

Pretext tasks

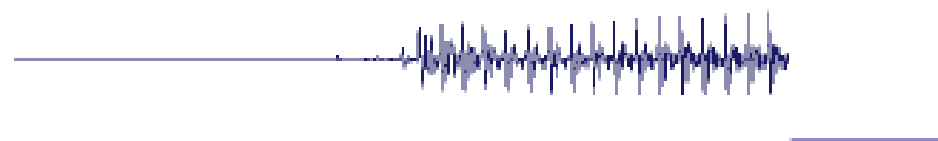
4. CPC

|| AR / LPC / VAR / RNN / APC / CPC ||

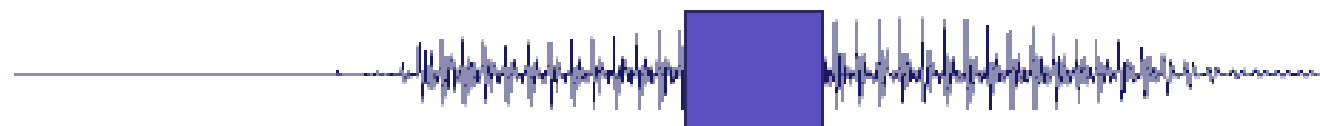
1	2	SELF-PREDICTION	INNATE RELATIONSHIP (Context-based)	1. ROTATION 2. RELATIVE POSITION	IMAGE
3		CONTRASTIVE LEARNING	INTER-SAMPLE CLASSIFICATION	1. Instance Discrimination 2. SimCLR [Contrastive Loss] 3. Theory – Guarantees / Bounds	IMAGE
4		CONTRASTIVE LEARNING	INTER-SAMPLE CLASSIFICATION	Contrastive Predictive Coding (CPC), [NCE, InfoNCE Loss]	AUDIO/ SPEECH
5		SELF-PREDICTION	GENERATIVE (VAE)	1. AE – Variational Bayes 2. VQ-VAE + AR	IMAGE AUDIO/ SPEECH
6		SELF-PREDICTION	GENERATIVE (AR)	1. AR-LM – GPT 2. Masked-LM – BERT	LANGUAGE
7		SELF-PREDICTION	MASKED-GEN (Masked LM for ASR)	1. Wav2Vec / 2.0 2. HuBERT	AUDIO/ SPEECH

