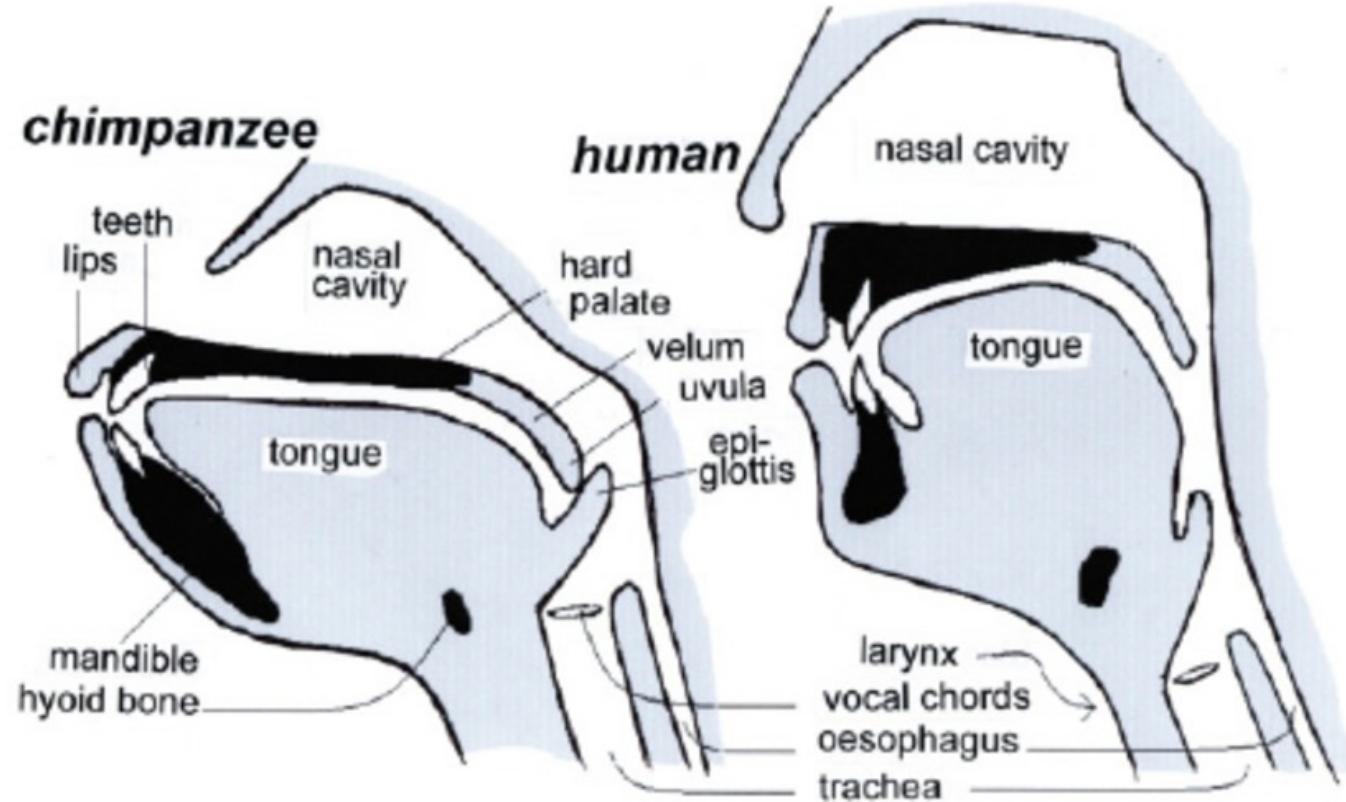


# Speech Production and Perception

ECE 448, Spring  
2023

Mark Hasegawa-Johnson,  
5/2023



# Speech

(Slide: Scharenborg, 2017)

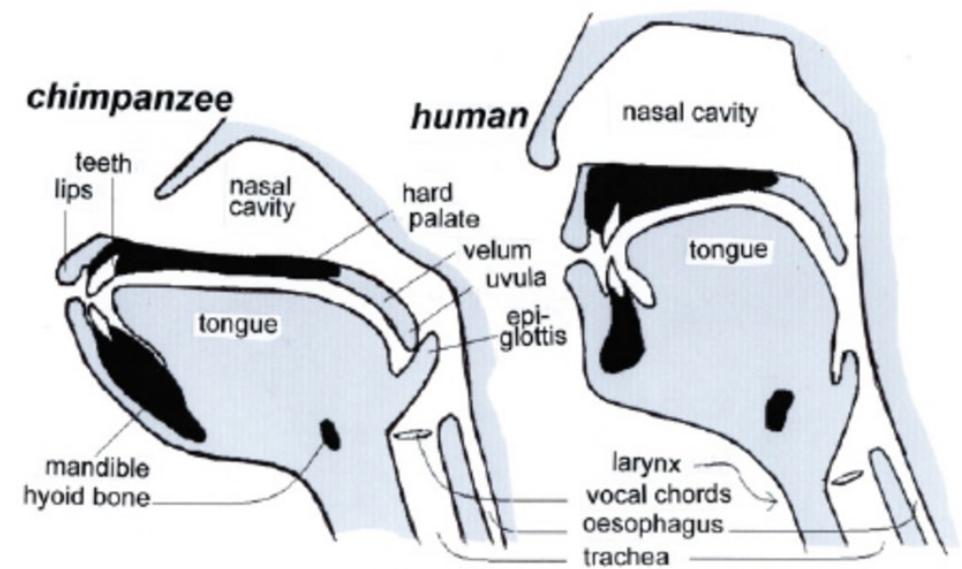
- Specific to humans
- Allows us to convey information very fast
- Central role in many other language-related processes
- One of the most complex skills humans perform:
  - <https://www.youtube.com/watch?v=DcNMCB-Gsn8>
  - <https://www.youtube.com/watch?v=KtN-FCOeWjl>

# Evolution of the vocal tract

(Slide: Scharenborg, 2017)

- Lowering of the tongue into the pharynx → lowering of the larynx
- Lengthening of the neck
- At the cost of an increase in the risk of choking on food

- Neanderthals were not capable of human speech
- Modern human vocal tract: since 50,000 years

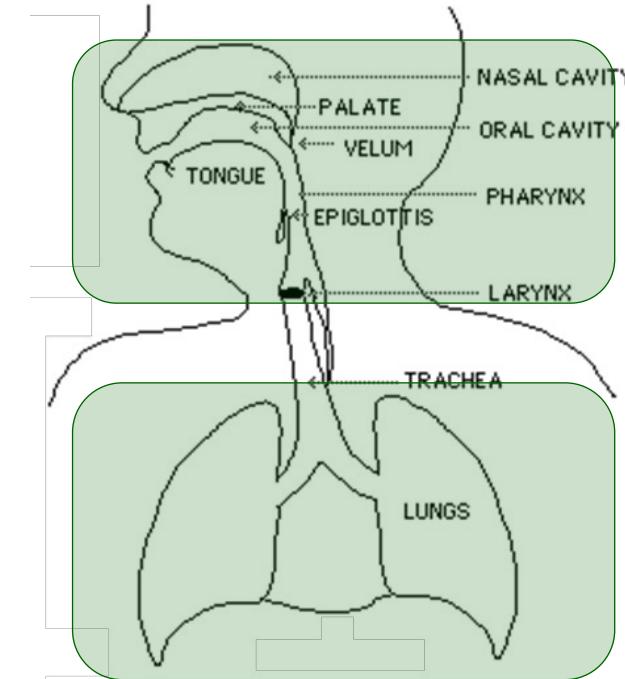


# The anatomy and physiology of speech

(slide: Scharenborg, 2017)

## Vocal tract

- Area between vocal cords and lips
- Pharynx + nasal cavity  
+ oral cavity



and lungs

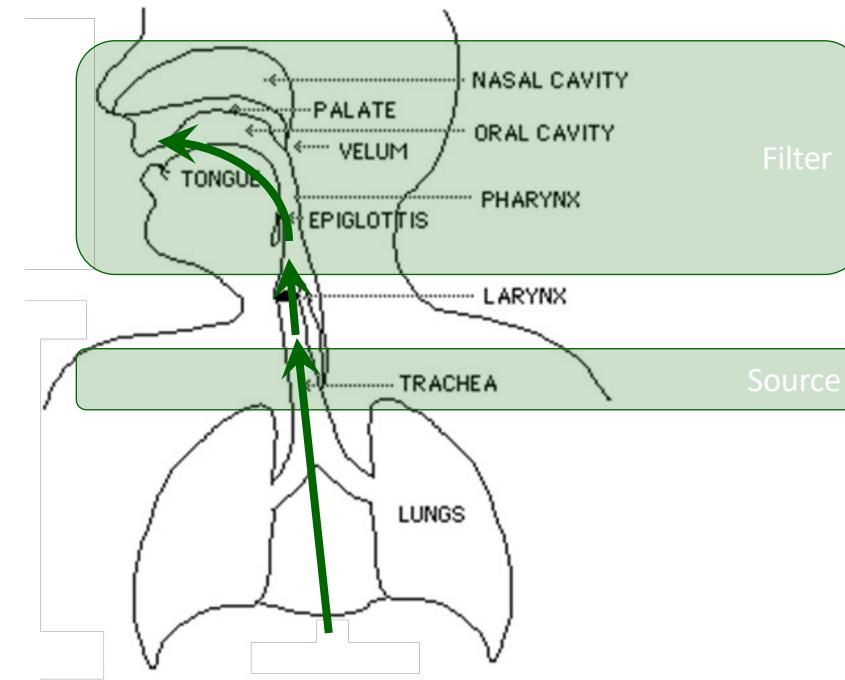
# 3 steps to produce sounds

(Slide: Scharenborg, 2017)

step 3: *articulation* =  
distortion of air  
→ time-varying formant-frequency  
pattern  
= speech

step 2: *phonation*

step 1: *initiation*



# The Source-Filter Model of Speech Production

(Chiba & Kajiyama, 1940)

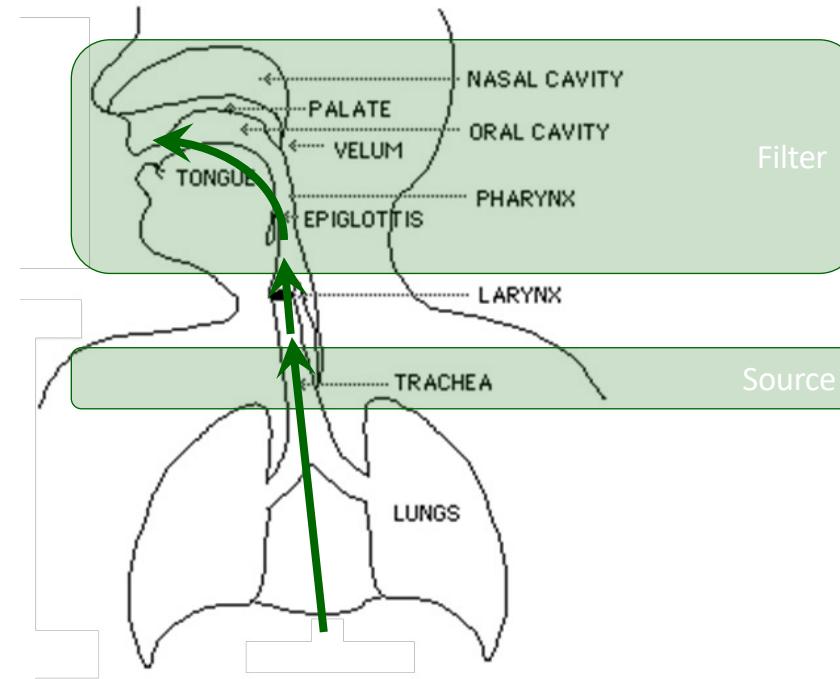
- Sources: there are only three, all of them have wideband spectrum
  - Voicing: vibration of the vocal folds, same type of aerodynamic mechanism as a flag flapping in the wind.
  - Frication or Aspiration: turbulence created when air passes through a narrow aperture
  - Burst: the “pop” that occurs when high air pressure is suddenly released
- Filter:
  - Vocal tract = the air cavity between glottis and lips
  - Just like a flute or a shower stall, it has resonances
  - The excitation has energy at all frequencies; excitation at the resonant frequencies is enhanced

# 3 steps to produce sounds

step 3: *articulation* =  
distortion of air  
→ time-varying formant-frequency  
pattern  
= speech

step 2: *phonation*

step 1: *initiation*



# The Source-Filter Model of Speech Production

A picture from Martin Rothenberg's website

The screenshot shows a web browser window with the title bar "Source-Filter-Lives-pap" and the URL "rothenberg.org/source-filter-lives/Source-Filter-Lives-paper-as-presented5.pdf". The page content is as follows:

**THE SOURCE-FILTER MODEL FOR VOICE PRODUCTION**

GLOTTAL AIRFLOW → **VOCAL TRACT** → RADIATED ACOUSTIC PRESSURE

Possible Assumptions

**LINEARITY** — The vocal tract is a linear acoustic system.

**INDEPENDENCE** — The properties of the glottal voice source and the supraglottal vocal tract are not dependent.

# The Source-Filter Model

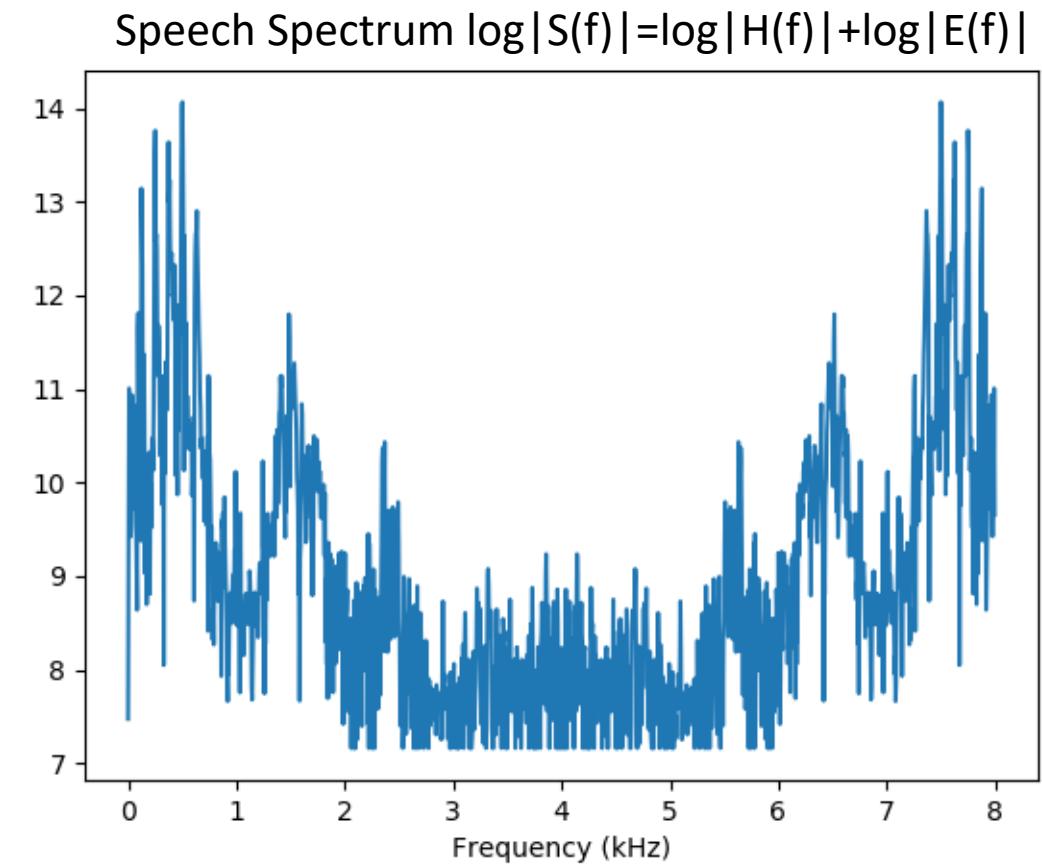
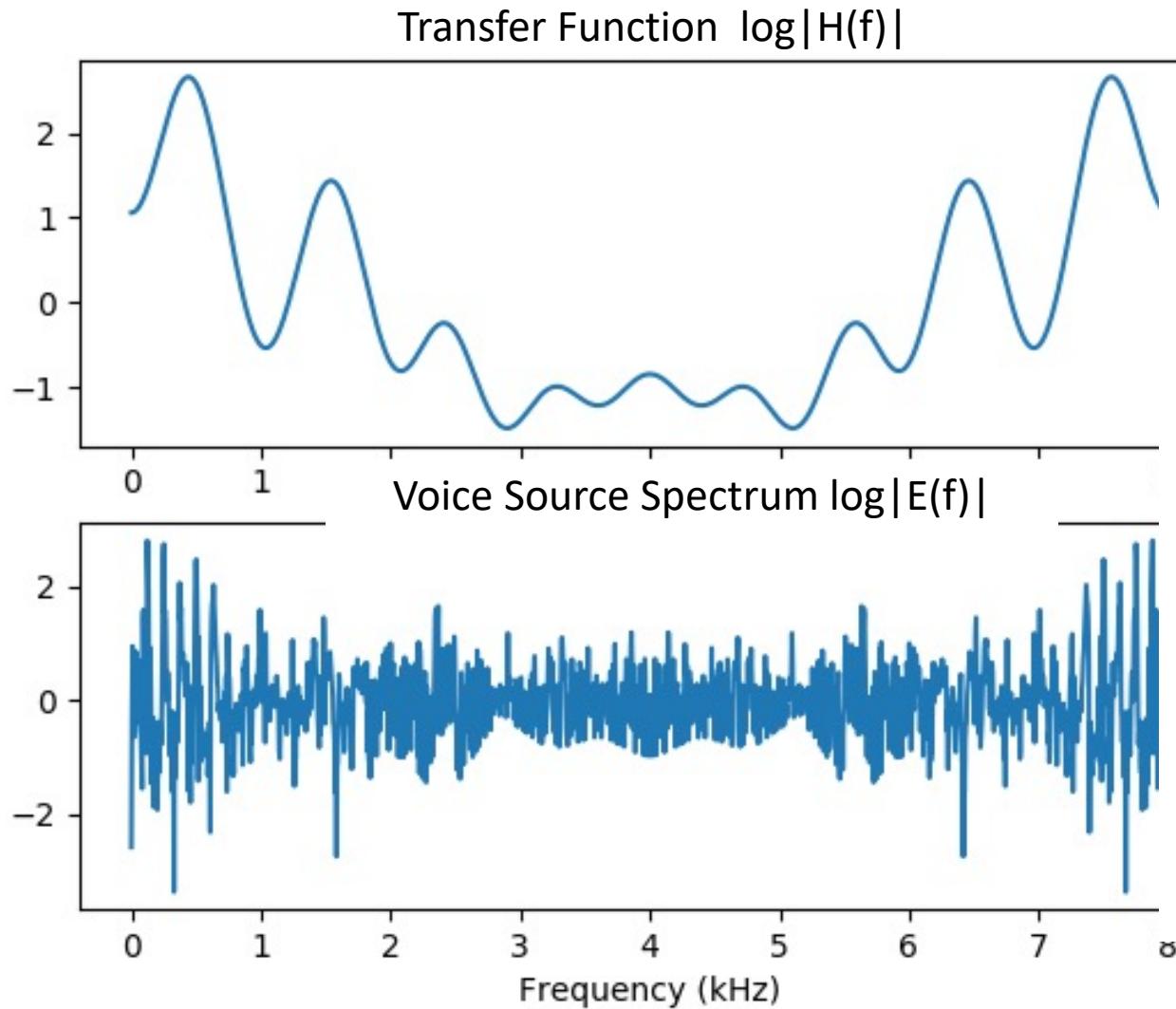
- The speech signal,  $s(t)$ , is created by convolving (\*) an excitation signal  $e(t)$  through a vocal tract transfer function  $h(t)$   
$$s(t) = h(t) * e(t)$$
- The Fourier transform of speech is therefore the product of excitation times transfer function:

$$S(f) = H(f)E(f)$$

...engineers usually compute Fourier transform using  $\Omega = 2\pi f$  rather than  $f$ . You can get one from the other if you remember that  $d\Omega = 2\pi df$ .

- Excitation includes all of the information about voicing, frication, or burst. Transfer function includes all of the information about the vocal tract resonances, which are called “formants.”

# The Source-Filter Model



# Source-Filter Model: Voice Source

- The most important thing about voiced excitation is that it is periodic, with a period called the “pitch period,”  $T_0$
- It’s reasonable to model voiced excitation as a simple sequence of impulses, one impulse every  $T_0$  seconds:

$$e(t) = \sum_{m=-\infty}^{\infty} \delta(t - mT_0)$$

- The Fourier transform of an impulse train is an impulse train (to prove this: use Fourier series):

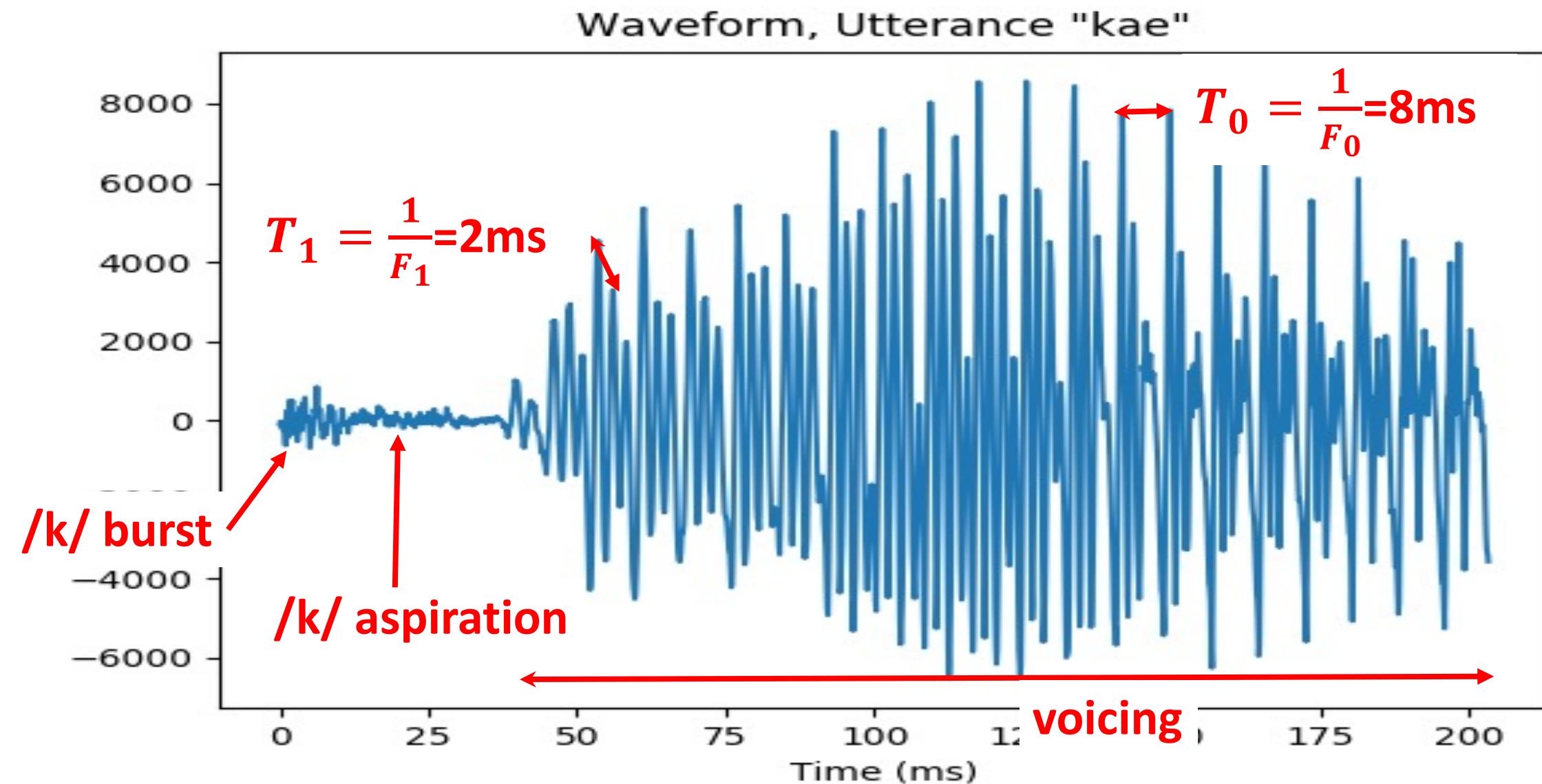
$$E(f) = \frac{1}{T_0} \sum_{k=-\infty}^{\infty} \delta(f - kF_0)$$

...where  $F_0 = \frac{1}{T_0}$  is the pitch frequency. It’s the number of times per second that the vocal folds slap together.

