


Speech Recognition and Synthesis




Berkeley
N L P

Dan Klein
UC Berkeley

1

Language Models



2

Noisy Channel Model: ASR

- We want to predict a sentence given acoustics:

$$w^* = \arg \max_w P(w|a)$$

- The noisy-channel approach:


$$w^* = \arg \max_w P(w|a)$$

$$= \arg \max_w P(a|w)P(w)/P(a)$$

$$\propto \arg \max_w P(a|w)P(w)$$

Acoustic model: score fit
between sounds and words

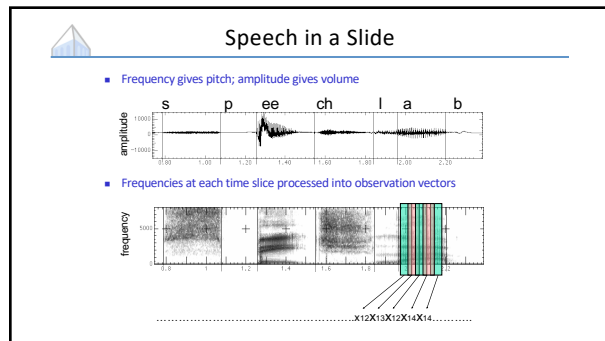
Language model: score
plausibility of word sequences



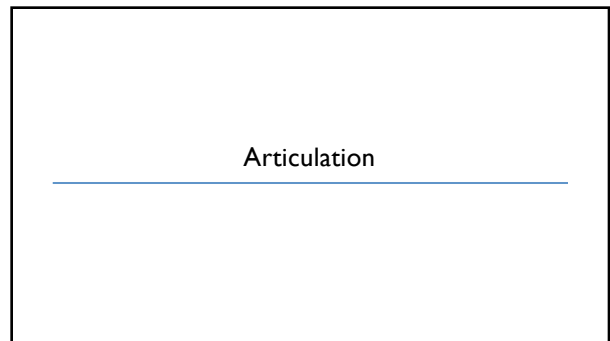
3

The Speech Signal

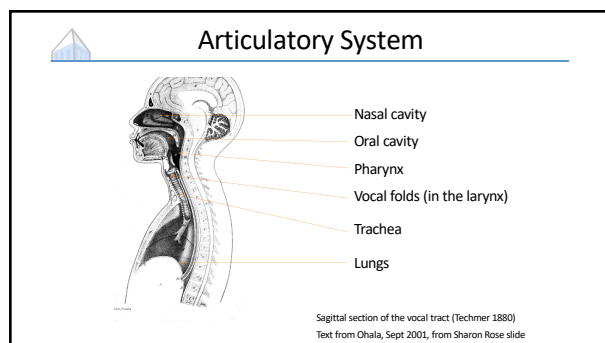
4



5



6



7

Space of Phonemes

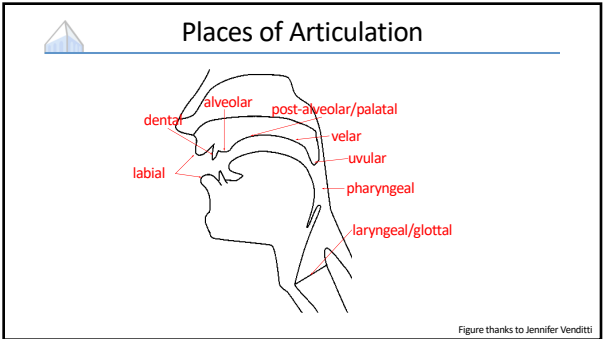
- Standard international phonetic alphabet (IPA) chart of consonants

	LABIAL		CORONAL				DORSAL			RAUCAL		LARYNGAL
	Bilabial	Labio-dental	Dental	Alveolar	Palato-alveolar	Retroflex	Palatal	Velar	Uvular	Pharyngeal	Epi-glottal	
Nasal	m	ɱ		n		ɳ	ɲ	ŋ	ɴ			
Plosive	p b	ɸ β	t d			ʈ ɖ	c ɟ	k ɡ	q ɢ		ʔ	ʔ
Fricative		f v	θ ð	s z	ʃ ʒ	ʂ ʐ	ç ʝ	x ɣ	χ ʁ	ħ	ħ	h ɦ
Approximant		ʋ		ɹ		ɻ	j	ɰ				ɹ̥ ɹ̥̄
Trill		ʙ		r					ʀ			ʀ̄
Tap, Flap		ɸ		ɾ		ɽ						
Lateral fricative				ɬ ɮ		ɮ̥ ɮ̄						
Lateral approximant				l		ɭ	ʎ	ʟ				
Lateral flap				ɭ		ɮ̣						

