

## LSA 352: Speech Recognition and Synthesis

Dan Jurafsky

### Lecture 5: Intro to ASR+HMMs: Forward, Viterbi, Baum-Welch

IP Notice:

LSA 352 Summer 2007

## Outline for Today

- Speech Recognition Architectural Overview
- Hidden Markov Models in general
  - Forward
  - Viterbi Decoding
  - Baum-Welch
- Applying HMMs to speech
- How this fits into the ASR component of course
  - July 6: Language Modeling
  - July 19 (today): HMMs, Forward, Viterbi, Start of Baum-Welch (EM) training
  - July 23: Feature Extraction, MFCCs, and Gaussian Acoustic modeling
  - July 26: Evaluation, Decoding, Advanced Topics

LSA 352 Summer 2007

## LVCSR

- Large Vocabulary Continuous Speech Recognition
- ~20,000-64,000 words
- Speaker independent (vs. speaker-dependent)
- Continuous speech (vs isolated-word)

LSA 352 Summer 2007

## Current error rates

Ballpark numbers; exact numbers depend very much on the specific corpus

Task	Vocabulary	Error Rate%
Digits	11	0.5
WSJ read speech	5K	3
WSJ read speech	20K	3
Broadcast news	64,000+	10
Conversational Telephone	64,000+	20

LSA 352 Summer 2007

## HSR versus ASR

Task	Vocab	ASR	Hum SR
Continuous digits	11	.5	.009
WSJ 1995 clean	5K	3	0.9
WSJ 1995 w/noise	5K	9	1.1
SWBD 2004	65K	20	4

- Conclusions:
  - Machines about 5 times worse than humans
  - Gap increases with noisy speech
  - These numbers are rough, take with grain of salt

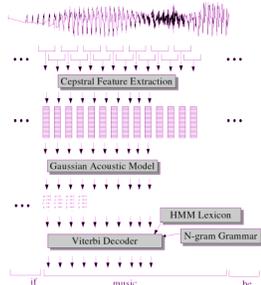
LSA 352 Summer 2007

## LVCSR Design Intuition

- Build a statistical model of the speech-to-words process
- Collect lots and lots of speech, and transcribe all the words.
- Train the model on the labeled speech
- Paradigm: Supervised Machine Learning + Search

LSA 352 Summer 2007

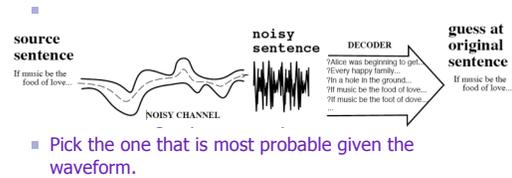
## Speech Recognition Architecture



LSA 352 Summer 2007

7

## The Noisy Channel Model



- Pick the one that is most probable given the waveform.

LSA 352 Summer 2007

8

## The Noisy Channel Model (II)

- What is the most likely sentence out of all sentences in the language  $L$  given some acoustic input  $O$ ?
- Treat acoustic input  $O$  as sequence of individual observations
  - $O = o_1, o_2, o_3, \dots, o_t$
- Define a sentence as a sequence of words:
  - $W = w_1, w_2, w_3, \dots, w_n$

LSA 352 Summer 2007

9

## Noisy Channel Model (III)

- Probabilistic implication: Pick the highest prob  $S$ :

$$\hat{W} = \arg \max_{W \in L} P(W | O)$$

- We can use Bayes rule to rewrite this:

$$\hat{W} = \arg \max_{W \in L} \frac{P(O|W)P(W)}{P(O)}$$

- Since denominator is the same for each candidate sentence  $W$ , we can ignore it for the argmax:

$$\hat{W} = \arg \max_{W \in L} P(O|W)P(W)$$

LSA 352 Summer 2007

10

## Noisy channel model

$$\hat{W} = \arg \max_{W \in L} P(O|W)P(W)$$

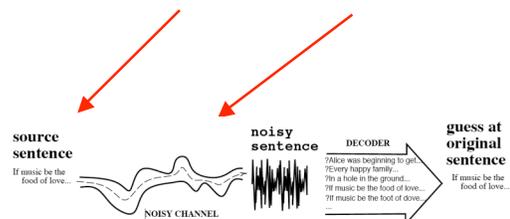
likelihood      prior  
↓                    ↓

