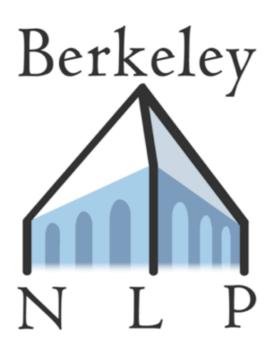
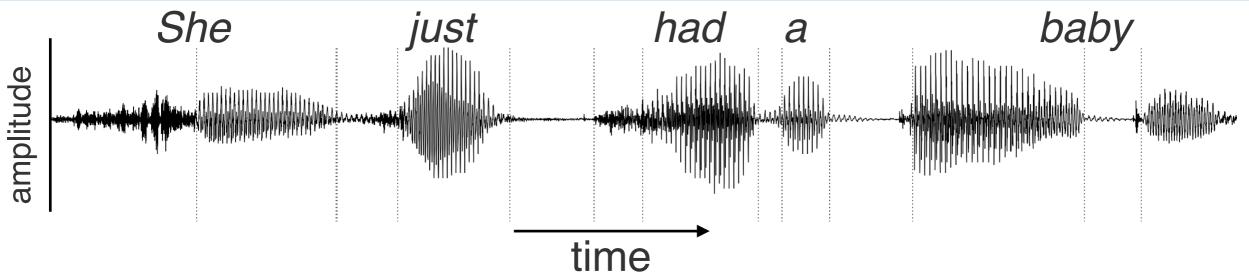
Speech Preliminaries + Noisy Channel Model



EECS 183/283a: Natural Language Processing

Recap:Speech Representations

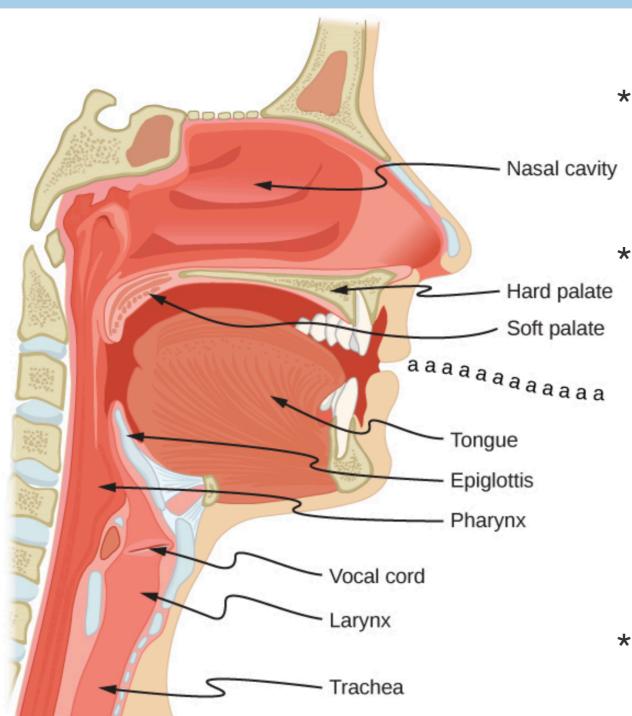




- * The result of the vocal articulation is an acoustic pressure wave
- * Speech can thus be represented as an acoustic waveform
- * Waveforms are continuous time series cannot be easily analyzed or interpreted, or computed with
- * Signal processing can give more interpretable information

Speech Production

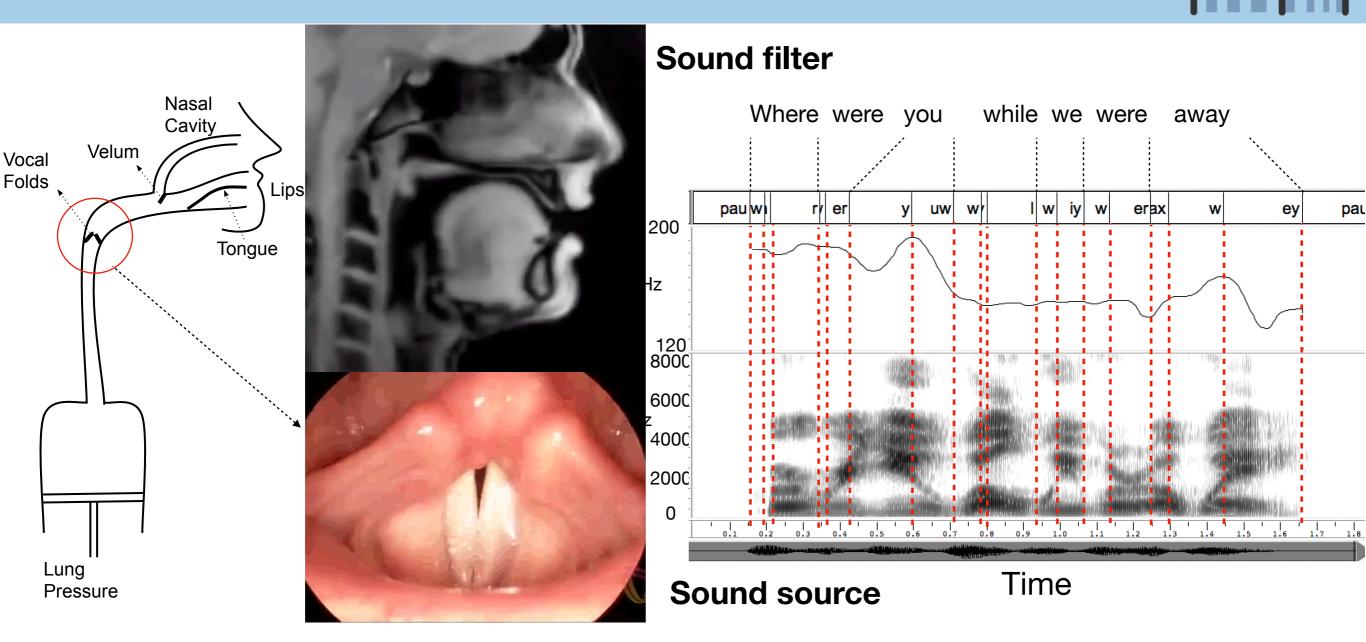




Vocal articulators that produce speech

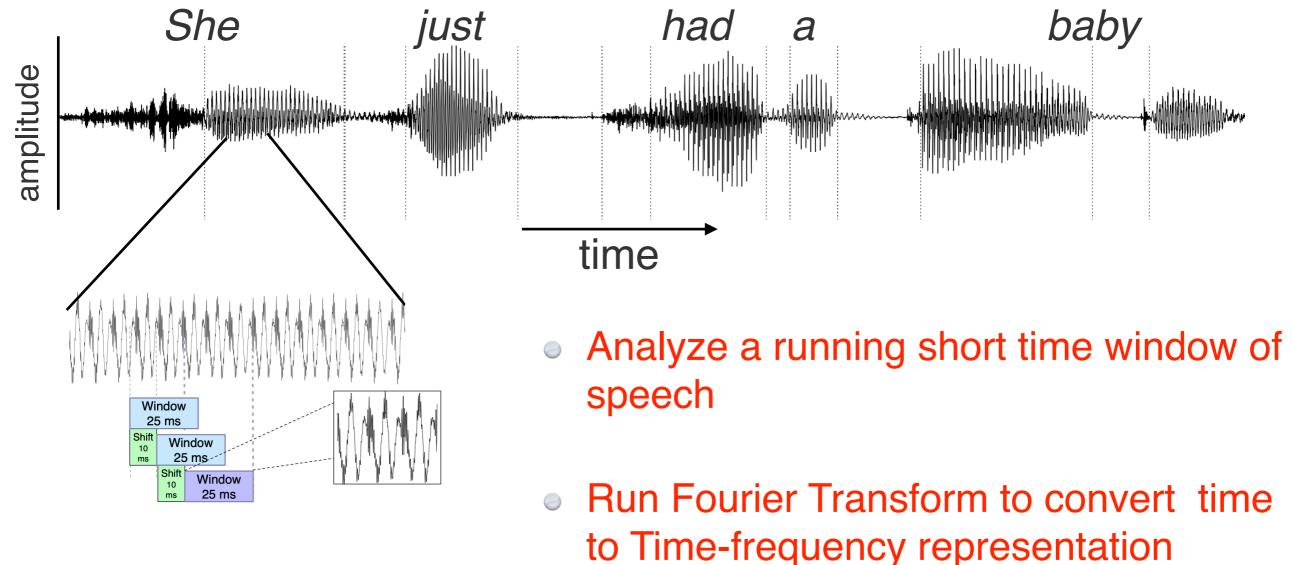
- Air passing through vocal articulators produces speech
- Vocal folds, tongue, jaw, lips, velum are both independently and jointly controlled to produce different sounds
- * eg. Vocal fold vibration causes voicing.
- The output of vocal articulation is an acoustic pressure wave