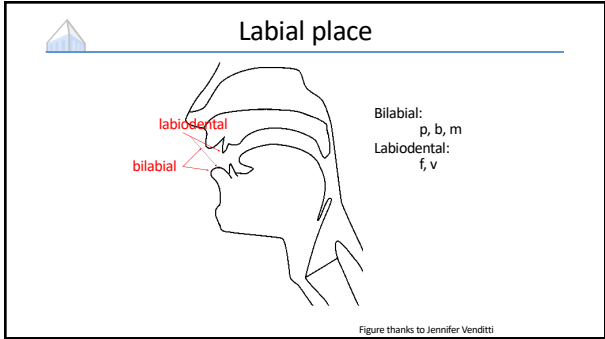
8

Articulation: Place

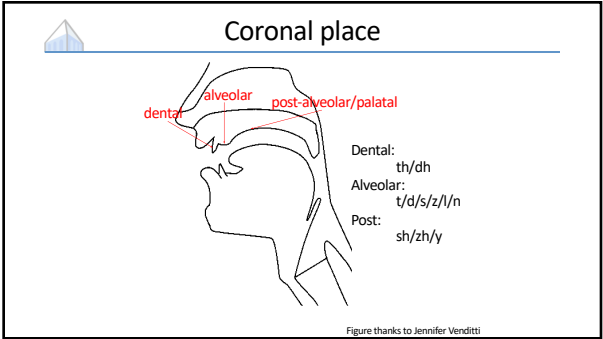
9



10



11



12

Dorsal Place

Velar:
k/g/ng

— velar
— uvular
— pharyngeal

Figure thanks to Jennifer Venditti

13

Space of Phonemes

- Standard international phonetic alphabet (IPA) chart of consonants

	LABIAL		CORONAL				DORSAL			RACIAL		GLOTTAL
	Bilabial	Labio-dental	Dental	Alveolar	Palato-alveolar	Retroflex	Palatal	Velar	Uvular	Pharyngeal	Epi-glottal	Glottal
Nasal	m	ɱ		n	ɲ	ɳ	ɲ	ŋ	ɴ			
Plosive	p b	ɸ β		t d	ʈ ɖ	ʈ ɖ	c ɟ	k ɡ	q ɢ		ʔ	ʔ
Fricative		f v	θ ð	s z	ʃ ʒ		ç ʝ	x ɣ	ħ	ħ	ʕ	h ɦ
Approximant				ɹ		ɻ	j	ɰ				
Trill	ʙ			ʀ					ʀ			ʀ
Tap, Flap		ⱱ		ɾ		ɽ						
Lateral fricative				ɬ ɮ	ɮ	ɬ						
Lateral approximant				l		ɭ	ʎ	ʟ				
Lateral flap				ɭ		ɻ						

14

Articulation: Manner

15

Manner of Articulation

- In addition to varying by place, sounds vary by manner
- Stop: complete closure of articulators, no air escapes via mouth
 - Oral stop: palate is raised (p, t, k, b, d, g)
 - Nasal stop: oral closure, but palate is lowered (m, n, ŋ)
- Fricatives: substantial closure, turbulent: (f, v, s, z)
- Approximants: slight closure, sonorant: (l, r, w)
- Vowels: no closure, sonorant: (i, e, a)

16

Space of Phonemes

- Standard international phonetic alphabet (IPA) chart of consonants

	LABIAL		CORONAL				DORSAL			RHYNOAL		LARYNGEAL
	Bilabial	Labio-dental	Dental	Alveolar	Palato-alveolar	Retroflex	Palatal	Velar	Uvular	Pharyngeal	Epi-glottal	Glottal
Nasal	m	ɱ		n		ɳ	ɲ	ŋ	ɴ			
Plosive	p b	ɸ β		t d		ʈ ɖ	c ɟ	k ɡ	q ɢ		ʔ	ʔ̰
Fricative	ɸ β	f v	θ ð	s z	ʃ ʒ	ʂ ʐ	ç ʝ	x ɣ	χ ʁ	ħ	ʕ	h ɦ
Approximant		ʋ		ɹ		ɻ	j	ɰ				
Trill	ʙ			ʀ					ʀ			ʀ̠
Tap, Flap		ⱱ		ɾ		ɽ						
Lateral fricative				ɬ ɮ		ɮ̰ ɮ̱		ɮ̥ ɮ̦				
Lateral approximant				l		ɭ	ʎ	ʟ				
Lateral flap				ɭ		ɮ̰						

17

Articulation: Vowels

18

Vowel Space

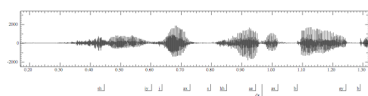
Vowels at right & left of bullets are rounded & unrounded.