LSA 352 Summer 2007

11

## The noisy channel model

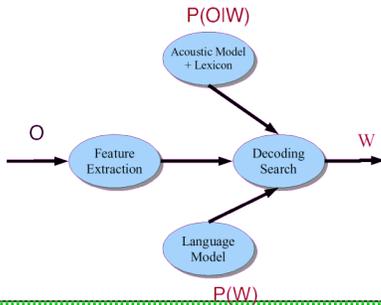
- Ignoring the denominator leaves us with two factors:  $P(\text{Source})$  and  $P(\text{Signal}|\text{Source})$



LSA 352 Summer 2007

12

## Speech Architecture meets Noisy Channel



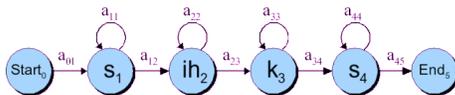
LSA 352 Summer 2007 14

## Architecture: Five easy pieces (only 2 for today)

- Feature extraction
- Acoustic Modeling
- HMMs, Lexicons, and Pronunciation
- Decoding
- Language Modeling

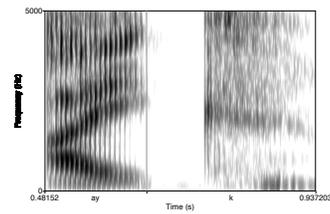
LSA 352 Summer 2007 14

## HMMs for speech



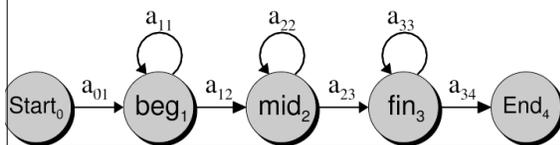
LSA 352 Summer 2007 15

## Phones are not homogeneous!



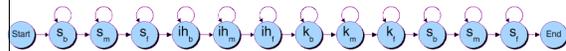
LSA 352 Summer 2007 16

## Each phone has 3 subphones



LSA 352 Summer 2007 17

## Resulting HMM word model for "six"



LSA 352 Summer 2007 18

## HMMs more formally

- Markov chains
- A kind of weighted finite-state automaton

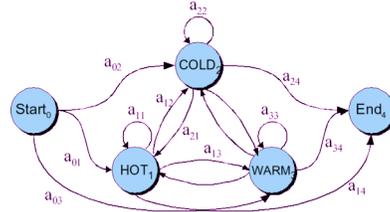
$Q = q_1 q_2 \dots q_N$  a set of **states**  
 $A = a_{01} a_{02} \dots a_{n1} \dots a_{nm}$  a **transition probability matrix**  $A$ , each  $a_{ij}$  representing the probability of moving from state  $i$  to state  $j$ , s.t.  $\sum_{j=1}^n a_{ij} = 1 \quad \forall i$   
 $q_0, q_{end}$  a special **start and end state** which are not associated with observations.

**Markov Assumption:**  $P(q_i | q_1 \dots q_{i-1}) = P(q_i | q_{i-1})$

LISA 352 Summer 2007 19

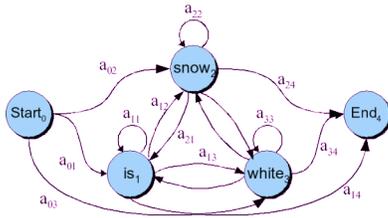
## HMMs more formally

- Markov chains
- A kind of weighted finite-state automaton



LISA 352 Summer 2007 20

## Another Markov chain

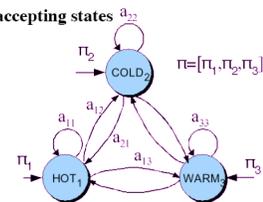


LISA 352 Summer 2007 21

## Another view of Markov chains

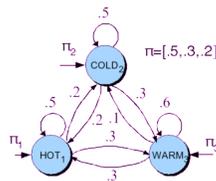
$\pi = \pi_1, \pi_2, \dots, \pi_N$  an **initial probability distribution** over states.  $\pi_i$  is the probability that the Markov chain will start in state  $i$ . Some states  $j$  may have  $\pi_j = 0$ , meaning that they cannot be initial states. Also,  $\sum_{i=1}^n \pi_i = 1$