Learning with or without supervision – speech and audio

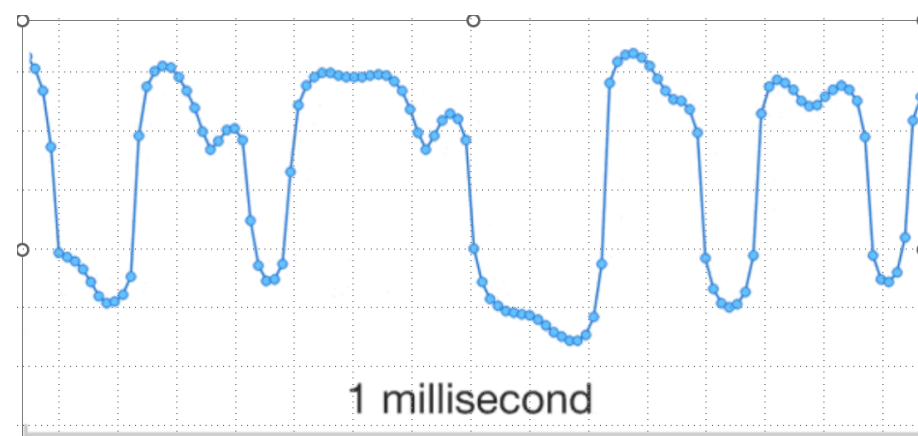
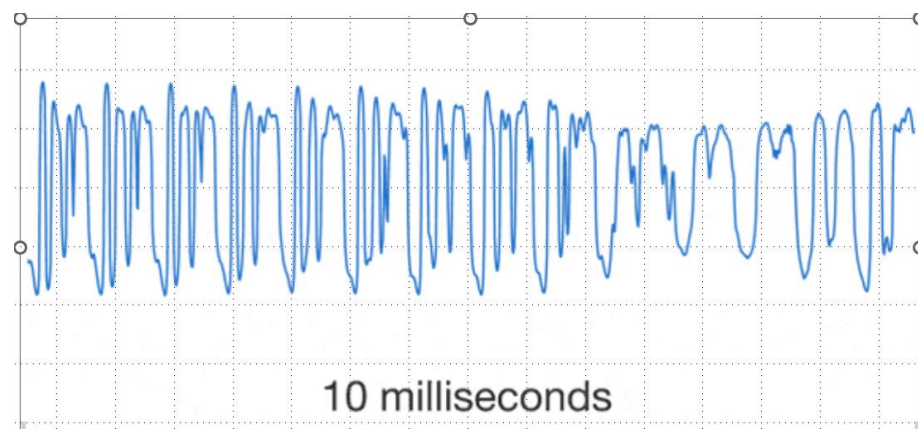
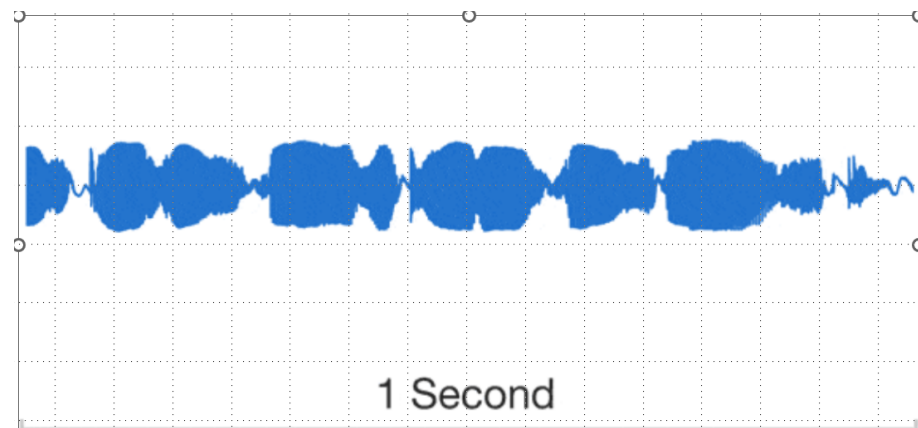
- Next frame prediction