The most complex action we do is speaking



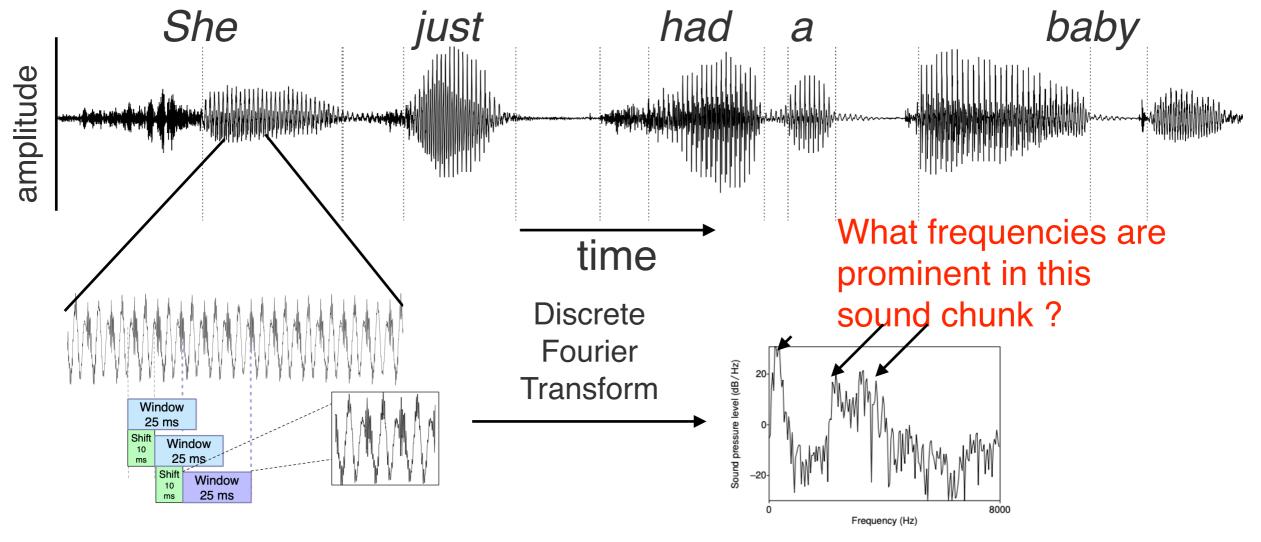
Speech Waveform





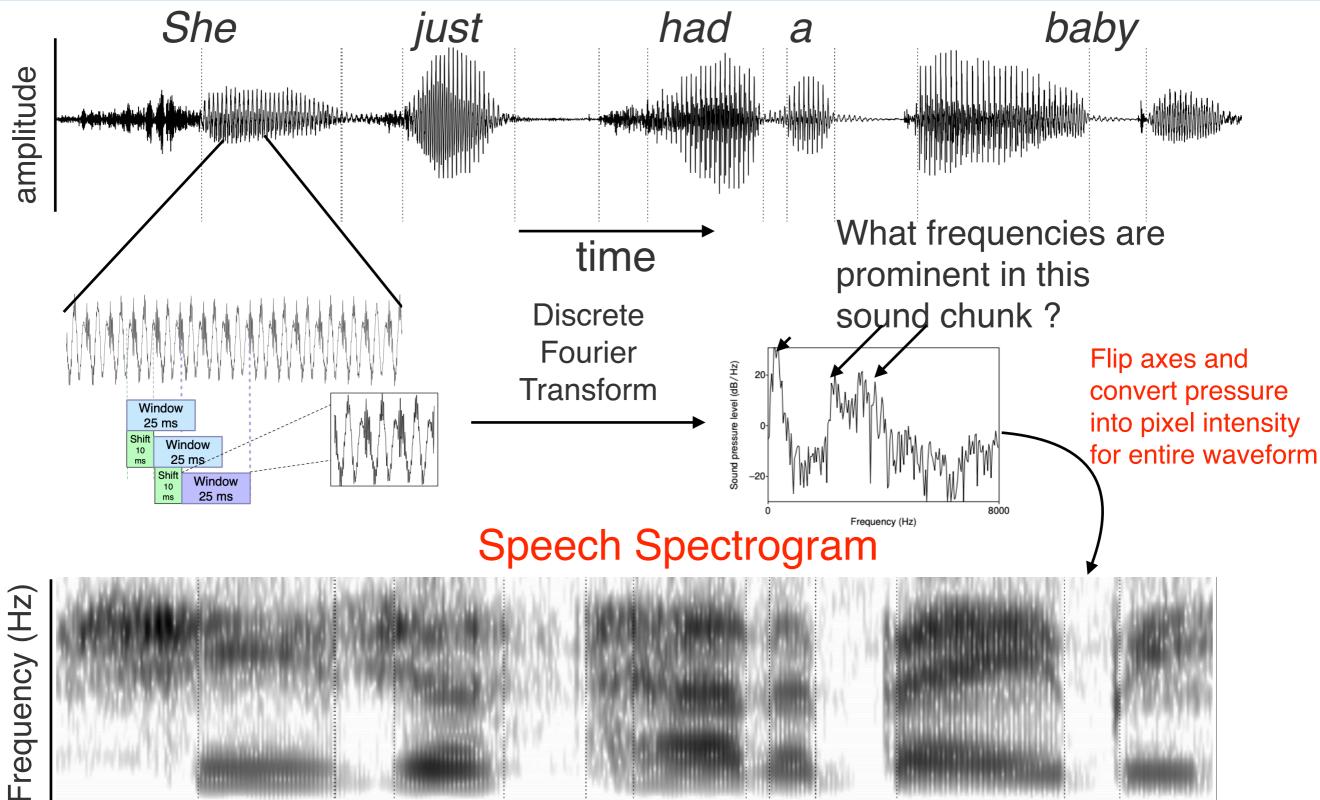
Speech Waveform





Speech Spectrogram

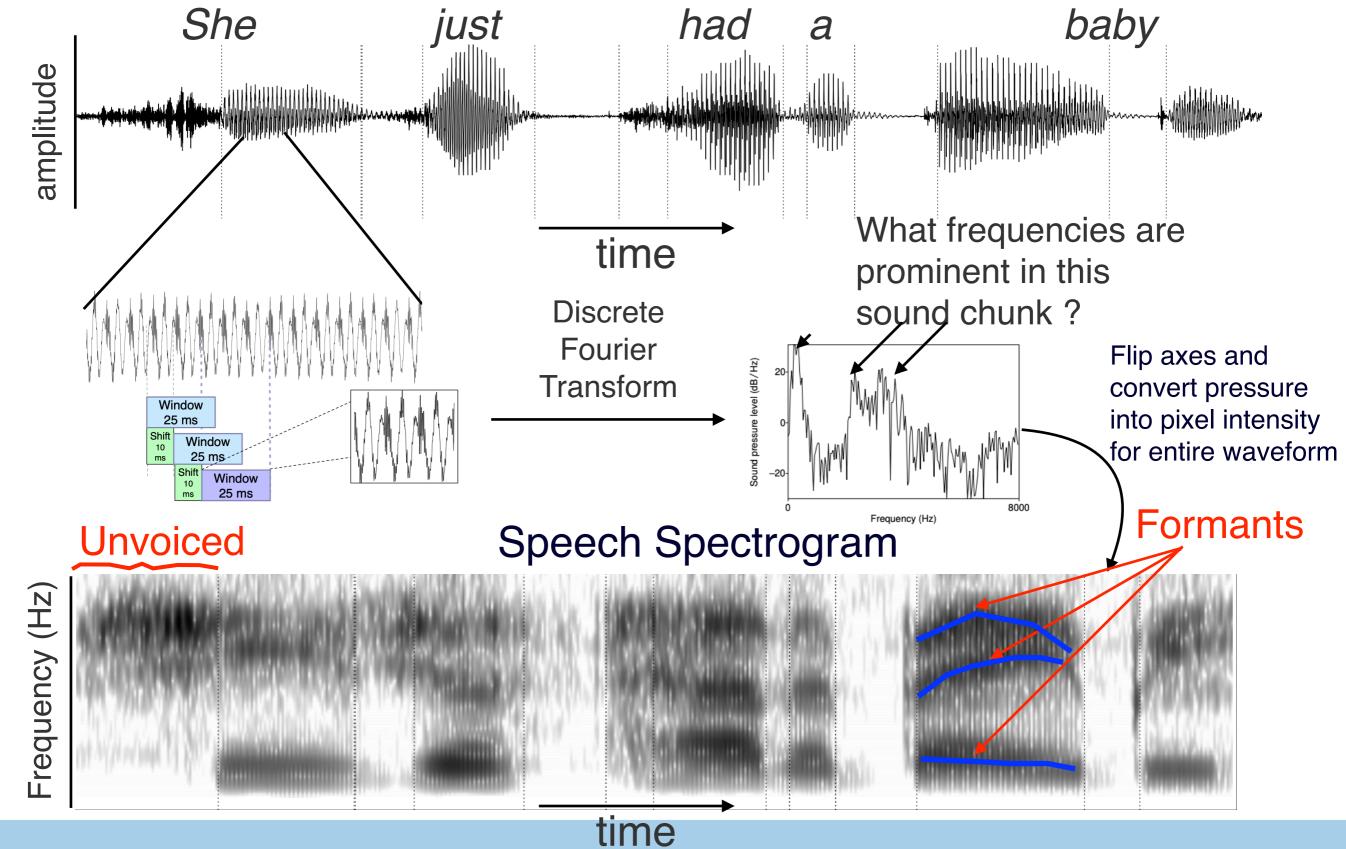




time

Speech Spectrogram

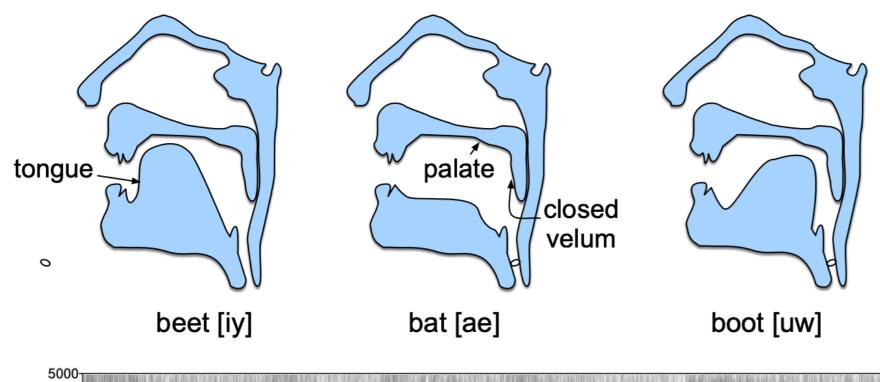




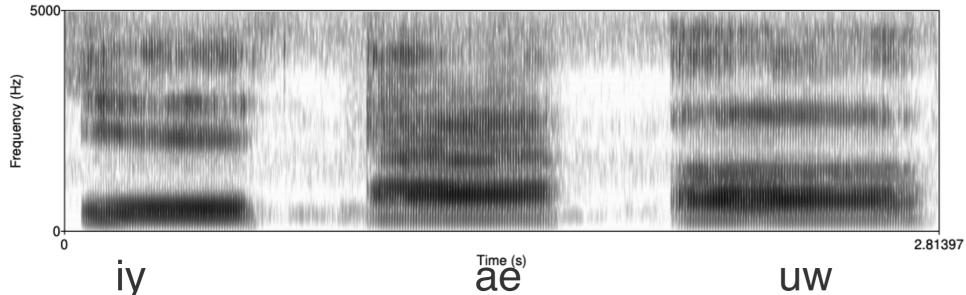
Phonetics



Study of speech sounds — their physical production, spectral and perceptual properties



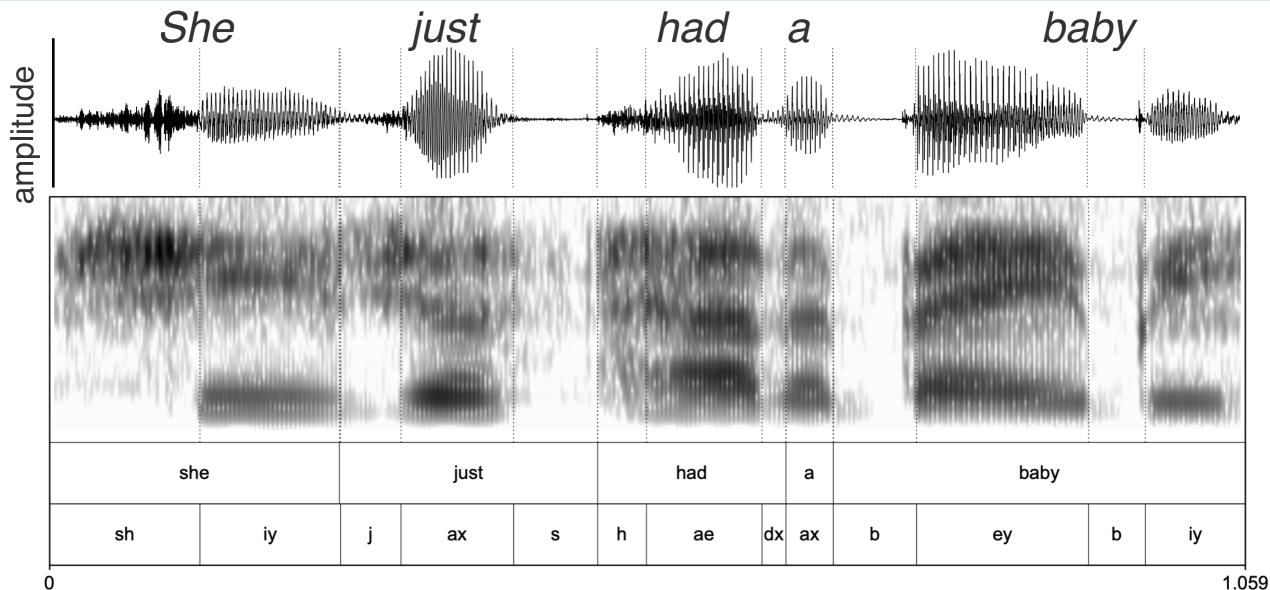
Articulatory Phonetics



Acoustic Phonetics

Acoustic Phonetics

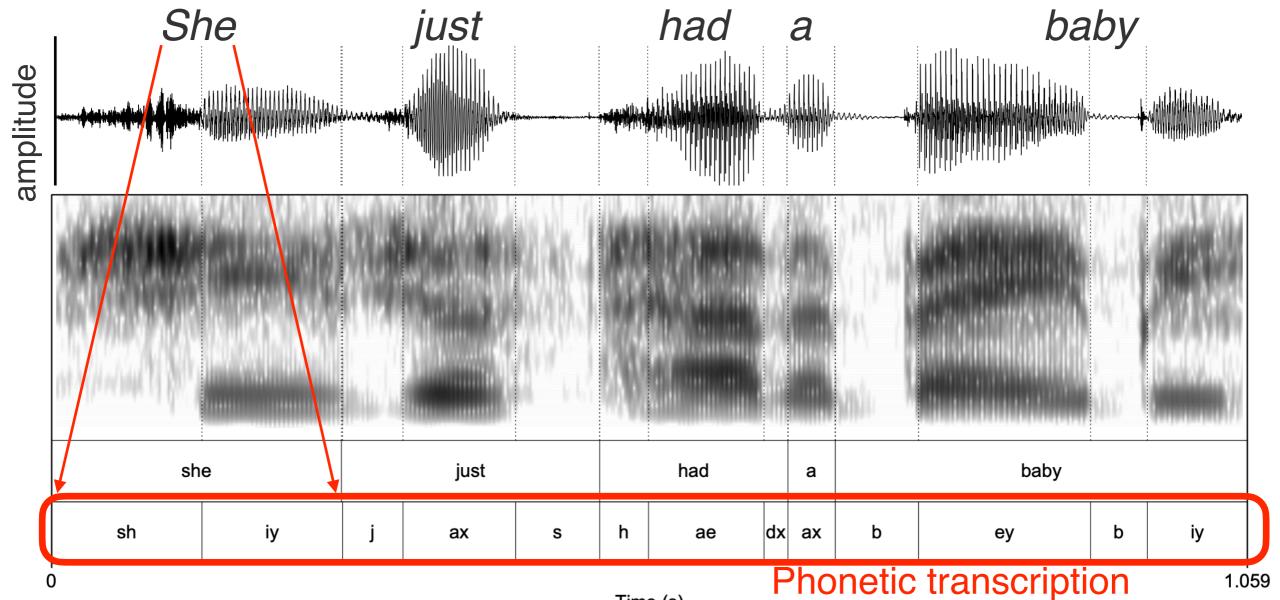




- Spectrogram reveals some segmental structure with distinct properties
- These are Phonemes perceptually distinct speech sounds

Acoustic Phonetics



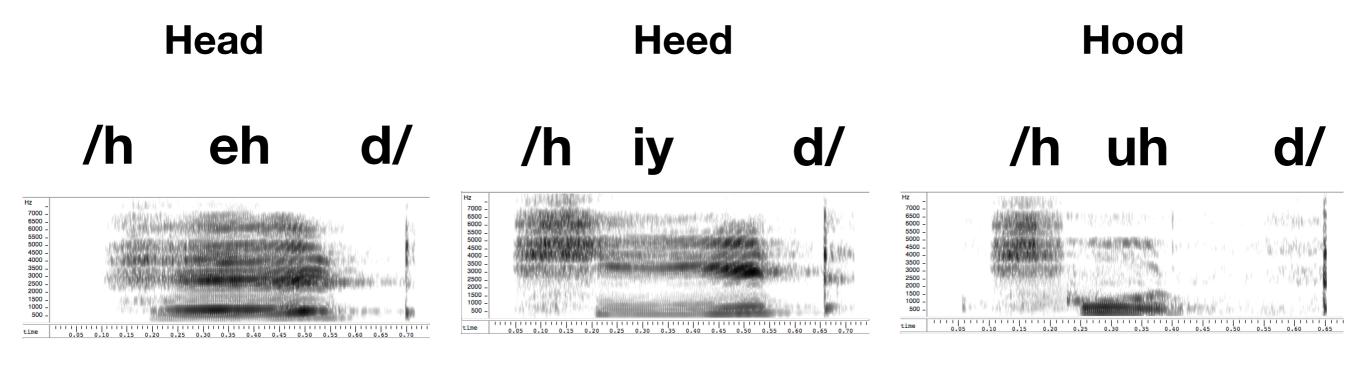


- Spectrogram reveals some segmental structure with distinct properties
- These are phonemes perceptually distinct speech sounds

 But, Speech is not just the phonemes

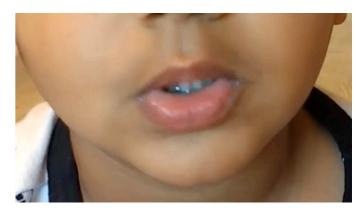
 This is a short example

Coarticulation: speech is not segmental or linear





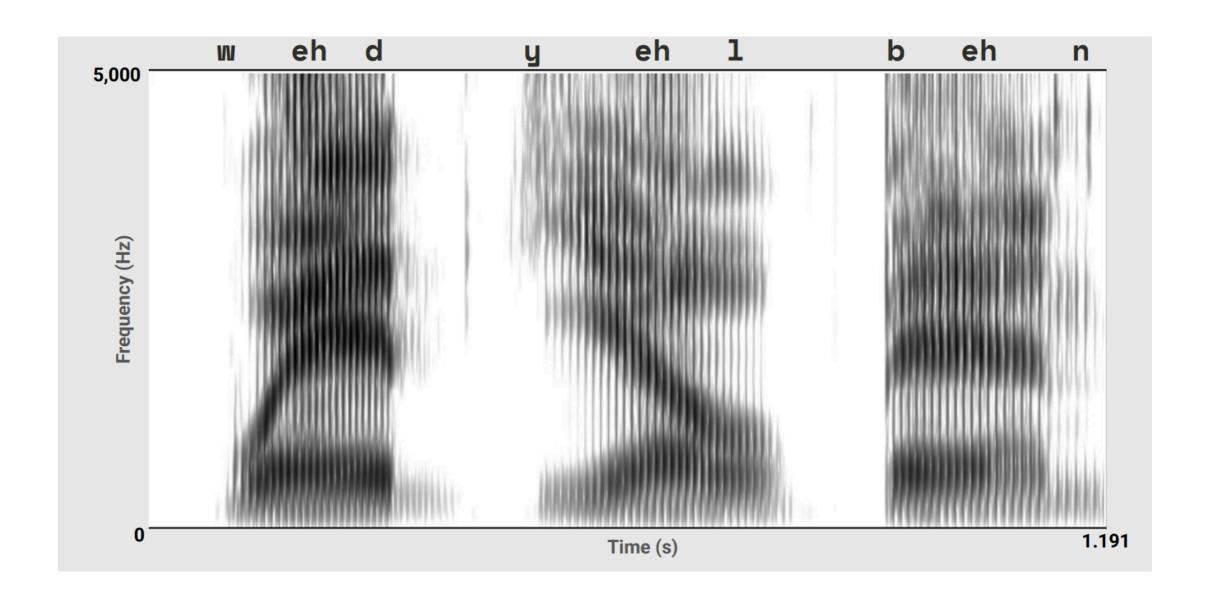




Co-articulation: An elegant phonological-articulatory transformation to facilitate rapid communication

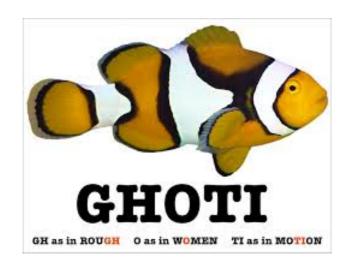
Coarticulation = Bad for pattern matching/

3 different "eh" sounds



International Phonetic Alphabet (IPA)

- THE INTERNATIONAL PHONETIC ALPHABET (revised to 2020)
- Phoneticians compiled a common set of sounds used to codify different speech sounds (across languages)
- English spelling is not phonetic
 - It has about 40 distinct phonemes represented by 26 graphemes
 - About 16 vowel sounds
 - About 24 consonant sounds



| CONSONANTS (PULMONIC) ⊕ ⊕ ⊕ 2020 IPA | | | | | | | | | | | | | | | | | |
|--------------------------------------|----------|-------------|--------|--------------|-------------------|---|----------------|---|---------|---|----|--------|--------------|------------|---|---------|---|
| | Bilabial | Labiodental | Dental | Alveolar | olar Postalveolar | | Retroflex | | Palatal | | ar | Uvular | | Pharyngeal | | Glottal | |
| Plosive | рb | | | t d | | t | d | С | J | k | g | q | G | | | 3 | |
| Nasal | m | m | | n | | | η | | n | | ŋ | | N | | | | |
| Trill | В | | | \mathbf{r} | | | | | | | | | \mathbf{R} | | | | |
| Tap or Flap | | V | | \mathbf{r} | | | \mathfrak{r} | | | | | | | | | | |
| Fricative | φβ | f v | θ ð | s z | \int 3 | ş | Z, | ç | j | X | γ | χ | \mathbf{R} | ħ | ? | h | ĥ |
| Lateral fricative | | | | łţ | | | | | | | | | | | | | |
| Approximant | | υ | | J | | | J | | j | | щ | | | | | | |
| Lateral approximant | | | | 1 | | | l | | Λ | | L | | | | | | |

Symbols to the right in a cell are voiced, to the left are voiceless. Shaded areas denote articulations judged impossible.