$QA = \{q_1, q_2, \dots\}$  a set  $QA \subset Q$  of legal **accepting states**



LISA 352 Summer 2007 22

## An example with numbers:



- What is probability of:
  - Hot hot hot hot
  - Cold hot cold hot

LISA 352 Summer 2007 23

## Hidden Markov Models

$Q = q_1 q_2 \dots q_N$  a set of  $N$  **states**  
 $A = a_{11} a_{12} \dots a_{n1} \dots a_{nm}$  a **transition probability matrix**  $A$ , each  $a_{ij}$  representing the probability of moving from state  $i$  to state  $j$ , s.t.  $\sum_{j=1}^n a_{ij} = 1 \quad \forall i$

$O = o_1 o_2 \dots o_T$  a sequence of  $T$  **observations**, each one drawn from a vocabulary  $V = v_1, v_2, \dots, v_V$ .

$B = b_1(o_t)$  A sequence of **observation likelihoods**, also called **emission probabilities**, each expressing the probability of an observation  $o_t$  being generated from a state  $i$ .

$q_0, q_F$  a special **start state** and **end (final) state** which are not associated with observations, together with transition probabilities  $a_{01} a_{02} \dots a_{0n}$  out of the start state and  $a_{1F} a_{2F} \dots a_{nF}$  into the end state.

LISA 352 Summer 2007 24

## Hidden Markov Models

**Markov Assumption:**  $P(q_i|q_1...q_{i-1}) = P(q_i|q_{i-1})$

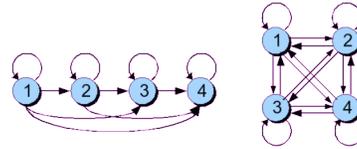
**Output Independence Assumption:**  $P(o_i|q_1...q_i, o_1, \dots, o_{i-1}, \dots, o_n) = P(o_i|q_i)$

LSA 352 Summer 2007

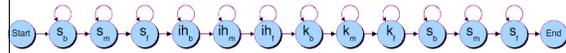
24

## Hidden Markov Models

- **Bakis network** Ergodic (fully-connected) network



- **Left-to-right network**



LSA 352 Summer 2007

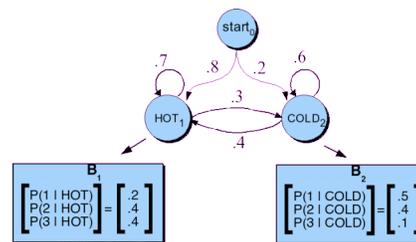
25

## The Jason Eisner task

- You are a climatologist in 2799 studying the history of global warming
- YOU can't find records of the weather in Baltimore for summer 2006
- But you do find Jason Eisner's diary
- Which records how many ice creams he ate each day.
- Can we use this to figure out the weather?
  - Given a sequence of observations  $O$ ,
    - each observation an integer = number of ice creams eaten
    - Figure out correct hidden sequence  $Q$  of weather states (H or C) which caused Jason to eat the ice cream

LSA 352 Summer 2007

27



LSA 352 Summer 2007

28

## HMMs more formally

- **Three fundamental problems**
  - Jack Ferguson at IDA in the 1960s
- 1) Given a specific HMM, determine likelihood of observation sequence.
- 2) Given an observation sequence and an HMM, discover the best (most probable) hidden state sequence
- 3) Given only an observation sequence, learn the HMM parameters (A, B matrix)

LSA 352 Summer 2007

29

## The Three Basic Problems for HMMs

- **Problem 1 (Evaluation):** Given the observation sequence  $O=(o_1o_2...o_T)$ , and an HMM model  $\Phi = (A,B)$ , how do we efficiently compute  $P(O | \Phi)$ , the probability of the observation sequence, given the model
- **Problem 2 (Decoding):** Given the observation sequence  $O=(o_1o_2...o_T)$ , and an HMM model  $\Phi = (A,B)$ , how do we choose a corresponding state sequence  $Q=(q_1q_2...q_T)$  that is optimal in some sense (i.e., best explains the observations)
- **Problem 3 (Learning):** How do we adjust the model parameters  $\Phi = (A,B)$  to maximize  $P(O | \Phi)$ ?