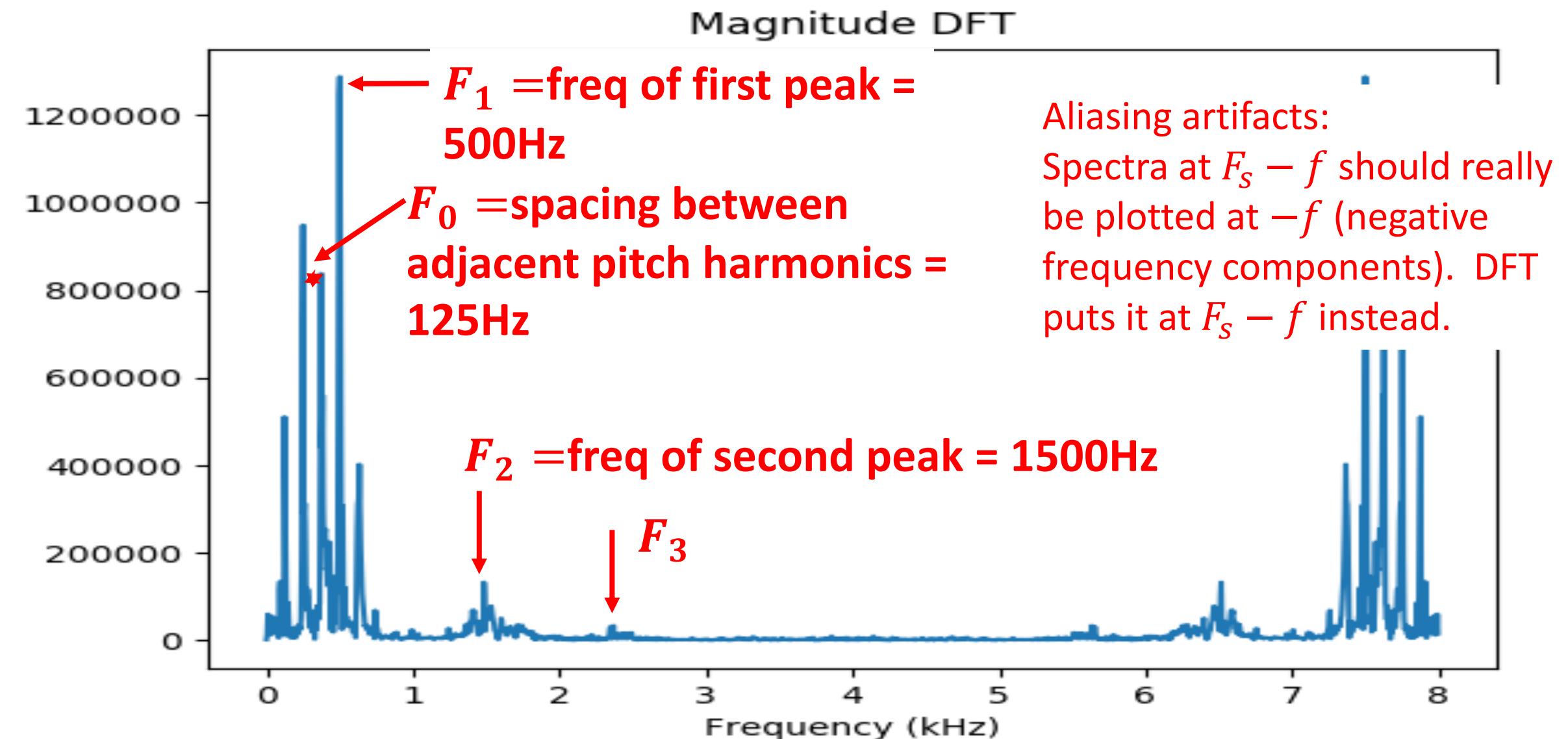
# Source-Filter Model: Filter

- The vocal tract is just a tube. At most frequencies, it just passes the excitation signal with no modification at all ( $H(f) = 1$ ).
- The important exception: the vocal tract has resonances, like a clarinet or a shower stall. These resonances are called “formant frequencies,” numbered in order:  $F_1 < F_2 < F_3 < \dots$ . Typically  $0 < F_1 < 1000 < F_2 < 2000 < F_3 < 3000\text{Hz}$  and so on, but there are some exceptions.
- At the resonant frequencies, the resonance enhances the energy of the excitation, so the transfer function  $H(f)$  is large at those frequencies, and small at other frequencies.

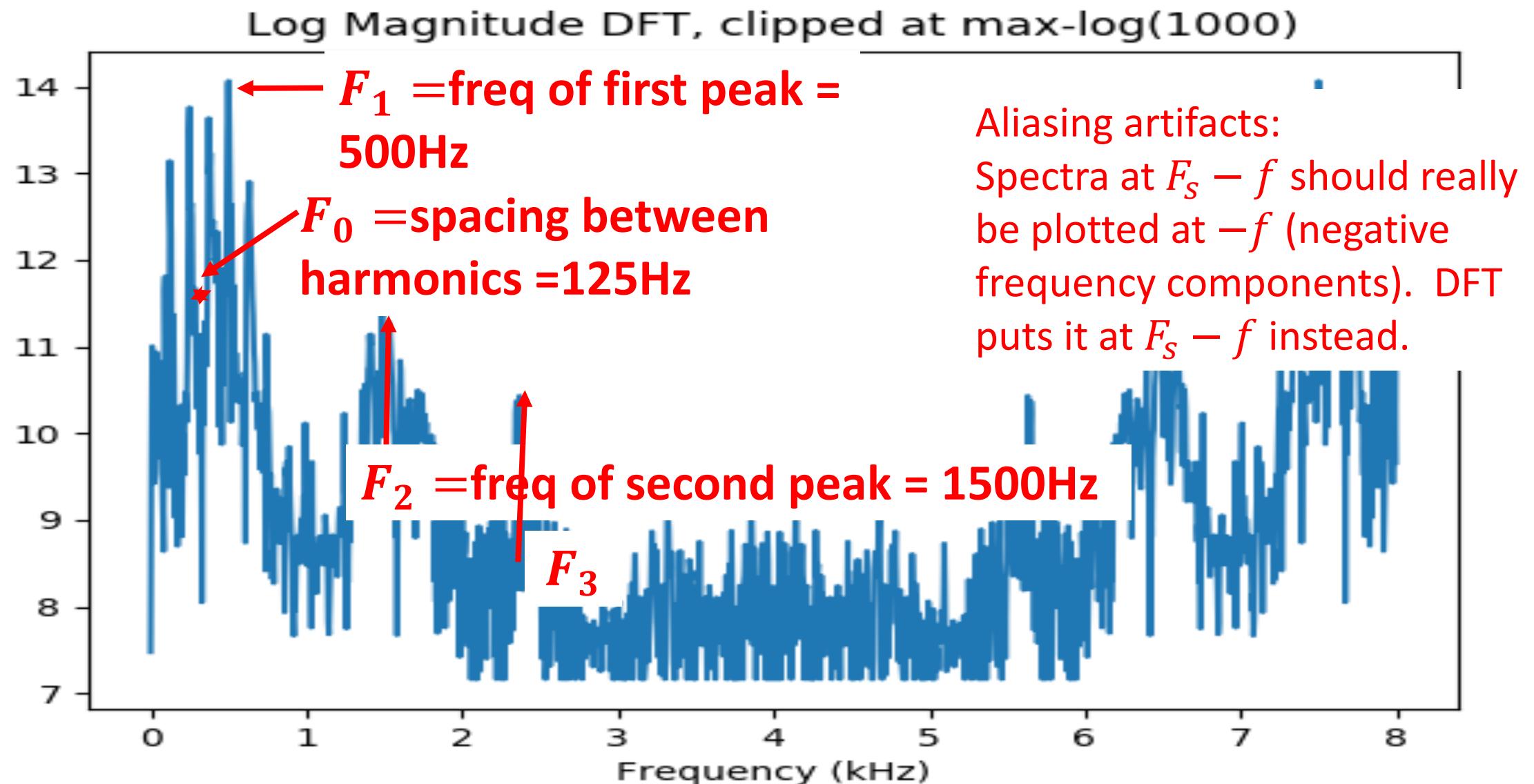
# Speech signal: Time domain



# Speech signal: Magnitude Fourier Transform



# Speech signal: Log Magnitude Transform



# Part 2: Linguistic units

Scharenborg, 2017

- Speech signal

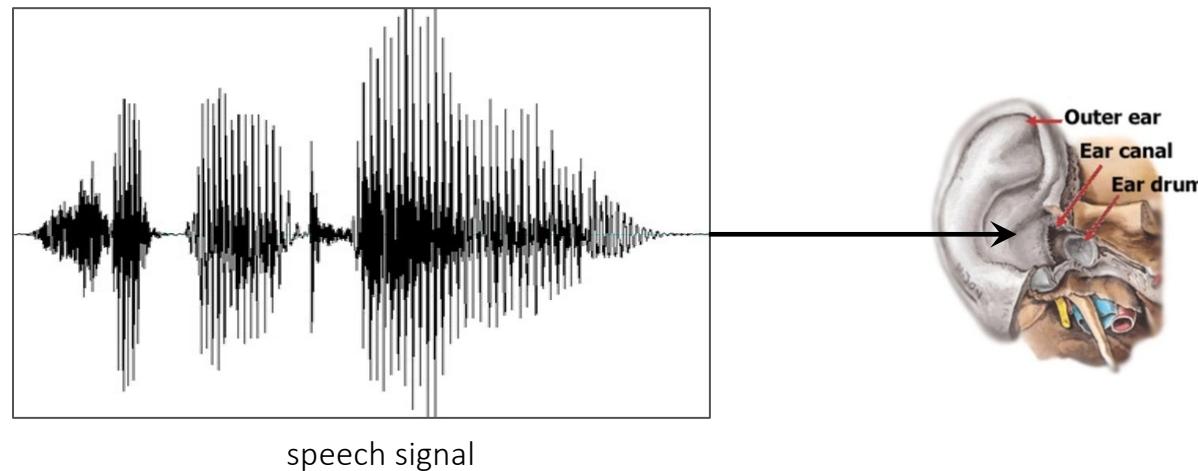
Linguistic units are:

- Phone(me)s
- Words

# Linguistic units

Scharenborg, 2017

- Speech = sound
- Sound = differences in air pressure
- Air pressure waves perceived as different phone(me)s, phone(me) sequences, and (partial or multi) words
- Via eardrum, cochlea, and auditory nerve to brain



# Some terminology

Scharenborg, 2017

- Phoneme: the smallest contrastive linguistic unit that distinguishes meaning, e.g., *tip* vs. *dip*
- Allophone: a variation of a phoneme, eg., *p<sup>h</sup>ot* vs. *spot*
- Phone: a distinct speech sound
- Word: the smallest distinct unit that can be uttered in isolation which has meaning

# Speech sounds

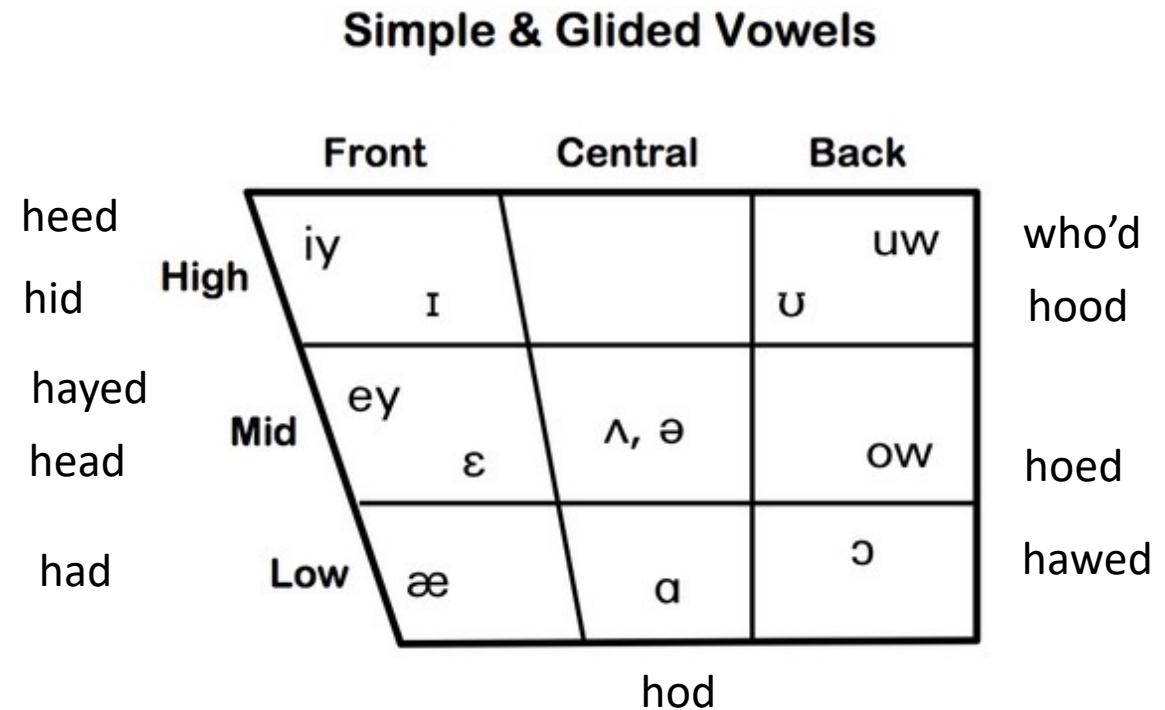
Scharenborg, 2017

- Vowels: unblocked air stream
- Consonants: constricted or blocked air stream

# Different sounds: Vowels

Scharenborg, 2017

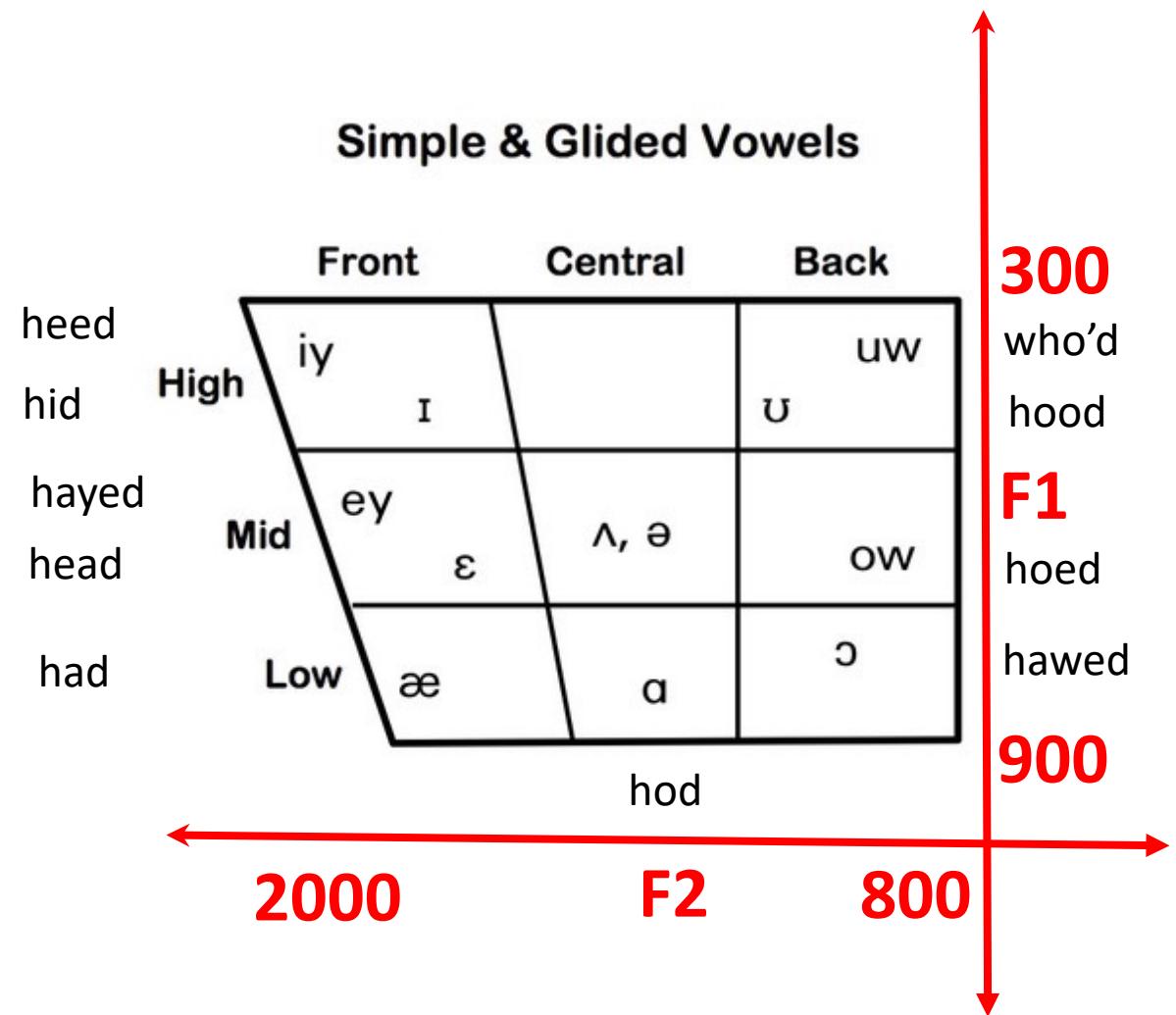
- Tongue height:
  - Low: e.g., /a/
  - Mid: e.g., /e/
  - High: e.g., /i/
- Tongue advancement:
  - Front : e.g., /i/
  - Central : e.g., /ə/
  - Back : e.g., /u/
- Lip rounding:
  - Unrounded: e.g., /ɪ, ɛ, e, ə/
  - Rounded: e.g., /u, o, ɔ/
- Tense/lax:
  - Tense: e.g., /i, e, u, o, ɔ, a/
  - Lax: e.g., /ɪ, ɛ, æ, ə/



# Different sounds: Vowels

Scharenborg, 2017

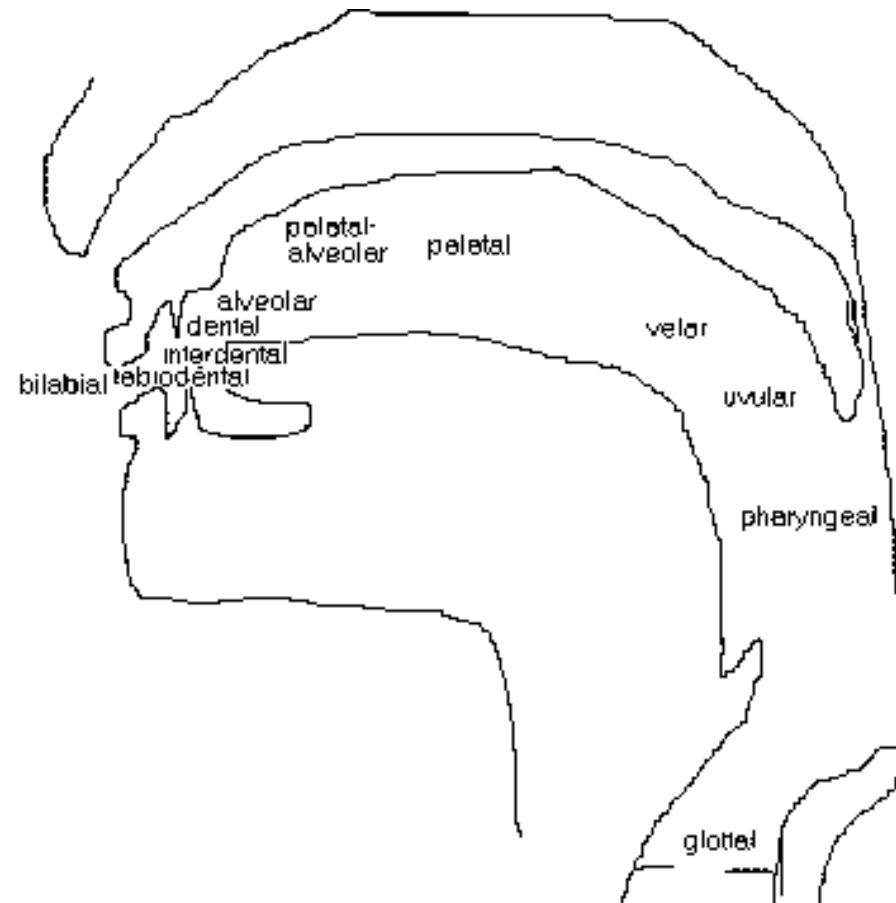
- Tongue height:
  - Low: e.g., /a/
  - Mid: e.g., /e/
  - High: e.g., /i/
- Tongue advancement:
  - Front : e.g., /i/
  - Central : e.g., /ə/
  - Back : e.g., /u/
- Lip rounding:
  - Unrounded: e.g., /ɪ, ɛ, e, ə/
  - Rounded: e.g., /u, o, ɔ/
- Tense/lax:
  - Tense: e.g., /i, e, u, o, ɔ, ɑ/
  - Lax: e.g., /ɪ, ɛ, æ, ə/