CONSONANTS (NON-PULMONIC)

| Clicks | Voiced implosives | Ejectives |
|------------------|-------------------|-----------------------|
| O Bilabial | 6 Bilabial | , Examples: |
| Dental | d Dental/alveolar | p' Bilabial |
| ! (Post)alveolar | f Palatal | t' Dental/alveolar |
| + Palatoalveolar | g Velar | k' Velar |
| Alveolar lateral | G Uvular | S' Alveolar fricative |

OTHER SYMBOLS

M Voiceless labial-velar fricative

W Voiced labial-velar approximant

U Voiced labial-palatal approximant
 H Voiceless epiglottal fricative

Yoiced epiglottal fricative

P Epiglottal plosive

| € Z Alveolo-palatal fricatives | |
|--------------------------------|--|
| J Voiced alveolar lateral flap | |
| | |

 ${\mathfrak f}$ Simultaneous ${\mathfrak f}$ and ${\mathfrak X}$

Affricates and double articulations can be represented by two symbols joined by a tie bar if necessary.

DIACRITICS

| 0 | Voiceless | ņ d | Breathy voiced b. a. Dental t. d. |
|----|-----------------|----------------------|---|
| ~ | Voiced | ş ţ | \sim Creaky voiced $\stackrel{b}{b}$ $\stackrel{a}{a}$ Apical $\stackrel{t}{b}$ $\stackrel{d}{d}$ |
| h | Aspirated | $t^{ m h} d^{ m h}$ | Linguolabial t d Laminal t d |
| , | More rounded | Ş | w Labialized t^{w} d^{w} $^{\sim}$ Nasalized $	ilde{e}$ |
| c | Less rounded | ò | $^{ m j}$ Palatalized ${ m t}^{ m j}$ ${ m d}^{ m j}$ $^{ m n}$ Nasal release ${ m d}^{ m n}$ |
| + | Advanced | ų | $^{\gamma}$ Velarized t^{γ} d^{γ} l Lateral release d^{l} |
| _ | Retracted | ė | $^{\Gamma}$ Pharyngealized $\ t^{\Gamma} \ d^{\Gamma}$ $^{\gamma}$ No audible release $\ d^{\gamma}$ |
| | Centralized | ë | ~ Velarized or pharyngealized } |
| × | Mid-centralized | ě | Raised $\underbrace{\mathbf{P}}_{\mathbf{L}}$ ($\underbrace{\mathbf{I}}_{\mathbf{L}}$ = voiced alveolar fricative) |
| , | Syllabic | ņ | Lowered $\underbrace{\mathbf{e}}_{T}$ ($\underbrace{\mathbf{F}}_{T}$ = voiced bilabial approximant) |
| ^ | Non-syllabic | ě | Advanced Tongue Root $\stackrel{\longleftarrow}{\mathbf{e}}$ |
| N. | Rhoticity | or ar | Retracted Tongue Root P |

Some diacritics may be placed above a symbol with a descender, e.g. $\mathring{\Pi}$



Close i • y

Close-mid e •

Open-mid

Where symbols appear in pairs, the to the right represents a rounded vov

w•u

SUPRASEGMENTALS

| 1 | Primary stress | founa |
|----------|-------------------------|------------------|
| 1 | Secondary stress | |
| I | Long | er |
| • | Half-long | e^{\centerdot} |
| | Extra-short | ĕ |
| | Minor (foot) group | |
| | Major (intonation) grou | p |
| | Syllable break | лі.ækt |
| \smile | Linking (absence of a b | reak) |
| | TONES AND WORD | COENT |

TONES AND WORD ACCENT LEVEL CONTOUR

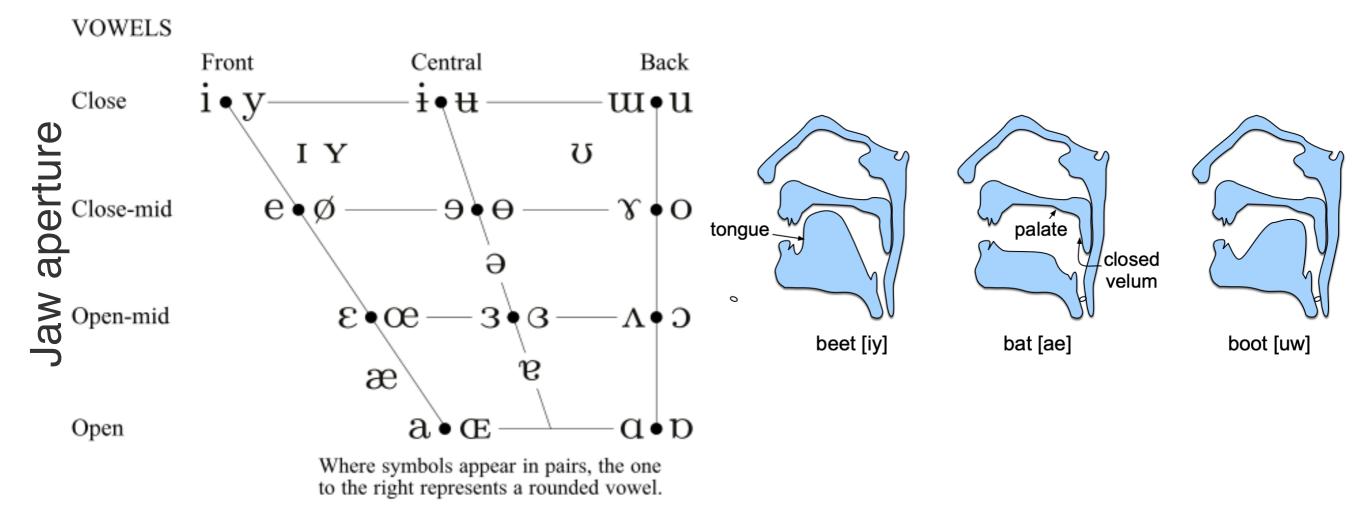
| e | or | high | e $^{\circ}$ | r / Risin |
|--------------------|-------|--------------|----------------|-------------|
| é | Ⅎ | High | ê | \ Falli |
| $\bar{\mathbf{e}}$ | 4 | Mid | é | 1 High |
| è | 4 | Low | ĕ | Low |
| ë | | Extra low | è | ↑ Risin |
| \downarrow | Doume | ten | 7.0 | Flobal rice |

↑ Upstep \(\square\) Globa

International Phonetic Alphabet (IPA)

- Vowels are characterized by jaw position and tongue shape
- Some vowels also use lips (eg. sound uw in cool)

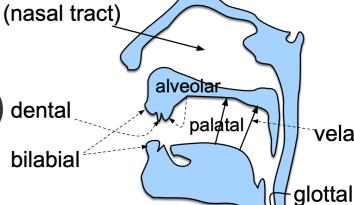
Tongue frontness



Manner of articulation

International Phonetic Alphabet (IPA)

- Consonants are characterized by place and manner of articulation
 - /p/ is caused by constriction at lips (labial)
 - /p/ is caused by sudden release of air (plosive) dental



velar

Place of articulation CONSONANTS (PULMONIC)

| | Bila | bial | Labio | dental | Dental Alveolar Postalveolar R | | Retr | oflex | Pal | atal | Velar | | Uvular | | Pharyngeal | | Glottal | | | | | |
|------------------------|------|------|-------|--------------|--------------------------------|---|------|--------------|-----|------|-------|----|--------|---|------------|---|---------|--------------|---|---|---|---|
| Plosive | р | b | | | | | t | d | | | t | d | С | J | k | g | q | G | | | 3 | |
| Nasal | | m | | m | | | | n | | | | η | | n | | ŋ | | N | | | | |
| Trill | | В | | | | | | r | | | | | | | | | | R | | | | |
| Tap or Flap | | | | \mathbf{V} | | | | ſ | | | | τ | | | | | | | | | | |
| Fricative | ф | β | f | V | θ | ð | S | \mathbf{Z} | ſ | 3 | ş | Z, | ç | j | x | γ | χ | \mathbf{R} | ħ | ? | h | ĥ |
| Lateral fricative | | | | | | | 4 | ৳ | | | | | | | | | | | | | | |
| Approximant | | | | υ | | | | J | | | | J | | j | | щ | | | | | | |
| Lateral approximant | | | | | | | | 1 | | | | 1. | | Λ | | L | | | | | | |

Symbols to the right in a cell are voiced, to the left are voiceless. Shaded areas denote articulations judged impossible.

What is phone and phoneme??? GO TO: "g oʊ t u" or "G OW T UW"



- Phone: g oʊ t u
 - Devised by International Phonetic Association
 - Physical categorization of speech sound
 - Not applicable to all languages, needs special characters, too many variations
- Phoneme: one of the units that distinguish one word from another in a particular language
 - /r/ and /l/ are degenerated in some languages (e.g., "rice" and "lice" sounds same for me!). Then, we don't have to distinguish them.
 - ARPAbet: G OW T UW
 - Proposed by ARPA for the development of speech recognition of only "American English"
 - Represented by ASCII characters

Pronunciation dictionary



- CMU dictionary
 - http://www.speech.cs.cmu.edu/cgi-bin/cmudict

"I want to go to campus"

→AY W AA N T T UW G OW T UW K AE M P AH S

- Powerful, but limited
- Out of vocabulary issue, especially new word
 - → Grapheme2Phoneme mapping based on machine learning

From letters to sounds



- Pronunciation dictionaries (often made by linguists) give the syllables and phonemes within each word in vocabulary
 - CMU Phonetic Dictionary gives the syllabic and phonetic spellings for >110K words in English
 - ML based phonetizers are built on such phonetic dictionaries

```
Graphemes She just had a baby

IPA ∫ix dʒʌst hæd ə 'beɪbi

Arpabet sh iy jh ah s t h ae d ah b ey b iy
```

Arpabet is an ASCII friendly representation of IPA

Phonology



- Phonology are the grammatical rules that phonemes of a language follow
- Lexical Phonology: Study of rules that govern the organization of sounds in a language (Phonemes —> Syllables —> Words)
 - Lexical Stress (project (noun) vs project (verb))
 - Allophones
 - (r and I in Japanese; p and b in Arabic; t and k in Hawaiian)
 - Phonological changes in continuous speech
 - Westside vs Westend
- Intonational Phonology: Study of the Fundamental Frequency (F0) in relation to the intended meaning of an utterance
 - I never said she stole my money (Emphasize each word for different meanings)

Lexical Phonology



 Phonology is the study of rules that govern the organization of sounds in a language (Phonemes —> Syllables —> Words)

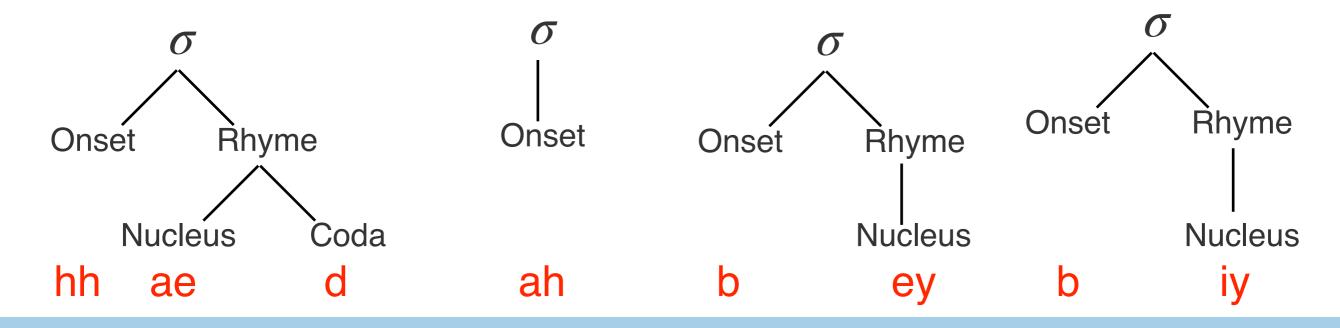
Syllabic constituency

 σ (denotes a syllable)

Onset Rhyme

Nucleus Coda

How is "had a baby" composed?



Speech in the Wild



Audio is neither clean nor just restricted to spoken language

- What technologies fall under spoken language research?
 - Speech recognition
 - Speech synthesis
 - Voice conversion
 - Speaker recognition
 - Language recognition
 - Speech emotion recognition
 - Speaker diarization
 - Speech coding
 - Speech perception

- Speech enhancement
- Microphone array processing
- Audio event classification and detection
- Speech separation
- Spoken language understanding
- Spoken dialogue systems
- Speech translation
- Multimodal processing

Speech in the Wild



Audio is neither clean nor just restricted to spoken language

- What technologies fall under spoken language research?
 - Speech recognition
 - Speech synthesis
 - Voice conversion
 - Speaker recognition
 - Language recognition
 - Speech emotion recognition
 - Speaker diarization
 - Speech coding
 - Speech perception

- Speech enhancement
- Microphone array processing
- Audio event classification and detection
- Speech separation
- Spoken language understanding
- Spoken dialogue systems
- Speech translation
- Multimodal processing

Speech Recognition



aka

ASR: Automatic Speech Recognition



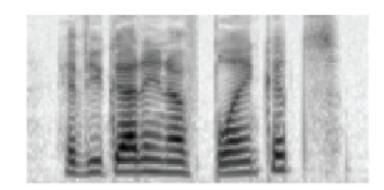
Speech Perception: Transforming Sound to Meaning



Have you got enough blankets?

hiv y gar if Af blæn kits

hivygaritnfblæŋkits





sentence

meaning word

syllable

phoneme

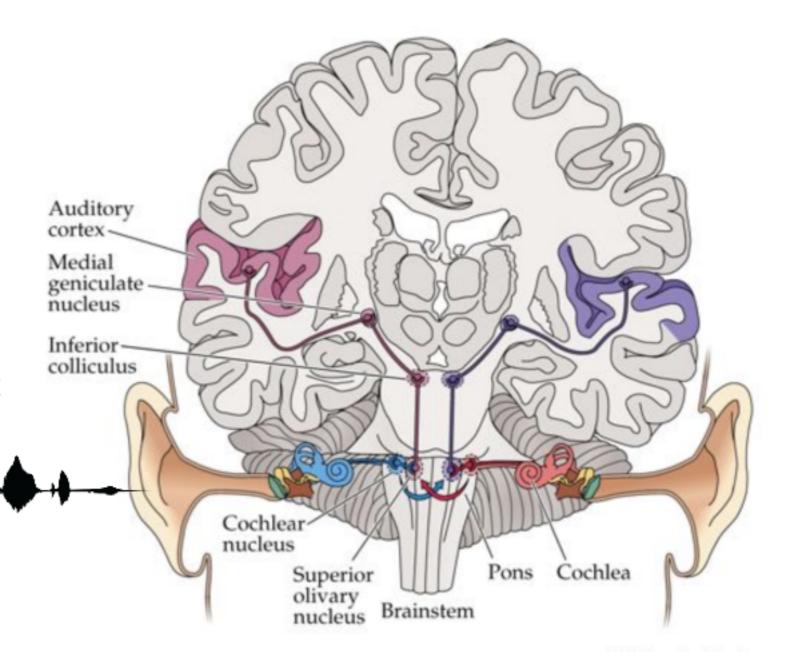
spectrotemporal (frequency decomposition)

acoustic vibration

Auditory Pathway



- ➤ Sounds reach the ear
- Sound pressure wave vibrates tympanic membrane
- Vibration converted to electrical signal
 - sent to auditory nerve, brain stem, thalamus, cerebral cortex



e 2001 Sinauer Associates, Inc.