- Masked prediction

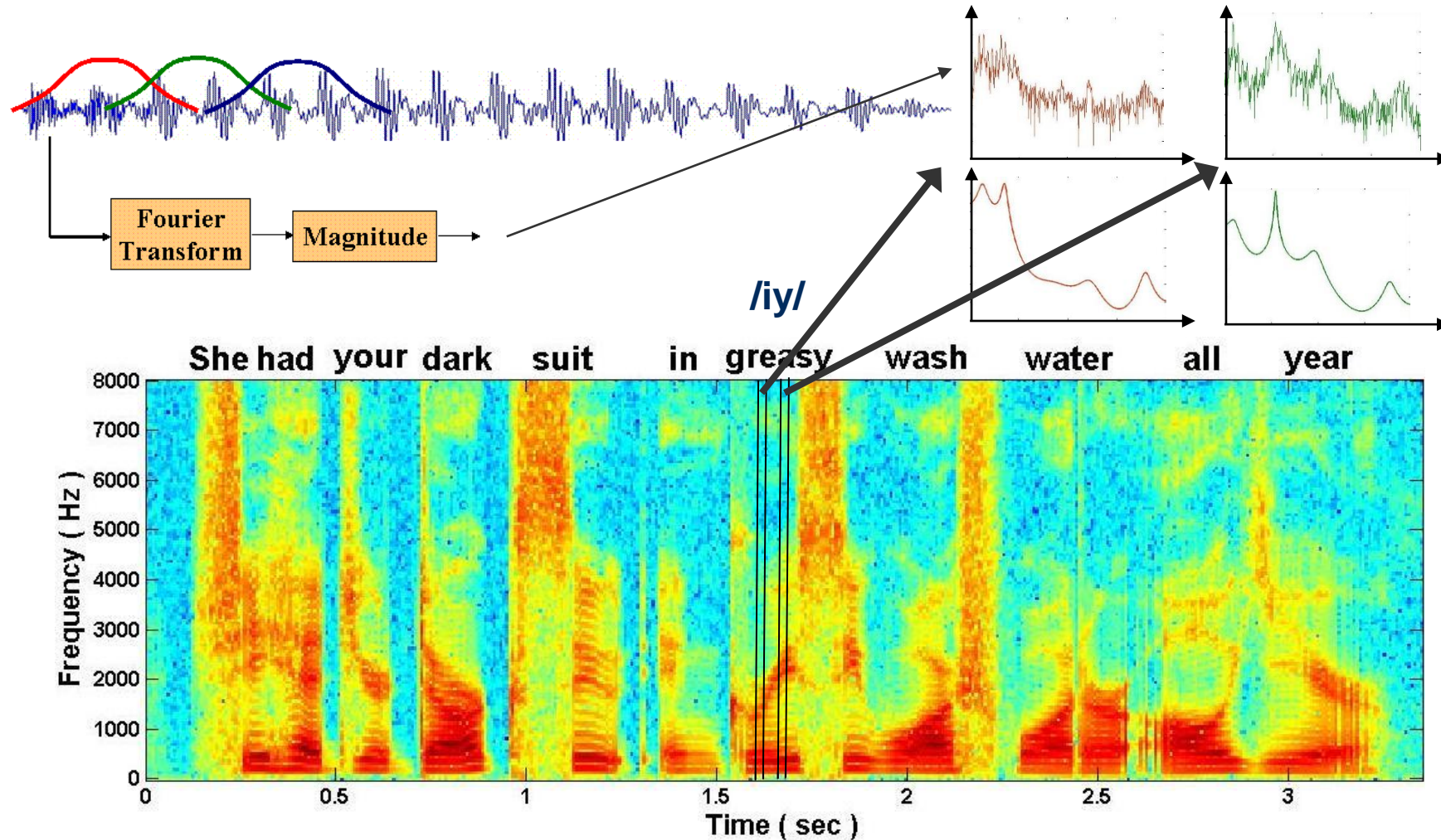


Auto-regressive Model and Speech Spectrum

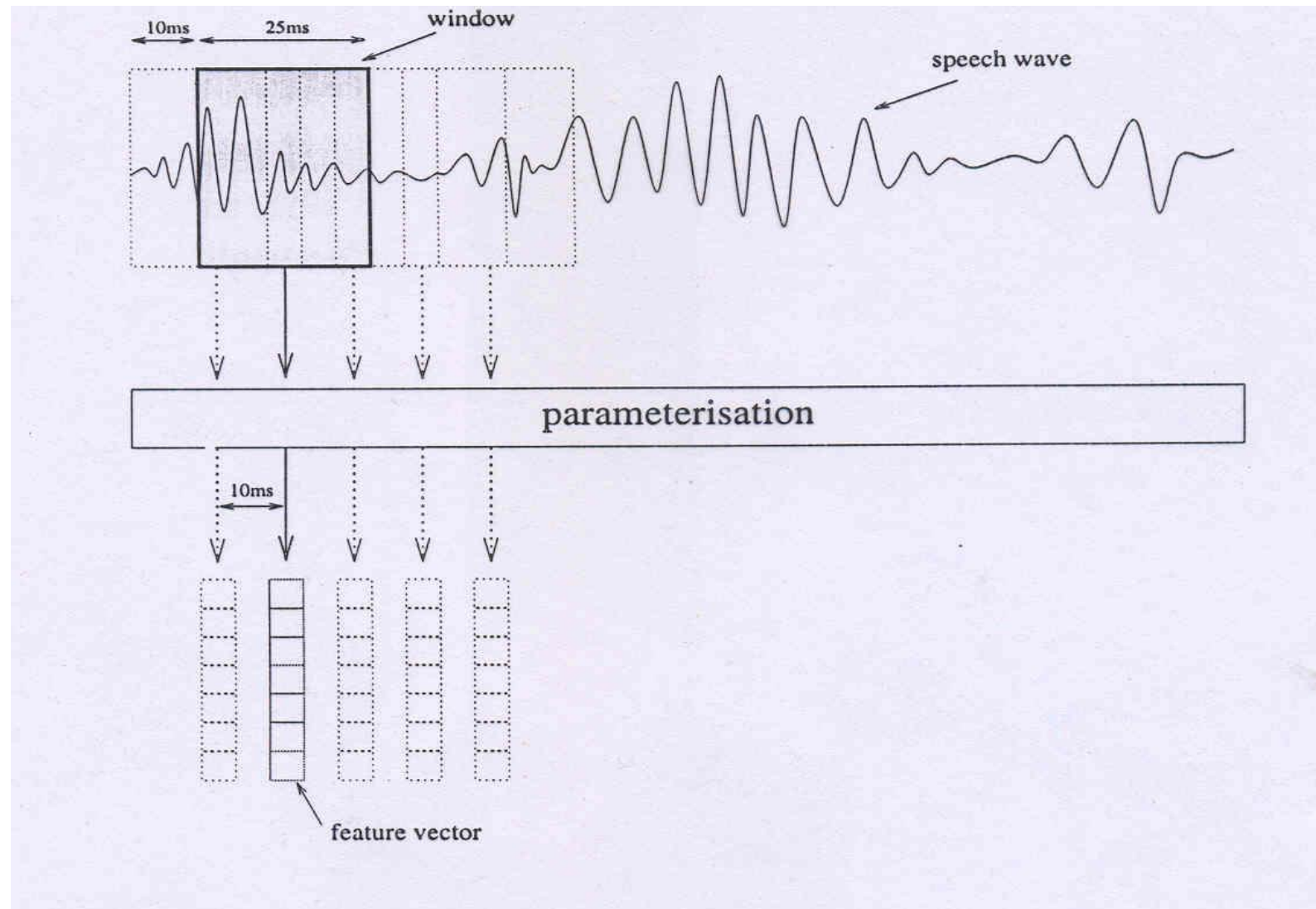


Spectrogram

- Speech is a continuous evolution of the vocal tract
- Spectrogram shows time-frequency evolution
- Represented as a time-series of short-time spectra



Short-time Analysis and Parameterization



