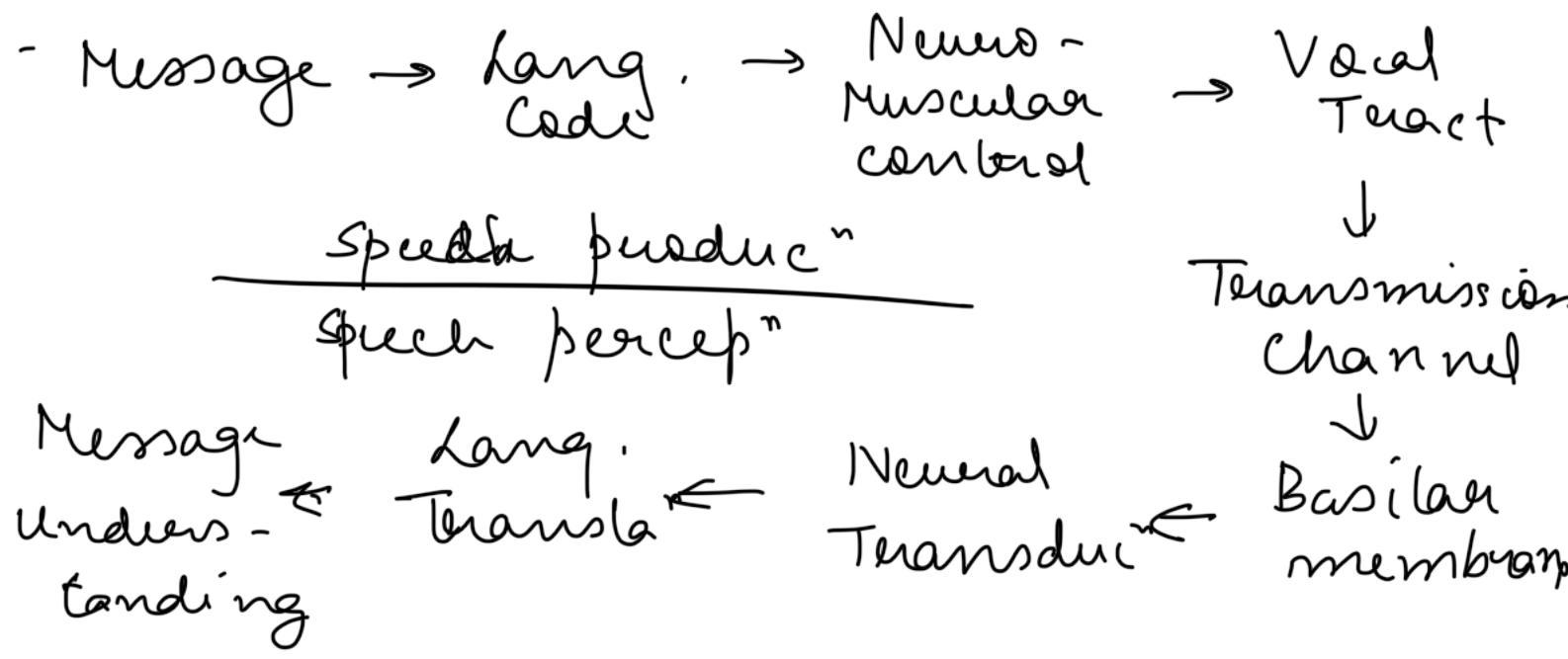


Speech Tech

Week - 1

L1 Course Introduct "