# Different sounds: Consonants

Scharenborg, 2017

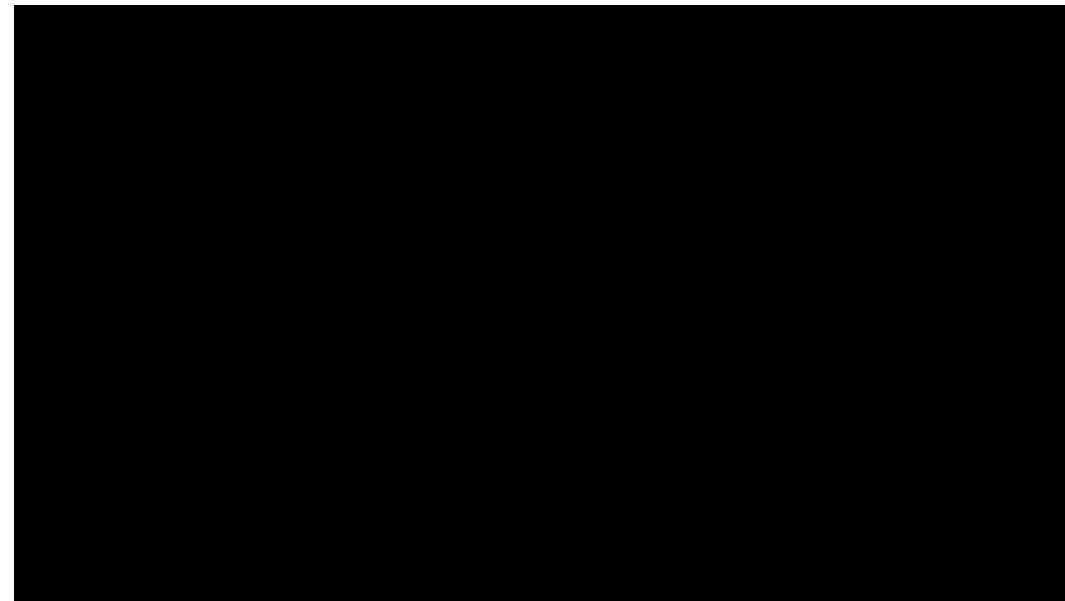
- Place of articulation
  - Where is the constriction/blocking of the air stream?
- Manner of articulation
  - Stops: /p, t, k, b, d, g/
  - Fricatives: /f, s, S, v, z, Z/
  - Affricates: /tS, dZ/
  - Approximants/Liquids: /l, r, w, j/
  - Nasals: /m, n, ng/
- Voicing



# Speech sound production

Scharenborg, 2017

- <https://www.youtube.com/watch?v=DcNMCB-Gsn8>



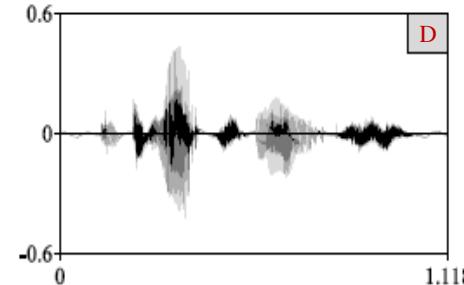
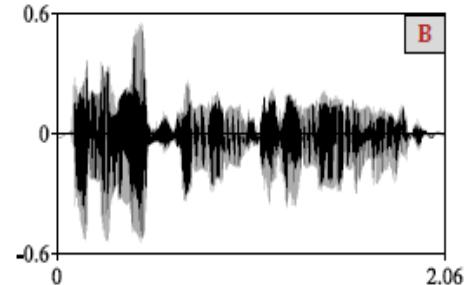
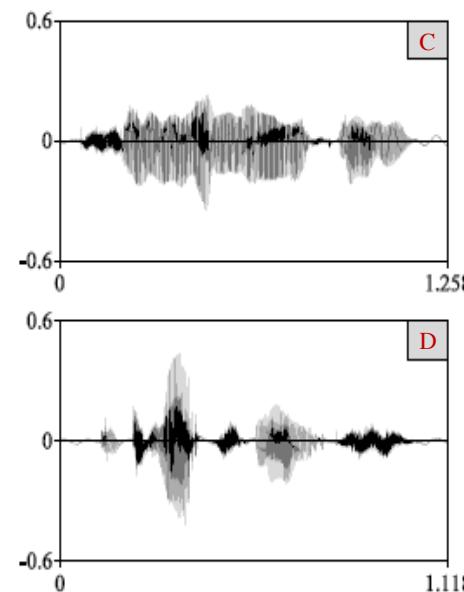
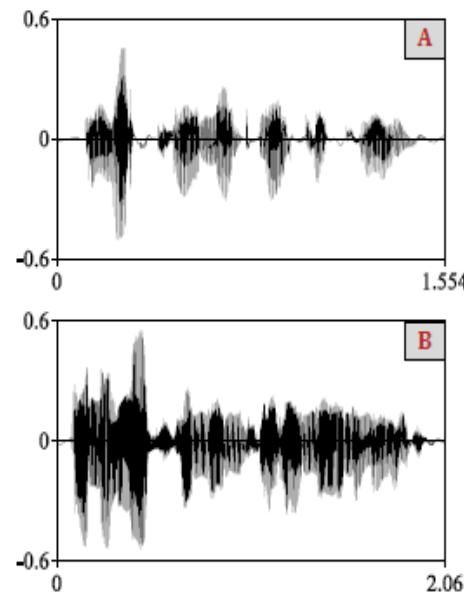
Recorded in 1962, Ken Stevens

Source: YouTube

# Quiz 1: How many words are there?

Scharenborg, 2017

Each picture shows a waveform of a short stretch of speech:

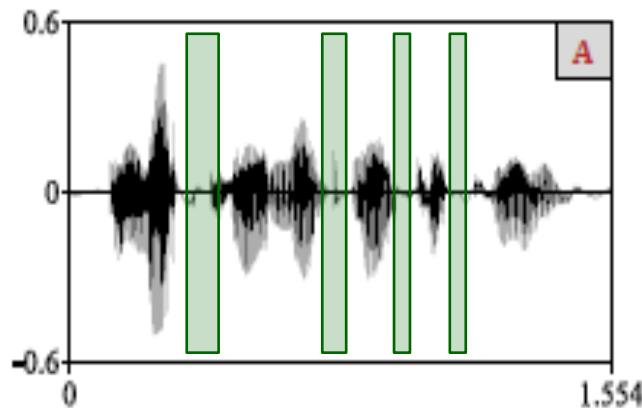


- A: Electromagnetically (1)
- B: Emma loves her mum's yellow marmelade (6)
- C: See you in the evening (5)
- D: Attachment (1)

# Electromagnetically

Scharenborg, 2017

Why is it so hard to determine the number of words?

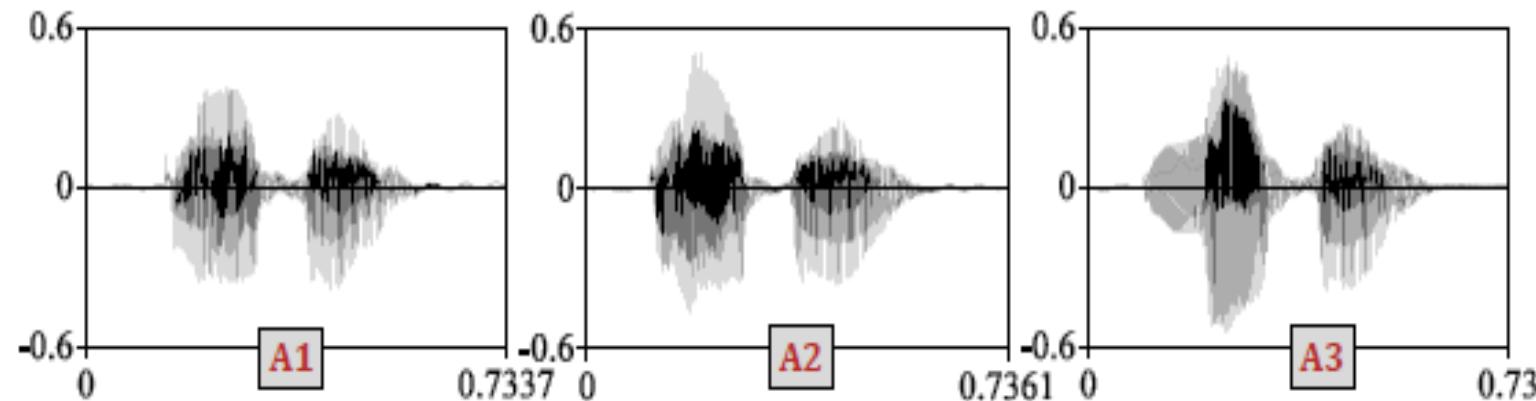


/i l ε kt romæ g nε t ɪ k ə l i/  
silence ≠ word boundary

# Quiz 2: Can you spot the odd one out?

Scharenborg, 2017

- Below are three waveforms each containing a single word:



*Every time you produce a word it sounds differently*

A3 (brother, brother, mother)

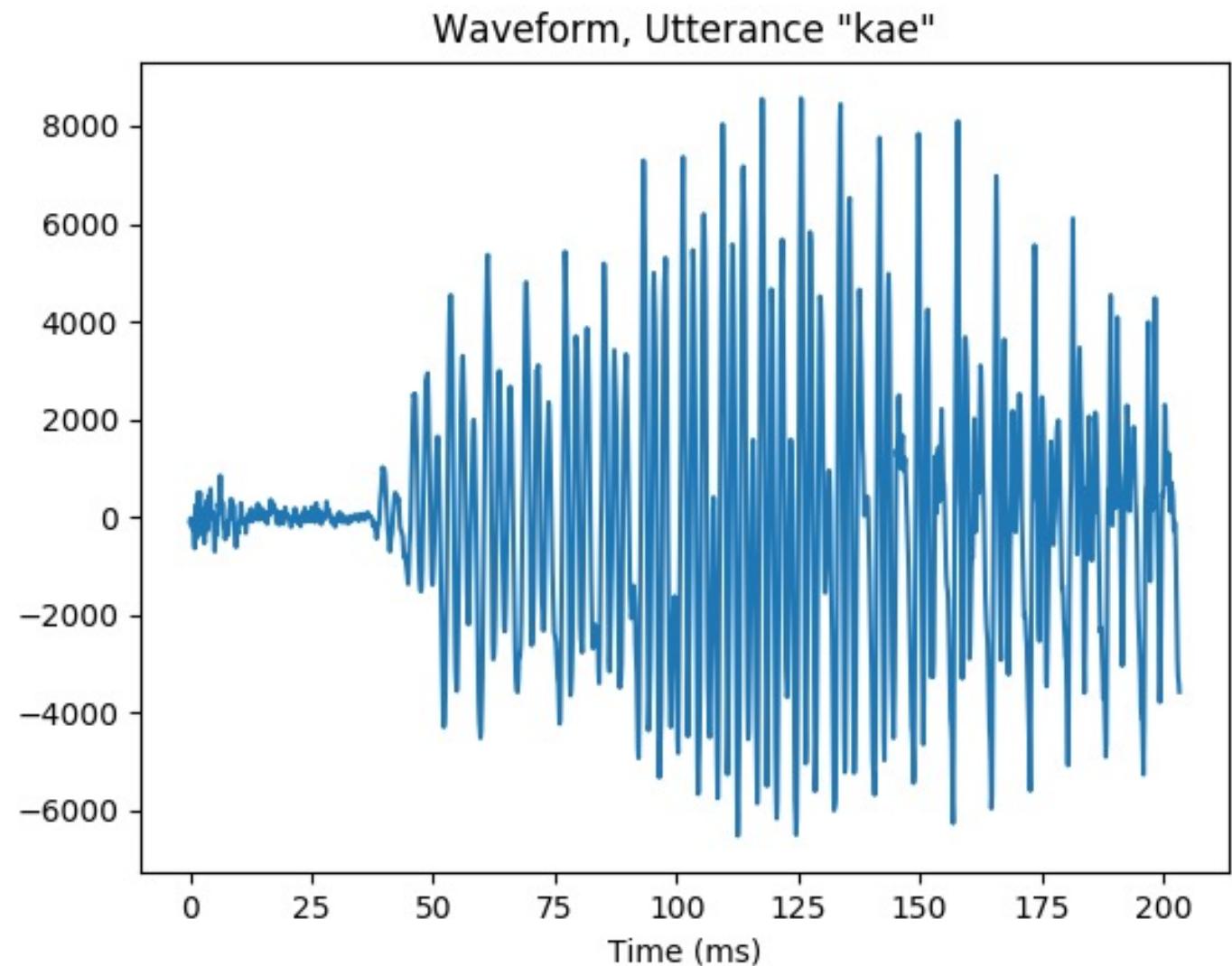
# Enormous variability

Scharenborg, 2017

- Speaker differences, e.g., gender, vocal tract length, age
- Speaker idiosyncracies , e.g., lisp, creaky voice
- Accent: dialects, non-nativeness
- Coarticulation: production of a speech sound becomes more like that of a preceding/following speech sound
- Speaking style → reductions

# Time domain signal: Hard to tell what he was saying

$$s(t) = h(t) * e(t)$$



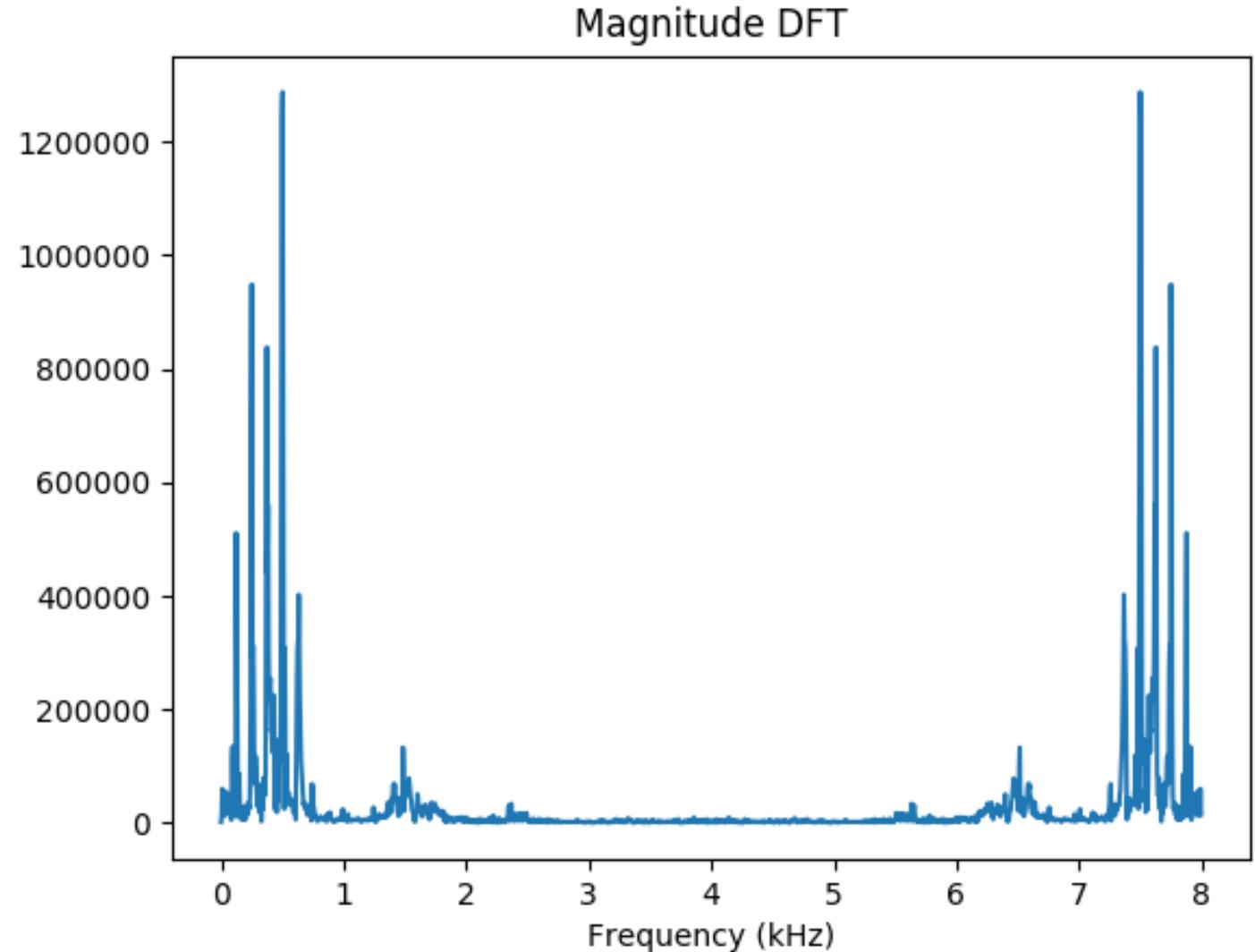
# Magnitude spectrum: A little easier

$$S(f) = H(f)E(f)$$

Easier to measure  
formants → easier to guess  
what he's saying.

Still easy to measure  
F0 → can still guess who he  
is.

(Formants ≈ phone-  
dependent, F0 ≈ person-  
dependent, though there's  
a lot of cross-talk)



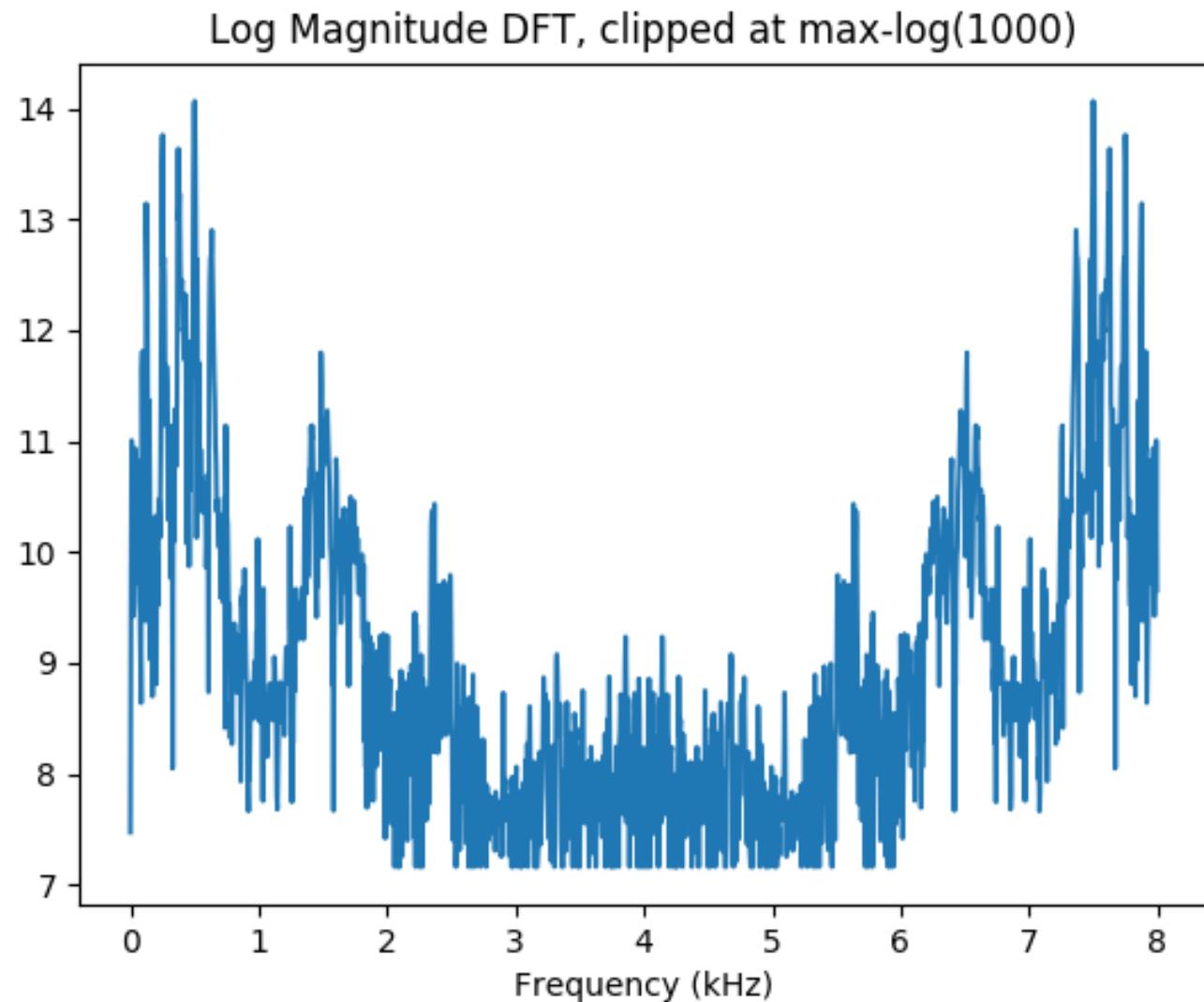
# Log magnitude spectrum: A lot easier

$$\begin{aligned}\ln |S(f)| \\ = \ln |H(f)| + \ln |E(f)|\end{aligned}$$

Easier to measure  
formants → easier to guess what  
he's saying.

Still easy to measure F0 → can  
still guess who he is.

