

# SPEECH PROCESSING

15 Sep 2020

## SPEECH SOUNDS AND FEATURES :

The number of linguistically distinct speech sounds are variant to different people and are often subject to a matter of judgement.

However, phonetic symbols of American English can be generalized into some standard representations, which are as follows:-

- Vowels
- Diphthongs
- Semi Vowels
- Nasal Consonants
- Unvoiced Fricatives
- Voiced Fricatives
- Voiced and Unvoiced stops.

## VOWELS :

Even though vowels are considered less important for the classification and representation of written text, most speech-recognition rely heavily on vowel recognition to achieve high performance.

For ex: This is an apple.

Th\_s -s -n -ppl-

Avg reader can fill in vowels to decode written text.

--i- i- a- a---e

The presence of vowels alone is not enough to decode the text.

Vowels are usually long in duration as compared to other types of sounds ~~at~~ and are spectrally well-defined. Hence, they are usually easily and reliably recognized and therefore contribute a lot to our ability to recognize speech.

### PRODUCTION OF VOWELS :

Vowels are produced by exciting the fixed vocal tract shape with quasi-periodic pulses of air caused by vibration of vocal cords.

The vowel sound produced is determined primarily by the following :-

- the position of tongue (primary point)
- the positions of jaw, lips, velum (slightly influence the resulting sound)

### CHARACTERIZATION, CLASSIFICATION OF VOWELS :

- Articulatory configuration : This refers to the tongue hump position (front, back, middle) tongue hump height (high, mid, low).
- Waveform plots : Front vowels show pronounced, high-frequency resonance, ~~back~~ <sup>mid</sup> vowels show balance of energy over broad freq. range and back
- Spectrogram plots : Front vowels show relatively high second, third formant frequency (resonance) whereas mid vowels show well-separated and balanced locations of formants. Back vowels show almost no energy beyond low frequency region with low first and second formant frequencies.

Vowels show mostly low-frequency spectral info



However, the concept of "typical" vowel sound is unreasonable because:-

- vowel pronunciation is variable among people with different regional accents

# This leads to a wide range of variability in first and second formant frequencies for a given vowel ~~son~~ sound.

# In addition to this, different vowel sounds by different speakers may have overlapping formant frequencies.

Hence, measuring formant frequencies or spectral peaks is not enough to accurately classify vowel sounds.

Note - the <sup>primary</sup> frequency at which the vocal chords vibrate is also called the fundamental frequency ( $f_0$ )

Diphthongs : It is a gliding monosyllabic speech sound that starts at or near the articulatory position of one vowel and moves to or towards the position of another vowel.

There are 6 diphthongs in American English:-

- $a^{\gamma}$  : EX buy
- $a^w$  : EX down
- $e^{\gamma}$  : EX bait
- $o^{\gamma}$  : EX boy boy
- $o$  : EX boat
- $ju$  : EX you

PRODUCTION OF DIPHTHONG : Diphthongs are produced by varying the vocal tract smoothly between vowel configurations appropriate to the diphthong. This is highly prominent in  $a^{\gamma}$ ,  $a^w$  etc but weaker for  $e^{\gamma}$

because of closeness in the two vowel sounds.

# Diphthong can also be thought of as a time-varying spectral characteristics, i.e. a plot of values of second formant versus first formant, as a func<sup>n</sup> of time.

∴ diphthong can be characterized by a time-varying vocal tract area function that varies between two vowel configurations. (It can be seen in a spectrogram : graphical plot with time on x axis and ~~amplitude~~ frequency of wave form on the y-axis)

SEMI-VOWELS : The sounds consisting of w, l, r, y are called semi-vowels because of their vowel-like nature. They are ~~not~~ generally characterized by a gliding transition in vocal tract area function between two adjacent phonemes.

Semi vowels are very difficult to characterize as the acoustic characteristics of w, l, r, y are strongly influenced by the context in which they occur.

Hence, semi-vowels are best described as vowel-like, transitional sounds, hence they are similar to both vowels and diphthongs.

NASAL CONSONANTS : The nasal consonants m, n, ŋ are produced with glottal excitation and the vocal tract is totally constricted at some point along the oral passageway.



## PRODUCTION OF NASAL CONSONANTS :

- The velum is lowered so that air flows through the nasal tract, with sound being radiated at the nostrils.
- The oral cavity, though restricted, is still acoustically coupled to the pharynx.
- Mouth serves as resonant cavity that traps acoustic energy at natural frequencies.

for m : constriction is at lips  
n : constriction is just behind the teeth  
ŋ : constriction is just forward of the velum.

WAVEFORMS : The waveforms of m and n are very similar.

SPECTROGRAMS : Nasal consonants show a concentration of low-frequency energy with a mid-range of frequencies that contain no prominent peaks. This is because of the particular combination of resonances and anti-resonances that result from coupling of nasal and oral tracts.

UNVOICED FRICATIVES : The sounds f, θ, s, sh are produced by exciting the vocal tract by a steady air flow, which becomes turbulent in the region of constriction in the vocal tract.

The location of constriction determines the type of sound:-

f : constriction near the lips  
θ : near the teeth  
s : middle of oral tract  
sh : back of the oral tract

Therefore, unvoiced fricatives consists of a source of noise at a constriction, which separates the vocal tract into two cavities.

- Sound is radiated from first cavity
- Back cavity (like in case of nasal consonants) traps energy and introduces antiresonances into vocal output

WAVEFORMS: Unvoiced fricatives have non-periodic waveforms