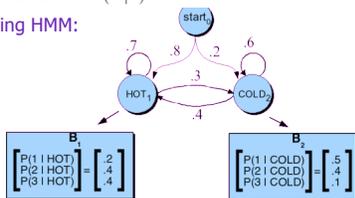
LSA 352 Summer 2007

30

## Problem 1: computing the observation likelihood

**Computing Likelihood:** Given an HMM  $\lambda = (A, B)$  and an observation sequence  $O$ , determine the likelihood  $P(O|\lambda)$ .

- Given the following HMM:



- How likely is the sequence 3 1 3?

LSA 352 Summer 2007

## How to compute likelihood

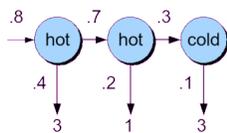
- For a Markov chain, we just follow the states 3 1 3 and multiply the probabilities
- But for an HMM, we don't know what the states are!
- So let's start with a simpler situation.
- Computing the observation likelihood for a **given** hidden state sequence
  - Suppose we knew the weather and wanted to predict how much ice cream Jason would eat.
  - I.e.  $P(3\ 1\ 3 | H\ H\ C)$

LSA 352 Summer 2007

## Computing likelihood for 1 given hidden state sequence

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

$$P(3\ 1\ 3|\text{hot hot cold}) = P(\text{hot}|\text{start}) \times P(\text{hot}|\text{hot}) \times P(\text{cold}|\text{hot}) \times P(3|\text{hot}) \times P(1|\text{hot}) \times P(3|\text{cold})$$



LSA 352 Summer 2007

## Computing total likelihood of 3 1 3

- We would need to sum over
  - Hot hot cold
  - Hot hot hot
  - Hot cold hot
  - ....
- How many possible hidden state sequences are there for this sequence?
- How about in general for an HMM with N hidden states and a sequence of T observations?
  - $N^T$
- So we can't just do separate computation for each hidden state sequence.

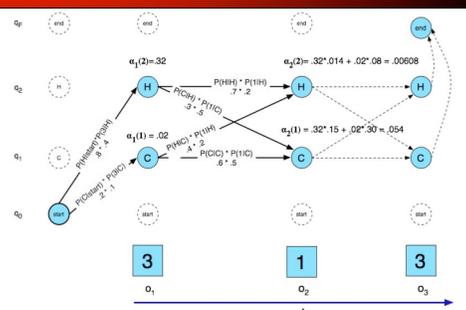
LSA 352 Summer 2007

## Instead: the Forward algorithm

- A kind of **dynamic programming** algorithm
  - Uses a table to store intermediate values
- Idea:
  - Compute the likelihood of the observation sequence
  - By summing over all possible hidden state sequences
  - But doing this efficiently
    - By folding all the sequences into a single **trellis**

LSA 352 Summer 2007

## The Forward Trellis



LSA 352 Summer 2007

## The forward algorithm

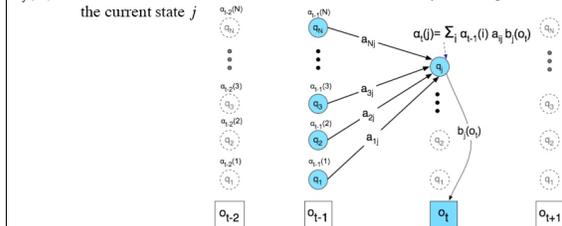
- Each cell of the forward algorithm trellis  $\alpha_t(j)$ 
  - Represents the probability of being in state  $j$
  - After seeing the first  $t$  observations
  - Given the automaton
- Each cell thus expresses the following probability