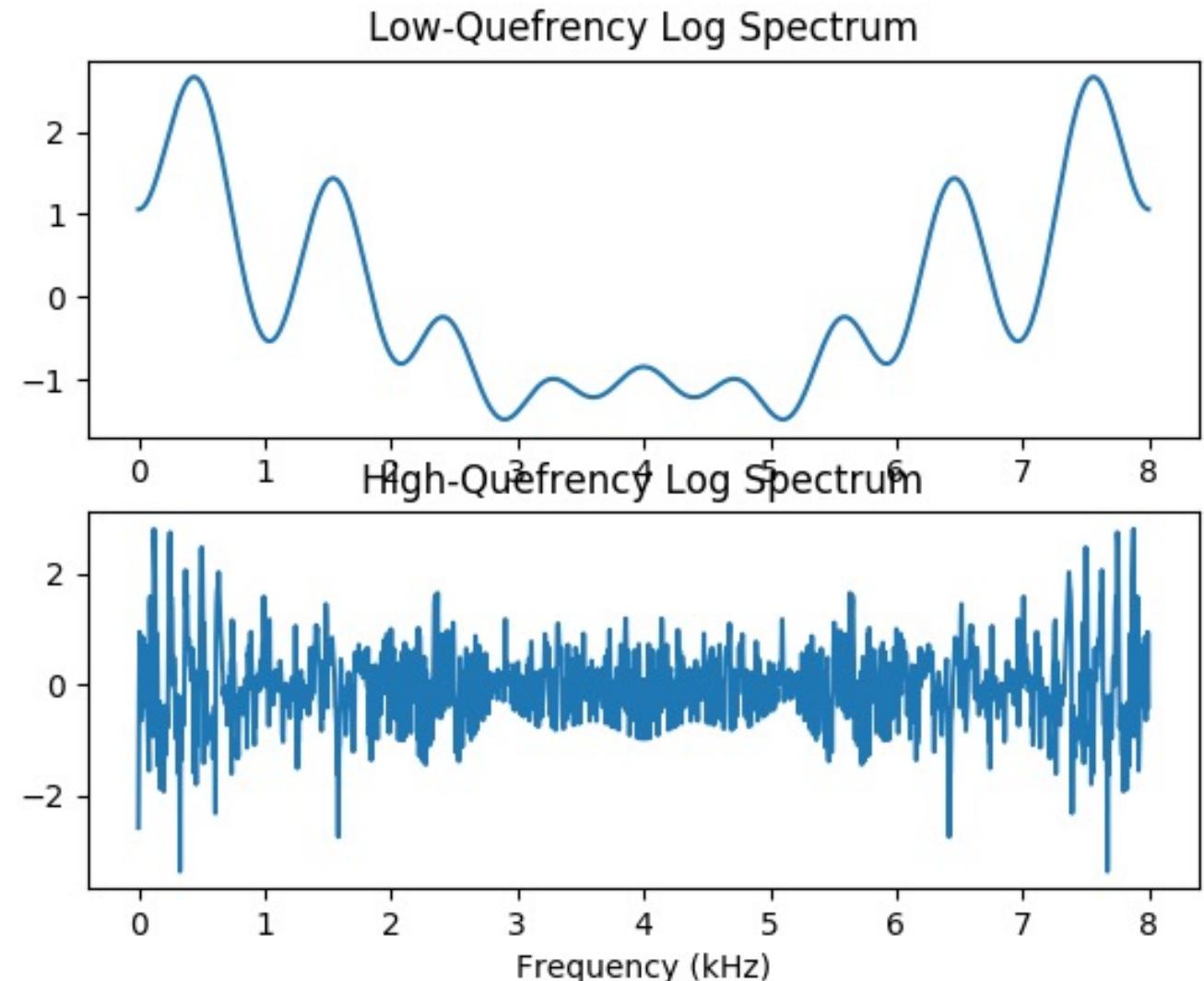
(Formants ≈ phone-dependent,  
F0 ≈ person-dependent, though  
there's a lot of cross-talk)



# Log spectrum = log filter + log excitation

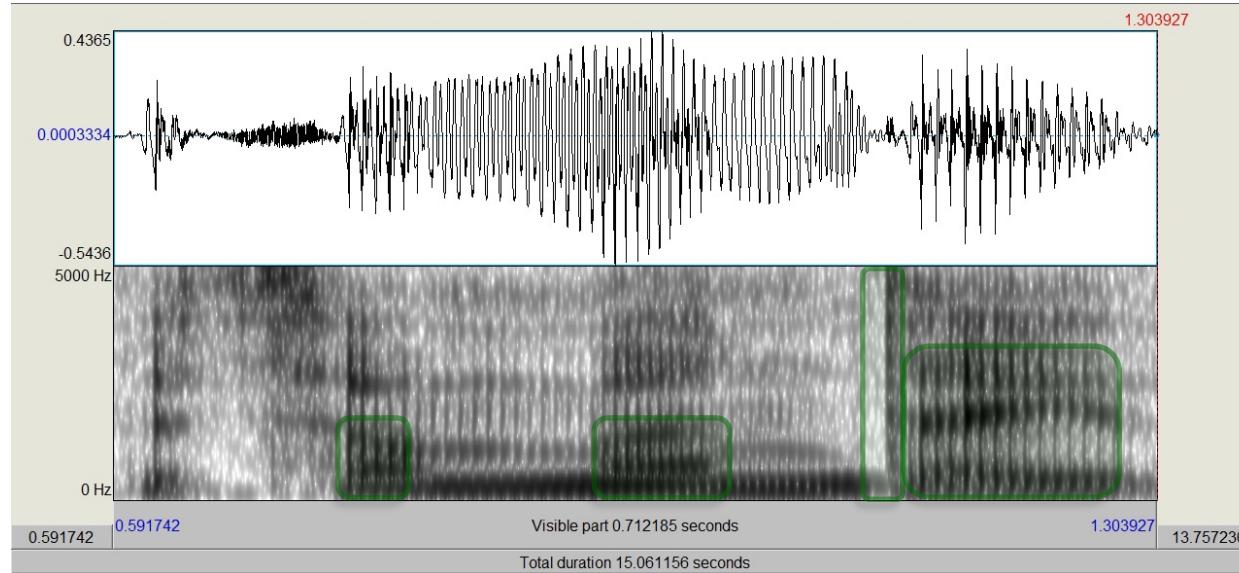
$$\begin{aligned}\ln |S(f)| \\ = \ln |H(f)| + \ln |E(f)|\end{aligned}$$

- But how can we separate the speech spectrum into the transfer function part, and the excitation part?
- Bogert, Healy & Tukey:
  - Excitation is high “quefrency” (varies rapidly as a function of frequency)
  - Transfer function is low “quefrency” (varies slowly as a function of frequency)



# Spectrogram: $\ln(\text{energy}(\text{frequency}, \text{time}))$

Scharenborg, 2017



bu t o nM o n da y

Spectrum lets you measure formants, so it gives some information about vowels.  
Timing is important to know about consonants.

**Spectrogram** = time on the horizontal axis, frequency on vertical axis.

# Summary: Speech Production

- Source-filter model:  $S(f) = H(f)E(f)$ 
  - Voiced excitation is an impulse train in time (with period = the pitch period  $T_0$ ), whose Fourier transform is an impulse train in frequency (with inter-harmonic spacing equal to the pitch frequency  $F_0$ )
  - Transfer function is nearly  $H(f) = 1$  at most frequencies, but with big peaks near the resonant frequencies, which are called formants
- Phones, phonemes, and allophones
- Estimating the transfer function and excitation
  - $\ln |S(f)| = \ln |H(f)| + \ln |E(f)|$
  - The transfer function is low-quefrency, excitation is high-quefrency

# Speech Perception: Outline

- Parseval's Theorem: Cepstral Distance = Spectral Distance
- What spectrum do people hear? The basilar membrane
- Frequency scales for hearing: mel, ERB
- Filterbank coefficients and MFCC

# Parseval's Theorem

L2 norm of a signal equals the L2 norm of its Fourier transform.

# Parseval's Theorem: Examples

- Fourier Series:

$$\frac{1}{T} \int_0^T |x(t)|^2 dt = \sum_{k=-\infty}^{\infty} |X_k|^2$$

- DTFT:

$$\sum_{n=-\infty}^{\infty} |x[n]|^2 = \frac{1}{2\pi} \int_{-\pi}^{\pi} |X(\omega)|^2 d\omega$$

- DFT:

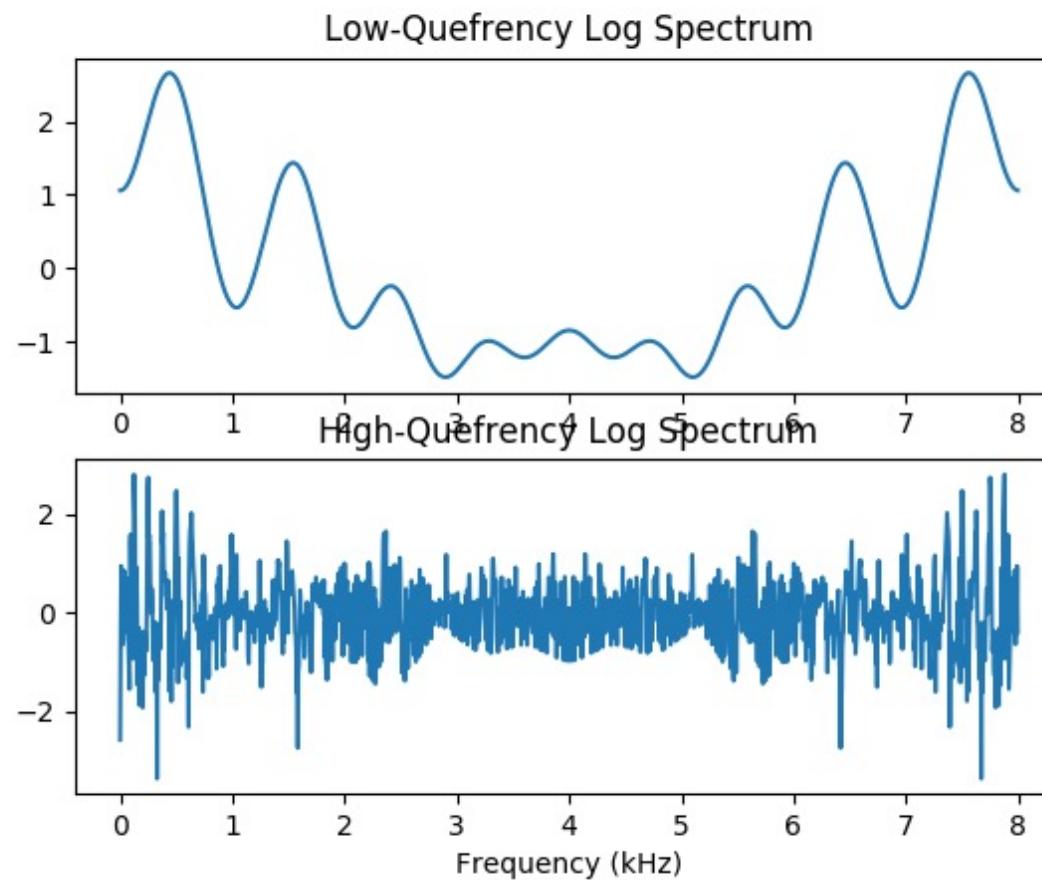
$$\sum_{n=0}^{N-1} |x[n]|^2 = \frac{1}{N} \sum_{k=0}^{N-1} |X[k]|^2$$

# What it means for deep learning

Suppose we have two acoustic signals  $x(t)$  and  $y(t)$ , and we want to find out how different they sound. If they have static spectra, then a good measure of their difference is the L2 difference between their log spectra:

$$\begin{aligned} D &= \sum_{k=0}^M \left( \ln \left| X \left( \frac{(k + 0.5)F_s}{N} \right) \right| - \ln \left| Y \left( \frac{(k + 0.5)F_s}{N} \right) \right| \right)^2 \\ &= \sum_{k=0}^{M-1} (X_k - Y_k)^2 = \sum_{n=0}^{M-1} (x_n - y_n)^2 = \|\vec{x} - \vec{y}\|^2 = \|\vec{X} - \vec{Y}\|^2 \end{aligned}$$

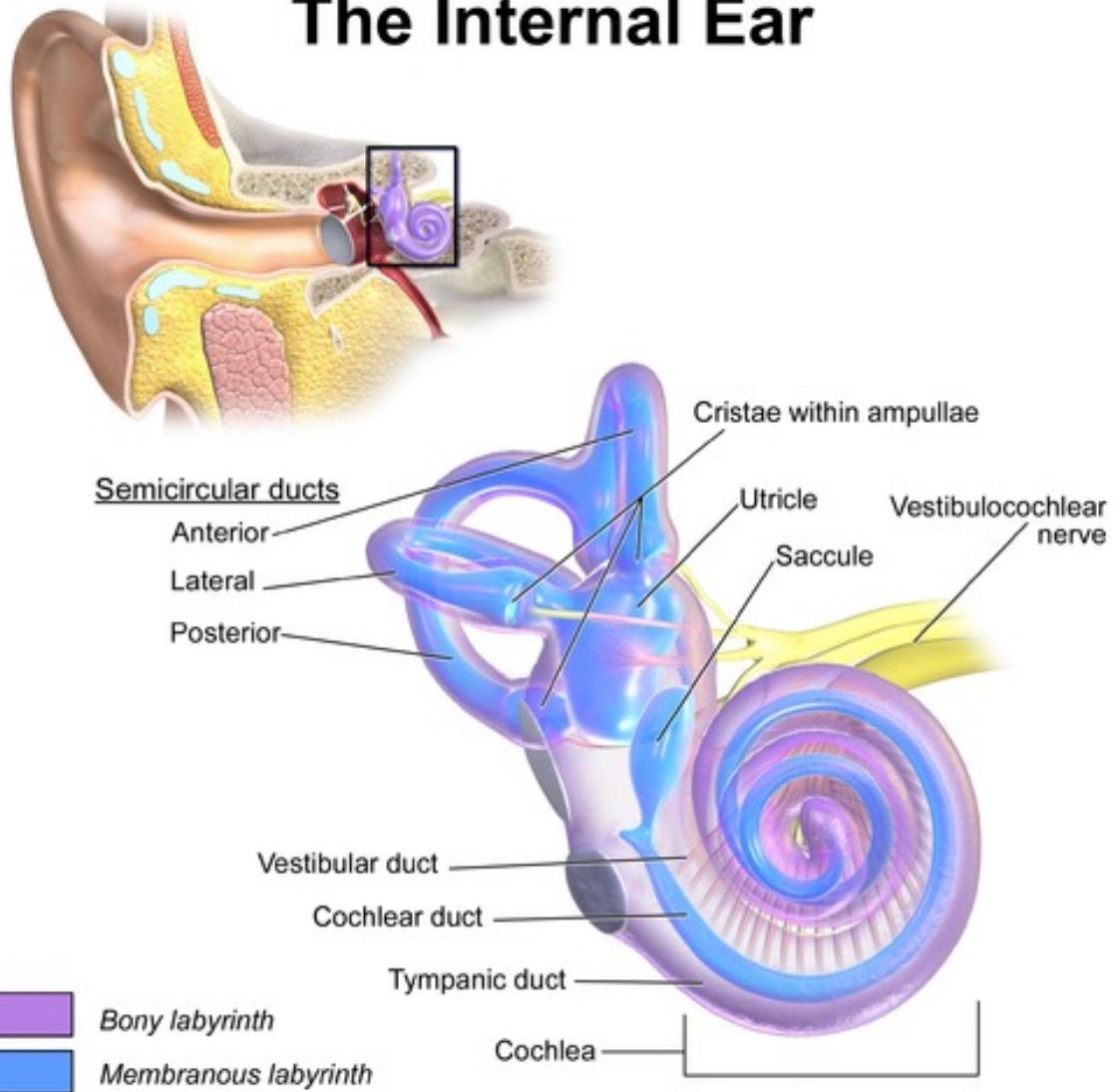
# Low-pass filtering smooths the spectrum



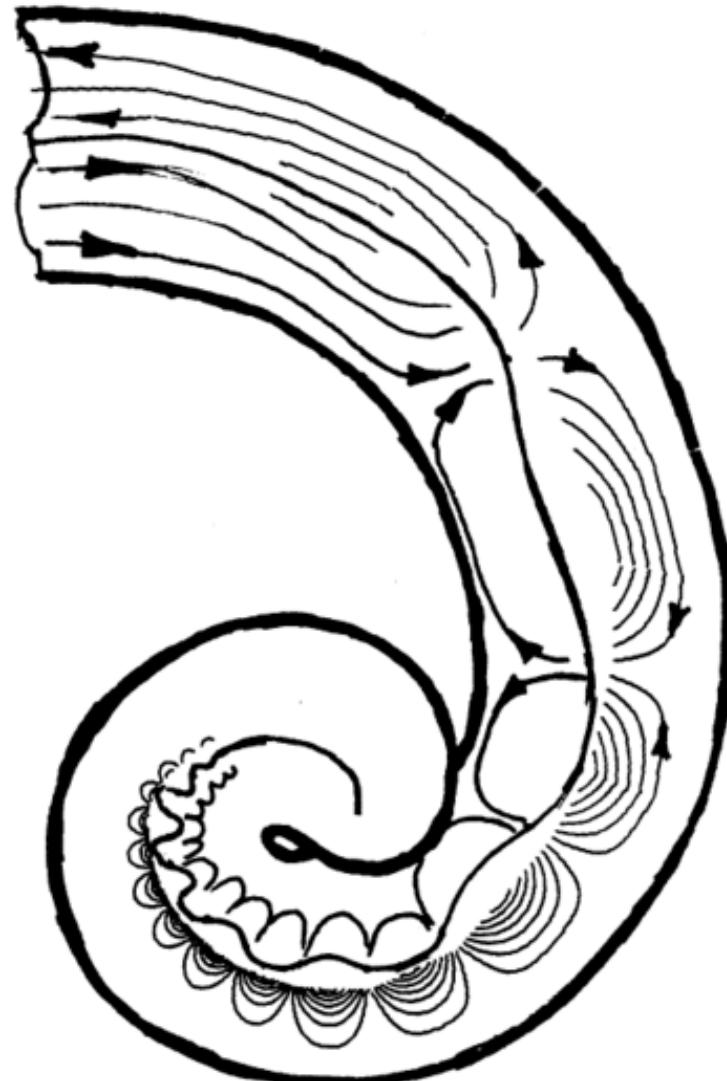
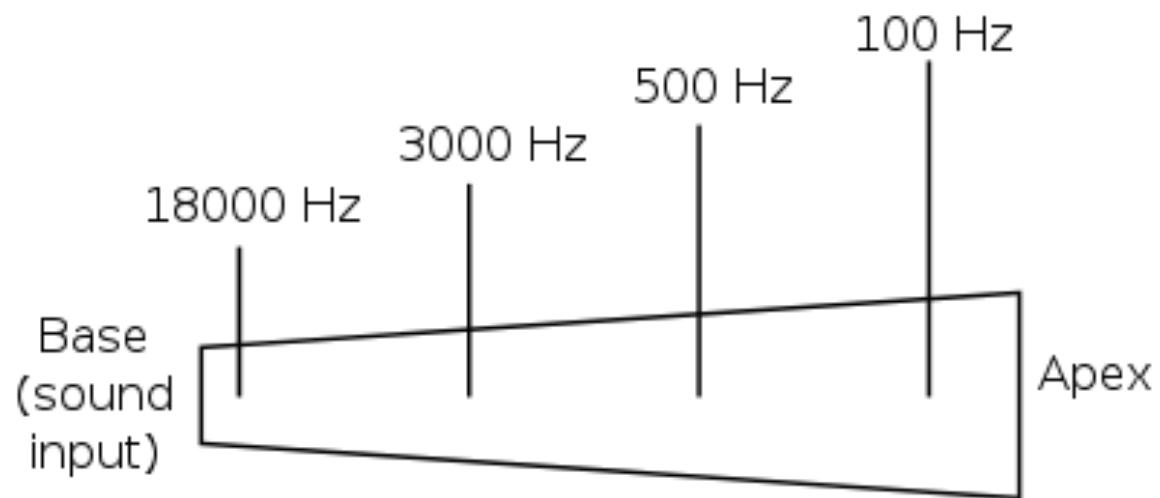
What spectrum do people  
hear? Basilar membrane

# Inner ear

## The Internal Ear



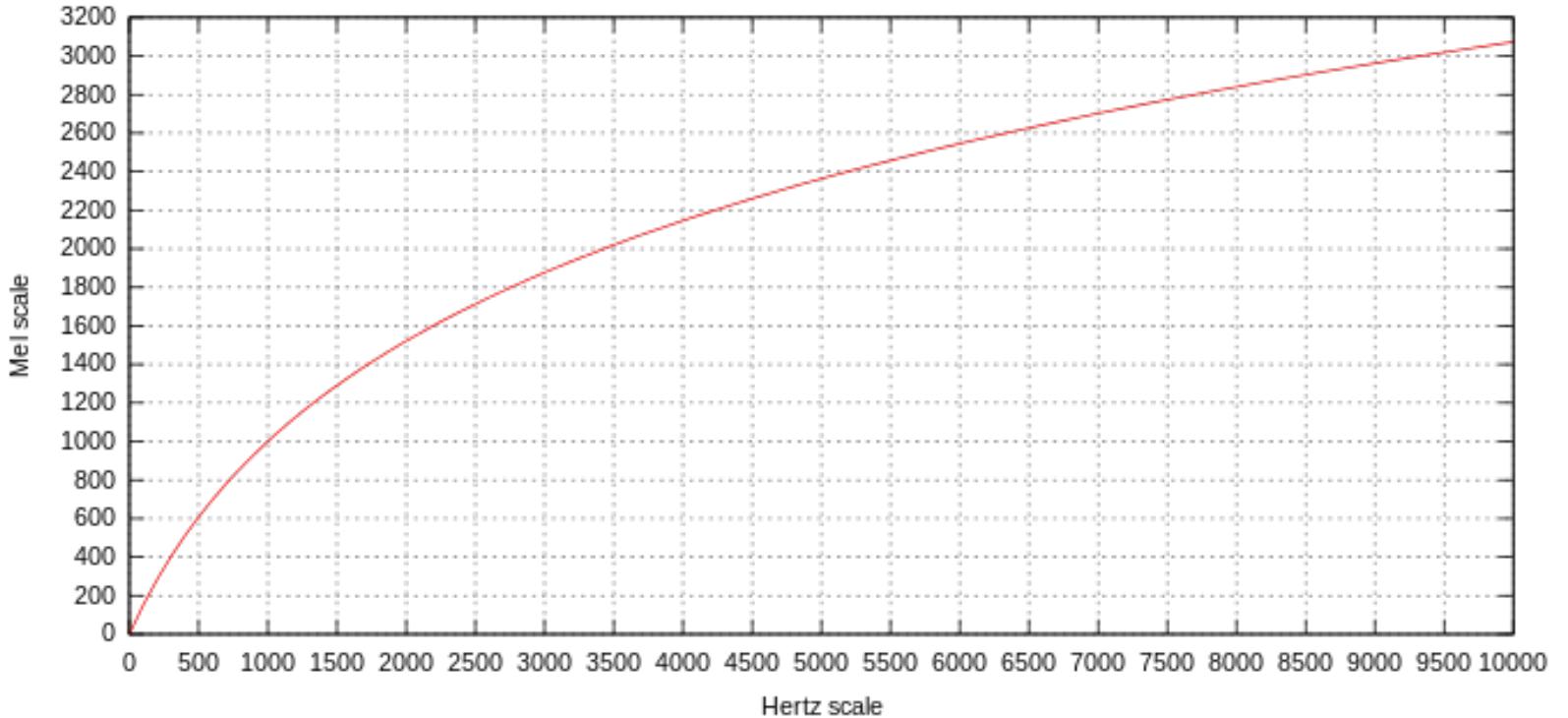
Basilar membrane  
of the cochlea = a  
bank of mechanical  
bandpass filters



Frequency scales for hearing:  
mel scale, ERB scale

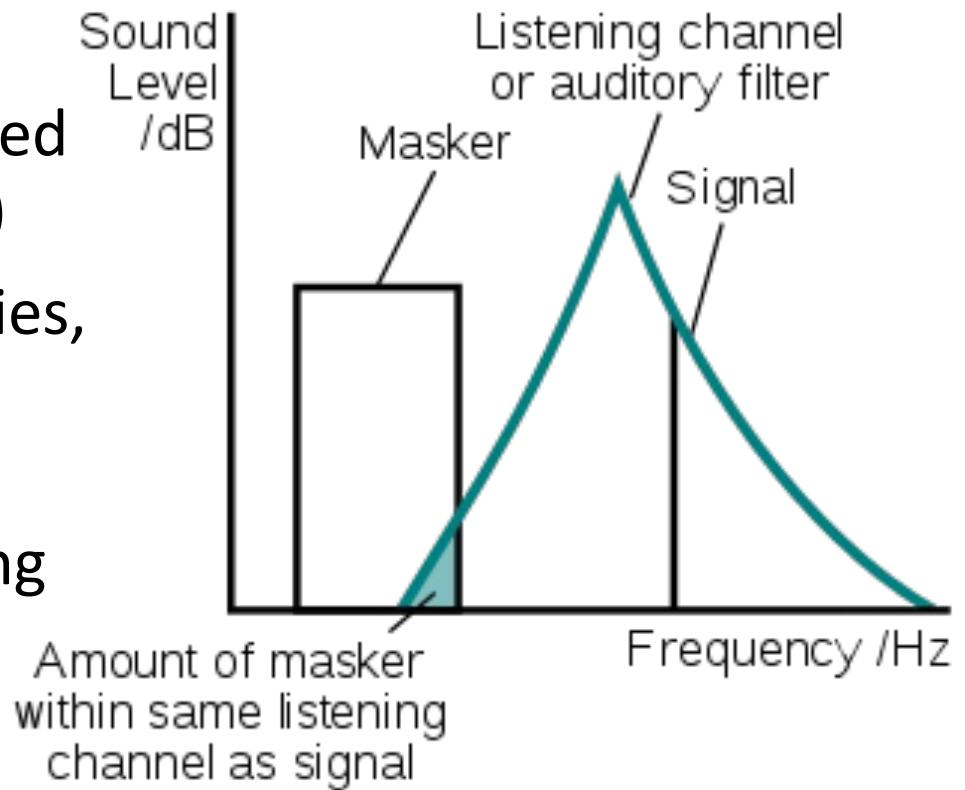
# Mel-scale

- The experiment:
  - Play tones A, B, C
  - Let the user adjust tone D until pitch(D)-pitch(C) sounds the same as pitch(B)-pitch(A)
- Analysis: create a frequency scale  $m(f)$  such that  $m(D)-m(C) = m(B)-m(A)$
- Result:  $m(f) = \frac{1}{2595} \log_{10} \left( 1 + \frac{f}{700} \right)$



# Critical bands

- When two tones play at exactly the same frequency, users can't tell the difference between  $x(t)$  versus  $x(t)+y(t)$  if  $y(t)$  is about 14dB below  $x(t)$  (in other words, the summed power is 1.03 times the power of  $x(t)$  alone)
- When  $x(t)$  and  $y(t)$  are at different frequencies, the masking power of  $x(t)$  is reduced
- Model: assume that the reduced masking power of  $x(t)$  is caused because  $x(t)$  is coming in through the tails of the bandpass filter centered at  $y(t)$ .