AR / LPC / VAR / RNN / APC / CPC

①

AR — Auto-Regressive

LPC — Linear Predictive Coding (Analysis)

VAR — Vector Auto-Regression

RNN — Recurrent Neural Networks

APC — Auto-regressive Predictive Coding

CPC — Contrastive Predictive Coding

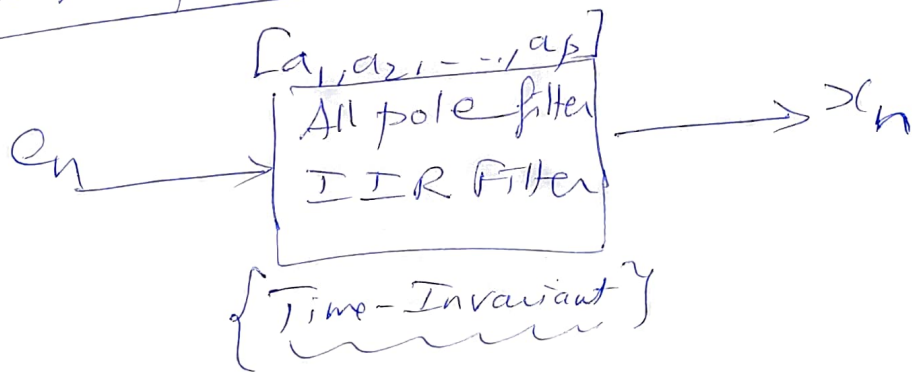
AR(p) : Auto-regressive model of order 'p' (2)