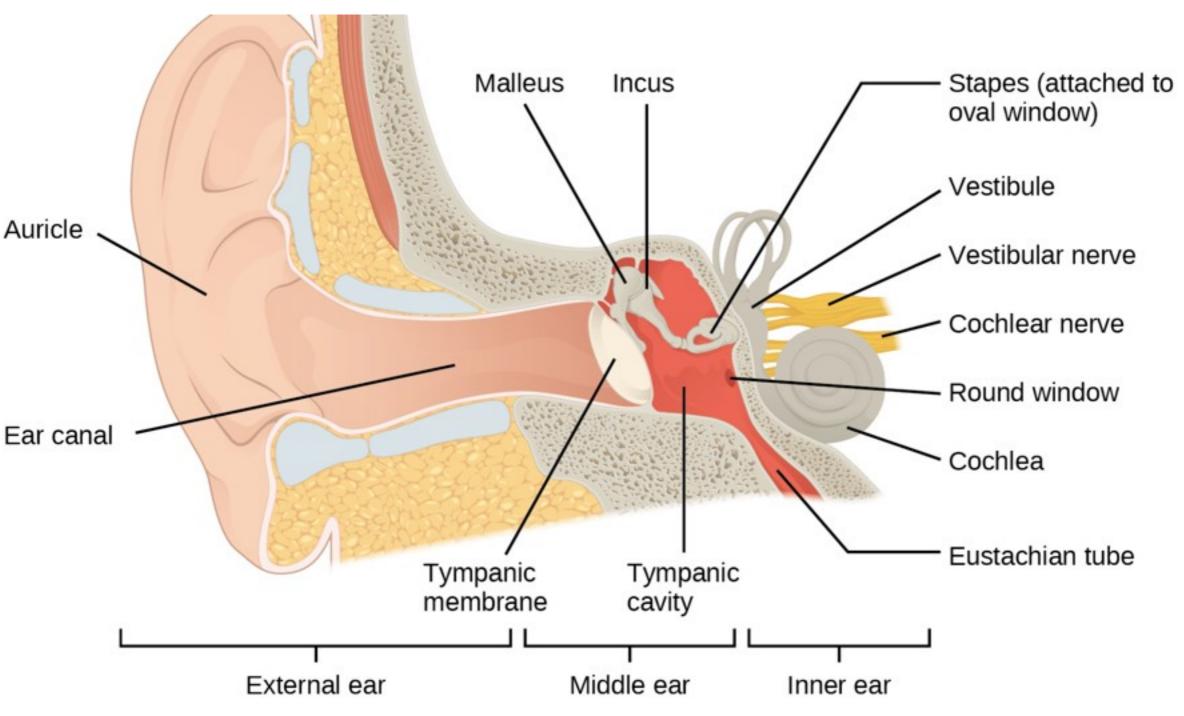
Computations along the auditory pathway



| Anatomy | Proposed Computation |
|-----------------|---|
| Cochlea | Frequency decomposition of sounds (FFT) |
| Brain stem | Sound localization (azimuth and elevation) |
| Cerebral cortex | Acoustic to phonetic transformation, meaning of words (semantics), higher cognitive functions (decision making) |

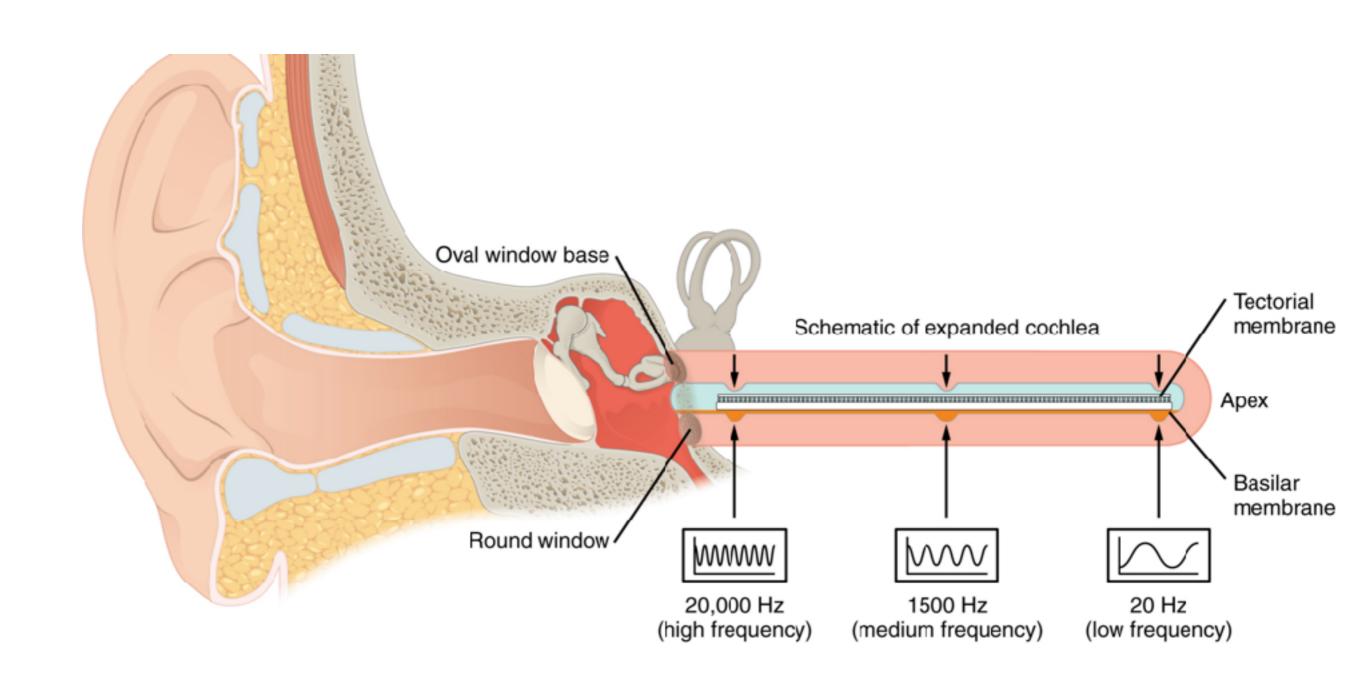
The ear





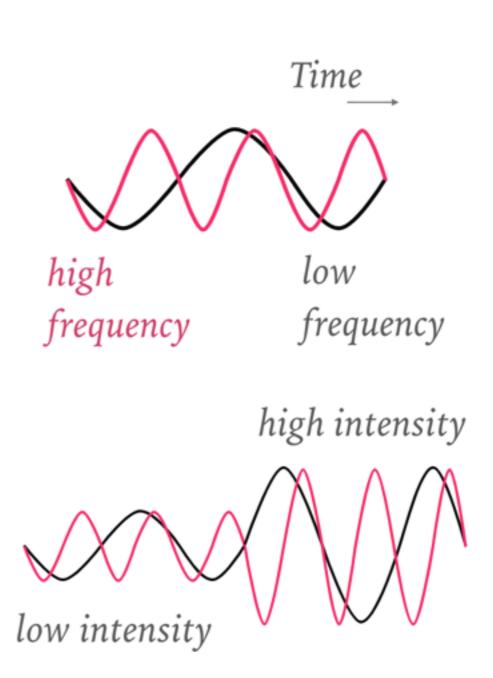
FFT in the ear





Important parameters of sound waves

- ➤ Frequency: # of waves that pass any point in one second. Measured in cycles/sec (Hz)
 - subjective correlate: pitch
- ➤ Intensity/amplitude:
 Magnitude of the
 movements produced
 (measured in dB)
 - subjective correlate: loudness



Perception of Sound

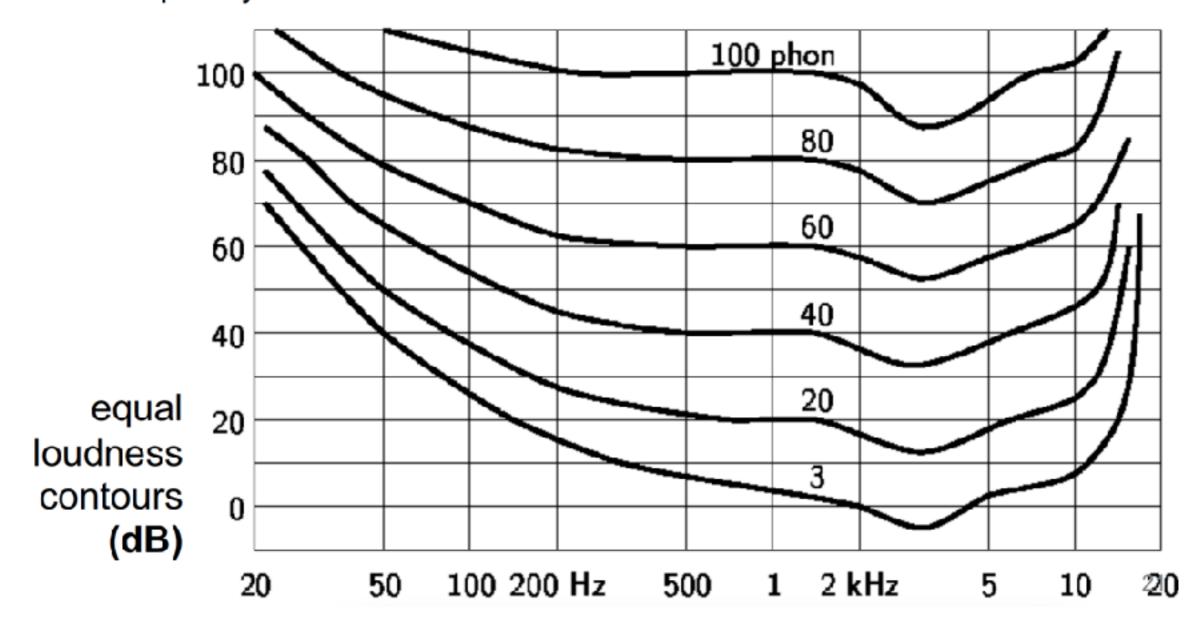


| Physical Quantity | Perceptual Quality |
|----------------------------------|--------------------|
| Intensity | Loudness |
| Fundamental frequency (f_0) | Pitch |
| Spectral envelope (formants) | Timbre |
| Onset/offset time | Timing |
| Phase difference in the two ears | Location |

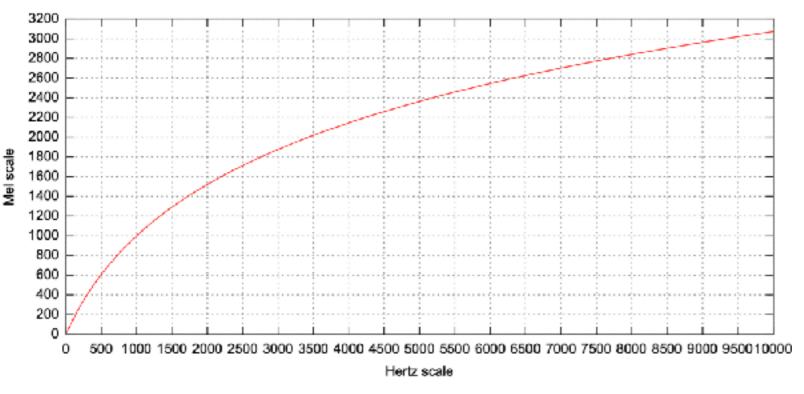
Perception of Sound



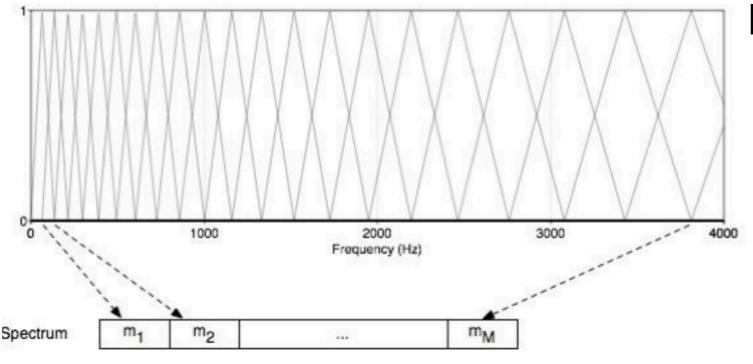
One divergence between perceptual and physical properties is the non-uniform equal loudness contours. The ears are most sensitive to sounds with a frequency of 3 ~ 4 kHz.



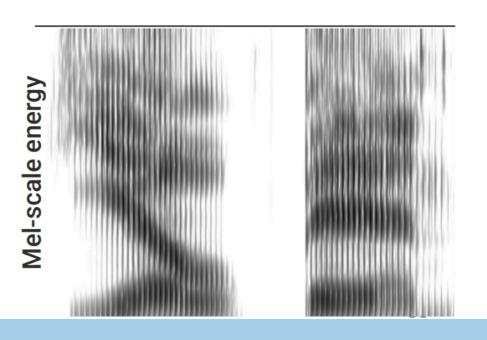
Mel scale: A Logarithmic filter bank for Perception



Mel Filter bank for MFCC

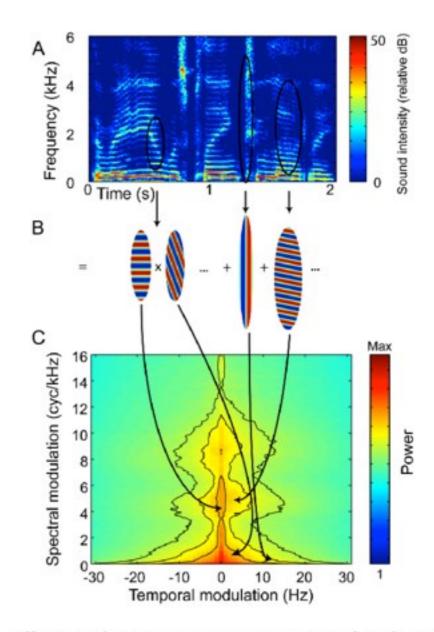


Mel Frequency Spectrogram



Speech is a small part of all possible sound

- Sounds can be decomposed into spectral and temporal modulations
- Natural stimuli (including speech) can be made up of many types
 - ➤ If high on one dimension, usually low on the other



Elliott & Theunissen PLoS Computational Biology 2011

Speech Recognition



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

| English Tasks | WER% |
|---|------|
| LibriSpeech audiobooks 960hour clean | 1.4 |
| LibriSpeech audiobooks 960hour other | 2.6 |
| Switchboard telephone conversations between strangers | 5.8 |
| CALLHOME telephone conversations between family | 11.0 |
| Sociolinguistic interviews, CORAAL (AAVE) | 27.0 |
| CHiMe5 dinner parties with body-worn microphones | 47.9 |
| CHiMe5 dinner parties with distant microphones | 81.3 |
| Chinese (Mandarin) Tasks | CER% |
| AISHELL-1 Mandarin read speech corpus | 6.7 |
| HKUST Mandarin Chinese telephone conversations | 23.5 |

Figure 27.1 Rough Word Error Rates (WER = % of words misrecognized) reported around 2020 for ASR on various American English recognition tasks, and character error rates (CER) for two Chinese recognition tasks.