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

LSA 352 Summer 2007 37

## We update each cell

- $\alpha_{t-1}(i)$  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  $j$



LSA 352 Summer 2007 38

## The Forward Recursion

### 1. Initialization:

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

### 2. Recursion (since states 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:

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

LSA 352 Summer 2007 39

## The Forward Algorithm

function FORWARD(observations of len  $T$ , state-graph of len  $N$ ) returns forward-prob

```

create a probability matrix forward[N+2,T]
for each state  $s$  from 1 to  $N$  do           ;initialization step
    forward[s,1] ←  $a_{0,s} * b_s(o_1)$ 
for each time step  $t$  from 2 to  $T$  do     ;recursion step
    for each state  $s$  from 1 to  $N$  do
        forward[s,t] ←  $\sum_{i=1}^N$  forward[s',t-1] *  $a_{s',s} * b_s(o_t)$ 
forward[qF,T] ←  $\sum_{i=1}^N$  forward[s,T] *  $a_{s,q_F}$  ; termination step
return forward[qF,T]
    
```

LSA 352 Summer 2007 40

## Decoding

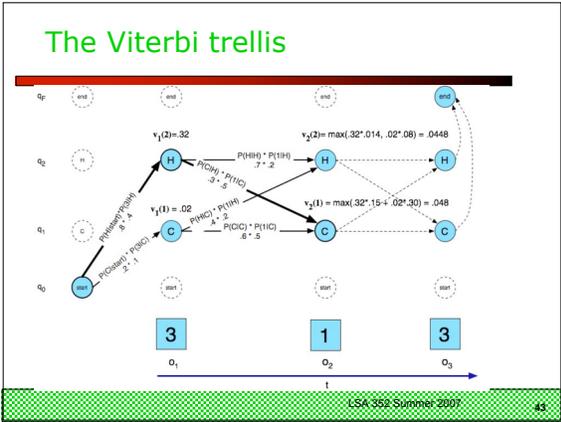
- Given an observation sequence
  - 3 1 3
- And an HMM
- The task of the **decoder**
  - To find the best **hidden** state sequence
- Given the observation sequence  $O=(o_1, o_2, \dots, o_T)$ , and an HMM model  $\Phi = (A, B)$ , how do we choose a corresponding state sequence  $Q=(q_1, q_2, \dots, q_T)$  that is optimal in some sense (i.e., best explains the observations)

LSA 352 Summer 2007 41

## Decoding

- One possibility:
  - For each hidden state sequence
    - HHH, HHC, HCH,
  - Run the forward algorithm to compute  $P(\Phi | O)$
- Why not?
  - $N^T$
- Instead:
  - The Viterbi algorithm
  - Is again a **dynamic programming** algorithm
  - Uses a similar trellis to the Forward algorithm

LSA 352 Summer 2007 42



### Viterbi intuition

- Process observation sequence left to right
- Filling out the trellis
- Each cell:

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

$$v_t(j) = \max_{i=1}^N v_{t-1}(i) a_{ij} b_j(o_t)$$

$v_{t-1}(i)$  the previous Viterbi 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  $j$