$$x_n = \sum_{i=1}^p a_i x_{n-i} + e_n \quad \text{--- (1)}$$

$\{a_1, a_2, \dots, a_i, \dots, a_p\}$: Parameters of the model

e_n : white noise

x_n as an AR process

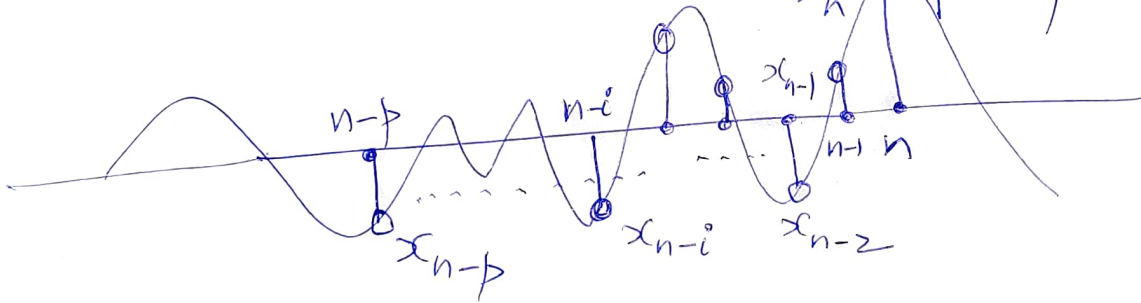


$$x_n = \underbrace{\sum_{i=1}^p a_i x_{n-i}}_{\hat{x}_n} + e_n \quad - (1) \quad (3)$$

$$\hat{x}_n = \sum_{i=1}^p c_i x_{n-i}$$

\hat{x}_n \times \downarrow e_n
 \uparrow
 x_n

Predict x_n
 as \hat{x}_n
 as a linear
 combination of
 past 'p' samples.



$$\hat{x}_n = a_1 x_{n-1} + a_2 x_{n-2} + \dots + a_i x_{n-i} + \dots + a_p x_{n-p}$$