19

Acoustics

20

“She just had a baby”

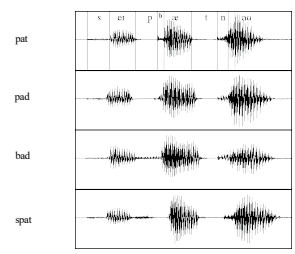


What can we learn from a wavefile?

- No gaps between words (l)
- Vowels are voiced, long, loud
- Length in time = length in space in waveform picture
- Voicing: regular peaks in amplitude
- When stops closed: no peaks, silence
- Peaks = voicing: .46 to .58 (vowel [iy], from second .65 to .74 (vowel [ax]) and so on
- Silence of stop closure (1.06 to 1.08 for first [b], or 1.26 to 1.28 for second [b])
- Fricatives like [sh]: intense irregular pattern; see .33 to .46

21

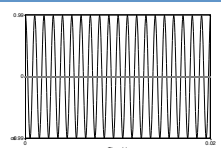
Time-Domain Information



Example from Ladefoged

22

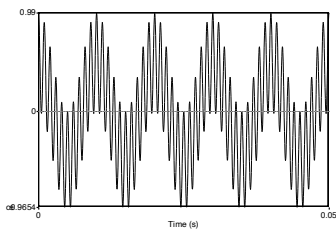
Simple Periodic Waves of Sound



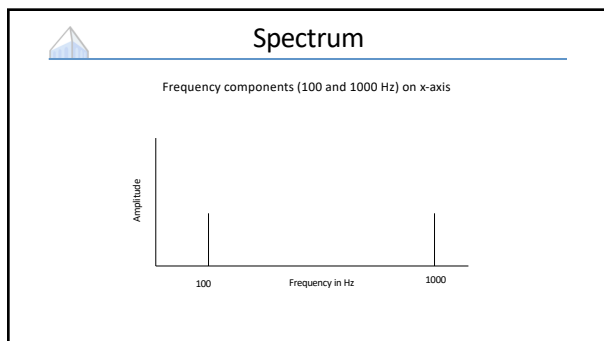
- Y axis: Amplitude = amount of air pressure at that point in time
 - Zero is normal air pressure, negative is rarefaction
- X axis: Time
- Frequency = number of cycles per second
- 20 cycles in .02 seconds = 1000 cycles/second = 1000 Hz

23

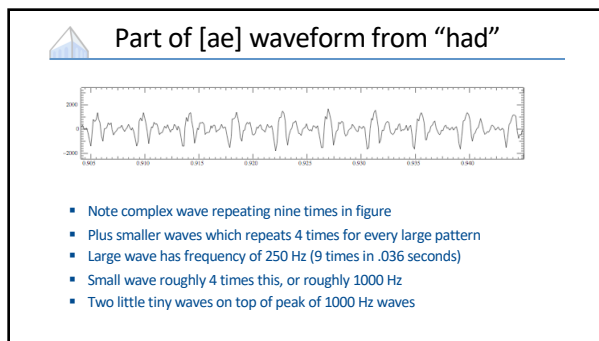
Complex Waves: 100Hz+1000Hz



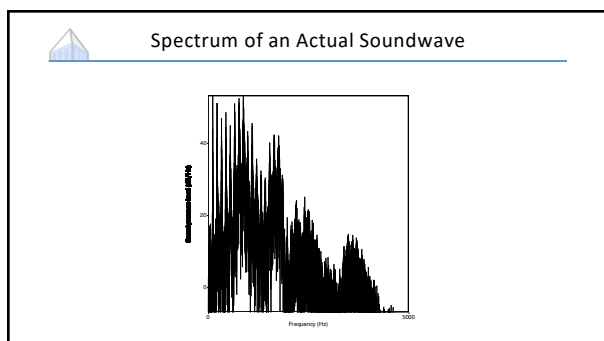
24



25



26



27

Source / Channel

28

Why these Peaks?

- Articulation process:**
 - The vocal cord vibrations create harmonics
 - The mouth is an amplifier
 - Depending on shape of mouth, some harmonics are amplified more than others

29

Vowel [i] at increasing pitches

Figures from Ratree Wayland

30

Resonances of the Vocal Tract

- The human vocal tract as an open tube:
 - Closed end
 - Open end
 - Length 17.5 cm.
- Air in a tube of a given length will tend to vibrate at resonance frequency of tube.
- Constraint: Pressure differential should be maximal at (closed) glottal end and minimal at (open) lip end.