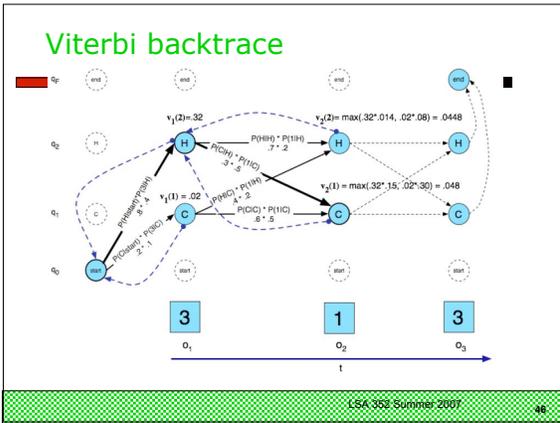
LSA 352 Summer 2007 44

### Viterbi Algorithm

```

function VITERBI(observations of len T, state-graph of len N) returns best-path
  create a path probability matrix viterbi[N+2,T]
  for each state s from 1 to N do           : initialization step
    viterbi[s,1] ← a0,s * bs(o1)
    backpointer[s,1] ← 0
  for each time step t from 2 to T do      : recursion step
    for each state s from 1 to N do
      viterbi[s,t] ← max_{s'=1}^N viterbi[s',t-1] * a_{s',s} * b_s(o_t)
      backpointer[s,t] ← argmax_{s'=1}^N viterbi[s',t-1] * a_{s',s}
  viterbi[qF,T] ← max_{s=1}^N viterbi[s,T] * a_{s,qF} ; termination step
  backpointer[qF,T] ← argmax_{s=1}^N viterbi[s,T] * a_{s,qF} ; termination step
  return the backtrace path by following backpointers to states back in time from backpointer[qF,T]
  
```

LSA 352 Summer 2007 45



### Viterbi Recursion

- Initialization:
 
$$v_1(j) = a_{0,j} b_j(o_1) \quad 1 \leq j \leq N$$

$$b_{t_1}(j) = 0$$
- Recursion (recall states 0 and q<sub>F</sub> are non-emitting):
 
$$v_t(j) = \max_{i=1}^N v_{t-1}(i) a_{ij} b_j(o_t); \quad 1 \leq j \leq N, 1 < t \leq T$$

$$b_{t_t}(j) = \max_{i=1}^N v_{t-1}(i) a_{ij} b_j(o_t); \quad 1 \leq j \leq N, 1 < t \leq T$$
- Termination:
 

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

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

LSA 352 Summer 2007 47

### Why "Dynamic Programming"

"I spent the Fall quarter (of 1950) at RAND. My first task was to find a name for multistage decision processes. An interesting question is, Where did the name, dynamic programming, come from? The 1950s were not good years for mathematical research. We had a very interesting gentleman in Washington named Wilson. He was Secretary of Defense, and he actually had a pathological fear and hatred of the word, research. I'm not using the term lightly; I'm using it precisely. His face would suffuse, he would turn red, and he would get violent if people used the term, research, in his presence. You can imagine how he felt, then, about the term, mathematical. The RAND Corporation was employed by the Air Force, and the Air Force had Wilson as its boss, essentially. Hence, I felt I had to do something to shield Wilson and the Air Force from the fact that I was really doing mathematics inside the RAND Corporation. What title, what name, could I choose? In the first place I was interested in planning, in decision making, in thinking. But planning, is not a good word for various reasons. I decided therefore to use the word, "programming" I wanted to get across the idea that this was dynamic, this was multistage, this was time-varying I thought, lets kill two birds with one stone. Lets take a word that has an absolutely precise meaning, namely dynamic, in the classical physical sense. It also has a very interesting property as an adjective, and that is its impossible to use the word, dynamic, in a pejorative sense. Try thinking of some combination that will possibly give it a pejorative meaning. Its impossible. Thus, I thought dynamic programming was a good name. It was something not even a Congressman could object to. So I used it as an umbrella for my activities." Richard Bellman, "Eye of the Hurricane: an autobiography" 1984.

Thanks to Chen, Pichery, Eric, Noel

LSA 352 Summer 2007 48



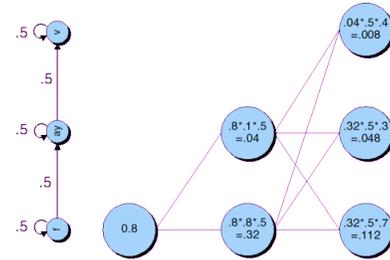
## The forward trellis for "five"