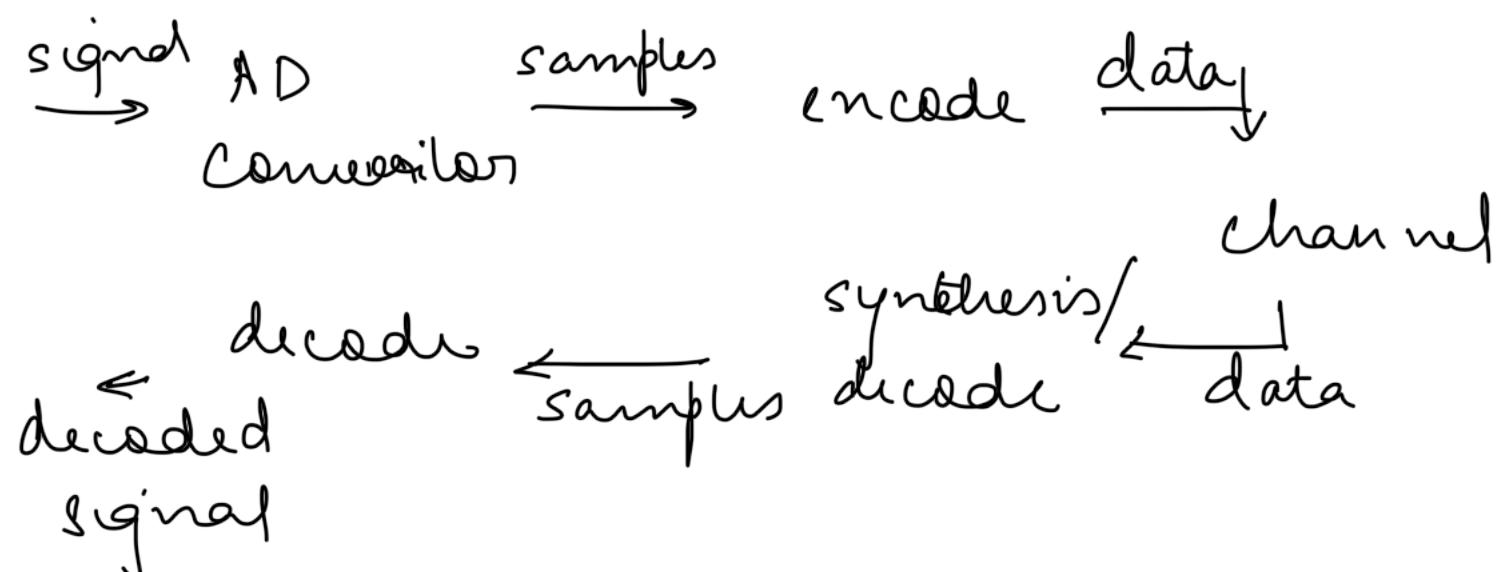
- basics → speech, DSP, stats
- ASR → HMM, DNN, TDNN, LM
- ASR DL → E2E, RNN, Transformers, LLM
- Self supervised models → BERT, GPT, APC, Wave2Vec, HuBERT
- Speaker verification → x, u, s vectors
- TTS tasks
 - primary purpose → communication.
Formulate the msg. in brain, then lang. & sentence forma" is selected.
The mouth articulators appropriate movement & finally spoken audio comes.



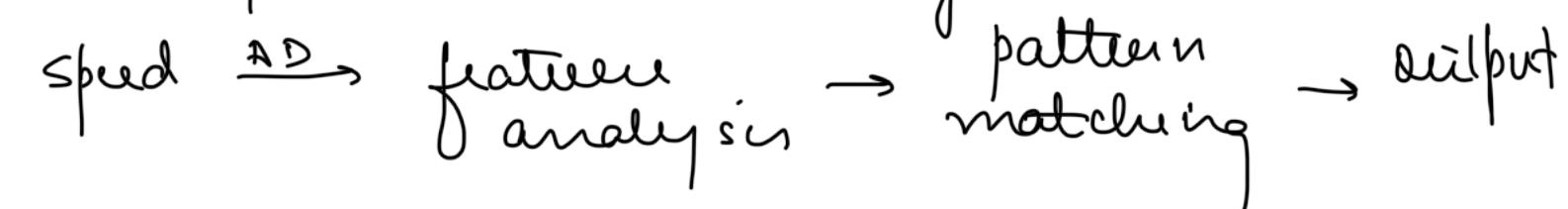
- speaker \rightarrow linguistic, Physiological, Acoustic

- listener \rightarrow Physiological, linguistic

- App. I \rightarrow store & transmit speech



- TTS, Pattern matching



- DSP Application - STT, TTS, Speech synthesis, speech enhancement, Aids for differently abled.

L2 DSP Fundamentals

- Sampling \rightarrow nos. only can be handled. Therefore, A-D (analog to digital conversion). They are equally spaced discrete pts.