Figure from W. Barry

31

From Sundberg

32

Computing the 3 Formants of Schwa

- Let the length of the tube be L
- $F_1 = c/\lambda_1 = c/(4L) = 35,000/4 \cdot 17.5 = 500\text{Hz}$
- $F_2 = c/\lambda_2 = c/(4/3L) = 3c/4L = 3 \cdot 35,000/4 \cdot 17.5 = 1500\text{Hz}$
- $F_3 = c/\lambda_3 = c/(4/5L) = 5c/4L = 5 \cdot 35,000/4 \cdot 17.5 = 2500\text{Hz}$

- So we expect a neutral vowel to have 3 resonances at 500, 1500, and 2500 Hz
- These vowel resonances are called **formants**

33

The diagram shows three rows of vocal tract models. Each row includes a cross-section of the vocal tract, a simplified acoustic tube model with labeled sections (Throat, Mouth, Lips), and an acoustic spectrum plot showing pressure amplitude versus frequency (0 to 4000 Hz). The spectra show three distinct peaks corresponding to the first three formants.

From Mark Liberman

34

Seeing Formants: the Spectrogram

The spectrograms are arranged in two rows. The top row shows [i], [ɪ], [e], and [æ]. The bottom row shows [ɑ], [ɔ], [o], and [u]. Each spectrogram plots frequency (0 to 4000 Hz) against time (0 to 400 ms). The dark horizontal bands represent the formants, with their positions varying according to the vowel's articulation.

35

Vowel Space

The vowel chart on the left plots tongue positions on a grid with axes: Front-Back and Close-Open. Vowels are marked with letters and symbols. A note states: "Vowels at right & left of bullets are rounded & unrounded." The plot on the right shows F2 Frequency (Hz) on the y-axis (500 to 3000) versus F1 Frequency (Hz) on the x-axis (200 to 1000). Vowels are plotted as points, showing their relative positions in the frequency space.

36

Spectrograms

37

How to Read Spectrograms

- [bab]: closure of lips lowers all formants: so rapid increase in all formants at beginning of "bab"
- [dad]: first formant increases, but F2 and F3 slight fall
- [gag]: F2 and F3 come together: this is a characteristic of velars. Formant transitions take longer in velars than in alveolars or labials

From Ladefoged "A Course in Phonetics"

38

"She came back and started again"

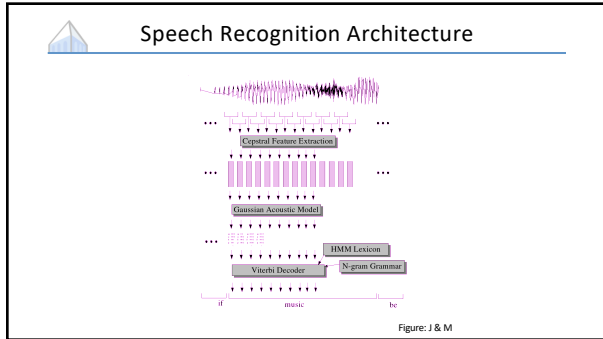
1. lots of high-freq energy
- 3.
- 4.
- 5.
- 6.
- 7.
- 8.
- 9.

From Ladefoged "A Course in Phonetics"