Conversational Speech



- Utterance without Context
- With context

HSR versus ASR

| nan | |
|------|--|
| is | |
| it | |
| and | |
| have | |
| a | |
| that | |
| i | |

| Deletions | | | Insertions | | | | |
|------------|----------|------------|------------|---------|----------|---------|----------|
| SWB CH | | SWB | | CH | | | |
| ASR | Human | ASR | Human | ASR | Human | ASR | Human |
| 30: it | 19: i | 46: i | 20: i | 13: i | 16: is | 23: a | 17: is |
| 20: i | 17: it | 46: it | 18: and | 10: a | 14: %hes | 14: is | 17: it |
| 17: that | 16: and | 39: and | 15: it | 7: and | 12: i | 11: i | 16: and |
| 16: a | 14: that | 32: is | 15: the | 7: of | 11: and | 10: are | 14: have |
| 14: and | 14: you | 26: oh | 14: is | 6: you | 9: it | 10: you | 13: a |
| 14: oh | 12: is | 25: a | 13: not | 5: do | 6: do | 9: the | 13: that |
| 14: you | 12: the | 20: to | 10: a | 5: the | 5: have | 8: have | 12: i |
| 12: %bcack | 11: a | 19: that | 10: in | 5: yeah | 5: yeah | 8: that | 11: %hes |
| 12: the | 10: of | 19: the | 10: that | 4: air | 5: you | 7: and | 10: not |
| 11: to | 9: have | 18: %bcack | 10: to | 4: in | 4: are | 7: it | 9: oh |

Table 3: Most frequent deletion and insertion errors for humans and ASR system on SWB and CH.

| SWB | | СН | | |
|---------------------|-----------------|----------------|---------------------|--|
| ASR | Human | ASR | Human | |
| 11: and / in | 16: (%hes) / oh | 21: was / is | 28: (%hes) / oh | |
| 9: was / is | 12: was / is | 16: him / them | 22: was / is | |
| 7: it / that | 7: (i-) / %hes | 15: in / and | 11: (%hes) / %bcack | |
| 6: (%hes) / oh | 5: (%hes) / a | 8: a / the | 10: bentsy / benji | |
| 6: him / them | 5: (%hes) / hmm | 8: and / in | 10: yeah / yep | |
| 6: to o / to | 5: (a-) / %hes | 8: is / was | 9: a / the | |
| 5: (%hes) / i | 5: could / can | 8: two / to | 8: is / was | |
| 5: then / and | 5: that/it | 7: the / a | 7: (%hes) / a | |
| 4: (%hes) / %bcack | 4: %bcack / oh | 7: too / to | 7: the / a | |
| 4: (%hes) / am | 4: and / in | 6: (%hes) / a | 7: well / oh | |

Table 2: Most frequent substitution errors for humans and ASR system on SWB and CH.

Why Study ASR?



- In the last ~5 years
 - Dramatic reduction in LVCSR error rates (16% to <3%)
 - Human level LVCSR performance on Switchboard
 - New class of recognizers (end to end neural network)

- Understanding how ASR works enables better ASR-enabled systems
 - What types of errors are easy to correct?
 - How can a downstream system make use of uncertain outputs?
 - How much would building our own improve on an API?

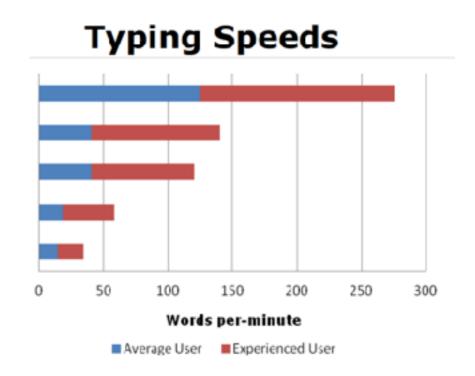
Next generation of ASR challenges as systems go live on phones and in homes

Why Study ASR?





With speech, interaction becomes independent of screen estate



If accurately recognized, speech is three times faster than QWERTY (Basapur et al. '07)



Only plausible interaction modality for 800 million non-literate users

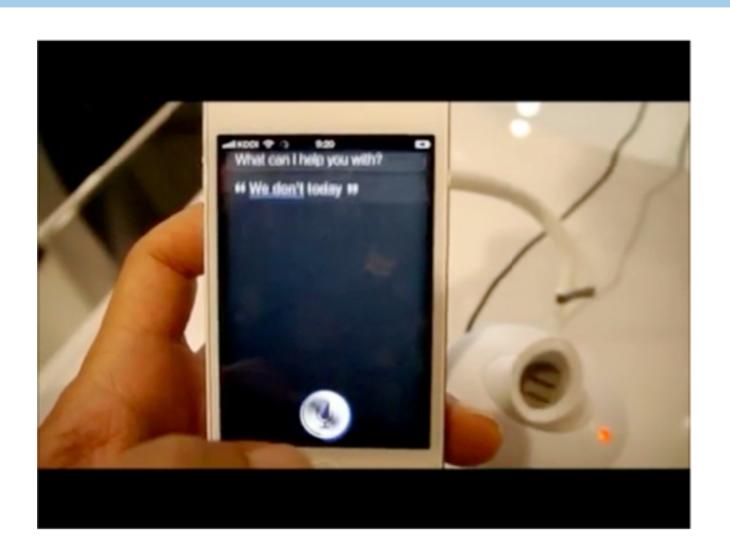
Design of Classical ASR Systems



- 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

Current ASR Systems





- ✓ Adult Native English Speakers
- ✓ Clean Lab Environment
- X Children
- X Non-native Accent
- X Noisy Environment
- X Variable Speaking Rate

No "universal" speech recognizer – one must be developed and optimized for each use case

Why is ASR Hard?



| written text: | Why is speech Recognition so Difficult? |
|----------------------------|---|
| spontaneous: | why's speech recognition so difficult |
| continuous: | whysspeechrecognitionsodifficult |
| pronunciation: | whazbeechregnizhnsadifcld |
| acoustic variability: | whazbrechnegnizhnsadifold |
| noise: | Character was a subject to |
| Cocktail party- Effect: | Land the second of the second |

What are the factors that determine difficulty

COMPLEXITY

amount of data: typically 32000 bytes per second (16khz)

class inventory: 50 phonemes, 5000 sounds, 100.000 words

combinatorial explosion: exponential growth of possible sentences

SEGMENTATION

our perception: Phones, syllables, words, sentences

actually there are: no boundary markers, continuous flow of samples

VARIABILITY

speaker: anatomy of vocal tract, speed, loudness, acoustic

stress, mood, dialect, speaking style, context

channel, environment: noise, microphones, channel conditions

AMBIGUITY

Homophones: two vs. too,

Word Boundaries: interface vs. in her face,

Semantics: He saw the Grand Canyon flying to New York,

Pragmatics: Time flies like an arrow.

Why is HSR easy?



"The main prerequisite of the uniquely human communication is that speaker and listener must have a common understanding that out of all possible sounds man can produce and hear, only a few have linguistic significance."

(Olli Aaltonen& Esa Uusipaikka: Why Speaking Is so Easy? – Because Talking Is Like Walking with a Mouth)

- Important feature of speech perception: we hear sounds either as speech or non-speech
- Once defined as speech, we hear a sequence of vowels and consonants not as buzzes and hisses, the segmentation into words happens on the fly
- Abstract away from sound variability we use an enormous database of background knowledge: phonotactics, morphology, syntax, semantics, pragmatic knowledge

Robustness of HSR



"The main prerequisite of the uniquely human communication is that speaker and listener must have a common understanding that out of all possible sounds man can produce and hear, only a few have linguistic significance."

(Olli Aaltonen& Esa Uusipaikka: Why Speaking Is so Easy? – Because Talking Is Like Walking with a Mouth)

- Important feature of speech perception: we hear sounds either as speech or non-speech
- Once defined as speech, we hear a sequence of vowels and consonants not as buzzes and hisses, the segmentation into words happens on the fly
- Abstract away from sound variability we use an enormous database of background knowledge: phonotactics, morphology, syntax, semantics, pragmatic knowledge

We also use multiple modalities - McGurk effect

Problems and Challenges



- Speech Recognition ("speech-to-text")
 - Finding Robust Acoustic Representations of Speech (how do things sound)
 - Dictionary Learning (how to decompose words into units)
 - Language Modeling (what is likely to be said)
 - Decoding (how to get an answer in finite time)
- Adaptation and robustness of models and techniques to changing conditions
 - Multi-modal and multi-task learning, audio-visual processing (deep learning)
- Language-universal (?) modeling
 - Can we port resources across languages, to enable processing of new, unwritten, low-resource languages?
 - How do Humans do it? Language acquisition?

Problems and Challenges



- Meta-data extraction (what is "not in text")
 - Speaker identification (age, gender, ...)
 - Emotions, personalities, ...
 - Languages, dialects, ...
- Optimality Criterion?
 - Speech-to-text useful? No, unless dictation
 - Optimize directly "speech-to-meaning" or "speech-to-action"
 - All neural architectures
 - Jointly optimized speech-to-X systems
 - Speech synthesis



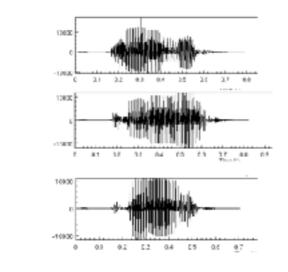
Template Matching

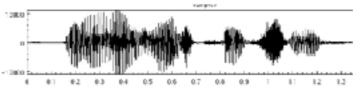
A Template based ASR System



Some applications only have limited vocabulary — Voice Dialing System

- Library
 - Mom
 - Dad
 - Bob
 - Mario's Pizza



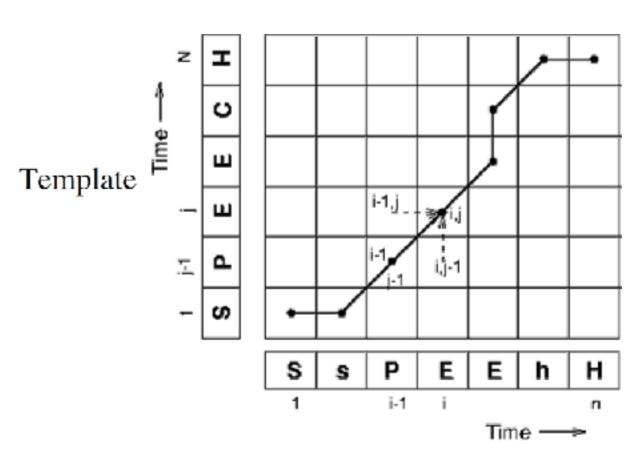


The matching problem



- We change durations
 - two utterances are never the same

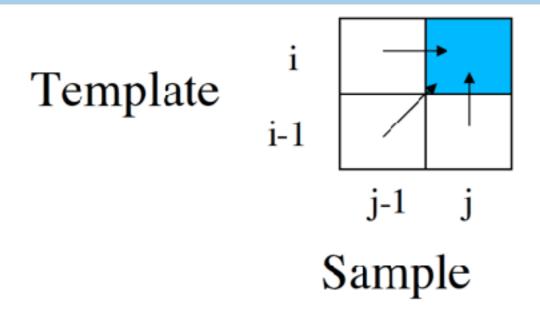
Dynamic Time Warping



Sample Speech

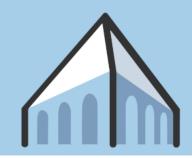
Dynamic Time Warping





- For each square

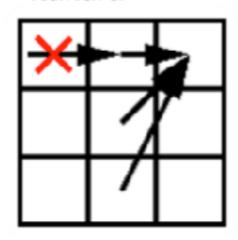
Dynamic Time Warping Paths

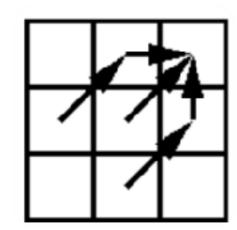


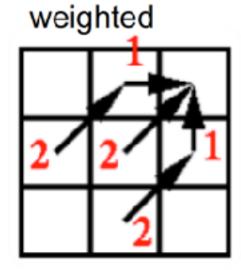
Many different warping steps are possible and have been used.

Examples:

Itakura





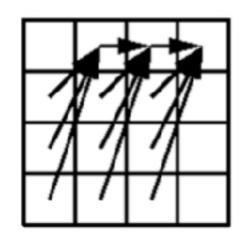


General rule is:

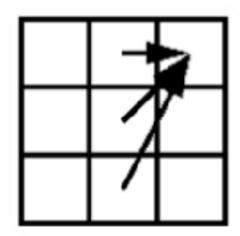
Cumulative cost of
destination = best-of
(cumulative cost of source
+ cost of step + distance in
destination)

symmetric (editing distance)



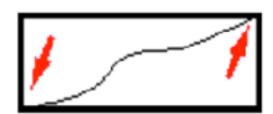


Bakis



Dynamic Time Warping Path Constraints

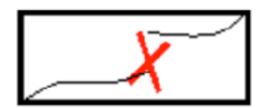
Endpoint constraints: we want the path to not skip a part at the beginning or end of the utterance



Monotonicity conditions: we can't go back in time (for any utterance)



Local continuity: no jumps etc



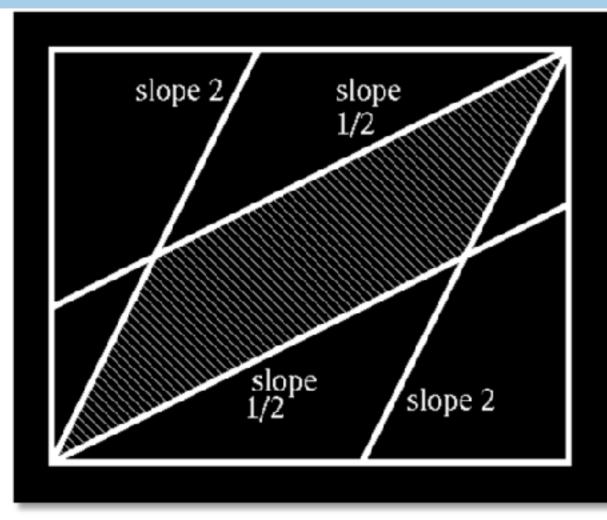
Global path constraints: path should be close to diagonal

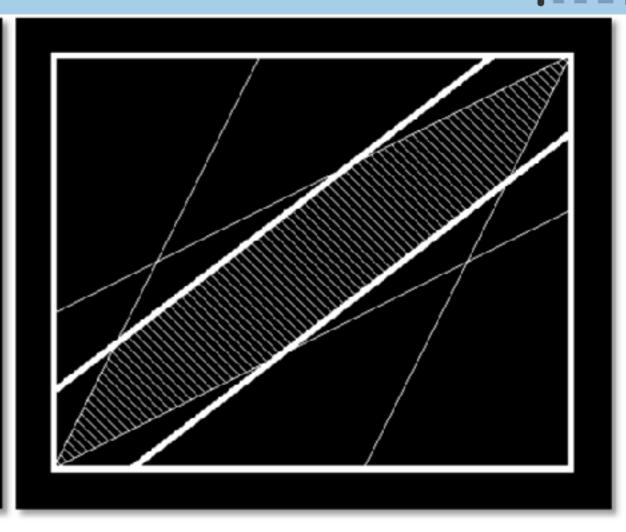


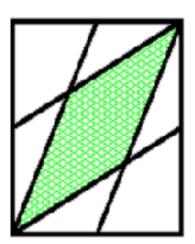
Slope weighting: we believe the DTW path sjould be somewhat "smooth"