- Sampling Th. \rightarrow Under conditions, we can reconstruct the original cont. time signals from samples.
- By Nyquist - Shannon sampling th. Sampling Freq. $\geq 2 \times$ Max. freq. in signal.
- Incorrect sampling freq. \rightarrow loss of info. (aliasing)
sampled signal look like a resultant of a slowly varying signal.
- The discrete sampled signals are then quantised (to a finite precision) to get the digital signals.

Signals

-

 Narrow Band
 \rightarrow telephone
 \rightarrow channel acts as band pass filter (300 to 3.4 kHz)
 \rightarrow sampling freq. = 8 kHz
 \hookrightarrow 8 k samples every second.

-
 Wide Band
 \rightarrow current std.
 \rightarrow HD voice (50 - 7 kHz)
 \rightarrow 16 kHz
 \downarrow
 16 k samples every second.

- Telephone \rightarrow A-law, mu law (64 kbit/s)
- Parameters \rightarrow based on code-excited
linear predictⁿ (mobil)
- Speech
- $$s[n] = \sum_{k=1}^P \alpha_k (s[n-k])$$
- GSM & CDMA \rightarrow 7-13 kbit/s

- Audio - 44.1 kHz & 48 kHz
- MP3 - are lossy & operate at lower bit rates but are still CD quality
- we are sensitive to frequencies \rightarrow e.g. colors (edges & corners) & speech.
- Recall -

$$\hat{i} = \begin{pmatrix} 1 \\ 0 \\ 0 \end{pmatrix} \quad \hat{j} = \begin{pmatrix} 0 \\ 1 \\ 0 \end{pmatrix} \quad \hat{k} = \begin{pmatrix} 0 \\ 0 \\ 1 \end{pmatrix}$$

$$\begin{pmatrix} 3 \\ 4 \\ 2 \end{pmatrix} = 3 \begin{pmatrix} 1 \\ 0 \\ 0 \end{pmatrix} + 4 \begin{pmatrix} 0 \\ 1 \\ 0 \end{pmatrix} + 2 \begin{pmatrix} 0 \\ 0 \\ 1 \end{pmatrix}$$

these are the basis funcⁿ.
Any pt. in 3D can be expressed
as linear combinaⁿ of unit vectors.

- A signal can be expressed as a linear combination of complex sinusoids of diff. freq. \Rightarrow Fourier series

$$f(x) = \sum_{n=-\infty}^{\infty} c_n e^{\frac{inx}{L}}$$

\hookrightarrow transformation from time-variable space to freq. variable space. 2 diff. views of same signal.

- Any periodic signal can be represented as a sum of complex sinusoids with diff. freq. Both conti. & discrete is known.

- CTFS

$$f(n) = \sum_{n=-\infty}^{\infty} c_n e^{\frac{inx}{L}}$$

$$c_n = \frac{1}{2\pi} \int_{-\pi}^{\pi} f(x) e^{-inx}$$

- \rightarrow Input - Conti. & Periodic
- \rightarrow Output - Discrete & Aperiodic

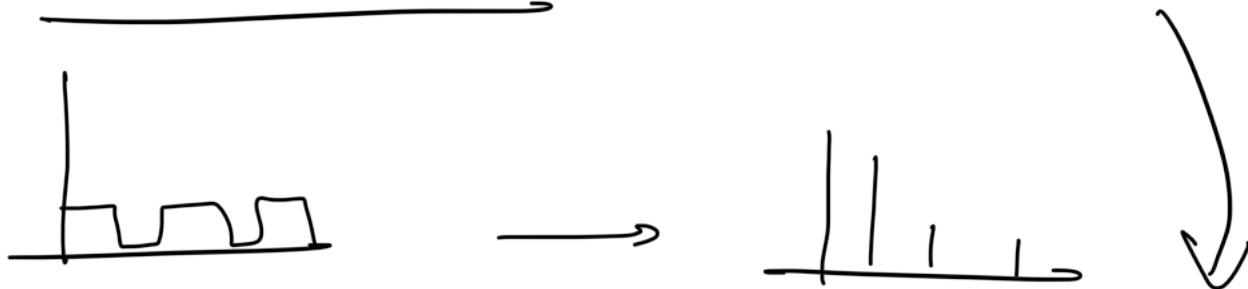
- DTFs

$$f([n]) = \sum_{k=0}^{N-1} c_k e^{i\omega_0 k n}$$

$$x^k = \frac{1}{N} \sum_{n=0}^{N-1} x(n) e^{-\frac{i 2 \pi k n}{N}}$$

Discrete &
Periodic
both sides

L3 Fourier Series



- DTFS leads to DFT which is usually the transform we use when analysing on computers.