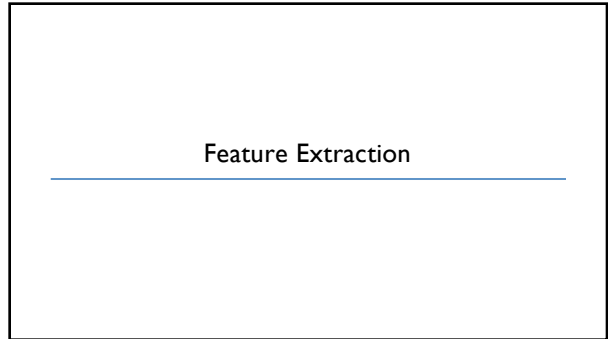
39

Speech Recognition

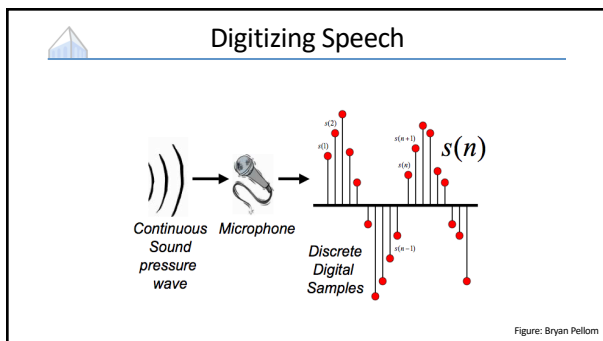
40



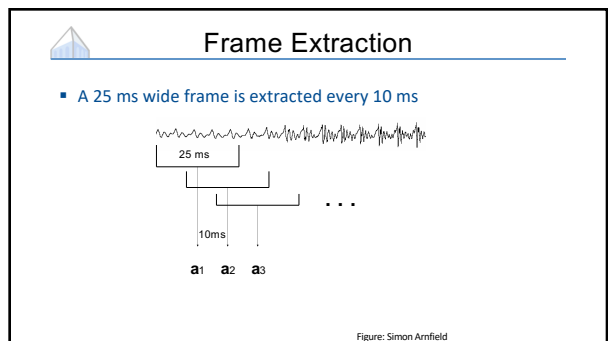
41



42



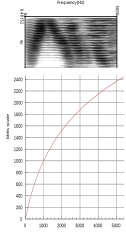
43



44

Mel Freq. Cepstral Coefficients

- Do FFT to get spectral information
 - Like the spectrogram we saw earlier
- Apply Mel scaling
 - Models human ear; more sensitivity in lower freqs
 - Approx linear below 1kHz, log above, equal samples above and below 1kHz
- Plus discrete cosine transform



[Graph: Wikipedia]

45

Final Feature Vector

- 39 (real) features per 10 ms frame:
 - 12 MFCC features
 - 12 delta MFCC features
 - 12 delta-delta MFCC features
 - 1 (log) frame energy
 - 1 delta (log) frame energy
 - 1 delta-delta (log frame energy)
- So each frame is represented by a 39D vector

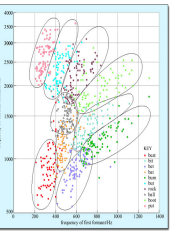
46

Emission Model

47

HMMs for Continuous Observations

- Solution 1: discretization
- Solution 2: continuous emission models
 - Gaussians
 - Multivariate Gaussians
 - Mixtures of multivariate Gaussians
- Solution 3: neural classifiers
- A state is progressively
 - Context independent subphone (~3 per phone)
 - Context dependent phone (triphones)
 - State tying of CD phone



48

Vector Quantization

- Idea: discretization
 - Map MFCC vectors onto discrete symbols
 - Compute probabilities just by counting
- This is called vector quantization or VQ
- Not used for ASR any more
- But: useful to consider as a starting point, and for understanding neural methods

The diagram shows an 'Input Feature Vector' being compared to a 'Codebook of 256' to produce an 'Output index of best vector' (144). Below this is a scatter plot of MFCC vectors with several Gaussian-like regions overlaid, representing the discretization process.

49

Gaussian Emissions

- VQ is insufficient for top-quality ASR
 - Hard to cover high-dimensional space with codebook
 - Moves ambiguity from the model to the preprocessing
- Instead: assume the possible values of the observation vectors are normally distributed.
 - Represent the observation likelihood function as a Gaussian?

The scatter plot shows MFCC vectors with several overlapping Gaussian regions, illustrating how Gaussian emissions can better model the continuous space of observation vectors compared to VQ.

From bartus.org/akustyk

50

But we're not there yet

- Single Gaussians may do a bad job of modeling a complex distribution in any dimension
- Even worse for diagonal covariances
- Classic solution: mixtures of Gaussians
- Modern solution: NN-based acoustic models map feature vectors to (sub)states

The scatter plot shows a complex distribution of MFCC vectors with many overlapping Gaussian regions, demonstrating the difficulty of modeling such a distribution with single Gaussians.

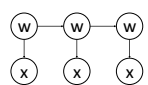
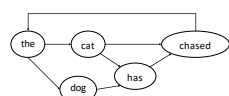
From opensteam.open.ac.uk