# ERB scale

- The experiment: find out the widths,  $B(f)$ , of the critical-band filters centered at every frequency  $f$ .
- Analysis: create a scale  $e(f)$  such that  $e(f+0.5B(f)) - e(f-0.5B(f)) = 1$ , for all frequencies
- Result:  $e(f) = 21.4 \log_{10}(1 + 0.00437f)$

# Mel-frequency filterbank coefficients

Mel filterbank coefficients: convert the spectrum from Hertz-frequency to mel-frequency

- Goal: instead of computing

$$C_k = \ln \left| S \left( \frac{(k+0.5)F_s}{N} \right) \right|$$

We want

$$C_k = \ln |S(f_k)|$$

Where the frequencies  $f_k$  are uniformly spaced on a mel-scale, i.e.,  $m(f_{k+1}) - m(f_k)$  is a constant across all  $k$ .

The problem with that idea: we don't want to just sample the spectrum. We want to summarize everything that's happening within a frequency band.

Mel filterbank coefficients: convert the spectrum from Hertz-frequency to mel-frequency

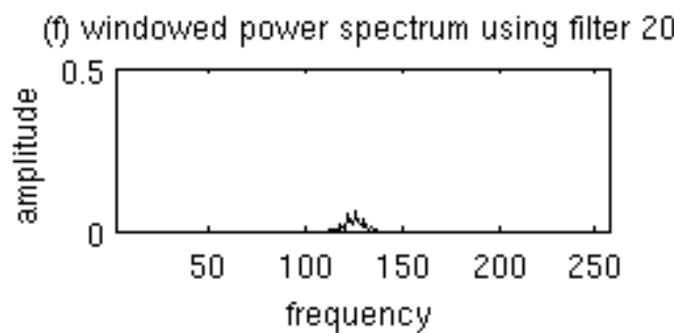
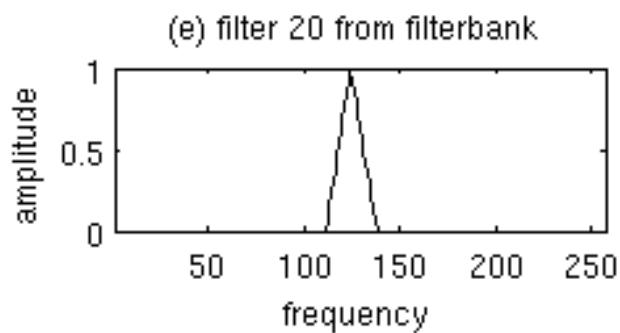
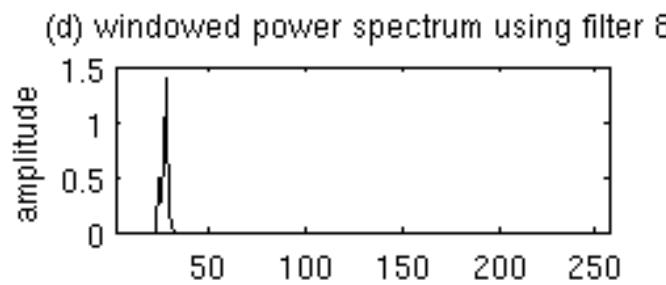
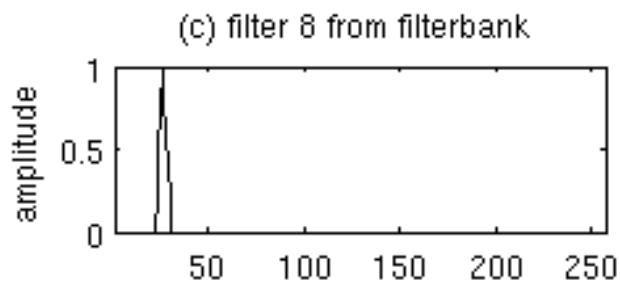
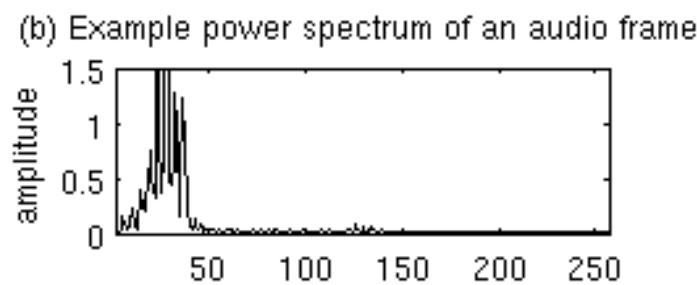
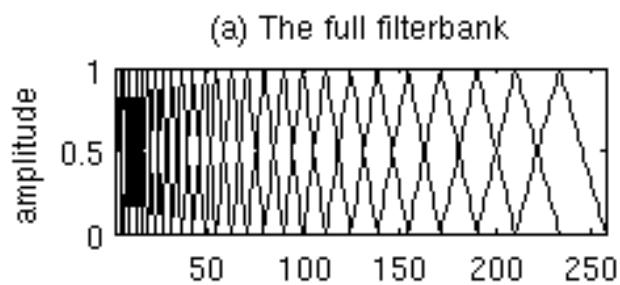
The solution:

$$C_m = \ln \sum_{k=0}^{\frac{N}{2}-1} W_m(k) \left| S\left(\frac{kF_s}{N}\right) \right|$$

Where

$$W_m(k) = \begin{cases} \frac{\frac{kF_s}{N} - f_{m-1}}{f_m - f_{m-1}} & f_m \geq \frac{kF_s}{N} \geq f_{m-1} \\ \frac{f_{m+1} - \frac{kF_s}{N}}{f_{m+1} - f_m} & f_{m+1} \geq \frac{kF_s}{N} \geq f_m \\ 0 & otherwise \end{cases}$$

# Mel filterbank coefficients: convert the spectrum from Hertz-frequency to mel-frequency



# Inclusive Speech Technology

# Outline

- The many dimensions of ability and disability
- Speech technology: the Transformer
- Speech technology: Pre-training and Fine-tuning
- Speech technology: Audiobooks
- Inclusive speech recognition
- The Speech Accessibility Project
- Disability in the public sphere

# Abilities that are financially useful

- Bench-press 225 pounds (weight of a washing machine)
- Rap 200 words/minute
- Differentiate a softmax w.r.t. the Fourier transform of its input
- Debug an operating system
- Notice when somebody is lying or unsure
- Conceptualize deliveries for a 200-component just-in-time assembly

# What is disability?

A disability is any condition of the body or mind (impairment) that makes it more difficult for the person with the condition to do certain activities (activity limitation) and interact with the world around them (participation restrictions). - CDC

## Why is disability a social construct?

The following things are socially constructed:

- How to do an activity
- How to interact with the world



...because sometimes comedy is the best way to describe the world...

<https://www.tinafriml.com>

...also...

<https://www.instagram.com/p/Ck6X6NCgc1c/?hl=en>

# Why aren't devices, buildings, and organizations accessible to everyone?

- Availability bias: Engineers, architects, and entrepreneurs overestimate the frequency of ability profiles that are most familiar to them personally

# Some neurological conditions that affect speech

- Progressive neurological conditions
  - Parkinson's disease (PD) and related cell-death conditions (~1m in USA)
  - Multiple sclerosis (MS) and related demyelination conditions (~800k in USA)
  - Amyotrophic lateral sclerosis (ALS)
  - Muscular dystrophies (MD: ~250k in USA)
- Trauma
  - Stroke/CVA (~4m in USA)
  - Traumatic brain injury/Concussion (~one quarter of people surveyed)
- Lifelong, non-progressive neurological conditions
  - Down Syndrome (~400k in USA)
  - Cerebral Palsy (~764k in USA)

# The engineering perspective

- If a device is usable by 20% of people, that's good.
- If it is usable by 80% of people, that's better.
- If it is usable by 95% of people, that's even better.
- If it is usable by 99% of people, that's even better.
- If everyone who wants to use the device can use it successfully, that's best of all.

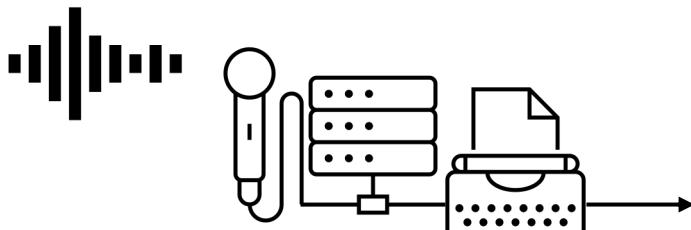
# Outline

- The many dimensions of ability and disability
- Speech technology: the Transformer
- Speech technology: Pre-training and Fine-tuning
- Speech technology: Audiobooks
- Inclusive speech recognition
- The Speech Accessibility Project
- Disability in the public sphere

# Automatic speech recognition (ASR)



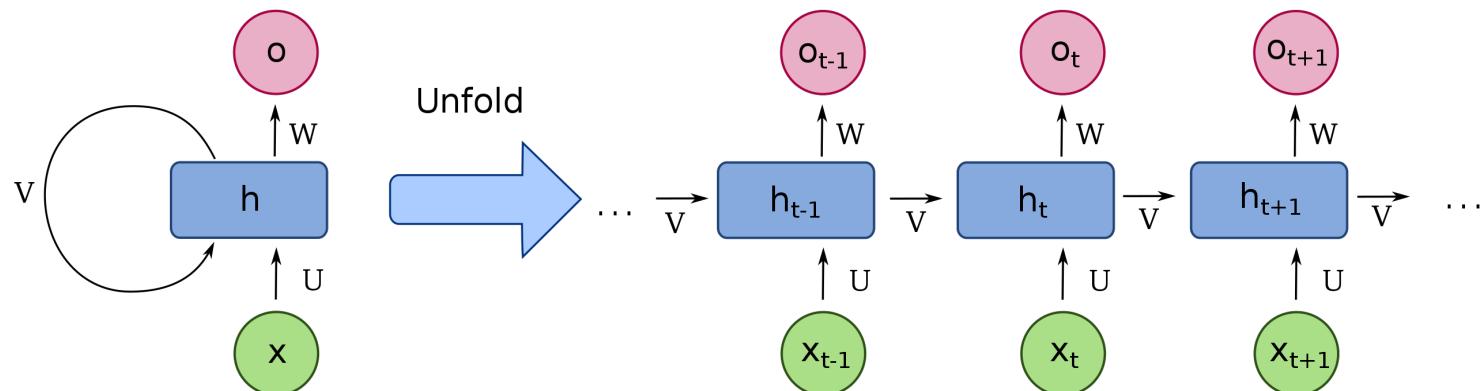
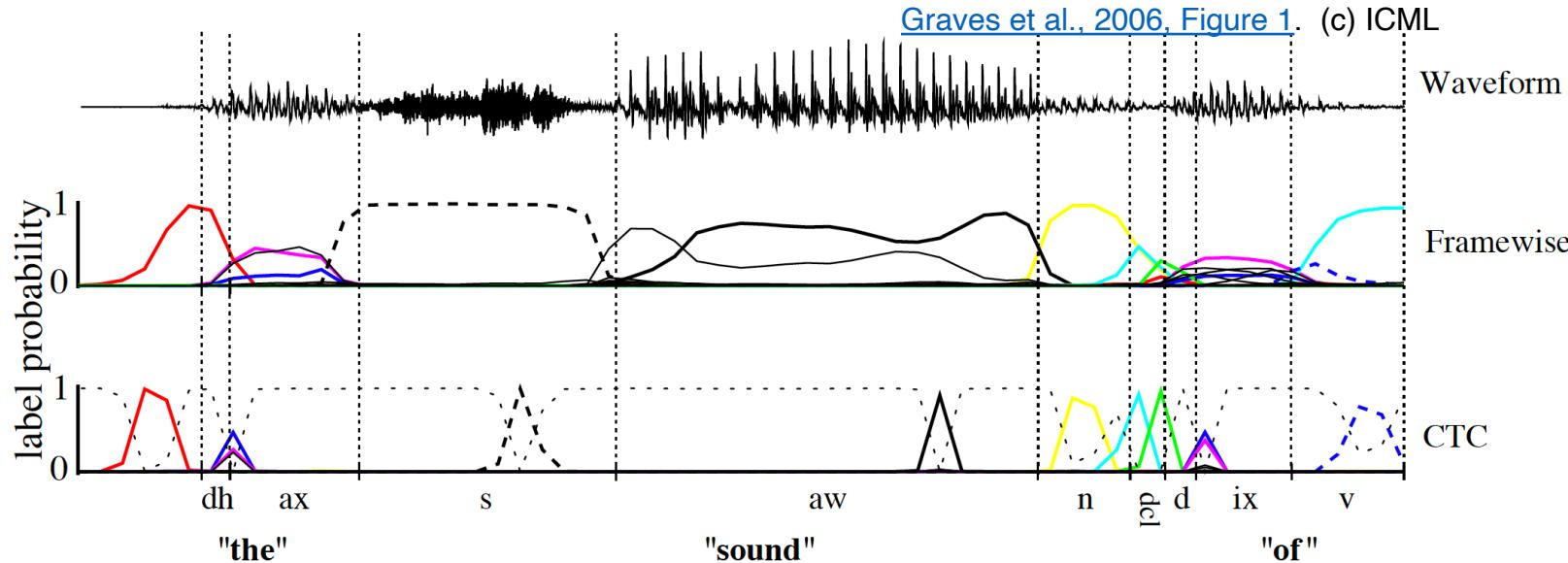
CC-BY 2.0, JCS, 2021



With every breath from my bronze-pounded chest, we will  
raise this wounded world into a wondrous one.  
Amanda Gorman, 2021

- Input: sampled audio waveform,  $x = [x[1], \dots, x[T]]$
- Output: characters,  $y = \{y_1, \dots, y_N\}$

# Recurrent neural network



- In a recurrent neural network (RNN), the hidden node activation vector,  $h_t$ , depends on the value of the same vector at time  $t - 1$ .
- From 2014-2017, the best speech recognizers used RNNs.
- The input is  $x_t$ =speech (waveform or spectrogram).
- The output is  $o_t$ =text characters.

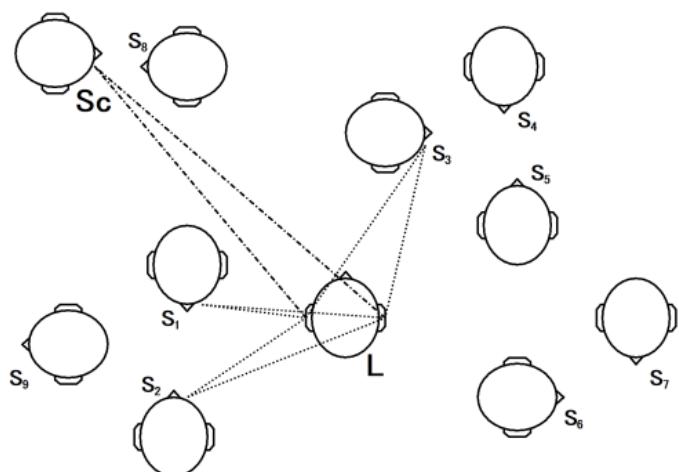
# The Cocktail-Party Effect

- If you are focusing on one person's voice, but hear your name spoken by another person, your attention immediately shifts to the second voice.
- This “cocktail-party effect” suggests a model of hearing in which all sounds are processed preconsciously. Trigger sounds in an unattended source will cause attention to re-orient to that source.

[https://commons.wikimedia.org/wiki/File:Cocktail\\_Party\\_At\\_The\\_Imperial\\_Hotel\\_March\\_13,\\_1961\\_\(Tokyo,\\_Japan\)\\_496610682.jpg](https://commons.wikimedia.org/wiki/File:Cocktail_Party_At_The_Imperial_Hotel_March_13,_1961_(Tokyo,_Japan)_496610682.jpg)



[https://commons.wikimedia.org/wiki/File:Cocktail\\_party\\_attendees\\_at\\_Fuller\\_Lodge,\\_1946.jpg](https://commons.wikimedia.org/wiki/File:Cocktail_party_attendees_at_Fuller_Lodge,_1946.jpg)

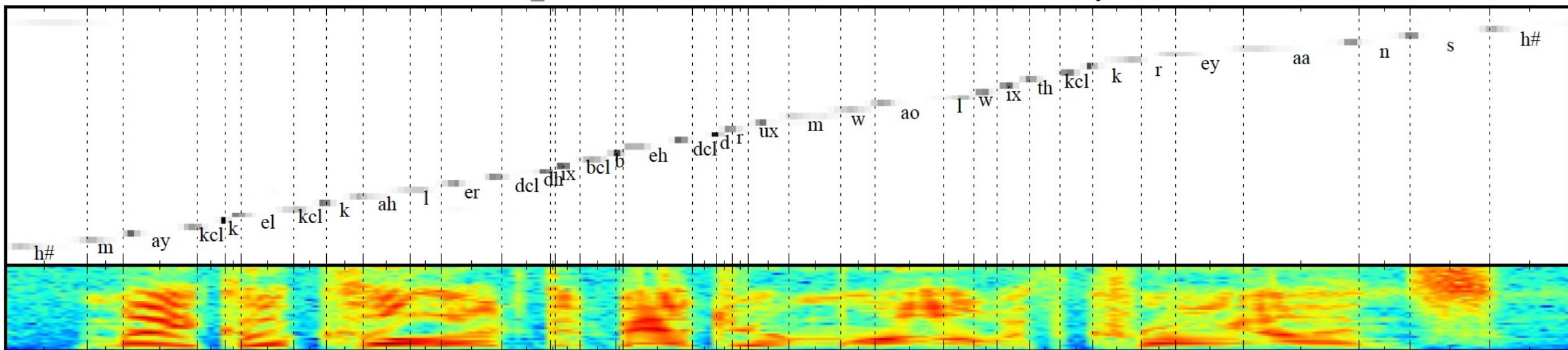


[https://commons.wikimedia.org/wiki/File:Cocktail-party\\_effect.svg](https://commons.wikimedia.org/wiki/File:Cocktail-party_effect.svg)

# Bottom-up attention as a strategy for machine listening

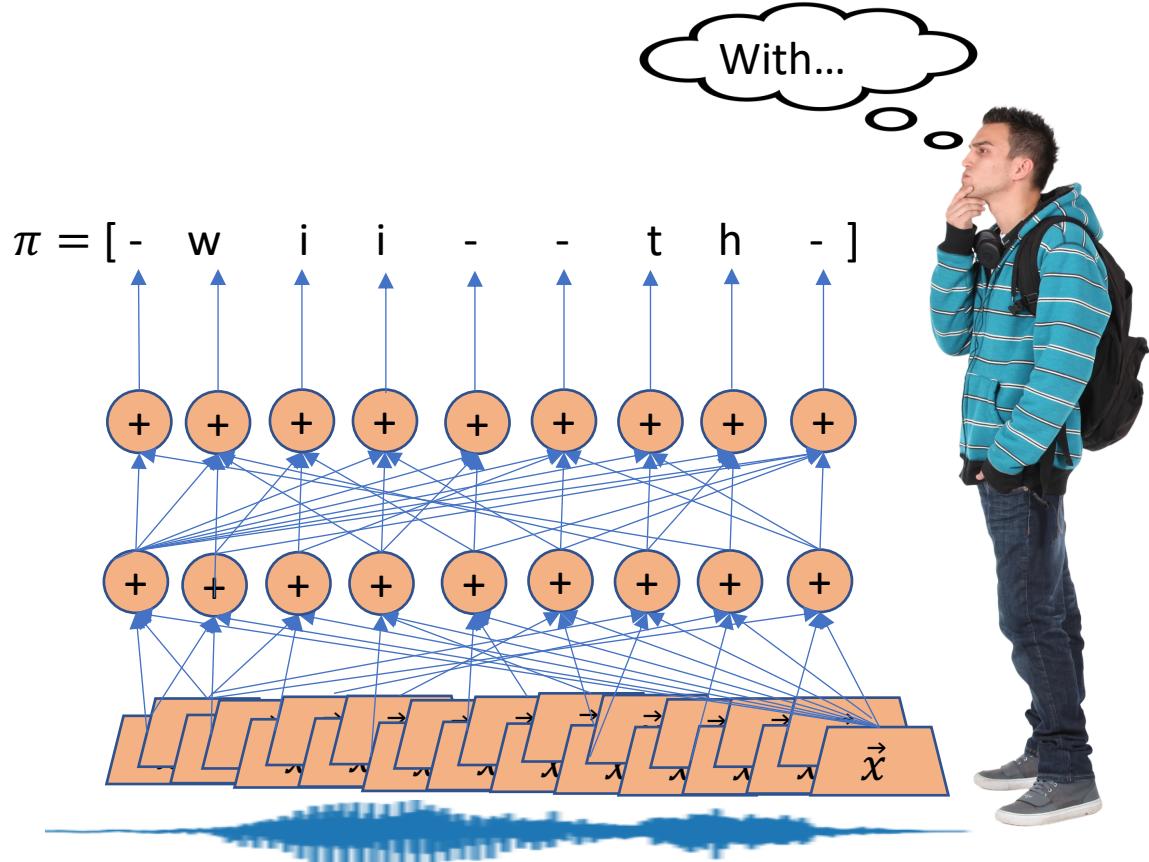
- In 2014, researchers proposed that the past 200ms of RNN state vectors should be stored in a “short-term memory buffer”
- A speech recognizer can attend to several centiseconds, all at one time, to decide what words it thinks it is hearing

FDHC0\_SX209: Michael colored the bedroom wall with crayons.