④

$$x_n = \hat{x}_n + e_n$$

$$\hat{x}_n = \sum_{i=1}^p a_i x_{n-i}$$

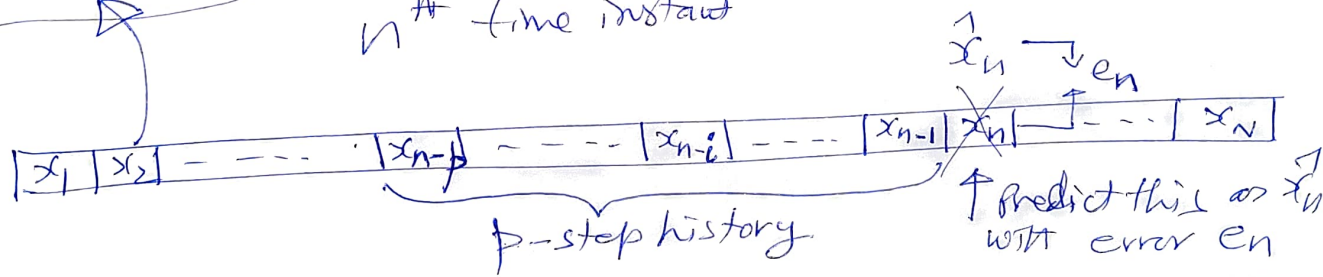
or

$$e_n = x_n - \hat{x}_n$$

True Signal
Value at
 n^{th} time instant

Predicted Signal
Value at
 n^{th} time instant

Frame/Record of
 N -Samples e.g. 10ms
of speech



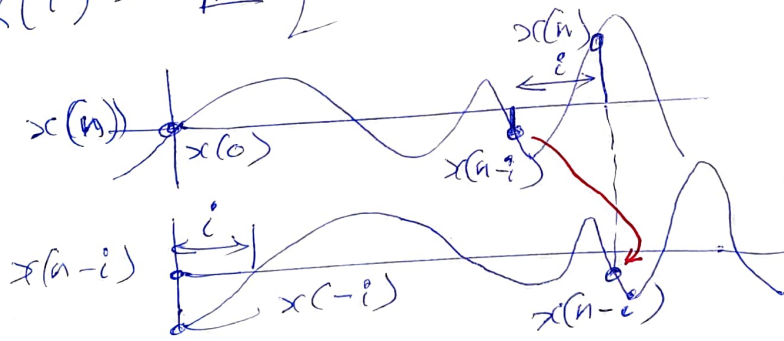
Optimization / Solution for $\{a_i\}_{i=1}^P, \left\{ \begin{matrix} \text{PREDICTION} \\ \text{CO-EFFICIENTS} \end{matrix} \right\}$ (5)

Minimize $E[e^2(n)]$ over a range of values

$$\Rightarrow \sum_{i=1}^P a_i R(j-i) = R(j) \quad \left| \begin{array}{l} \text{Normal} \\ \text{Equations} \\ \text{or} \\ \text{Yule-Walker} \\ \text{Equations} \end{array} \right. \quad j=1, \dots, P$$

$R(i)$: Auto Correlation Coeff of $x(n)$ for lag ' i '

$$R(i) = E[x(n) x(n-i)]$$



→ Right shifted by ' i ' samples

Solution for $\{a_i\}_{i=1}^p$

⑥

$$\sum_{i=1}^p a_i R(j-i) = R(j), \quad \underline{1 \leq j \leq p}$$

Matrix form

$$[a_1, a_2, \dots, a_{p-1}, a_p]^T \rightarrow [p \times 1]$$

Symmetric, Toeplitz

$$R \quad A = \gamma \rightarrow \begin{matrix} 1 \\ 2 \\ \vdots \\ j \\ \vdots \\ p \end{matrix} \begin{bmatrix} R(1) \\ R(2) \\ \vdots \\ R(j) \\ \vdots \\ R(p) \end{bmatrix}$$

$$\begin{matrix} 1 & 2 & \dots & j & \dots & p \\ \begin{bmatrix} 1 \\ 2 \\ \vdots \\ i \\ \vdots \\ p \end{bmatrix} & & & & & \end{bmatrix} \begin{matrix} [p \times p] & [p \times 1] & [p \times 1] \end{matrix}$$

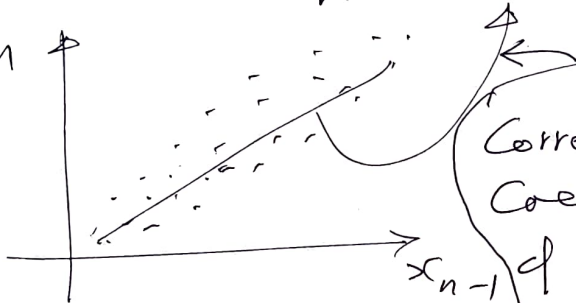
$\gamma_{ij} = R(i-j)$

Order 1 AR: AR(1)