51

HMM / State Model

52

State Transition Diagrams

- Bayes Net: HMM as a Graphical Model
 
- State Transition Diagram: Markov Model as a Weighted FSA
 

53

ASR Lexicon

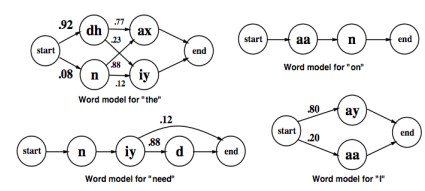


Figure: J & M

54

Lexical State Structure

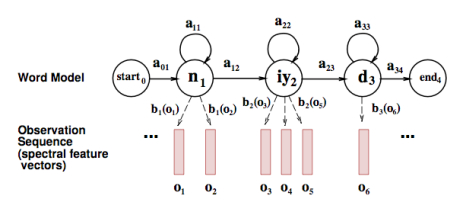


Figure: J & M

55

Adding an LM

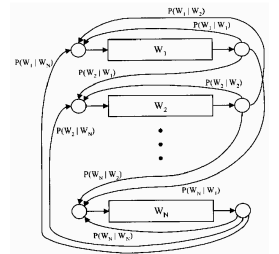


Figure from Huang et al page 618

56

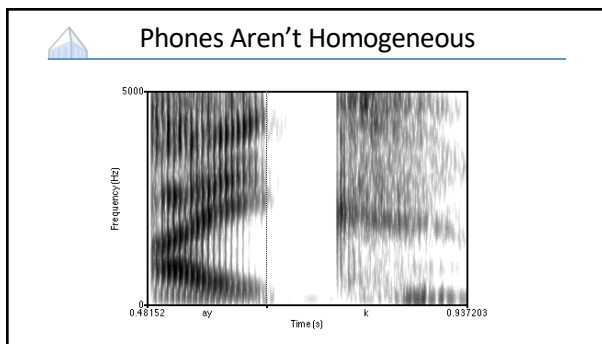
State Space

- State space must include
 - Current word ($|V|$ on order of 50K+)
 - Index within current word ($|L|$ on order of 5)
 - E.g. (lec[t]ure) (though not in orthography!)
- Acoustic probabilities only depend on (contextual) phone type
 - E.g. $P(x|lec[t]ure) = P(x|t)$
- From a state sequence, can read a word sequence

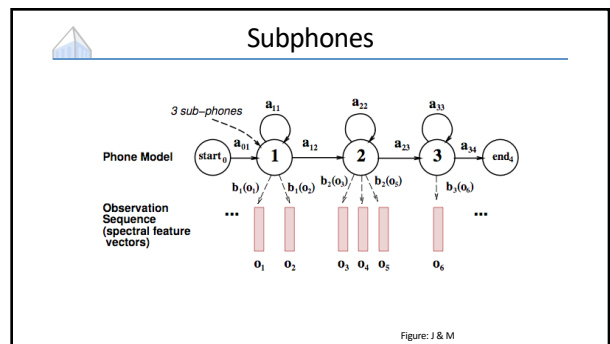
57

State Refinement

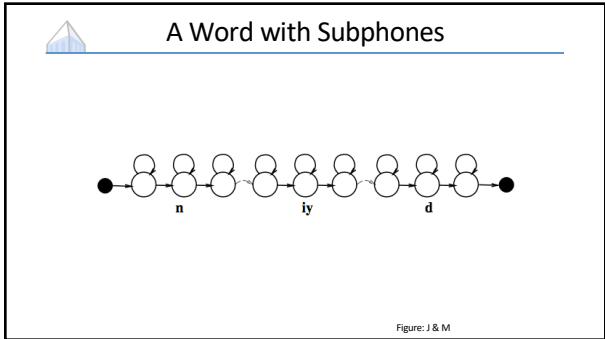
58



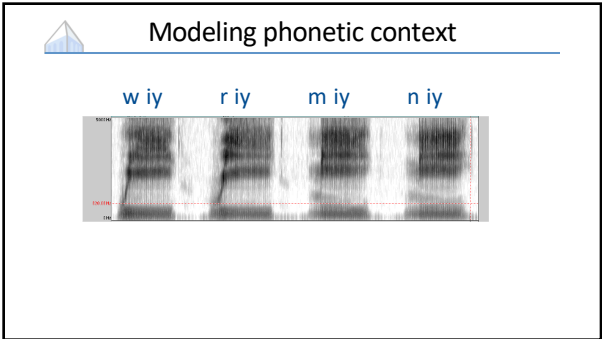
59



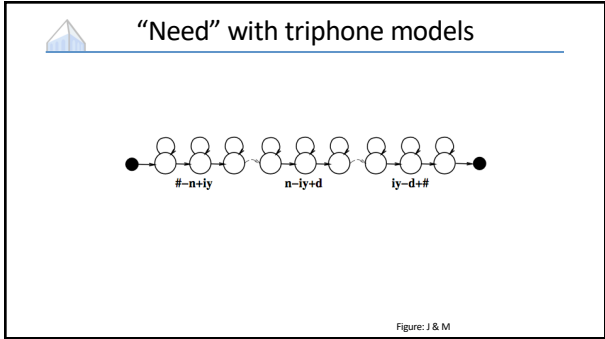
60



61



62



63

- ### Lots of Triphones
- Possible triphones: $50 \times 50 \times 50 = 125,000$
 - How many triphone types actually occur?
 - 20K word WSJ Task (from Bryan Pellom)
 - Word internal models: need 14,300 triphones
 - Cross word models: need 54,400 triphones
 - Need to generalize models, tie triphones