Chorowski, Bahdanau, Serdyk, Cho & Bengio, [Attention-Based Models for Speech Recognition](#), Fig. 1

# The Transformer: “Attention is all you need”



- In 2017, researchers proposed that the short-term memory buffer should contain raw signals, not processed signals.
- All processing is done using a model of bottom-up attention.
- The “context” for each part of the signal is computed by a model of perceptual similarity:
  - Perceptual similarity in the first layer
  - Semantic similarity in the top layer

# Word Error Rates using Transformers

By 9/2020, transformers had error rates of:

- 2%: English, quiet recording conditions
- 4%: Chinese or Japanese, quiet recording conditions
- 5-7%: if the reference transcript has errors
- 14%: 2-talker mixtures, synthetic reverberation
- 38%: actual in-home recordings in noisy households

**Table 1.** CER/WER results on various open source ASR corpora. Both Transformer and Conformer models are implemented based on ESPnet toolkit. \* marks ESPnet2 results. † and ‡ indicate only w/ speed or only w/ SpecAugment, respectively. § denotes w/o any data augmentation.

Dataset	Vocab	Metric	Evaluation Sets	Transformer	Conformer
AIDATATANG	Char	CER	dev / test	(†) 5.9 / 6.7	<b>4.3 / 5.0</b>
AISHELL-1	Char	CER	dev / test	(†) 6.0 / 6.7	(*) <b>4.4 / 4.7</b>
AISHELL-2	Char	CER	android / ios / mic	(†) 8.9 / 7.5 / 8.6	<b>7.6 / 6.8 / 7.4</b>
AURORA4	Char	WER	dev_0330 (A / B / C / D)	<b>3.3 / 6.0 / 4.5 / 10.6</b>	4.3 / 6.0 / 5.4 / <b>9.3</b>
CSJ	Char	CER	eval{1, 2, 3}	(*) 4.7 / 3.7 / 3.9	(*) <b>4.5 / 3.3 / 3.6</b>
CHiME4	Char	WER	{dt05, et05}_[simu, real]	(†) 9.6 / 8.2 / 15.7 / 14.5	<b>9.1 / 7.9 / 14.2 / 13.4</b>
Fisher-CallHome	BPE	WER	dev / dev2 / test / devtest / evltest	22.1 / 21.5 / 19.9 / 38.1 / 38.2	<b>21.5 / 21.1 / 19.4 / 37.4 / 37.5</b>
HKUST	Char	CER	dev	(†) 23.5	(†) <b>22.2</b>
JSUT	Char	CER	our split	(†) 18.7	<b>14.5</b>
LibriSpeech	BPE	WER	{dev, test}_[clean, other]	2.1 / 5.3 / 2.5 / 5.5	<b>1.9 / 4.9 / 2.1 / 4.9</b>
REVERB	Char	WER	et_{near, far}	(†) 13.1 / 15.4	(†) <b>10.5 / 13.9</b>
Switchboard	BPE	WER	eval2000 (callhm / swbd)	17.2 / 8.2	<b>14.0 / 6.8</b>
TEDLIUM2	BPE	WER	dev / test	9.3 / 8.1	<b>8.6 / 7.2</b>
TEDLIUM3	BPE	WER	dev / test	10.8 / 8.4	<b>9.6 / 7.6</b>
VoxForge	Char	CER	our split	(§) 9.4 / 9.1	(§) <b>8.7 / 8.2</b>
WSJ	BPE	WER	dev93 / eval92	(‡) <b>7.4 / 4.9</b>	(‡) 7.7 / 5.3
WSJ-2mix	Char	WER	tt	(§) 12.6	(§) <b>11.7</b>

# Outline

- The many dimensions of ability and disability
- Speech technology: the Transformer
- Speech technology: Pre-training and Fine-tuning
- Speech technology: Audiobooks
- Inclusive speech recognition
- The Speech Accessibility Project
- Disability in the public sphere

# Label quality

- Audiobooks on librivox.org are readings of texts from Gutenberg.org.
- Often, readers make mistakes, or read different versions, or change the title enough to make it hard to know which Gutenberg text they read.
- Until 2020, speech technology was trained using “labeled data”:
  - 400 hours of librivox books with verified transcripts (“Librispeech Clean”)
  - 600 hours of librivox books with acoustic noise or transcript errors (“Librispeech Other”)

**LibriVox**  
free public domain audiobooks

Search by Author, Title or Reader

Advanced search

**Free public domain audiobooks**  
Read by volunteers from around the world.

**Read**  
LibriVox audiobooks are read by volunteers from all over the world. Perhaps you would like to join us?  
[VOLUNTEER](#)

**Listen**  
LibriVox audiobooks are free for anyone to listen to, on their computers, iPods or other mobile device, or to burn onto a CD.  
[CATALOG](#)

**Welcome to Project Gutenberg**

**Project Gutenberg is a library of over 60,000 free eBooks**

Choose among free epub and Kindle eBooks, download them or read them online. You will find the world's great literature here, with focus on older works for which U.S. copyright has expired. Thousands of volunteers digitized and diligently proofread the eBooks, for you to enjoy.

Kapellendorf by Sophie Hoechstetter	Uncle Wiggily's Airship by Howard R. Garis	The first church's Christmas by George Michael Sterry	The old South A Monograph by H. H. Hamill, D.D.	Least Said, Soonest Mended by Edwin August Grover	X-mas Sketches from the Dartmouth Literary Monthly Edited by Edwin August Grover	X-mas Sketches from the Dartmouth Literary Monthly Edited by Edwin August Grover

Some of our latest eBooks [Click Here for more latest books!](#)

# What do babies hear?

How much unlabeled speech does a baby hear?

- 2000-15000 words/day = 600-4500 hours of speech by age 6 (Weisleder & Fernald, 2013)

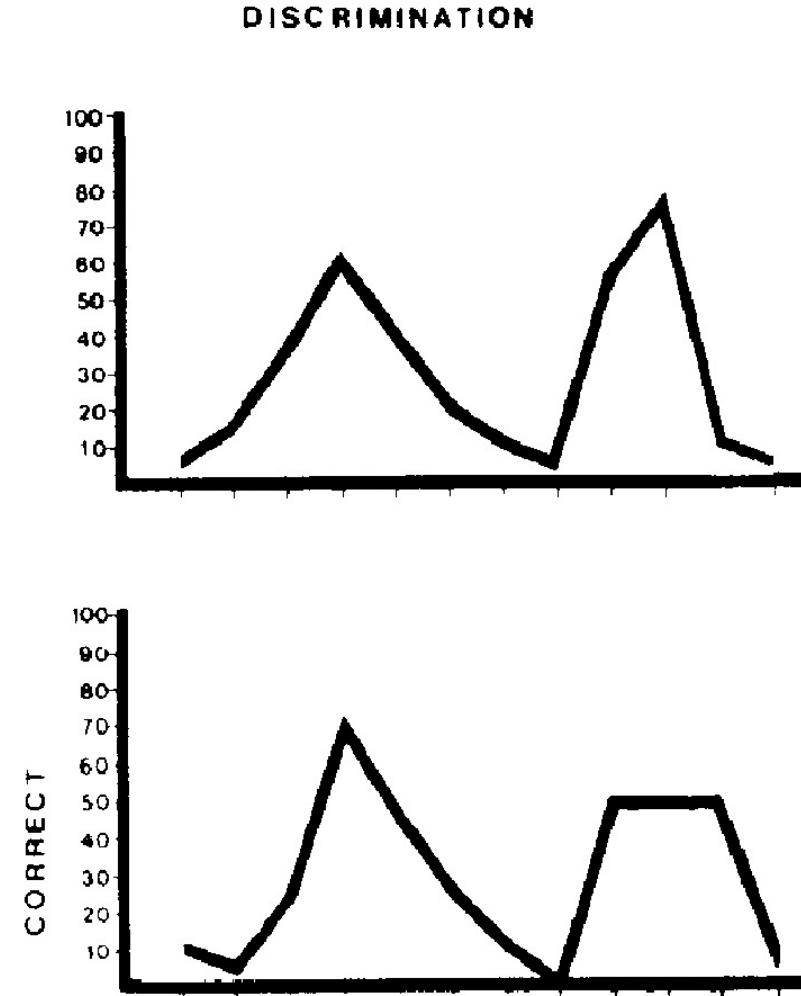
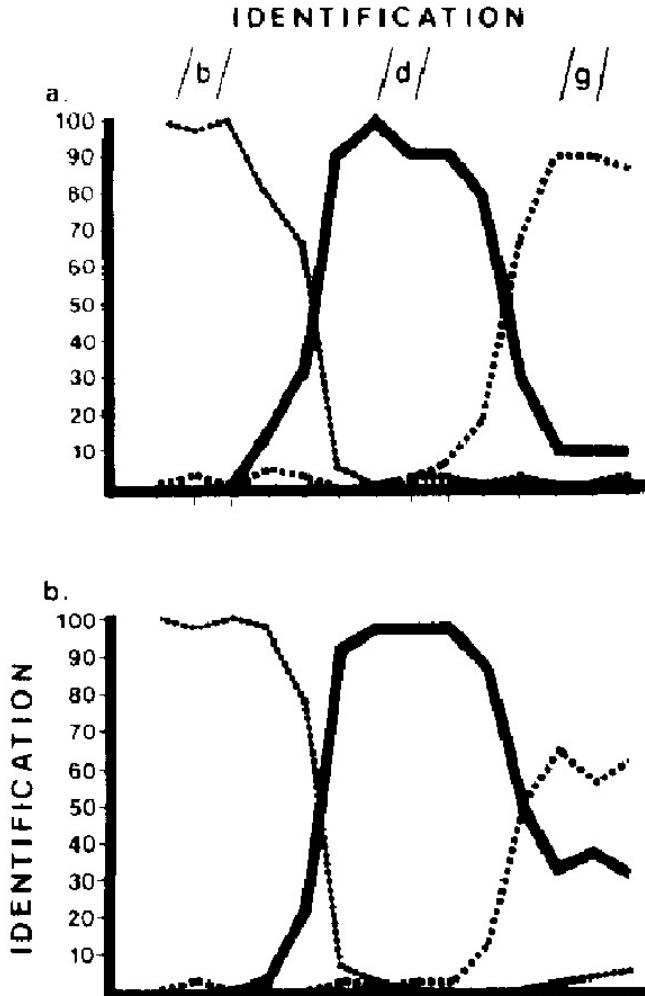
How much labeled speech does a baby hear?

- 30 (?) words/day accompanied by referential gestures = 9.1 hours of speech by age 6



By Steve Jurvetson from Los Altos, USA - A Proper Space Book for Babies, CC BY 2.0, <https://commons.wikimedia.org/w/index.php?curid=105132804>

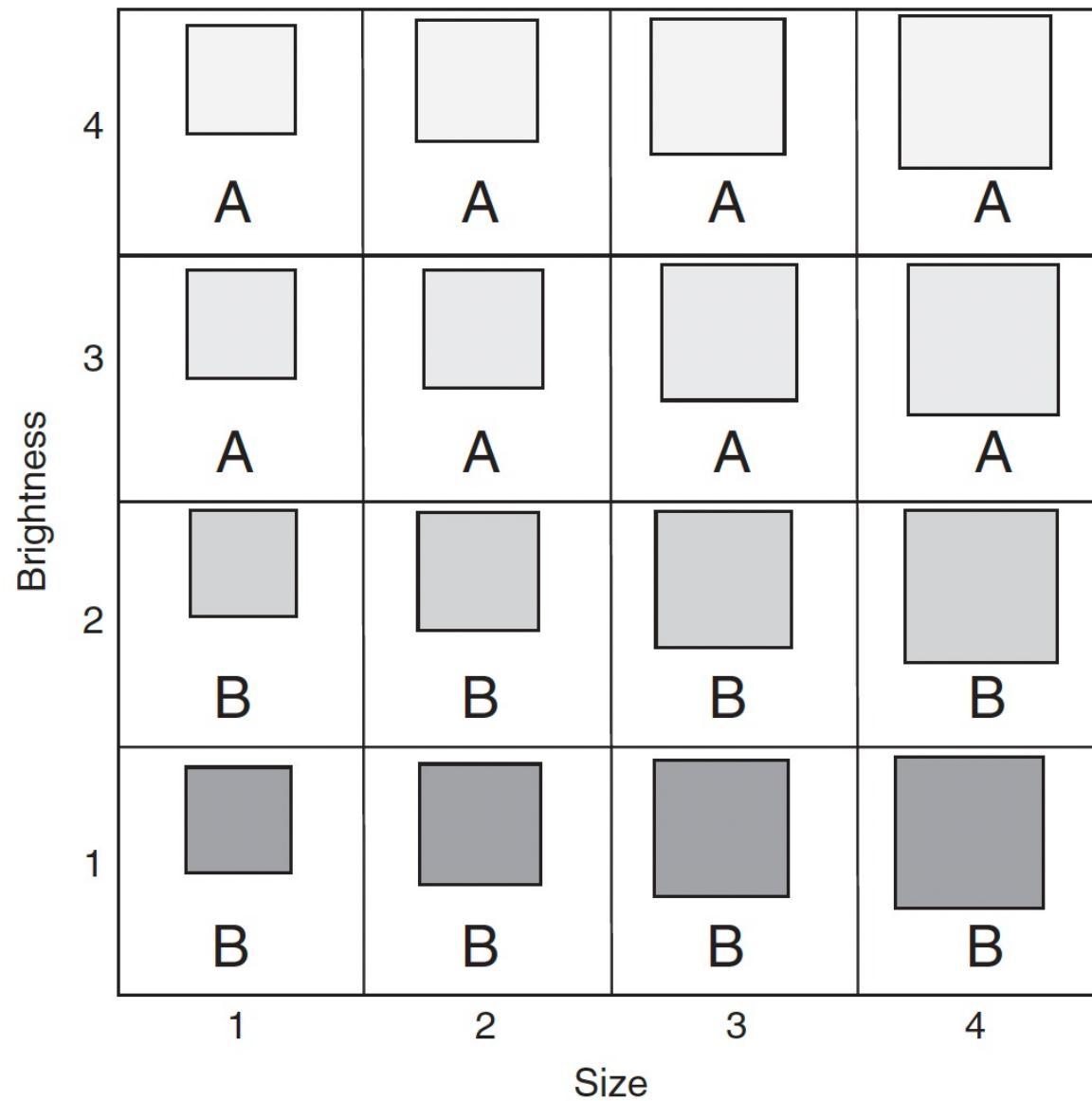
# Is speech learned, or innate? (Hint: it's a trick question)



- 15 synthetic syllables, continuous from /ba/ to /da/ to /ga/
- same label  $\Rightarrow$  hard to tell if they are same sound or different sounds
- different labels  $\Rightarrow$  heard as obviously different

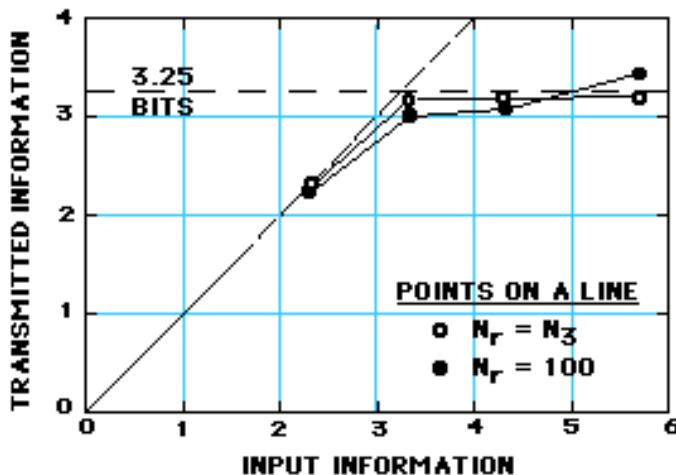
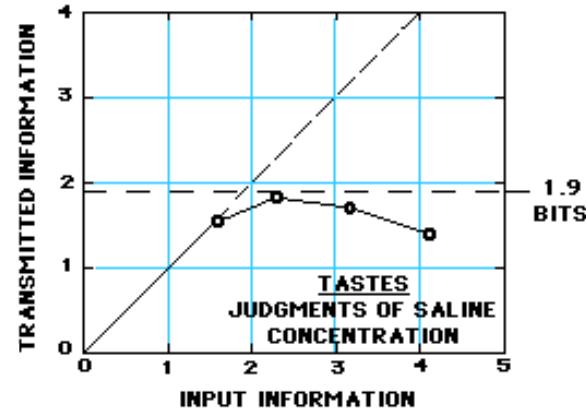
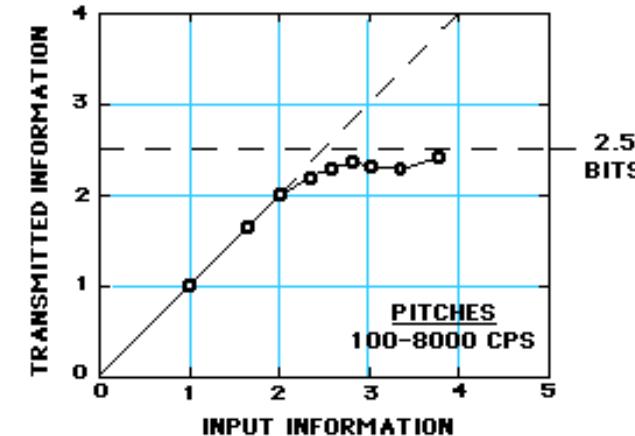
# Categorical perception can be learned!

- People trained to categorize based on **brightness** show reduced within-category perceptual memory, and greater across-category perceptual memory, for **brightness**, but **size** is perceived on a continuum.
- People trained to categorize based on **size** show reduced within-category perceptual memory, and greater across-category perceptual memory, for **size**, but **brightness** is perceived on a continuum.



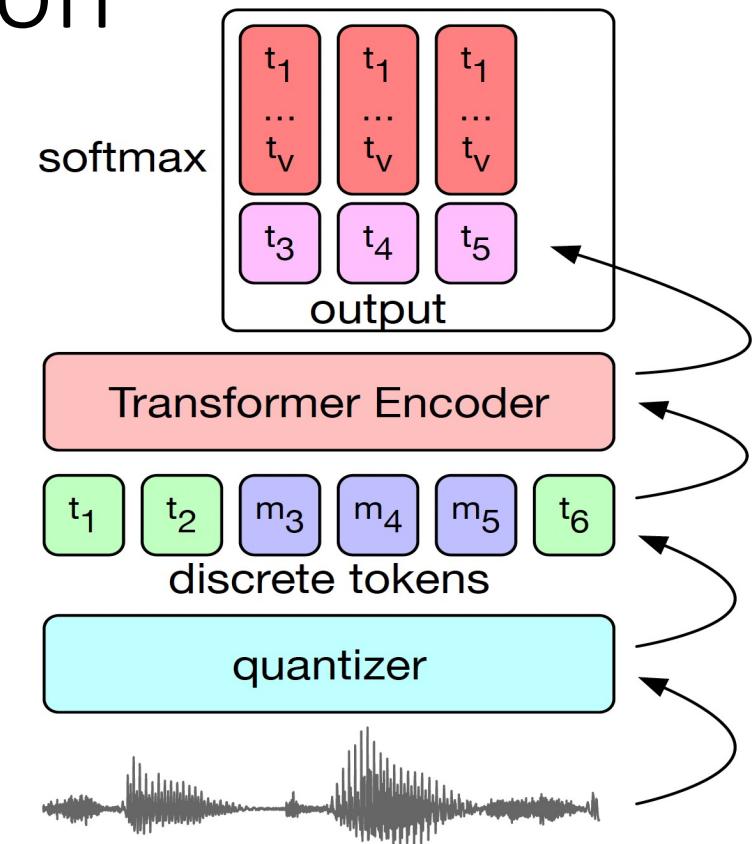
# Categorical perception as a cognitive bias

- If we categorize things, maybe we can remember them longer.
- In “The Magic Number Seven,” Miller argued that people can be taught to categorize any continuum (pitch, taste, position, size) into seven categories, but not more.



# Unsupervised pre-training of transformers based on categorical perception

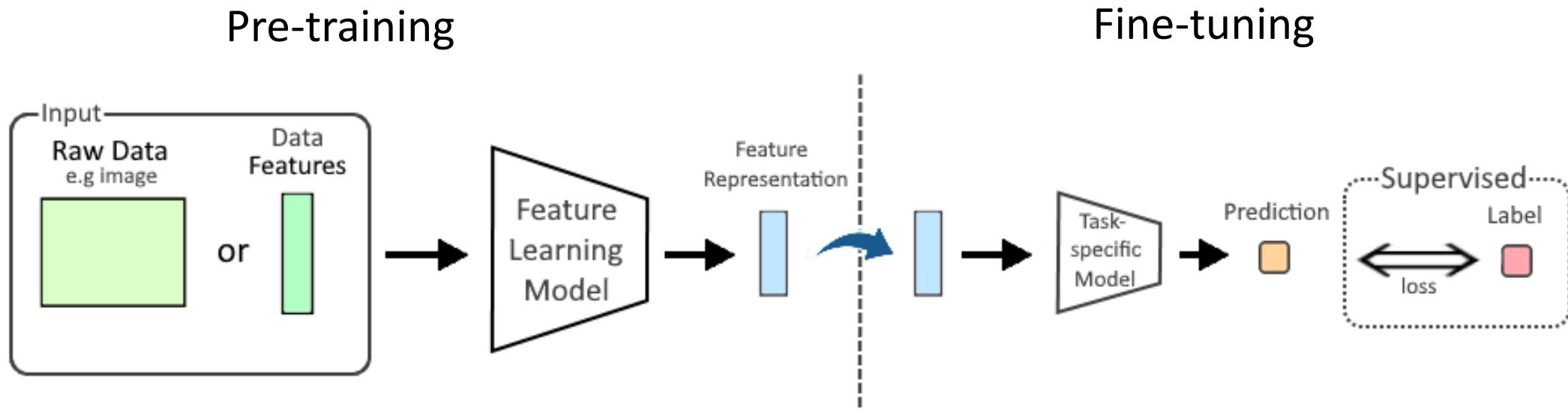
- Given: 60,000 hours of speech, with no associated text.
- Suppose we train the neural network to form its own categories. What would make those categories speech-like?
- **Context-Predictable Speech Categories:** given the context (the quantized units  $t_1$ ,  $t_2$ , and  $t_6$ ), it should be possible to figure out what phonemes were masked ( $t_3$ ,  $t_4$ ,  $t_5$ ).