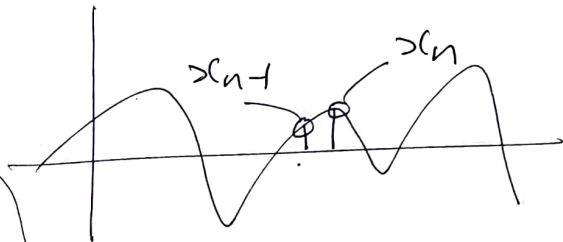
7

$$\hat{x}_n = a_1 x_{n-1}$$

x_n



Correlation Coefficient of lag 1



AR(0): Noise

$$x_n = \sum_{i=1}^p a_i x_{n-i} + e_n$$

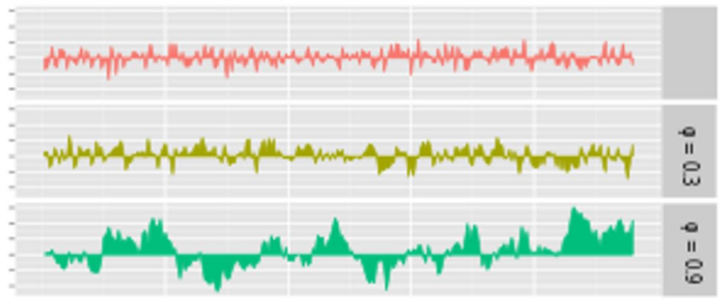
$$AR(1): x_n = a_1 x_{n-1} + e_n$$

For Small $a_1 = .3$

$$\Rightarrow x_n = .3 x_{n-1} + e_n \quad \left\| \begin{array}{l} x_n \text{ looks} \\ \text{like } e_n \end{array} \right.$$

$$a_1 \rightarrow 1 \text{ eg } a_1 = .9$$

$$x_n = .9 x_{n-1} + e_n \quad \left\| \begin{array}{l} \text{"smoothing"} \\ \text{of output (Low Pass} \\ \text{filter)} \end{array} \right.$$



$\Rightarrow \text{AR}(1) \Rightarrow a_1 \Rightarrow \text{Spectral Property}$

(8)

e.g.
$$S(f) = \frac{\sigma^2}{1 + a_1^2 + 2a_1 \cos 2\pi f}$$

In general for p^{th} order AR(p)

$$[a_1 \ a_2 \ \dots \ a_p] \Leftrightarrow [p_1 \ p_2 \ \dots \ p_p] \begin{array}{l} \text{Auto Correlation} \\ \text{Coefficients} \end{array}$$

Fourier
Transform
Pair

Power Spectral Density

$d[x_1 \ x_2 \ \dots \ x_N]$
e.g. 10ms or short-time frame/window
of speech/audio.

More precisely in z-transform

⑨

$$E(z) \left[\frac{1}{1 - a_1 z^{-1} - \dots - a_p z^{-p}} \right] = X(z)$$

$$\Rightarrow e(n) * h(n) = x(n)$$

↳ impulse response of filter.

$$E(z) \cdot H(z) = X(z)$$

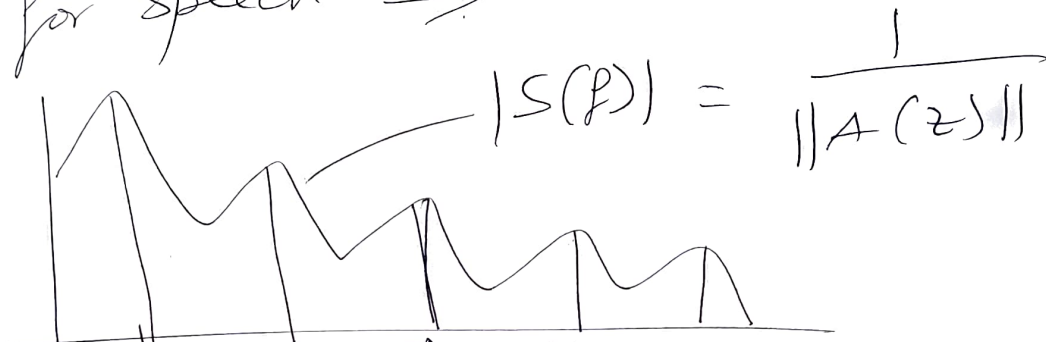
$$H(z) = \frac{1}{1 - a_1 z^{-1} - \dots - a_p z^{-p}} = \frac{1}{\prod_{i=1}^p (z - z_i)}$$

z_i : Roots of $A(z) = 0$

Spectral Envelope

(10)