- CTFT

$$X(\omega) = \int_{-\infty}^{\infty} x(t) e^{-i \omega t} dt$$

I O → Continuous & Aperiodic

- DTFT

$$X(\omega) = \sum_{n=-\infty}^{\infty} x(n) e^{-i \omega n}$$

I → Discrete & Aperiodic

O → Continuous & Periodic

- DFT

$$X(k) = \sum_{n=0}^{N-1} x(n) e^{-\frac{j 2 \pi k n}{N}}$$

→ discrete & Periodic

→ used in analysis of freq. components in a signal. Used in actual DSP with computers.

- Z-transform → system manipulates the signals. This is a tool for understanding the internal behavior of signals & systems. For discrete, a signal w/o DFT can have Z-transform.

$$X(z) = \sum_{n=-\infty}^{\infty} x(n) z^{-n}$$

- Sys. $x[n] \xrightarrow[\text{sys.}]{\quad} y[n]$

Properties

- linear / non
- time variant / non
- memory / less
- causal / anti - / non -
- stable / unstable

- Linear

$$x_1[n] \xrightarrow{\text{sys}} y_1[n]$$

$$x_2[n] \xrightarrow{\text{sys.}} y_2[n]$$

$$x_1[n] + x_2[n] \xrightarrow{\text{sys.}} y_1[n] + y_2[n]$$

- Time Variant

$$x[n] \xrightarrow{\text{sys.}} y[n]$$

$$x[n-k] \xrightarrow{\text{sys.}} y_2[n-k]$$

- memoryless $\rightarrow y[n]$ depends only on the current sample of $x[n]$.

- causal \rightarrow present & past $x[n]$

anti \rightarrow present & future $x[n]$

non. \rightarrow present, past, future $x[n]$

- Stable $\rightarrow y[n]$ is bounded for $x[n]$
BIBO

- LTI (linear & time variant) sys.

$$y[n] = x[n] * h[n]$$

\downarrow
convoluⁿ operaⁿ

$$y[n] = \sum_{-\infty}^{\infty} x[k] h[n-k]$$

Time domain \rightarrow convolution

Freq. domain \rightarrow multiplication

filters to remove the noise & get the voiced signal in a better way.

- Rational Digital filters

$\xleftarrow{\text{FIR}}$
Finite Impulse Response
 $a_k = 0$

$\xrightarrow{\text{IIR}}$
Infinite Impulse Response
 $b_n = 0$

\rightarrow does not take past outputs into account.

\rightarrow takes past inputs & outputs into account

\rightarrow causal (takes current & past inputs)

- Filter diff. eqⁿ

$$y[n] - \sum_{k=1}^N a_k y[n-k] = \sum_{n=1}^M b_n x[n-k]$$

- Filter transfer funcⁿ

$$H(z) = \frac{Y(z)}{X(z)} = \frac{\sum_{n=0}^M b_n z^{-n}}{1 - \sum_{k=1}^N a_k z^{-k}}$$

$$H(z) = \frac{k}{(s - z_1)(s - z_2)} \frac{(s - p_1)(s - p_2)(s - p_3)}{(s - \bar{p}_1)(s - \bar{p}_2)(s - \bar{p}_3)}$$

$z_1, z_2 \rightarrow$ zeroes
poles of dev. \rightarrow poles of filter

stable \rightarrow ROC must be within unit circle.

Causal \rightarrow ROC must be outside of the outermost pole.

stable & causal \rightarrow all poles must be within unit circle.

