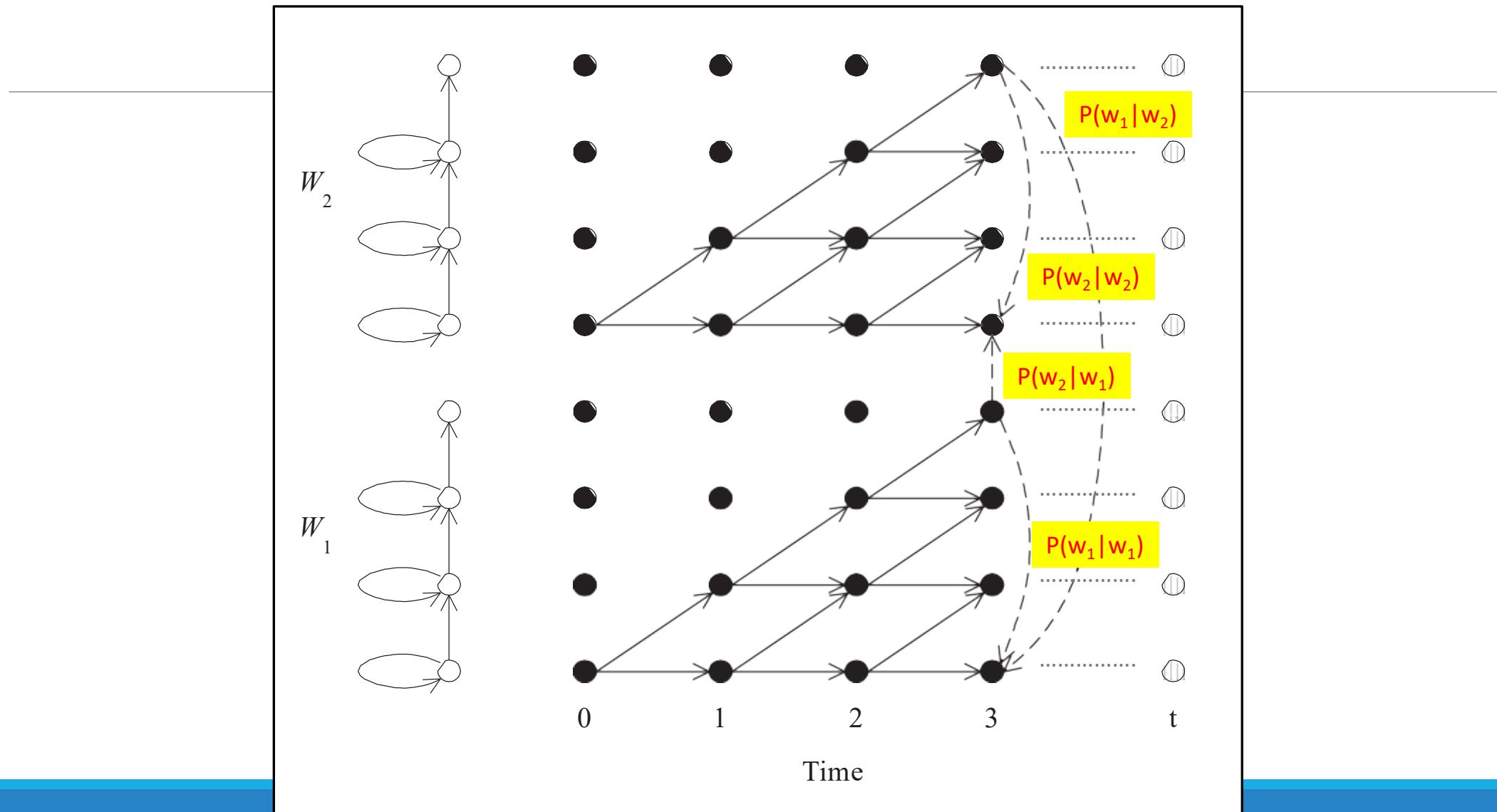
Viterbi Trellis



Viterbi Backtrace

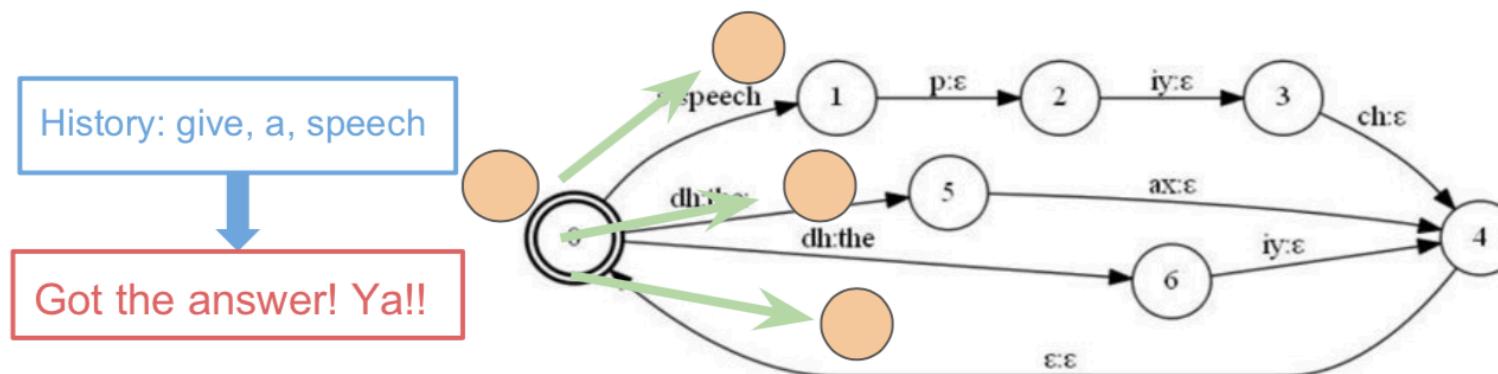


Viterbi Path with Complex Models



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