[Baevski, Auli & Mohamed, 2019](#)

# Pre-training and Fine-tuning

- A transformer is pre-trained to create its own context-predictable speech categories using, say, 60,000 hours of speech
- Then it is fine-tuned using a few hours, or a few hundred hours, or labeled speech



# Word Error Rates using Pre-Training

Pre-training makes it possible to achieve error rates of

- 4.4% using only 10 minutes of labeled data
- 2.6% using only 1 hour of labeled data

Model	Unlabeled Data	LM	dev-clean	dev-other	test-clean	test-other
<i>10-min labeled</i>						
DiscreteBERT [52]	LS-960	4-gram	15.7	24.1	16.3	25.2
wav2vec 2.0 BASE [7]	LS-960	4-gram	8.9	15.7	9.1	15.6
wav2vec 2.0 LARGE [7]	LL-60k	4-gram	6.3	9.8	6.6	10.3
wav2vec 2.0 LARGE [7]	LL-60k	Transformer	4.6	7.9	4.8	8.2
HUBERT BASE	LS-960	4-gram	9.1	15.0	9.7	15.3
HUBERT LARGE	LL-60k	4-gram	6.1	9.4	6.6	10.1
HUBERT LARGE	LL-60k	Transformer	<b>4.3</b>	7.0	4.7	7.6
HUBERT X-LARGE	LL-60k	Transformer	4.4	<b>6.1</b>	<b>4.6</b>	<b>6.8</b>
<i>1-hour labeled</i>						
DeCoAR 2.0 [51]	LS-960	4-gram	-	-	13.8	29.1
DiscreteBERT [52]	LS-960	4-gram	8.5	16.4	9.0	17.6
wav2vec 2.0 BASE [7]	LS-960	4-gram	5.0	10.8	5.5	11.3
wav2vec 2.0 LARGE [7]	LL-60k	Transformer	2.9	5.4	2.9	5.8
HUBERT BASE	LS-960	4-gram	5.6	10.9	6.1	11.3
HUBERT LARGE	LL-60k	Transformer	<b>2.6</b>	4.9	2.9	5.4
HUBERT X-LARGE	LL-60k	Transformer	<b>2.6</b>	<b>4.2</b>	<b>2.8</b>	<b>4.8</b>

# Outline

- The many dimensions of ability and disability
- Speech technology: the Transformer
- Speech technology: Pre-training and Fine-tuning
- Speech technology: Audiobooks
- Inclusive speech recognition
- The Speech Accessibility Project
- Disability in the public sphere

# A brief history of speech recognizer training corpora

---

- TIMIT (1993): 4 hours, phonetically transcribed, DARPA-sponsored
  - Broadcast News (1996): 104 hours, orthographically transcribed, DARPA-sponsored
  - Switchboard (1997): 300 hours, orthographically transcribed, DARPA-sponsored
- 
- Corporate data providers enter the picture. From 2000-2020, high-cost training datasets grow in both size and quality.
  - LibriSpeech (2015): 1000 hours, curated from public domain audiobooks recorded by contributors to librivox.org.
  - Multilingual LibriSpeech (2017): 3000 hours, curated from public domain audiobooks recorded by contributors to librivox.org.
-

# Love, not money

- Audiobooks on librivox.org are usually readings of texts from Gutenberg.org.
- In order to be posted on librivox.org, the audio must be released into the public domain.
  - Portfolio samples of professional audiobook narrators
  - Books contributed by people who love them, and want them to be available
  - Languages contributed by people who love them, and want them to be available
  - Books contributed by people who love audiobooks

**LibriVox**  
free public domain audiobooks

Search by Author, Title or Reader

Advanced search

**Free public domain audiobooks**  
Read by volunteers from around the world.

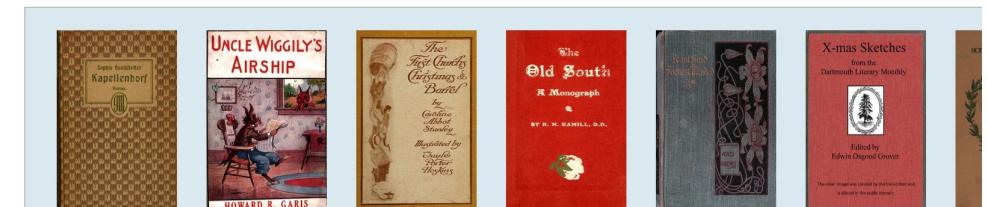
**Read**  
LibriVox audiobooks are read by volunteers from all over the world. Perhaps you would like to join us?  
[VOLUNTEER](#)

**Listen**  
LibriVox audiobooks are free for anyone to listen to, on their computers, iPods or other mobile device, or to burn onto a CD.  
[CATALOG](#)

## Welcome to Project Gutenberg

Project Gutenberg is a library of over 60,000 free eBooks

Choose among free epub and Kindle eBooks, download them or read them online. You will find the world's great literature here, with focus on older works for which U.S. copyright has expired. Thousands of volunteers digitized and diligently proofread the eBooks, for you to enjoy.



Kapellendorf  
by Sophie  
Hoechstetter

Uncle  
Wiggily's  
Airship  
by  
Howard  
Garis

The first  
church's  
Christmas  
by Howard  
Melanthon

The old South  
A Monograph  
by H. H. Hamill, D.D.

Least Said,  
Soonest  
Mended  
by Charles  
Dickens

X-mas  
Sketches from  
the Dartmouth  
Literary Monthly  
Edited by  
Edwin August Grover

The Old Testament  
by H. H. Hamill

Some of our latest eBooks [Click Here for more latest books!](#)

### LibriSpeech ASR corpus

**Identifier:** SLR12

**Summary:** Large-scale (1000 hours) corpus of read English speech

**Category:** Speech

**License:** CC BY 4.0

**Downloads (use a mirror closer to you):**

[dev-clean.tar.gz](#) [337M] (development set, "clean" speech )

Mirrors: [\[US\]](#) [\[EU\]](#) [\[CN\]](#)

[dev-other.tar.gz](#) [314M] (development set, "other", more challenging, speech ) Mirrors: [\[US\]](#) [\[EU\]](#) [\[CN\]](#)

[test-clean.tar.gz](#) [346M] (test set, "clean" speech ) Mirrors: [\[US\]](#) [\[EU\]](#) [\[CN\]](#)

[test-other.tar.gz](#) [328M] (test set, "other" speech ) Mirrors: [\[US\]](#) [\[EU\]](#) [\[CN\]](#)

[train-clean-100.tar.gz](#) [6.3G] (training set of 100 hours "clean" speech ) Mirrors: [\[US\]](#) [\[EU\]](#) [\[CN\]](#)

[train-clean-360.tar.gz](#) [23G] (training set of 360 hours "clean" speech ) Mirrors: [\[US\]](#) [\[EU\]](#) [\[CN\]](#)

[train-other-500.tar.gz](#) [30G] (training set of 500 hours "other" speech ) Mirrors: [\[US\]](#) [\[EU\]](#) [\[CN\]](#)

### Librispeech

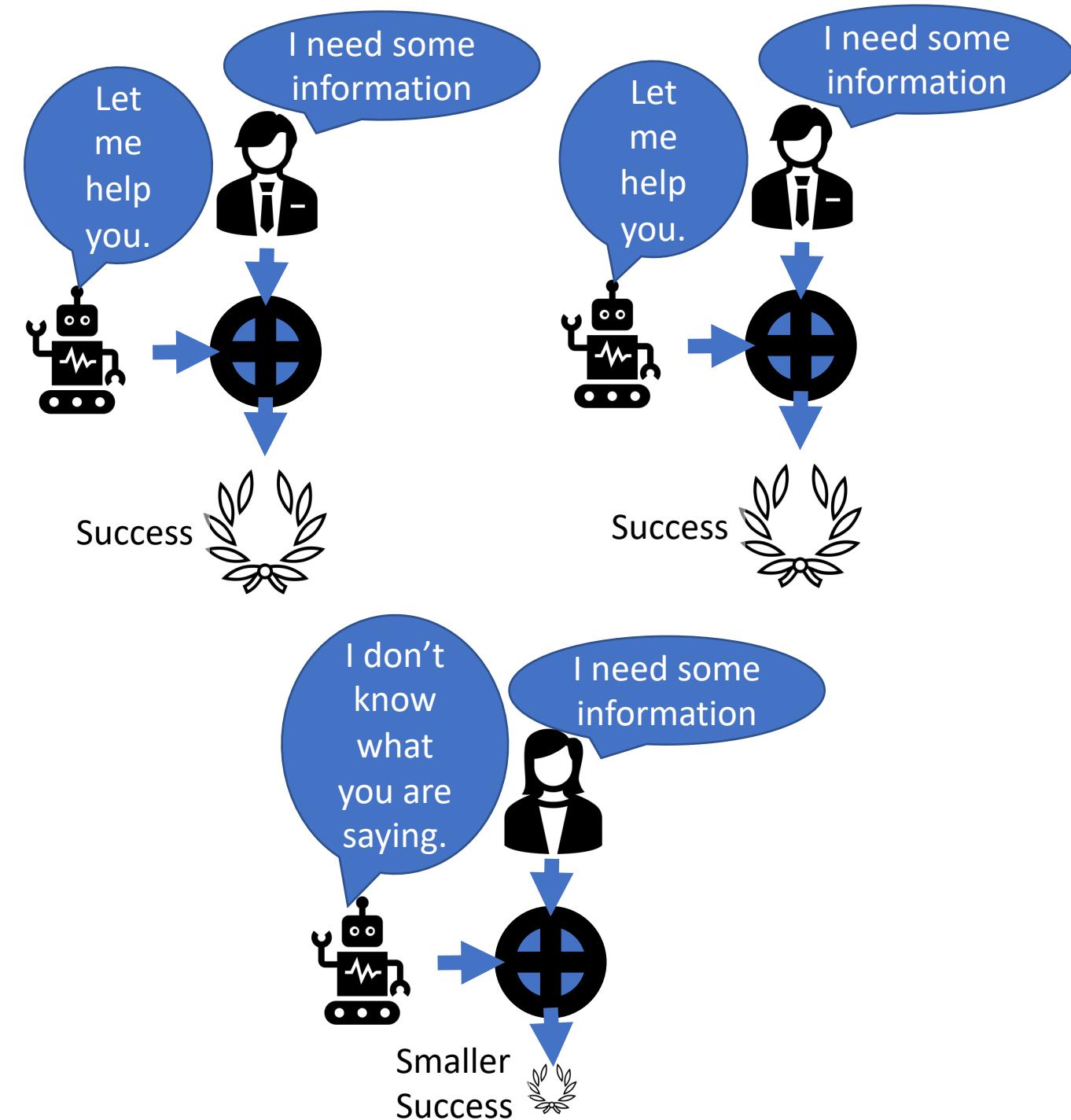
- Librispeech includes 360h clean speech, 500h other speech
- Curated from [librivox.org](#) and [Gutenberg.org](#)
- Free to download, free to use and redistribute
- Librispeech is the reason for the deep-learning revolution in automatic speech recognition.

# Outline

- The many dimensions of ability and disability
- Speech technology: the Transformer
- Speech technology: Pre-training and Fine-tuning
- Speech technology: Audiobooks
- Inclusive speech recognition
- The Speech Accessibility Project
- Disability in the public sphere

# Inclusive ASR: why?

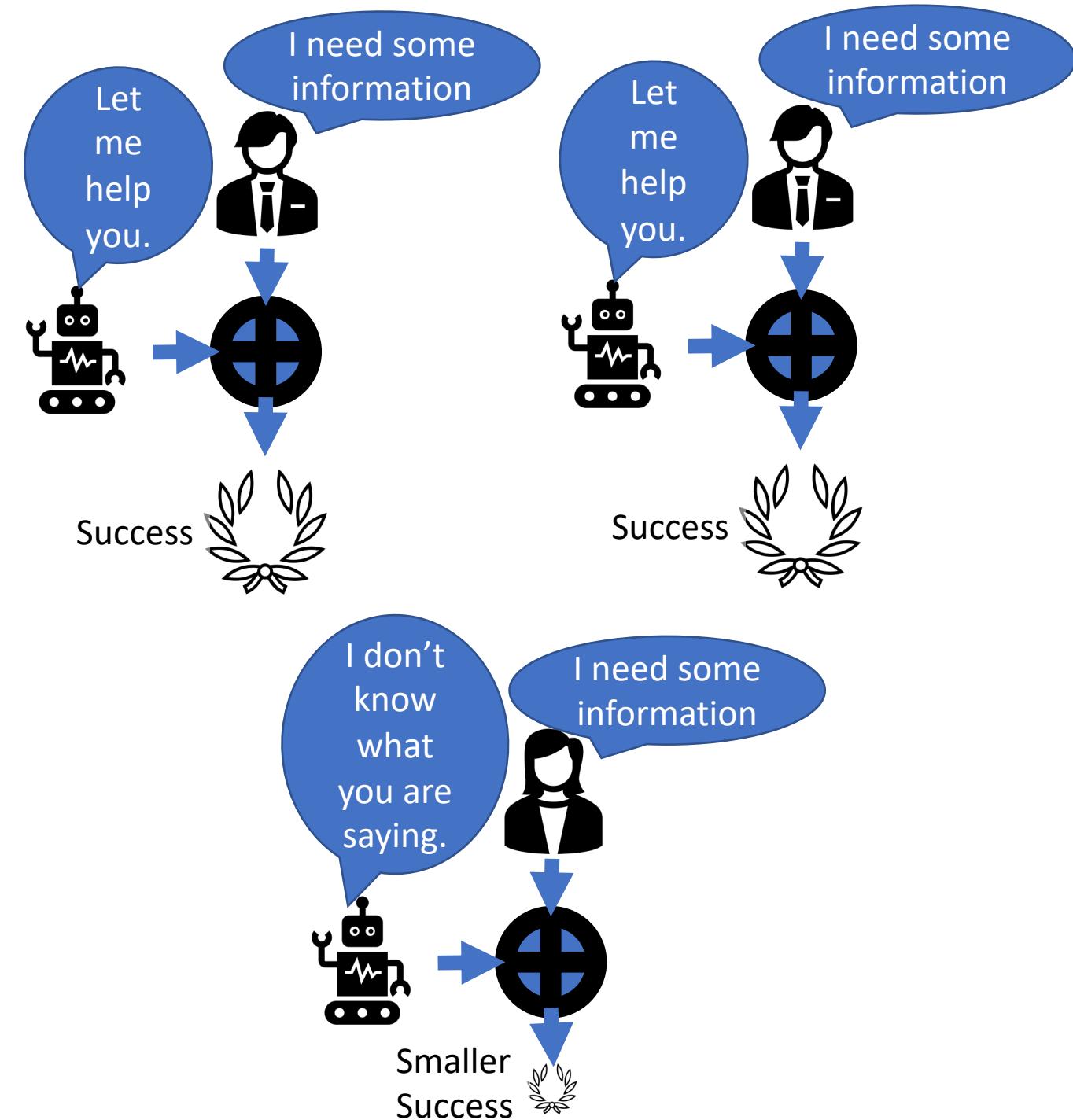
- ASR is a useful productivity tool
- If ASR works for you, it gives you a socioeconomic advantage, compared to somebody for whom ASR fails
- The people for whom ASR fails are often those who are *already* socioeconomically disadvantaged



# Why does ASR fail?

Reasons it may fail for a group:

- **Under-representation**: The training corpus doesn't have enough examples
- **Inter-group variance**: The group speaks differently from other groups
- **Intra-group variance**: Members of the group all speak differently from one another
- **Intra-individual variance**: Members of the group speak less precisely



# Automatic Speech Recognition Word Error Rates (WER)

## Gender:

- Women > Men (51% > 38%: Tatman, 2017, YouTube captioning)
- Women > Men (61% > 47%: Garnerin et al., 2019, European broadcast news)
- Black Men > Black Women (41% > 30%: Koenecke et al., 2020)

## Dialect:

- Scottish > American (53% > 42%: Tatman, 2017)
- American Deep South > General American (Picone 1991)

## Race:

- Black > White (35% > 19%: Koenecke: avg of Amazon, Apple, Google, IBM, Microsoft)

## Disability:

- People w/Cerebral Palsy > People w/o (41% > 33%: Issa et al., in review)

## Age:

- Teenage (<20)) > Old (>70) > Adult (Feng & Scharenborg, 2020)
- Old adult (50-70) < Young adult (20-30) (Sari et al., 2021)

# Database Collection Efforts

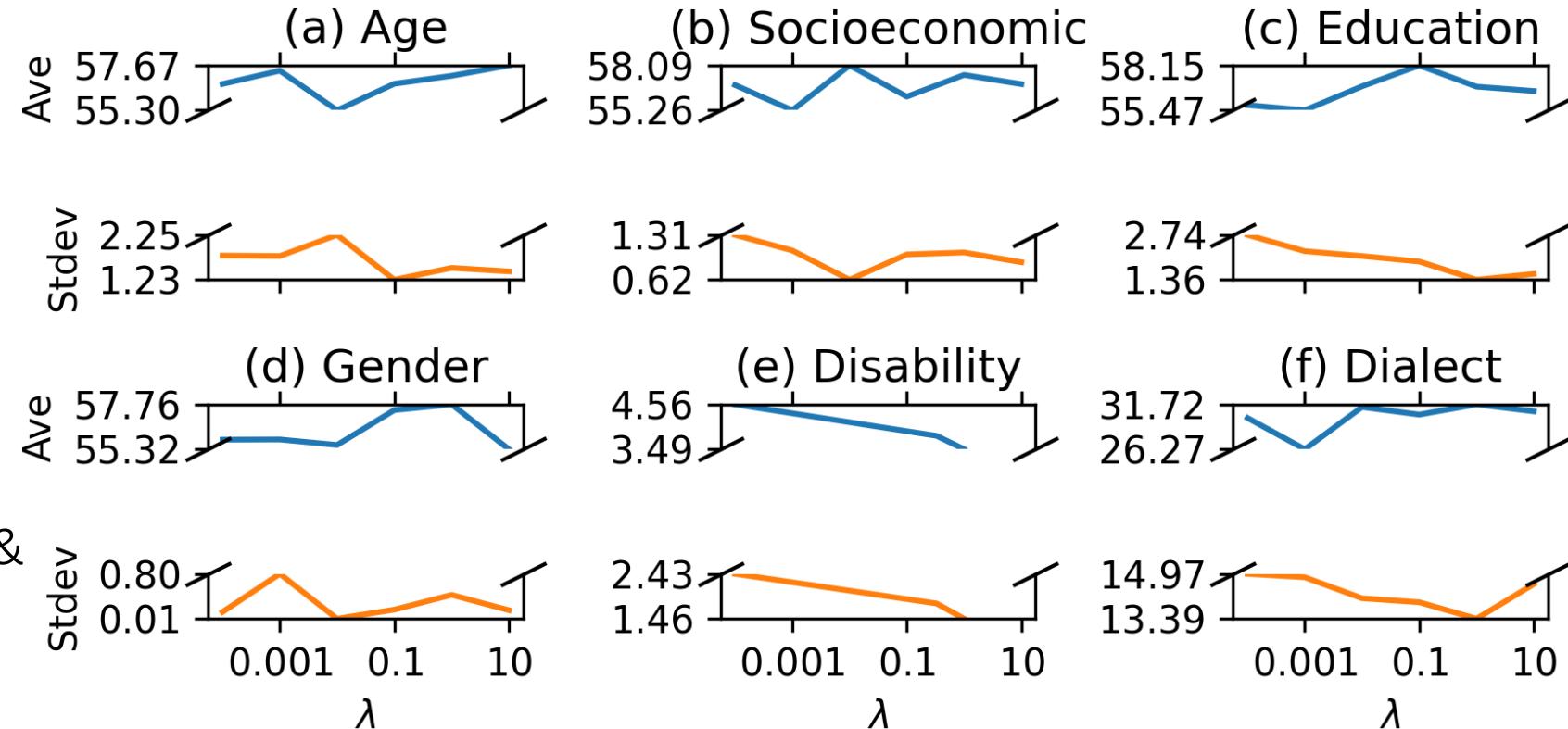
- UASpeech: dysarthria as a symptom of Cerebral Palsy (Kim et al., 2008)
  - Perhaps 3.2 hours
- CORAAL (Corpus of Regional African American Languages, Kendall et al., 2018)
  - Perhaps 200 hours

... but ...

- Librispeech: largest publicly available corpus of General American English speech:
  - 1000 hours
- LibriVox Complete (Baevski et al., 2020):
  - 10,000 hours

# Seamless equal accuracy ratio for inclusive CTC speech recognition

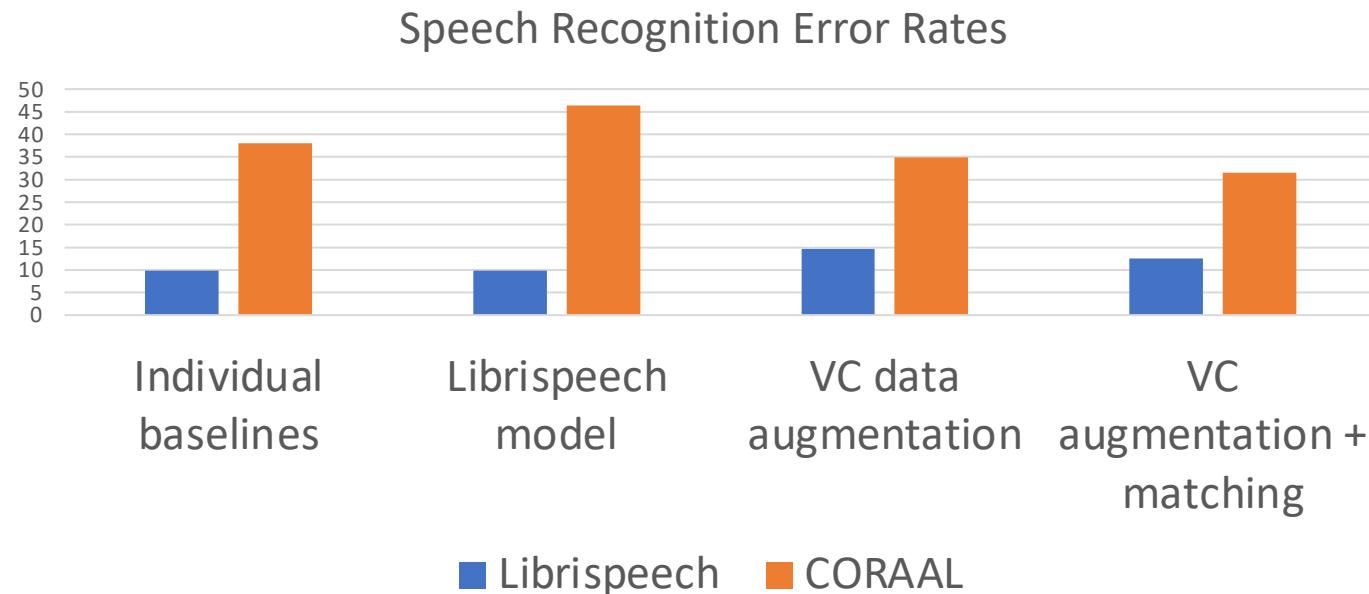
Gao, Wang, Kang, Mina, Issa,  
Harvill, Sari, Hasegawa-Johnson &  
Yoo, 2020



- $\lambda = 0$ : train speech recognizer to reduce average error rate
- $\lambda = 1$ : train it to reduce the maximum error rate incurred by any group
- As  $\lambda$  gets larger, standard deviation of error rate (among speakers) decreases. This was expected.
- As  $\lambda$  gets larger, the average error rate (averaged across speakers) also decreases. This was unexpected. Apparently, focusing on the most poorly-served group can help to improve the error rate for all groups.