V	0	0	0.008	0.0093	0.0114	0.00703	0.00345	0.00306	0.00206	0.00117	
AY	0	0.04	0.054	0.0664	0.0355	0.016	0.00676	0.00208	0.000532	0.000109	
F	0.8	0.32	0.112	0.0224	0.00448	0.000896	0.000179	4.48e-05	1.12e-05	2.8e-06	
Time	1	2	3	4	5	6	7	8	9	10	
B	f 0.8 ay 0.1 v 0.6 p 0.4 iy 0.1	f 0.8 ay 0.1 v 0.6 p 0.4 iy 0.1	f 0.7 ay 0.3 v 0.4 p 0.2 iy 0.3	f 0.4 ay 0.8 v 0.3 p 0.1 iy 0.6	f 0.5 ay 0.6 v 0.6 p 0.1 iy 0.5	f 0.5 ay 0.5 v 0.8 p 0.3 iy 0.5	f 0.5 ay 0.4 v 0.9 p 0.3 iy 0.4	f 0.5 ay 0.4 v 0.9 p 0.3 iy 0.4			

LSA 352 Summer 2007

44

## Viterbi trellis for "five"



LSA 352 Summer 2007

45

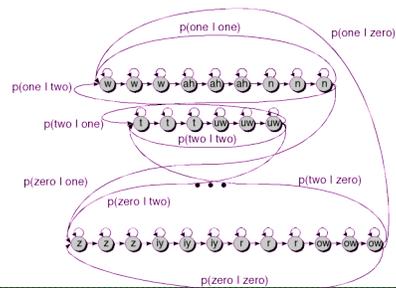
## Viterbi trellis for "five"

V	0	0	0.008	0.0072	0.00672	0.00403	0.00188	0.00161	0.000667	0.000493	
AY	0	0.04	0.048	0.0448	0.0269	0.0125	0.00538	0.00167	0.000428	8.78e-05	
F	0.8	0.32	0.112	0.0224	0.00448	0.000896	0.000179	4.48e-05	1.12e-05	2.8e-06	
Time	1	2	3	4	5	6	7	8	9	10	
B	f 0.8 ay 0.1 v 0.6 p 0.4 iy 0.1	f 0.8 ay 0.1 v 0.6 p 0.4 iy 0.1	f 0.7 ay 0.3 v 0.4 p 0.2 iy 0.3	f 0.4 ay 0.8 v 0.3 p 0.1 iy 0.6	f 0.5 ay 0.6 v 0.6 p 0.1 iy 0.5	f 0.5 ay 0.5 v 0.8 p 0.3 iy 0.5	f 0.5 ay 0.4 v 0.9 p 0.3 iy 0.4	f 0.5 ay 0.4 v 0.9 p 0.3 iy 0.4			

LSA 352 Summer 2007

47

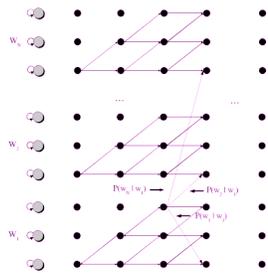
## Search space with bigrams



LSA 352 Summer 2007

48

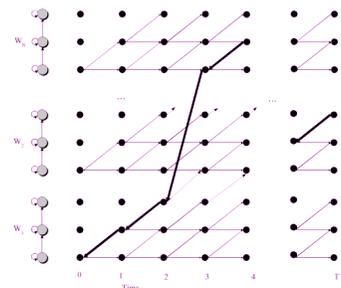
## Viterbi trellis with 2 words and uniform LM



LSA 352 Summer 2007

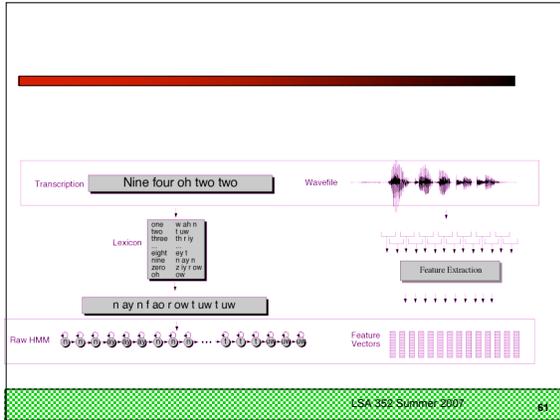
49

## Viterbi backtrace



LSA 352 Summer 2007

50



## Evaluation

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

LSA 352 Summer 2007 62

## Word Error Rate

- Word Error Rate =  $\frac{100 \times (\text{Insertions} + \text{Substitutions} + \text{Deletions})}{\text{Total Word in Correct Transcript}}$