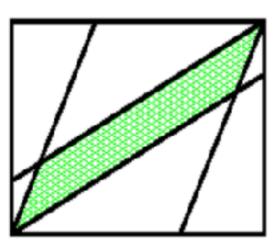


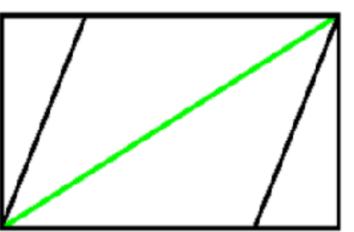
Dynamic Time Warping Path Constraints/











← only one path

DTW Search Space

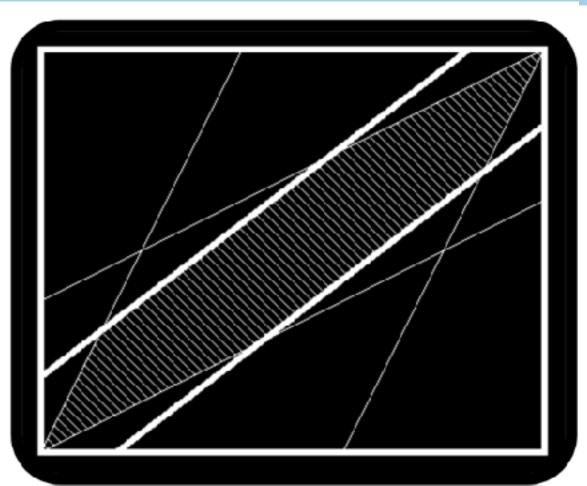


Already suggested: restrict search space by window around diagonal. Caveats:

- Silence period in one utterance can cause "edgy" path
- Search area becomes too restricted when utterance durations differ a lot

Other reason (besides global path constraints) for restricting search space:

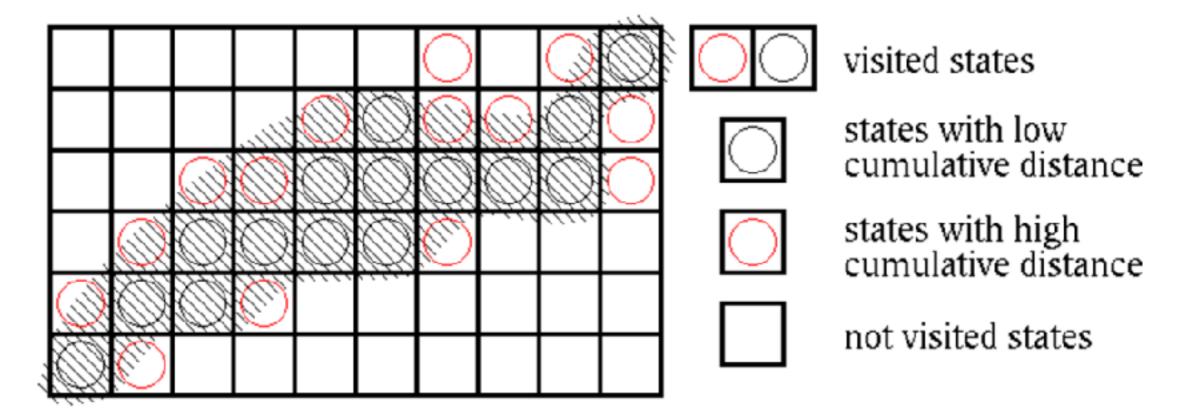
- Save time: A window that has a constant width, reduces the search effort from O(n²) to O(n)
- To overcome caveats of "diagonal window" restriction, use: beam search.



DTW with Beam Search



Idea: do not consider steps to be possible out of states that have "too high" cumulative distances.



Approaches:

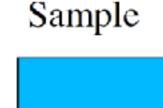
- "expand" only a fixed number of states per column of DTW matrix
- expand only states that have a cumulative distance less than a factor (the "beam") times the best distance so far

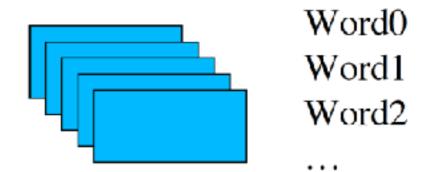
DTW and Multiple Templates



- Compare against each
- Find closest
- Need to normalize scores
 - (divide by length of matches)

Template Library





```
For Word in Templates

Score = dtw(Template[Word], Sample);

if (Score < BestScore)

BestWord = Word;

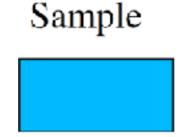
DoAction(Action[BestWord])
```

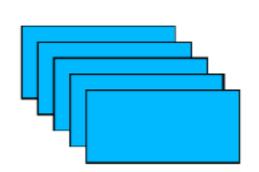
Also adapted for Speaker ID



- Compare against each
- Find closest
- Need to normalize scores
 - (divide by length of matches)

Template Library





Speaker0 Speaker1 Speaker2

For Speaker in Templates

Score = dtw(Template[Speaker], Sample);

if (Score < BestScore)

BestSpeaker = Speaker;

Template Matching



Advantages

- Works well for small number of templates (<20)
- Language independent
- Speaker specific
- Easy to train (end user controls it)

Disadvantages

- Limited number of templates
- Speaker specific
- Need actual training examples

Computing distances



- Distance metric
 - Euclidean

$$\sqrt{\sum_{i=0}^{N} (T_i - S_i)^2}$$

- But some distances are bigger than others
 - Silence is pretty similar
 - Fricatives are quite larger
 - A longer fricative might give large score
 - A longer vowel might give smaller score

Isolated Word Recognition using Template Matching

- For each word in the vocabulary, store at least one reference pattern
- When multiple reference patterns are available, either use all of them or compute an average
- During recognition
 - Record a spoken word
 - Perform pattern matching with all stored patterns (or at least with those that can be used in the current context)
 - Compute a DTW score for every vocabulary word (when using multiple references, compute one score out of many, e.g. average or max)
 - Recognize the word with the best DTW score
- This approach works only for very small vocabularies and/ or for speakerdependent recognition
- Forms the basis of Hidden Markov Model based recognizers

More reliable distances



- Instead of Euclidean distance
 - Doesn't care about the standard deviation

$$\sqrt{\sum_{i=0}^{N} (T_i - S_i)^2}$$

- Use Mahalanobis distance
 - Care about means and standard deviation

$$\sqrt{\sum_{i=0}^{N} \left(\frac{(\mu_i - S_i)}{\sigma_i}\right)^2}$$

More reliable matching



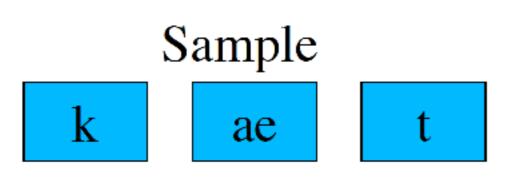
- Having multiple template examples
 - Individual matches or
 - Average them together
- DTW align all of the examples
- Collect statistics as a Gaussian
 - Mean and standard deviation for each coeff

$$\{\mu_o, \sigma_0, \mu_1, \sigma_1, \mu_2, \sigma_2, \ldots\}$$

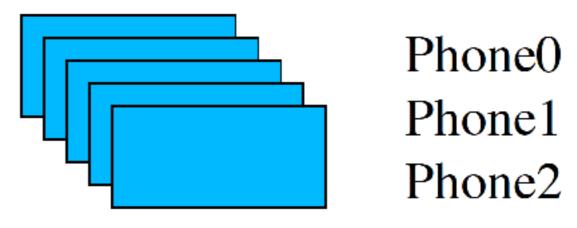
Beyond Template Matching



- String phoneme templates together
 - A template model for each phoneme



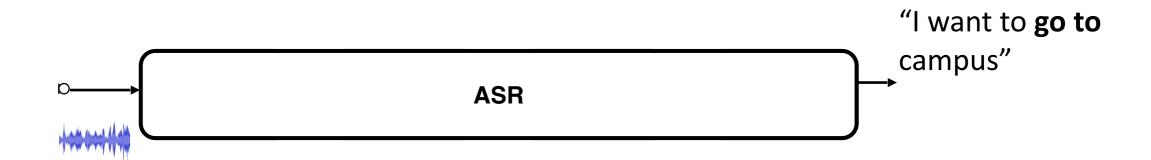
Phoneme Templates



. .

Automatic speech recognition





Automatic speech recognition





- Instead of starting from the waveform, we will often start from speech features (MFCC, etc.) through the feature extraction module
- Let's think of the conversion from speech feature o to text w

Speech recognition with a probabilistic formulation



• MAP decision theory: Estimate the most probable word sequence among all possible word sequences (I'll omit the domain sometimes) \hat{W}

$$\hat{W} = \underset{W \in \mathcal{W}}{\operatorname{argmax}} p(W|O)$$

Speech recognition with a probabilistic formulation



• MAP decision theory: Estimate the most probable W word sequence W among all possible word sequences (I'll omit the domain sometimes)

$$\hat{W} = \operatorname*{argmax}_{W \in \mathcal{W}} p(W|O)$$

- The following parts will discuss how to make this equation tractable
- To do that we need to prepare some basic math

Notation



| Туре | Font, case | Latex command | Looks like |
|-----------------|-------------------------|------------------|------------------|
| Scalar variable | Italic font, lower case | \$x\$ | \boldsymbol{x} |
| Vector variable | Bold font, lower case | \$\mathbf{x}\$ | X |
| Matrix variable | Bold font, upper case | \$\mathbf{X}\$ | \mathbf{X} |

Notation



- Please specify the domain of variables
 - D-dimensional continuous vector: $\mathbf{o} \in \mathbb{R}^D$
 - $(D \times D)$ -dimensional matrix: $\mathbf{\Sigma} \in \mathbb{R}^{D \times D}$
 - ullet Word with vocabulary $\mathcal{V}:w\in\mathcal{V}$
- Set: calligraphic font, upper case, a set of elements are represented with curly brackets

$$\mathcal{V} = \{\text{"one"}, \text{"two"}, \text{"three"}, \cdots \}$$

 Sequence: italic font, upper case, a sequence of elements are represented with round brackets

$$O = (\mathbf{o}_1, \mathbf{o}_2, \cdots)$$
 $O = (\mathbf{o}_t \in \mathbb{R}^D | t = 1, \cdots, T)$

Notation



Subsequence

$$O=(\mathbf{o}_1,\cdots,\mathbf{o}_{t_1},\mathbf{o}_{t_1+1},\cdots,\mathbf{o}_{t_2},\cdots,\mathbf{o}_{T})$$
 $o_{t_1+1:t_2}=(\mathbf{o}_{t_1+1},\cdots,\mathbf{o}_{t_2})$ or $O_{t_1+1:t_2}=(\mathbf{o}_{t_1+1},\cdots,\mathbf{o}_{t_2})$

Speech recognition cases



T-length speech feature sequence (D-dimensional vector)

$$O = (\mathbf{o}_t \in \mathbb{R}^D | t = 1, \dots, T)$$

• N-length word sequence with vocabulary $\mathcal V$

$$W = (w_n \in \mathcal{V} | n = 1, \dots, N)$$

Notation cont'd



Operation: non-italic

| Operation type | Latex command |
|----------------|-------------------------------------|
| | \log() |
| | \arg \max (), \text{argmax}(), etc. |
| | \text{sigmoid}() |

• Index: subscript, italic

$$W_n$$
, \mathbf{o}_t

• Type of variables, functions: superscript, non-italic

$$x^{\text{HMM}}, x^{\text{DNN}}$$
 $p^{\text{HMM}}(x), p^{\text{DNN}}(x)$

Probabilistic rules



Product rule

$$p(x|y)p(y) = p(x,y)$$

Sum rule

$$p(y) = \sum_{x} p(x, y)$$

Conditional independence assumption

$$p(x|y,z) = p(x|z) \qquad p(x,y|z) = p(x|z)p(y|z)$$

Other rules



Bayes rule

$$p(x|y) = \frac{p(y|x)p(x)}{p(y)} = \frac{p(y|x)p(x)}{\sum_{x} p(y|x)p(x)}$$

Probabilistic chain rule

$$p(x_1, \dots, x_N) = \prod_{n=1}^N p(x_n | x_{1:n-1})$$
 where $p(x_1 | x_{1:0}) = p(x_1)$

Both are derived with a combination of the product and/or sum rules

Other approximation



Viterbi approximation

$$p(x|y) = \sum_{z} p(x, z|y) \approx \max_{z} p(x, z|y)$$

We often use this approximation to avoid \sum_{z}

- Set an actual distribution, e.g.,
 - p(x): Gaussian distribution, Gamma distribution, softmax probability obtained by a neural network etc.

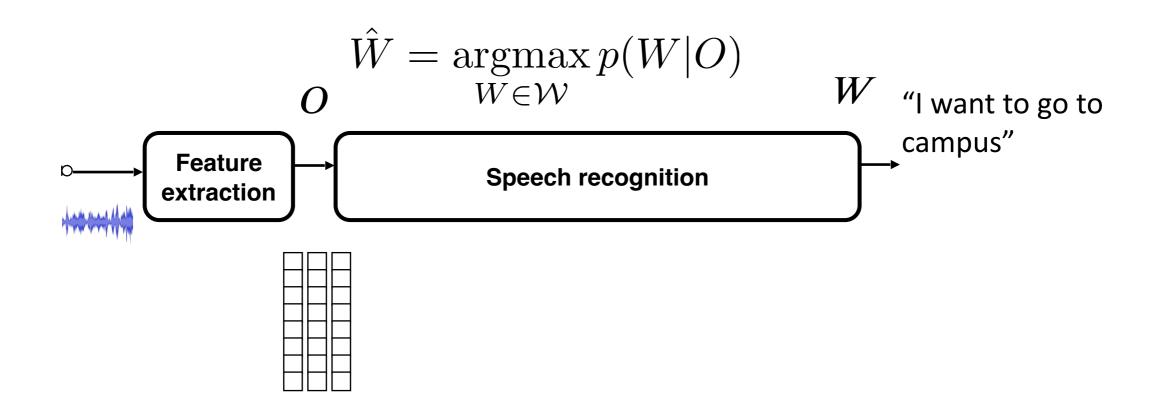


Now let us use product and sum rules and conditional independence assumption to formulate the speech recognition problem

- End-to-End Speech Recognition
- Classical speech recognition
 - Pipeline

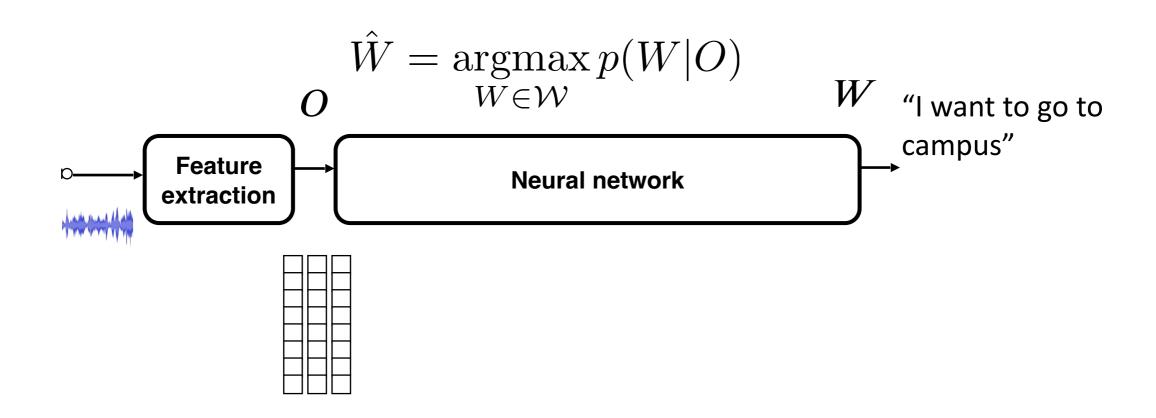
Speech recognition





End-to-end speech recognition





How to obtain the posterior p(W|O)