64

State Tying / Clustering

- [Young, Odell, Woodland 1994]
- How do we decide which triphones to cluster together?
- Use **phonetic features** (or 'broad phonetic classes')
- Stop
- Nasal
- Fricative
- Sibilant
- Vowel
- lateral

Initial set of untied states

Figure: J & M

65

State Space

- Full state space
(LM context, lexicon index, subphone)
- Details:
 - LM context is the past n-1 words
 - Lexicon index is a phone position within a word (or a trie of the lexicon)
 - Subphone is begin, middle, or end
 - E.g. (after the, lec[t-mid]ure)
- Acoustic model depends on clustered phone context
 - But this doesn't grow the state space

66

Learning Acoustic Models

67

What Needs to be Learned?

- Emissions: $P(x \mid \text{phone class})$
 - x is MFCC-valued
 - In neural methods, actually have $P(\text{phone} \mid \text{window around } x)$ and then coerce those scores into $P(x \mid \text{phone})$
- Transitions: $P(\text{state} \mid \text{prev state})$
 - If between words, this is $P(\text{word} \mid \text{history})$
 - If inside words, this is $P(\text{advance} \mid \text{phone class})$
 - (Really a hierarchical model)

68

Estimation from Aligned Data

- What if each time step were labeled with its (context-dependent sub) phone?

- Can estimate $P(x|/ae/)$ as empirical mean and (co-)variance of x 's with label $/ae/$, or mixture, etc/
- Problem: Don't know alignment at the frame and phone level

69

Forced Alignment

- What if the acoustic model $P(x|phone)$ were known (or approximately known)?
 - ...and also the correct sequences of words / phones
- Can predict the best alignment of frames to phones

"speech lab"

sssssspppppeeeeeetshshshllllaeaeabbbb

- Called "forced alignment"

70

Forced Alignment

- Create a new state space that forces the hidden variables to transition through phones in the (known) order

- Still have uncertainty about durations: this key uncertainty persists in neural models (and in some ways is worse now)
- In this HMM, all the parameters are known
 - Transitions determined by known utterance
 - Emissions assumed to be known
 - Minor detail: self-loop probabilities
- Just run Viterbi (or approximations) to get the best alignment

71

EM for Alignment

- Input: acoustic sequences with word-level transcriptions
- We don't know either the emission model or the frame alignments
- Expectation Maximization
 - Alternating optimization
 - Impute completions for unlabeled variables (here, the states at each time step)
 - Re-estimate model parameters (here, Gaussian means, variances, mixture ids)
 - Repeat
- One of the earliest uses of EM for structured problems

72

Staged Training and State Tying

- Creating CD phones:**
 - Start with monophone, do EM training
 - Clone Gaussians into triphones
 - Build decision tree and cluster Gaussians
 - Clone and train mixtures (GMMs)
- General idea:**
 - Introduce complexity gradually
 - Interleave constraint with flexibility

73

Neural Acoustic Models

- Given an input x , map to s ; this score coerced into generative $P(x|s)$ via Bayes rule (liberally ignoring terms)**
 - One major advantage of the neural net**
 - is that you can look at many x 's at once to capture dynamics (important!)

[Diagram from Hung-yi Li]

74

Decoding

75

State Trellis

$$\phi_t(s_{t-1}, s_t) = P(x_t|s_t)P(s_t|s_{t-1})$$

$$P(x, s) = \prod_t P(x_t|s_t)P(s_t|s_{t-1})$$

$$= \prod_t \phi_t(s_{t-1}, s_t)$$

Figure: Enrique Benimeli

76

Beam Search

- Lattice is not regular in structure! Dynamic vs static decoding
- At each time step
 - Start: Beam (collection) v_t of hypotheses s at time t
 - For each s in v_t
 - Compute all extensions s' at time $t+1$
 - Score s' from s
 - Put s' in v_{t+1} replacing existing s' if better
 - Advance to $t+1$
- Beams are priority queues of fixed size* k (e.g. 30) and retain only the top k hypotheses

77

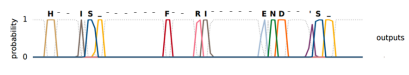
Dynamic vs Static Decoding

- Dynamic decoding
 - Build transitions on the fly based on model / grammar / etc
 - Very flexible, allows heterogeneous contexts easily (eg complex LMs)
- Static decoding
 - Compile entire subphone/vocabulary/LM into a huge weighted FST and use FST optimization methods (eg pushing, merging)
 - Much more common at scale, better eng and speed properties