Alignment example:  
 REF: portable \*\*\*\* PHONE UPSTAIRS last night so  
 HYP: portable FORM OF STORES last night so  
 Eval: I S S  
 WER =  $100 \times (1+2+0)/6 = 50\%$

LSA 352 Summer 2007 63

## NIST sctk-1.3 scoring software: Computing WER with sclite

- <http://www.nist.gov/speech/tools/>
- Sclite aligns a hypothesized text (HYP) (from the recognizer) with a correct or reference text (REF) (human transcribed)

```

id: (2347-a-013)
Scores: (#C #S #D #I) 9 3 1 2
REF: was an engineer SO I I was always with **** ** MEN UM and they
HYP: was an engineer ** AND I was always with THEM THEY ALL THAT and they
Eval: D S I I S S
  
```

LSA 352 Summer 2007 64

## Sclite output for error analysis

```

CONFUSSION PAIRS          Total          (972)
                          With >= 1 occurrences (972)

1: 6 -> ($hesitation) ==> on
2: 6 -> the ==> that
3: 5 -> but ==> that
4: 4 -> a ==> the
5: 4 -> four ==> for
6: 4 -> in ==> and
7: 4 -> there ==> that
8: 3 -> ($hesitation) ==> and
9: 3 -> ($hesitation) ==> the
10: 3 -> (a-) ==> i
11: 3 -> and ==> i
12: 3 -> and ==> is
13: 3 -> are ==> there
14: 3 -> as ==> is
15: 3 -> have ==> that
16: 3 -> is ==> this
  
```

LSA 352 Summer 2007 65

## Sclite output for error analysis

```

17: 3 -> it ==> that
18: 3 -> mouse ==> most
19: 3 -> was ==> is
20: 3 -> was ==> this
21: 3 -> you ==> we
22: 2 -> ($hesitation) ==> it
23: 2 -> ($hesitation) ==> that
24: 2 -> ($hesitation) ==> to
25: 2 -> ($hesitation) ==> yeah
26: 2 -> a ==> all
27: 2 -> a ==> know
28: 2 -> a ==> you
29: 2 -> along ==> well
30: 2 -> and ==> it
31: 2 -> and ==> we
32: 2 -> and ==> you
33: 2 -> are ==> i
34: 2 -> are ==> were
  
```

LSA 352 Summer 2007 66

## Better metrics than WER?

- WER has been useful
- But should we be more concerned with meaning ("semantic error rate")?
  - Good idea, but hard to agree on
  - Has been applied in dialogue systems, where desired semantic output is more clear

LSA 352 Summer 2007

67

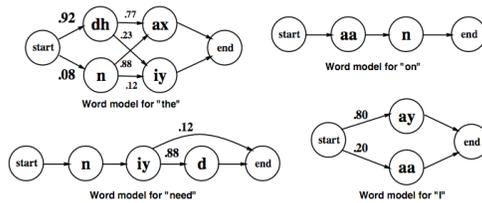
## Summary: ASR Architecture

- Five easy pieces: ASR Noisy Channel architecture
  - 1) Feature Extraction:
    - 39 "MFCC" features
  - 2) Acoustic Model:
    - Gaussians for computing  $p(o|q)$
  - 3) Lexicon/Pronunciation Model
    - HMM: what phones can follow each other
  - 4) Language Model
    - N-grams for computing  $p(w_i|w_{1:i})$
  - 5) Decoder
    - Viterbi algorithm: dynamic programming for combining all these to get word sequence from speech!

LSA 352 Summer 2007

68

## ASR Lexicon: Markov Models for pronunciation



LSA 352 Summer 2007

69

## Summary

- Speech Recognition Architectural Overview
- Hidden Markov Models in general
  - Forward
  - Viterbi Decoding
- Hidden Markov models for Speech
- Evaluation

LSA 352 Summer 2007

70