• We just replace it with a neural network-based function $f^{nn}(\cdot)$

$$\operatorname*{argmax}_{W} p(W|O) = \operatorname*{argmax}_{W} f^{\mathrm{nn}}(W|O)$$

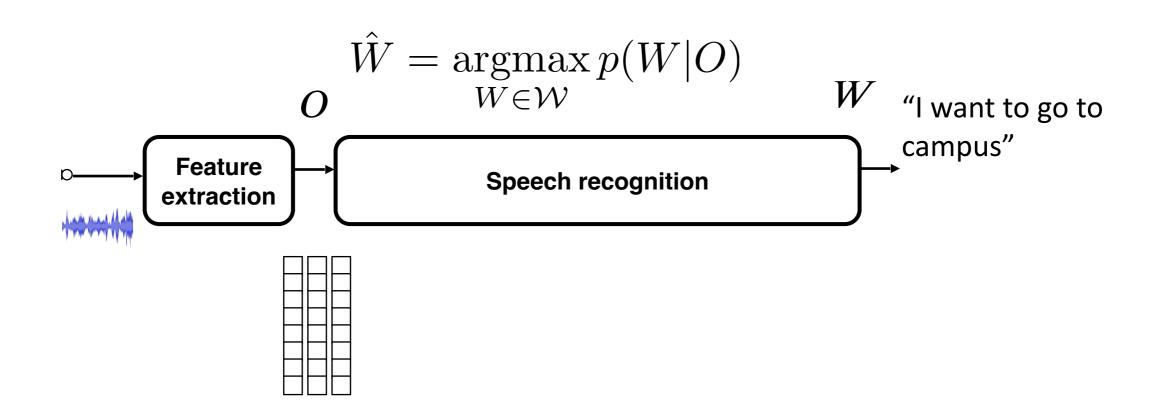
- Easy and simple, no math (in this level), however
- W is a sequence! $W = (w_n \in \mathcal{V} | n = 1, ..., N)$
 - Very difficult to deal with it
 - Say N = 10, $|\mathcal{V}| = 100$, we have to deal with 100^{10} possible sequences
 - Also, the length N is variable
 - We have to use a special neural network (e.g., attention, CTC, and RNN-transducer)



- Classical speech recognition
 - Pipeline

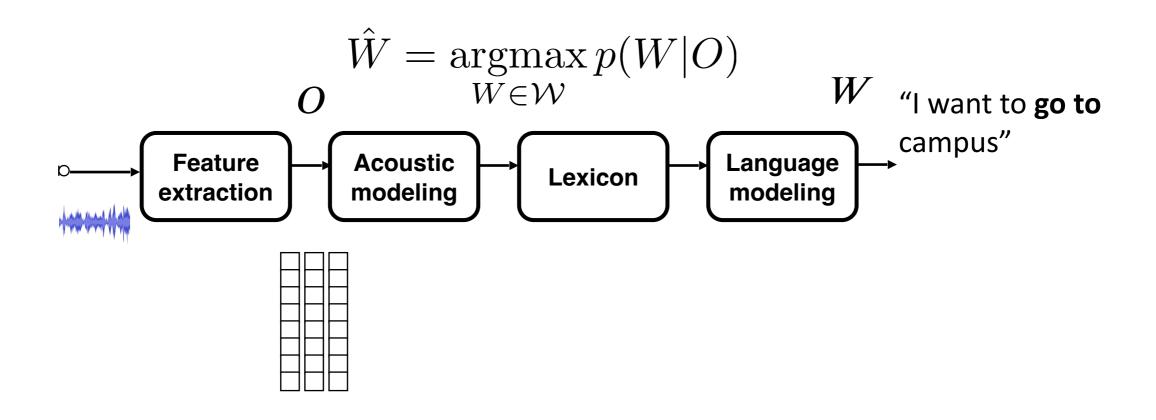
Speech recognition





Classical (non-end-to-end) speech recognition





Speech -> Text





Text W: I want to go to the campus

Speech -> Phoneme -> Text



Speech o:



Phoneme L: AY W AA N T T UW G OW T UW K AE M P AH S



Text W: I want to go to campus

How to obtain the posterior p(W|O)



- Factorize the model with phoneme
 - Let $L=(l_i\in\{/\mathrm{AA}/,\,/\mathrm{AE}/,\cdots\}|i=1,\cdots,J)$ be a phoneme sequence

$$\arg\max_{W} p(W|O) = \arg\max_{W} \sum_{L} p(W,L|O) \qquad \text{Sum rule}$$

$$= \arg\max_{W} \sum_{L} \frac{p(O|W,L)p(L|W)p(W)}{p(O)} \qquad \text{Bayes+ Product rule}$$

$$= \arg\max_{W} \sum_{L} p(O|W,L)p(L|W)p(W) \qquad \text{does not depend}$$
 on W
$$= \arg\max_{W} \sum_{L} p(O|L)p(L|W)p(W) \qquad \text{Conditional independence}$$
 assumption

Note: the right hand side is not the probability as it lacks a sum to one constraint

Noisy channel model



$$\underset{\mathbf{W}}{\operatorname{argmax}} p(W \mid O) = \underset{W}{\operatorname{argmax}} p(O \mid W) p(W) \approx \underset{W}{\operatorname{argmax}} \sum_{L} p(O \mid L) p(L \mid W) p(W)$$

Speech recognition

• p(O|L): Acoustic model (Hidden Markov model)

• p(L|W): Lexicon

• p(W): Language model (n-gram)

Noisy channel model



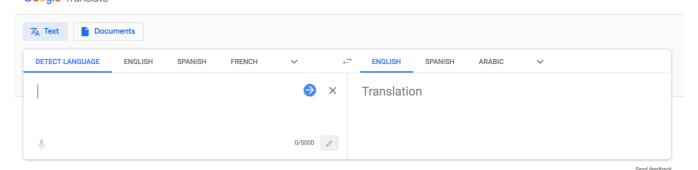
$$\underset{\mathbf{W}}{\operatorname{argmax}} p(W | Y) = \underset{\mathbf{W}}{\operatorname{argmax}} p(Y | W) p(W)$$

W: Targetlanguage textY: Source languagetext

Machine translation

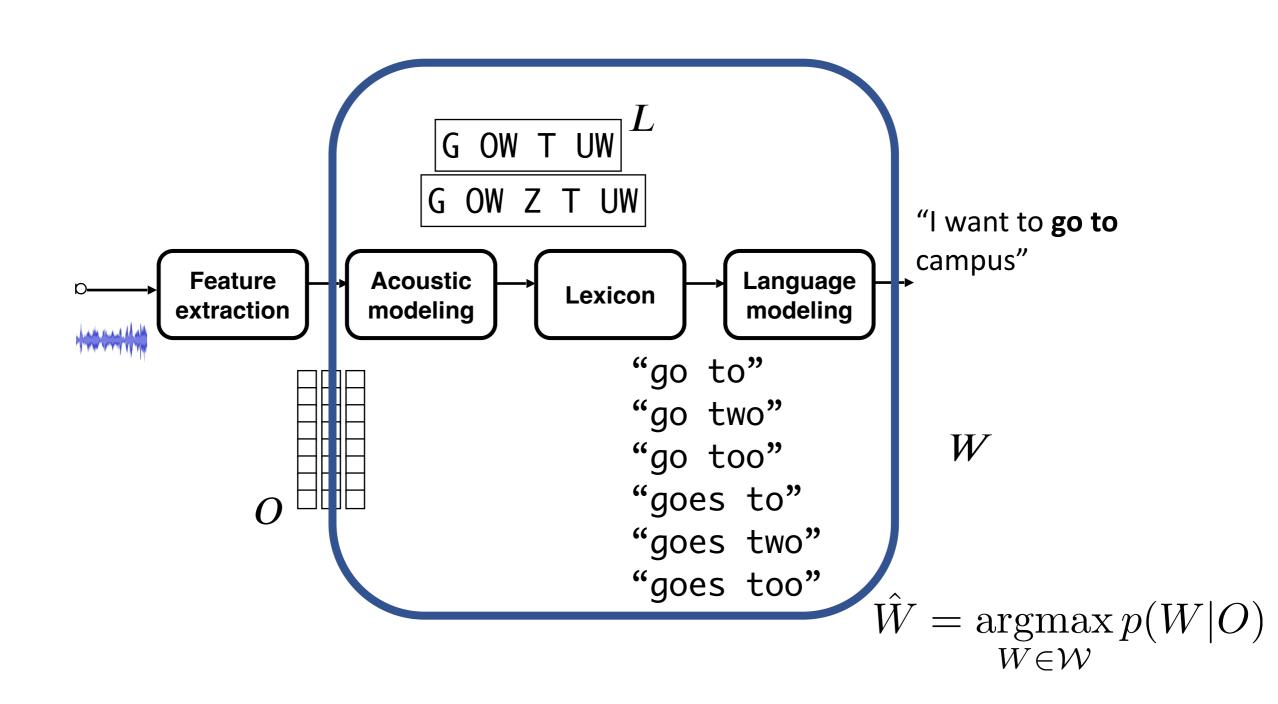
• p(Y|W): Translation model

• p(W): Language model



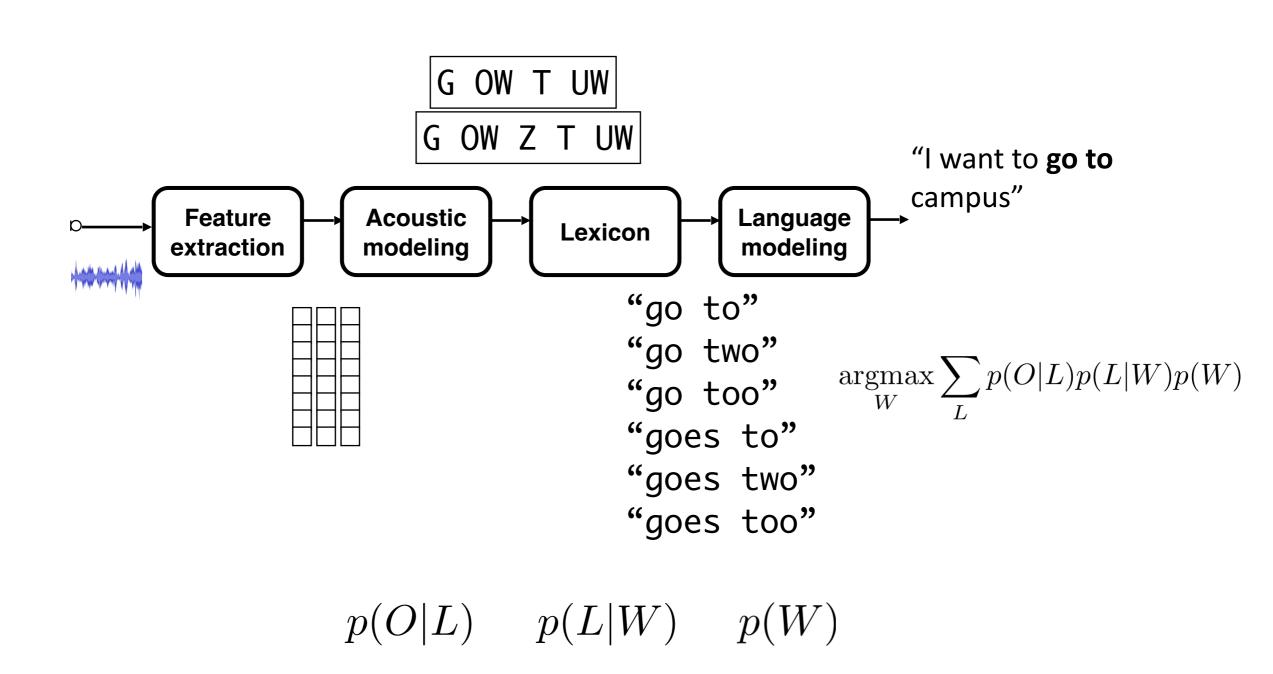
Speech recognition pipeline





Speech recognition pipeline





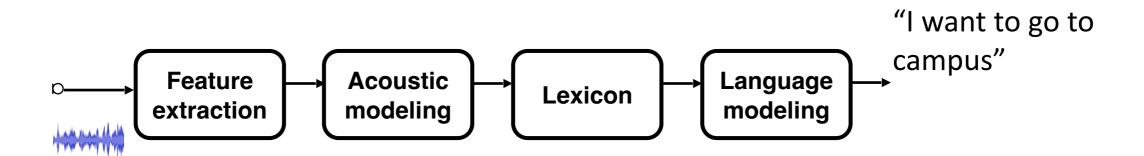
Please remember the noisy channel model

- Factorization
- Conditional independence (Markov) assumptions

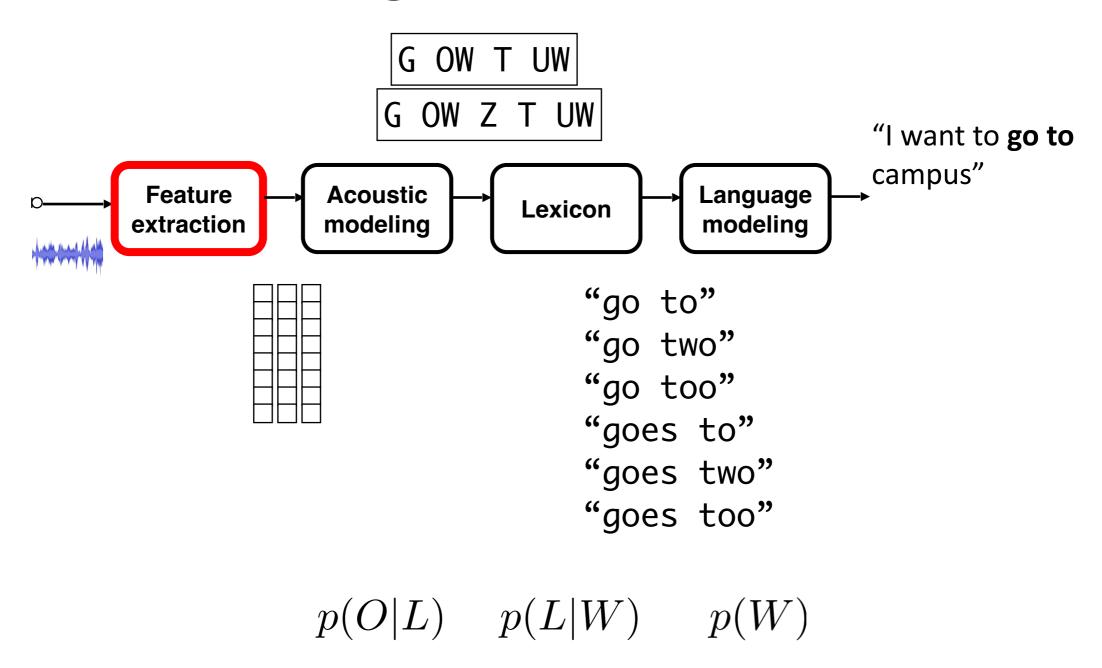
We can elegantly factorize the speech recognition problem with a tractable subproblem