78

Direct Neural Decoders

- Lots of work in decoders that skip explicit / discrete alignment
 - Decode to phone, or character, or word
 - Handle alignments softly (eg attention) or discretely (eg CTC)



- Catching up but not yet as good as structured systems

[Diagram from Graves 2014]

79

Speech Synthesis

[Many slides from Dan Jurafsky]

80

Early TTS

- Von Kempelen, 1791

81

The Voder

Developed by Homer Dudley at Bell Telephone Laboratories, 1939

82

Voder Architecture

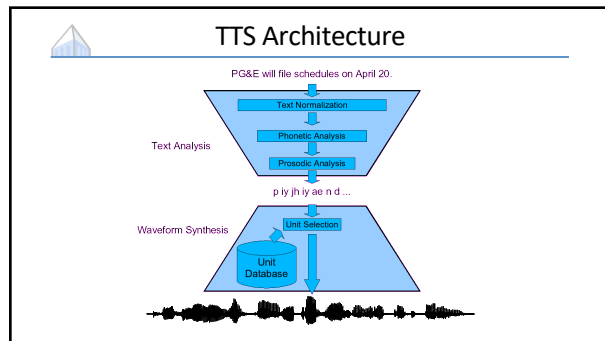
- An early hardware solution that already captured the flow of parametric synthesizers

83

Modern TTS

- 1960's first full TTS: Umeda et al (1968)
- 1970's
 - Joe Olive 1977 concatenation of linear-prediction diphones
 - Speak and Spell
- 1980's
 - 1979 MIT MITalk (Allen, Hunnicut, Klatt)
- 1990's – 2000's
 - Diphone synthesis
 - Unit selection synthesis
- Recent
 - Parametric synthesis returns!

84

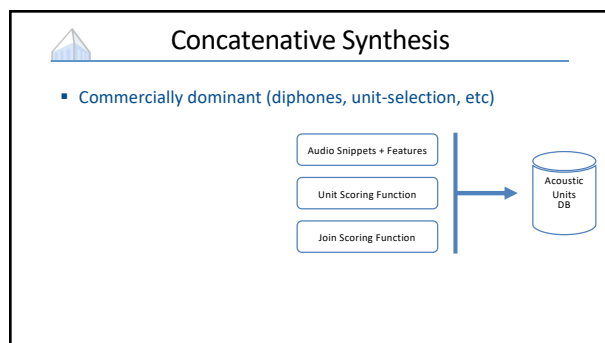


85

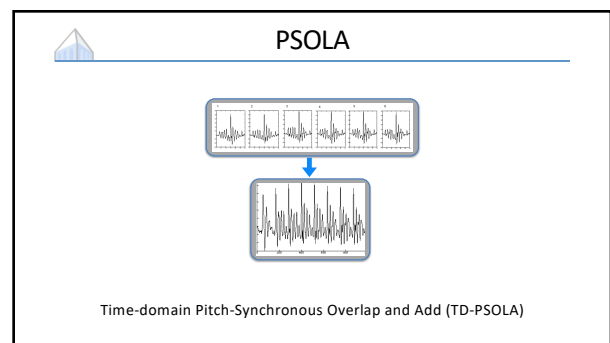
Typical Data for TTS

- Professional voice actor
- Carefully selected material
- High-quality recordings
 - 10-100 hours @ 44kHz
 - High signal-to-noise ratio
 - Consistent audio levels
 - No vocal issues (creaky voice)
 - Anechoic-like environment
- Usually lots of post-processing (alignments, pronunciations, ...)

86



87





88

Trade-offs

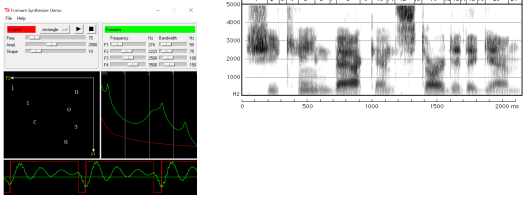
Concatenative: 1 Parametric: 2

"Beauty and the Beast is playing at AMC Burlington at 6:30pm, 9:30pm, and 10pm tomorrow night"

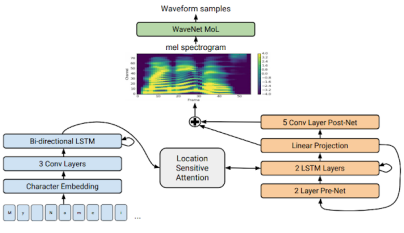
89

Formant Synthesis



90

Direct-to-Wave Synthesis



<https://ai.googleblog.com/2017/12/tracron-2-generating-human-like-speech.html>

91