# Counterfactually fair automatic speech recognition

Sari, Hasegawa-Johnson & Yoo, 2021

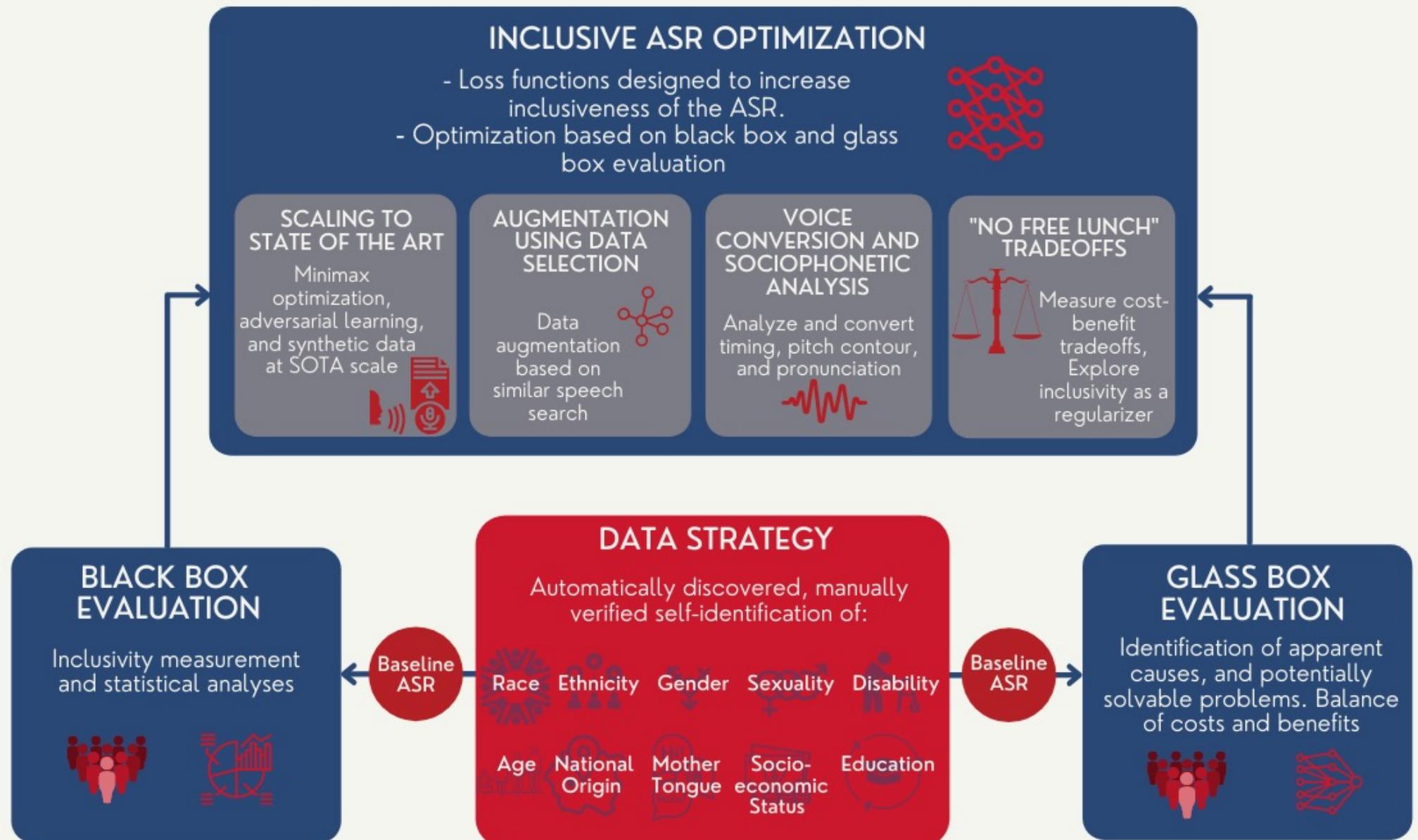


- Librispeech = audiobooks, 1000h; plurality are “General American English” (GAE)
- CORAAL = corpus of regional African American Languages (AAL), 200h, autobiographical narratives from Atlanta, DC, Detroit, Manhattan, Princeville, Rochester, Valdosta
- Applying the Librispeech model to CORAAL data is worse than using the CORAAL-only baseline
- Synthesizing new data using voice conversion (VC) reduces error rates on CORAAL
- Counterfactual matching forces the transcriptions for a given speaker to remain unchanged even after the voice is converted from AAL to GAE or vice versa

# FAI: A new paradigm for the evaluation and training of inclusive automatic speech recognition

Hasegawa-Johnson, Fagyal, Dehak, Zelasko, and Moro-Velazquez, 2021

## Evaluation and Training of Inclusive Automatic Speech Recognition



# Outline

- The many dimensions of ability and disability
- Speech technology: the Transformer
- Speech technology: Pre-training and Fine-tuning
- Speech technology: Audiobooks
- Inclusive speech recognition
- **The Speech Accessibility Project**
- Disability in the public sphere

# The Speech Accessibility Project

- Speech technology is now good enough to be useful for people speaking “General American English” without disabilities.
- To make it useful for people with disabilities, we need about 1000 hours of transcribed speech (~1.2 million sentences).



# SPEECH ACCESSIBILITY PROJECT

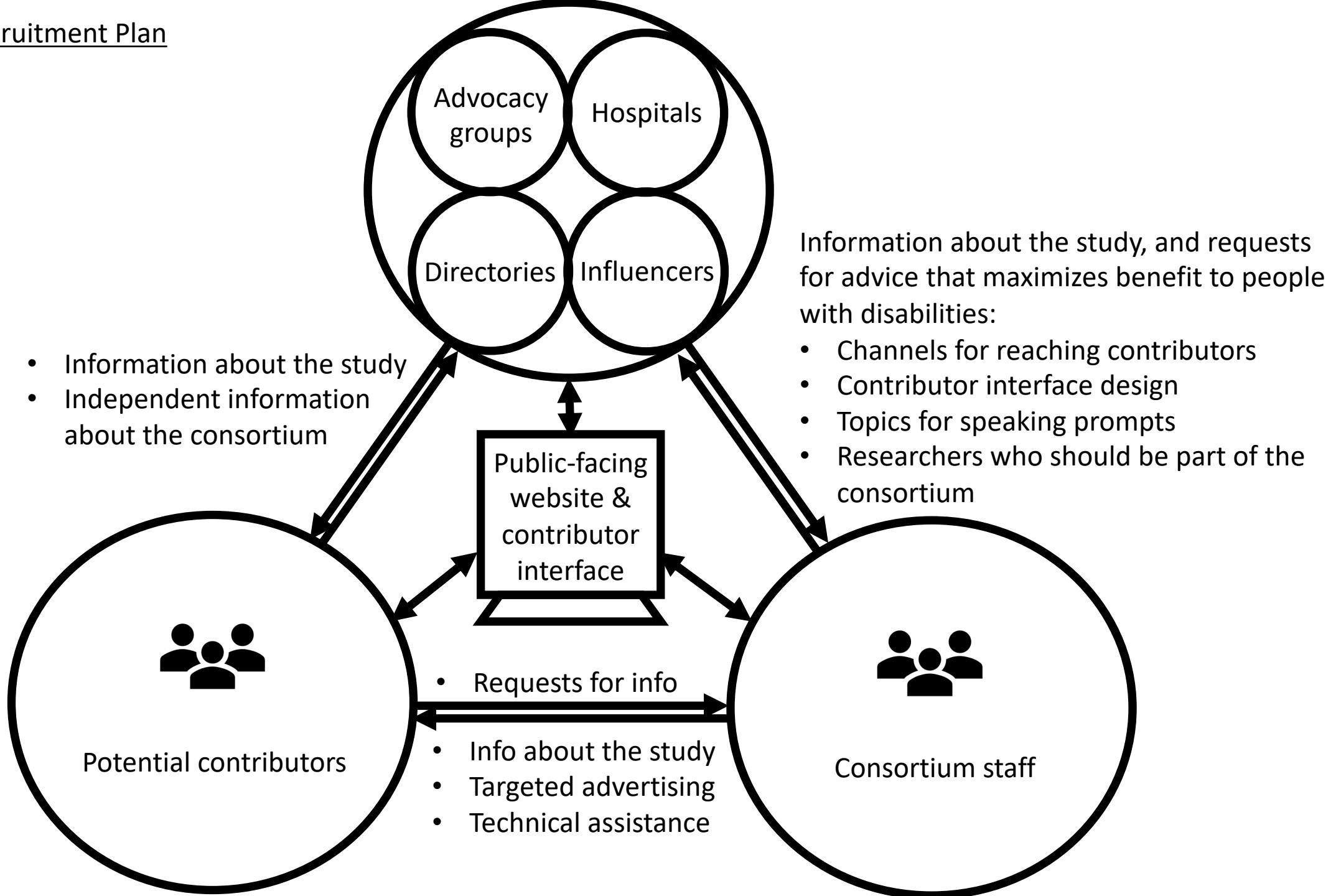
Beckman Institute for Advanced Science and Technology

The Speech Accessibility Project is intended as a communication channel that will connect people with disabilities with the engineers and computer scientists who create speech technology.

# Etiologies

- Etiology = cause of the disability
- Different etiologies have different speech patterns
- Having many different etiologies is desirable (corpus includes speech representative of different etiologies)
- Having many speakers for each etiology is also desirable (speech patterns within each group may be quite variable)
- Our compromise: we intend to recruit 400 speakers from each of 5 different etiologies, starting with Parkinson's, ALS, and Cerebral Palsy.

## Outreach & Recruitment Plan



# Human subject protection principles

Contributors should know how their data will be used;  
Researchers should commit to the same terms.

- The purpose of the study and permitted uses, as specified in the participant consent form,
- ...should be the same as the Permitted Purpose and permitted uses in the data use agreement.



University of Illinois Urbana-Champaign  
Online Consent Form

## Speech Accessibility Project: Individuals with disabilities helping researchers to improve technology

You are being asked to participate in a voluntary research study that is called the "Speech Accessibility Project." The purpose of this study is to help researchers at universities and companies to develop spoken-language user interfaces that work for people with atypical speech. Software programs that understand speech are developed using machine learning. Machine learning is a software development method that imitates the way human behavior whenever it makes a mistake mistakes in the future. By providing e possible for the software to learn how the future. Participating in this study your audio recordings to a secure end samples or about 5 hours of your time convenience. Risks associated with the expected benefit of this research to you and organizations using your speech with speech and motor disabilities, w

**Project Name:** The Speech Accessibility Project  
<https://speechaccessibilityproject.be>

### LICENSED SPEECH SIGNAL DATA DATA USE TERMS AND CONDITIONS

The licenses to Licensed Speech Signal Data granted to each Member pursuant to Section 9.b. of the Agreement are subject to the following Terms and Conditions. Each Member receiving such a license is referred to herein as the "**Data User**". Any capitalized terms used but not otherwise defined in this Exhibit C shall have the meanings given to them elsewhere in the Agreement.

#### 1. Provision and Use of Data

1.1 Effective as of the Effective Date, Data User may make, have made, use, load, access, store, copy, reproduce, destroy, modify, transmit, display, make derivative works of and otherwise use the Data made available to it by Organization under the Agreement on a non-exclusive, non-assignable, non-sublicensable (except as provided in Section 2.1 below) basis solely for the Permitted Purpose subject to the terms of this Exhibit C and the Agreement. Any such modifications or derivative works made by Data User are "**Modification(s)**".

1.2 This Agreement does not restrict Data User's use or modification of any portions of the Data that Organization makes publicly available under a more permissive license, if applicable, to the extent Data User uses or modifies the Data under the terms of such public license, or that otherwise become public.

1.3 The rights and restrictions set forth in this Exhibit C as it relates to the use, modification, distribution, disclosure, and privacy of the Data may only be modified by a written agreement between Data User and Organization specifically indicating it is amending these Terms and Conditions, e.g., these Terms and Conditions may not be amended by or superseded by any general terms and conditions that are presented or included as part of access to or use of data storage, data processing, platform, or computing resources.



IRB Number: 23183  
IRB Approval Date: 03/07/2023  
IRB Expiration Date: 03/02/2028

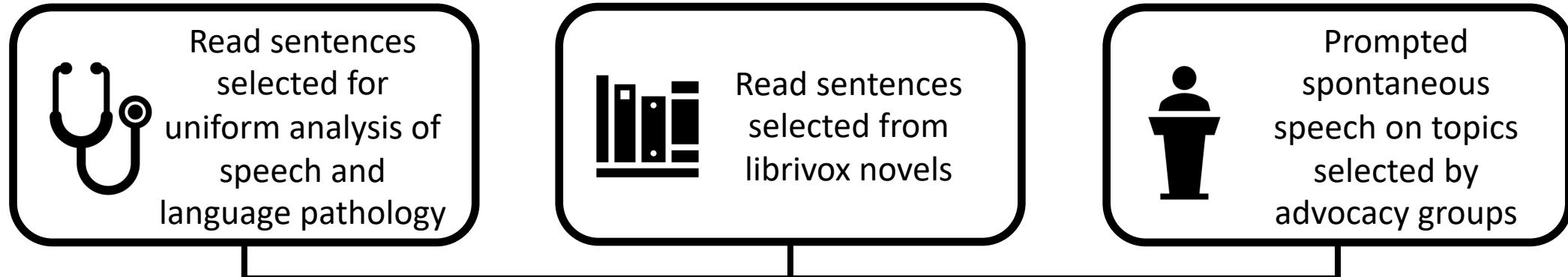
The screenshot shows a GitHub repository page for 'speechaccessibility/ AudioControls'. The repository is public and has 1 branch and 7 tags. The main file listed is 'AudioControls.ts' with a commit message 'dispatch error events in catch blocks' made yesterday. Other commits include 'LICENSE' added more specific copyright holder information 5 months ago, and 'README.md' modified Readme to remove references to various 'Play...' 2 weeks ago. The README.md file is expanded, showing its content: 'Typescript implementation of audio recording, playback and simulated waveform display based on the MediaRecorder API.' An example code snippet is provided:

```
let recordButton = document.getElementById('recordButton')
recordButton.addEventListener(
  'AudioControls.RecordingStarted',
```

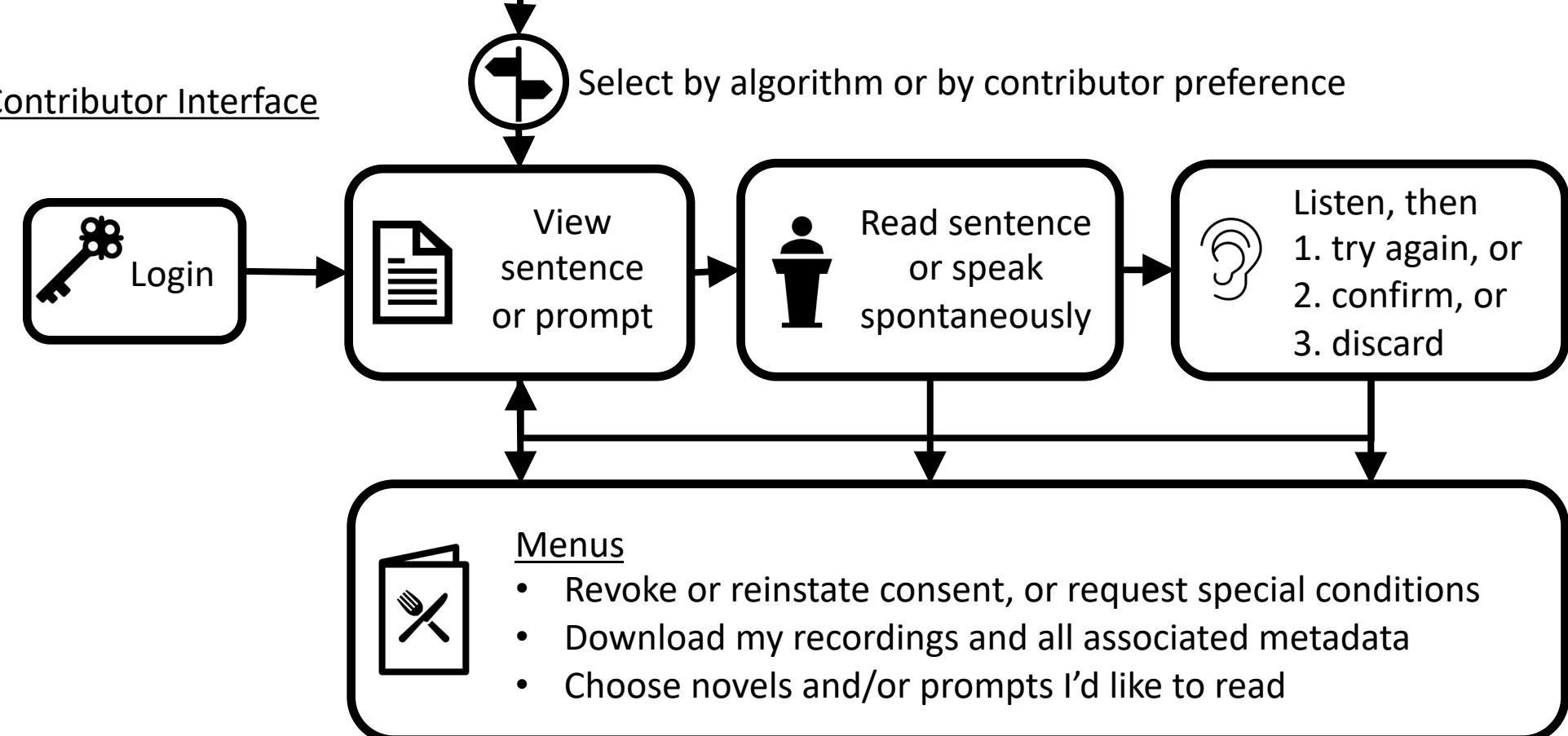
- Software design principles:
  - Contributor safety
  - Data recoverability
  - Contributor ease of use
  - Annotator ease of use
- Contributor and Annotator user interfaces go through dev and test environments before being pushed to production.
- All finished products are released open-source.

<https://github.com/speechaccessibility/ AudioControls>

## Prompt Sources



## Contributor Interface



# Types of speech, Transcription quality

- Speech recognition is most robust if training samples are as natural and spontaneous as possible.  
...however...
- Read speech is easier to transcribe: We know what the person was trying to say, so “transcribing” just means correcting the prompt text so it matches what they actually said.
- Our compromise:
  - 300 command & control sentences (“find my phone,” “What’s on my shopping list?”)
  - 100 phonetically diverse sentences (“When I was a young boy, father always said I was a born business manager.”)
  - 200 spontaneous sentences, in response to 50 prompts for spontaneous speech (“Tell us about one of your favorite bands or singers.”)

# Distribution to researchers

- The data distribution includes:
  - Contributed audio
  - Prompt text, Corrected text,
  - Annotations of dysarthria or dysphonia type and severity
  - Start time & End time of speech, signal to noise ratio, clipping severity
- Data is packaged into an encrypted tar file
- Each consortium member (each organization that has signed a data use agreement) can specify one person who gets a decryption key
- Consortium members will begin testing data quality & attempting to train ASR as soon as we begin distributing the data.

# Outline

- The many dimensions of ability and disability
- Speech technology: the Transformer
- Speech technology: Pre-training and Fine-tuning
- Speech technology: Audiobooks
- The Speech Accessibility Project
- Inclusive speech recognition
- **Disability in the public sphere**

# Disabled in public



## Grief 03: Chronic Pain and Disability in Motherhood—An Interview with Vaneetha Risner

Motherhood has a way of revealing our limits, be it physical, mental, or emotional. When we find ourselves at the end of our abilities, it can be easy to wonder if God has equipped us wit...



Mar 8 · 31 min 50 sec



## 4 ways to design a disability-friendly future | Meghan Hussey

Nearly fifteen percent of the world's population lives with a disability, yet this massive chunk of humanity is still routinely excluded from opportunities. Sharing her experience growing up wi...



Oct 2022 · 10 min 35 sec



## Uncovering The Layers of Disability w/ Tiffany Yu CEO & Founder of Diversability

Episode Notes On episode 319, I sit down with my friend and colleague, the multihyphenate, Tiffany Yu, CEO and Founder of Diversability. We talk about her journey to acceptance of...



Nov 2022 · 1 hr 34 min



## r/AmITheA--Hole My Teacher Mocked My Disability

<https://www.youtube.com/rslash>



Jun 2022 · 16 min 16 sec



## What to know and how to talk about disability

Do you find yourself avoiding conversations on disabilities? Worried you'll offend a disabled friend? A disability rights activist shares ways to be a better ally and to destigmatize disability...

<https://www.youtube.com/watch?v=0MLRcsjfLsE>



THE MICHAEL J. FOX FOUNDATION  
FOR PARKINSON'S RESEARCH

We Exist ▾ Understanding Parkinson's ▾ For Researchers ▾

<https://www.michaeljfox.org/our-promise>



TEAM<sup>37</sup> GLEASON

<https://teamgleason.org/>

[https://commons.wikimedia.org/wiki/File:Davis\\_Phinney\\_1991\\_Thrift\\_Drug\\_Classic.jpg](https://commons.wikimedia.org/wiki/File:Davis_Phinney_1991_Thrift_Drug_Classic.jpg)

# The growing population of Americans with disabilities

- The number of people diagnosed with Parkinson's Disease in the United States is expected to double by 2040. – [NINDS](#)
- The estimated number of people with Multiple Sclerosis ... is 30% higher than in 2013. – [Walton et al., 2020](#)
- By our estimates, the number of cases of ALS in the world will increase from 222,801 in 2015 to 376,674 in 2040 – [Arthur et al., 2016](#)
- In the United States, about 85% of stroke victims now survive. – [Donkor, 2018](#)
- In 1900 the life expectancy of people with Down Syndrome was 9 years. By 1984, [it] had increased to 28 years. Since then, the life expectancy of people with DS has increased to about 60 years. – [Chicoine, 2021](#)

# Summary

---

- Speech technology benefits from psychologically-inspired innovations in machine learning.
- Speech technology benefits from community involvement.
- Physical differences mean disability only if the technology fails to work.
- The Speech Accessibility Project aims to make the technology work.