Main blocks of Classical ASR



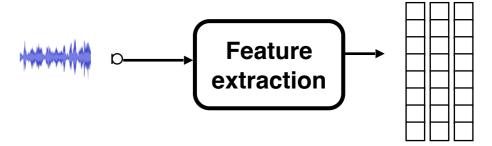


Speech recognition pipeline



Waveform to speech feature

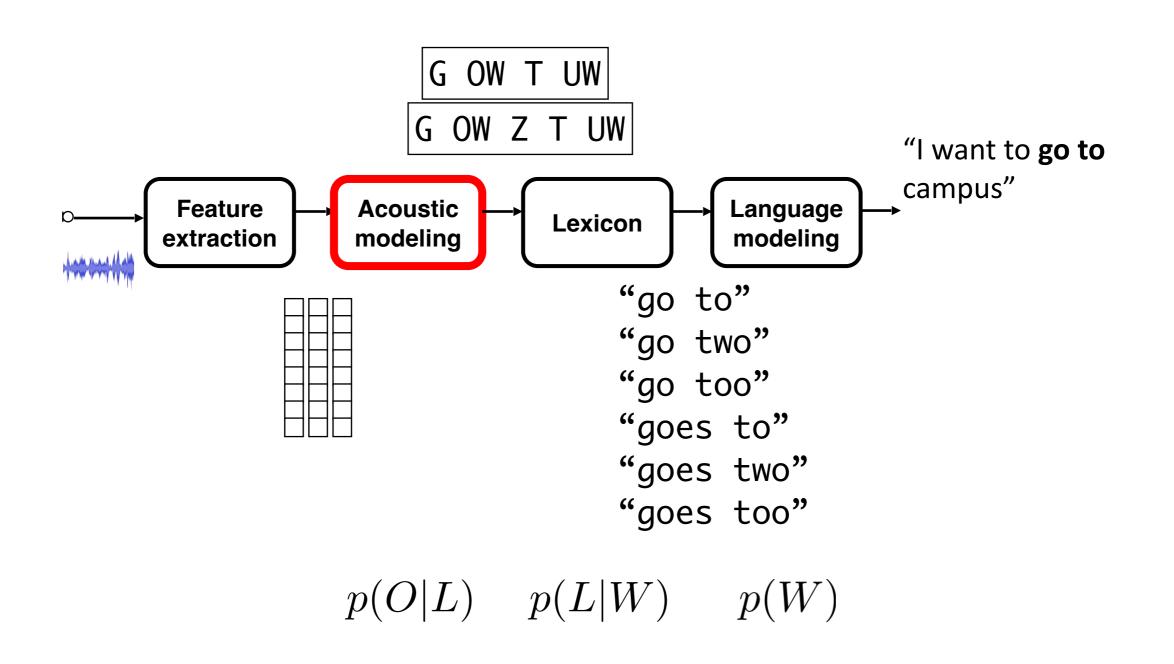




- Performed by so-called feature extraction module
 - Mel-frequency cepstral coefficient (MFCC), Perceptual Linear Prediction (PLP) used for Gaussian mixture model (GMM)
 - Log Mel filterbank used for deep neural network (DNN)
- Time scale
 - 0.0625 milliseconds (16kHz) to 10 milliseconds
- Type of values
 - Scalar (or discrete) to 12—40 dimensional vector

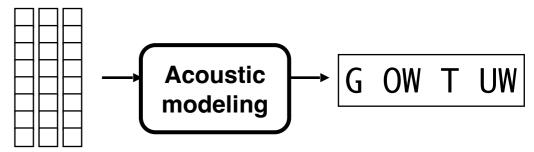
Speech recognition pipeline





Speech feature to phoneme





- Performed by so-called acoustic modeling module
 - Hidden Markov model (HMM) with GMM as an emission probability function
 - Hidden Markov model (HMM) with DNN as an emission probability function
- Time scale
 - 10 milliseconds to ~100 milliseconds (depending on a phoneme)
- Type of values
 - 12-dimensional continuous vector to 50 categorical value (~6bit)
- The most critical component to get the ASR performance
- It can be a probability of possible phoneme sequences, e.g.,

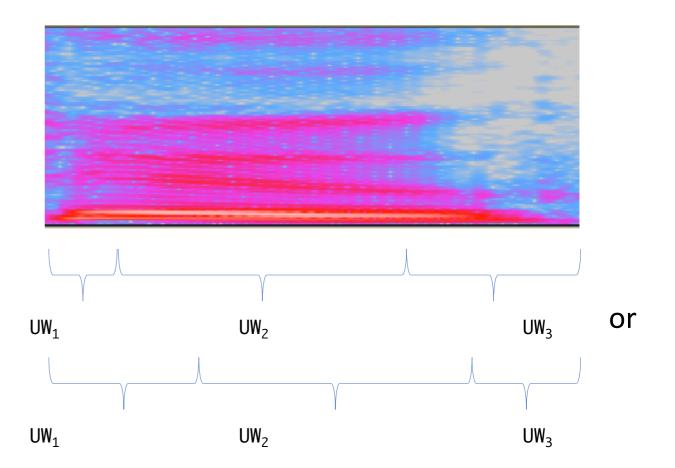
G OW T UW or G OW Z T UW with some scores

Acoustic model p(O|L)



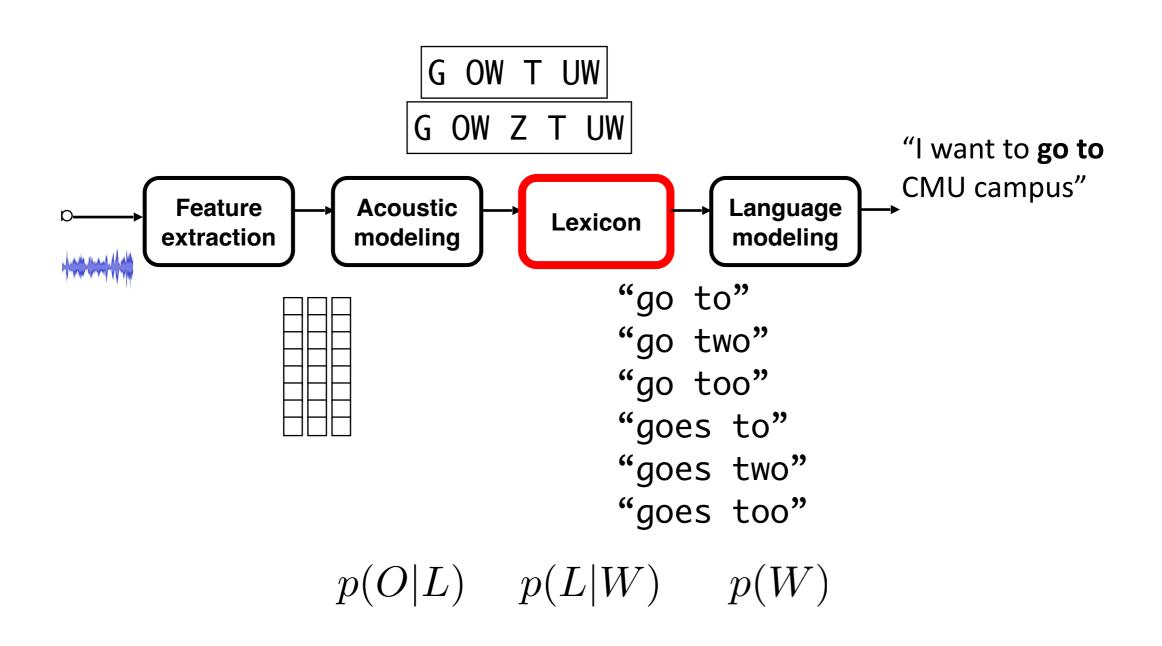
- O and L are different lengths
- Align speech features and phoneme sequences by using HMM

- Provide p(O|L) based on this alignment and model
- The most important problem in speech recognition



Speech recognition pipeline





Phoneme to word





- Performed by lexicon module
 - American English: CMU dictionary
- Time scale
 - 100 milliseconds (depending on a phoneme) to 1 second (depending on a word and also language)
- Type of values
 - 50 categorical value (~6bit) to 100K categorical value (~2Byte)
- We need a pronunciation dictionary
- It can be multiple word sequences (one to many)

Lexicon p(L | W)

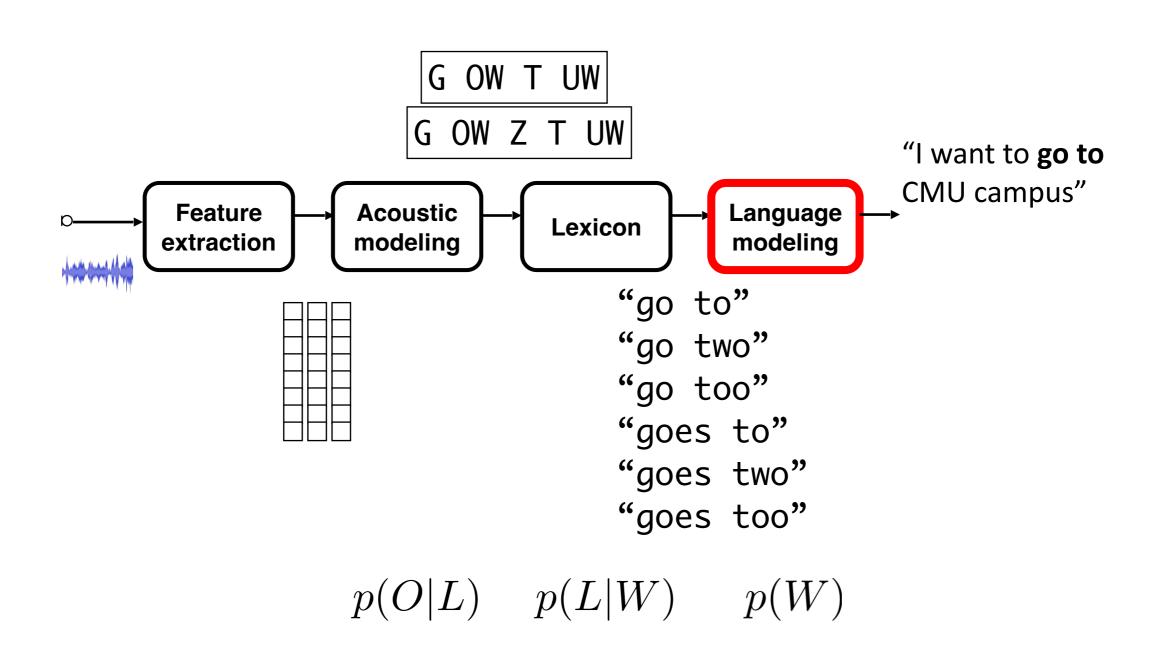


- Basically use a pronunciation dictionary, and map a word to the corresponding phoneme sequence
 - with the probability = 1.0 when single pronunciation
 - with the probability = 1.0/J when multiple (J) pronunciations

$$p(L|W) = p(/T/, /OW/|"two") = 1.0$$

Speech recognition pipeline







```
"go to"

"go two"

"go too"

Language modeling
```

- Performed by language modeling module p(W)
 - N-gram
 - Neural language model (recurrent neural network or transformer)
- From training data, we can basically find how possibly "to", "two", and "too" will be appeared after "go"
 - Part of WSJ training data, 37,416 utterances
 - "go to": **51** times
 - "go two":
 - "go too":

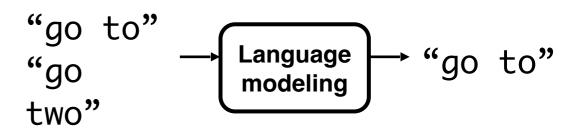
THE WALL STREET JOURNAL.



- Performed by language modeling module p(W)
 - N-gram
 - Neural language model (recurrent neural network or transformer)
- From training data, we can basically find how possibly "to", "two", and "too" will be appeared after "go"
 - Part of WSJ training data, 37,416 utterances
 - "go to": **51** times
 - "go two": **0** times
 - "go too": **0** times



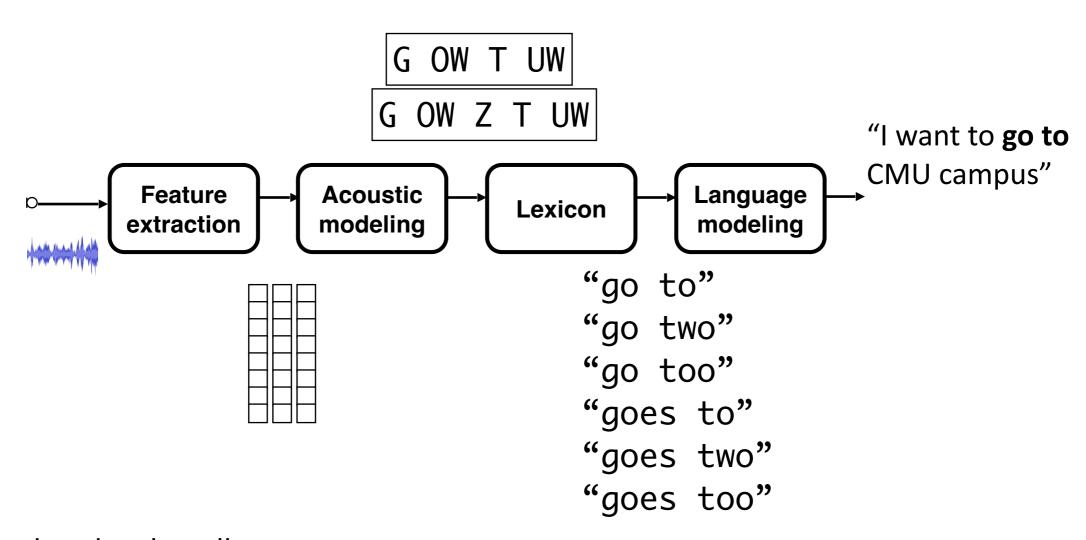
- Performed by language modeling module p(W)
 - N-gram
 - Neural language model (recurrent neural network or transformer)
- From training data, we can basically find how possibly "to", "two", and "too" will be appeared after "go"
 - WSJ all text data, 6,375,622 sentences
 - "go to": **2710** times
 - "go two":
 - "go too":



- Performed by language modeling mod "go too"
 - N-gram
 - Neural language model (recurrent neural network or transformer)
- From training data, we can basically find how possibly "to", "two", and "too" will be appeared after "go"
 - WSJ all text data, 6,375,622 sentences
 - "go to": **2710** times
 - "go two": 2 times, e.g., "those serving shore plants often go two hundred miles or more"
 - "go too":

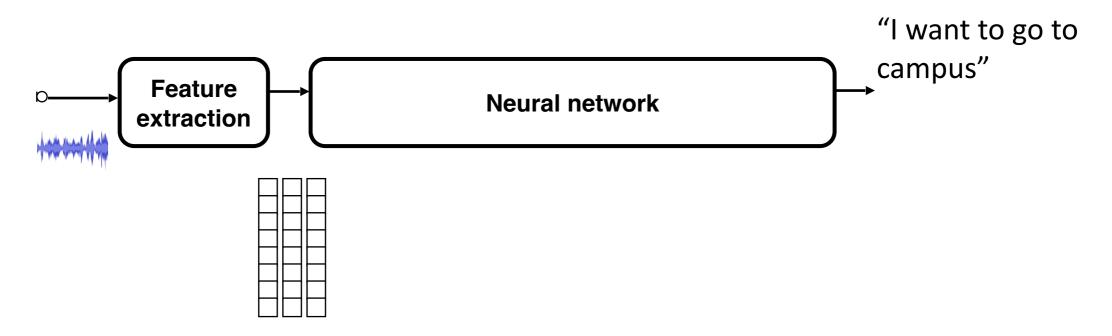
Building speech recognition was really difficult...





- We need to develop all components
- Each component requires a lot of background knowledge
- We need to tune hyper-parameters in each module

Next: End-to-end speech recognition



- We can simply the complicated models
- Optimize all components by using back propagation
- We still need some formulations to make a problem tractable