Speech Production

- Lungs \rightarrow Trachea \rightarrow Glottis \rightarrow Nasal cavity +
 exciteⁿ air open / close mouth +
 source flow quasi- pharynx
 path periodically

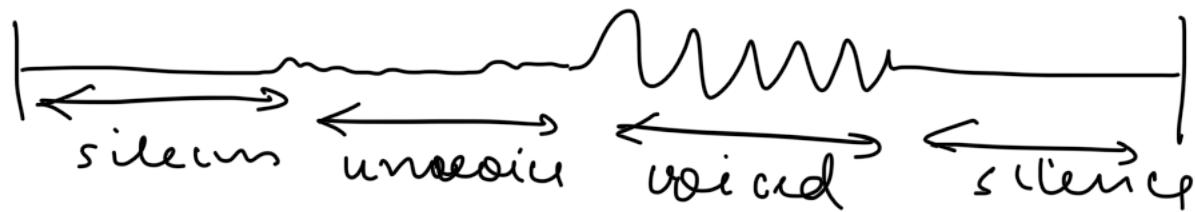
- time period of vibraⁿ cycle depends on the height & weight.

$T_{men} > T_{women} > T_{children}$

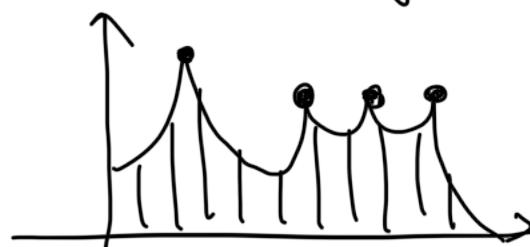
- different mouth configuraⁿs leads to diff voices (vowels)

- all spoken vowels are voiced. Phonemes are like alphabets for speech. They are the basic units from which speech is formed.
- Sound articulators → glide, nasal, stop, fricative, voicing, mixed source, whispered.
labial, labiodental, dental, alveolar, palatal, velar, pharyngeal

L6 Waveforms



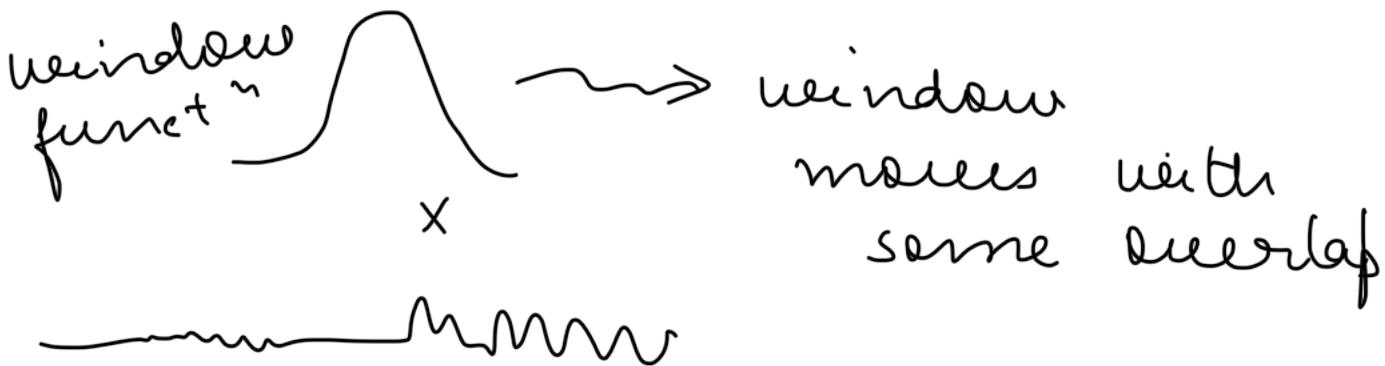
- Speech production model → mathematical equivalent of speech, also called as terminal analogies model



- formants
marks the sound

| vertical lines
↓
pitch harmonics

- short time fourier transform (STFT)



↓
Take fourier
transform

spectrogram → log magnitude of the resultant

wideband

→ small time
windows

→ good time resoluⁿ

→ captures formants

narrowband

→ larger time
windows

→ poor

→ good freq.
resoluⁿ- pitch
harmonics