$p=10$ for speech \Rightarrow 5 Resonances

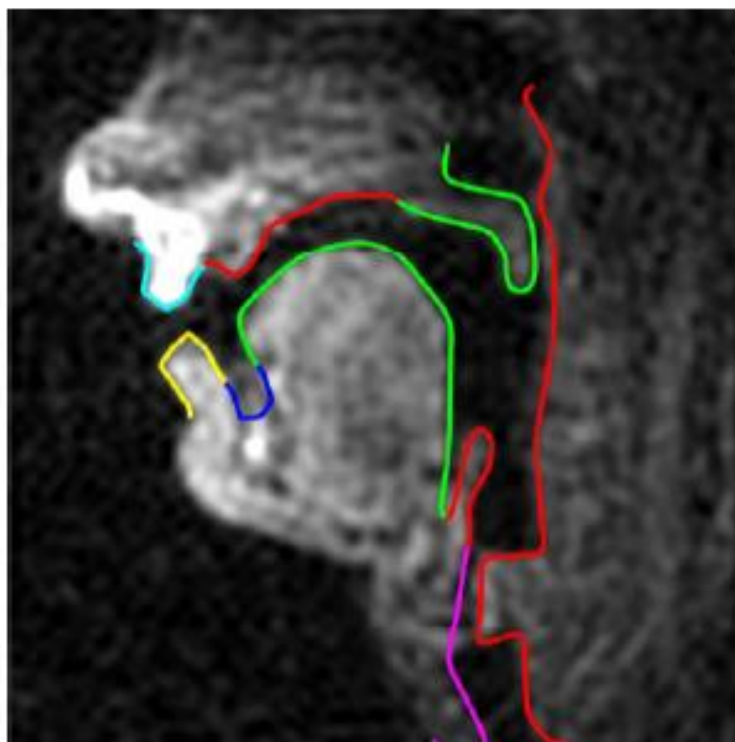


~~$\phi_1, \phi_2, \dots, \phi_N = 1/\sqrt{1-A^2}$~~

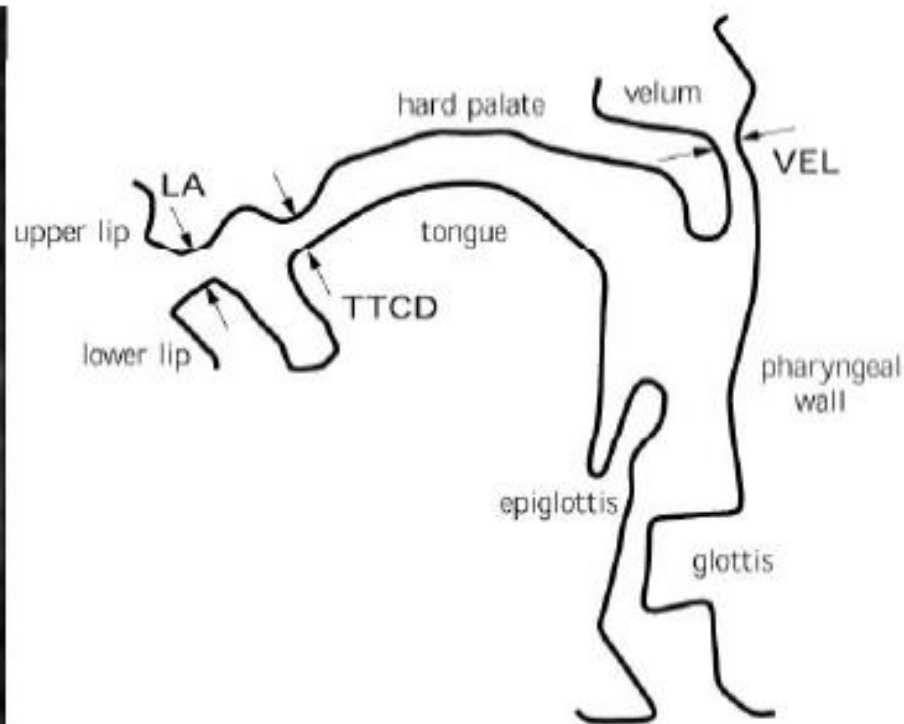
$[a_1, a_2, \dots, a_p] \Rightarrow A(z) \Rightarrow S(f)$

\hookrightarrow Feature Vector of $[x_1, x_2, \dots, x_N]$

MRI of Speech (Prof. Shri Narayanan, USC)



(a)



(b)

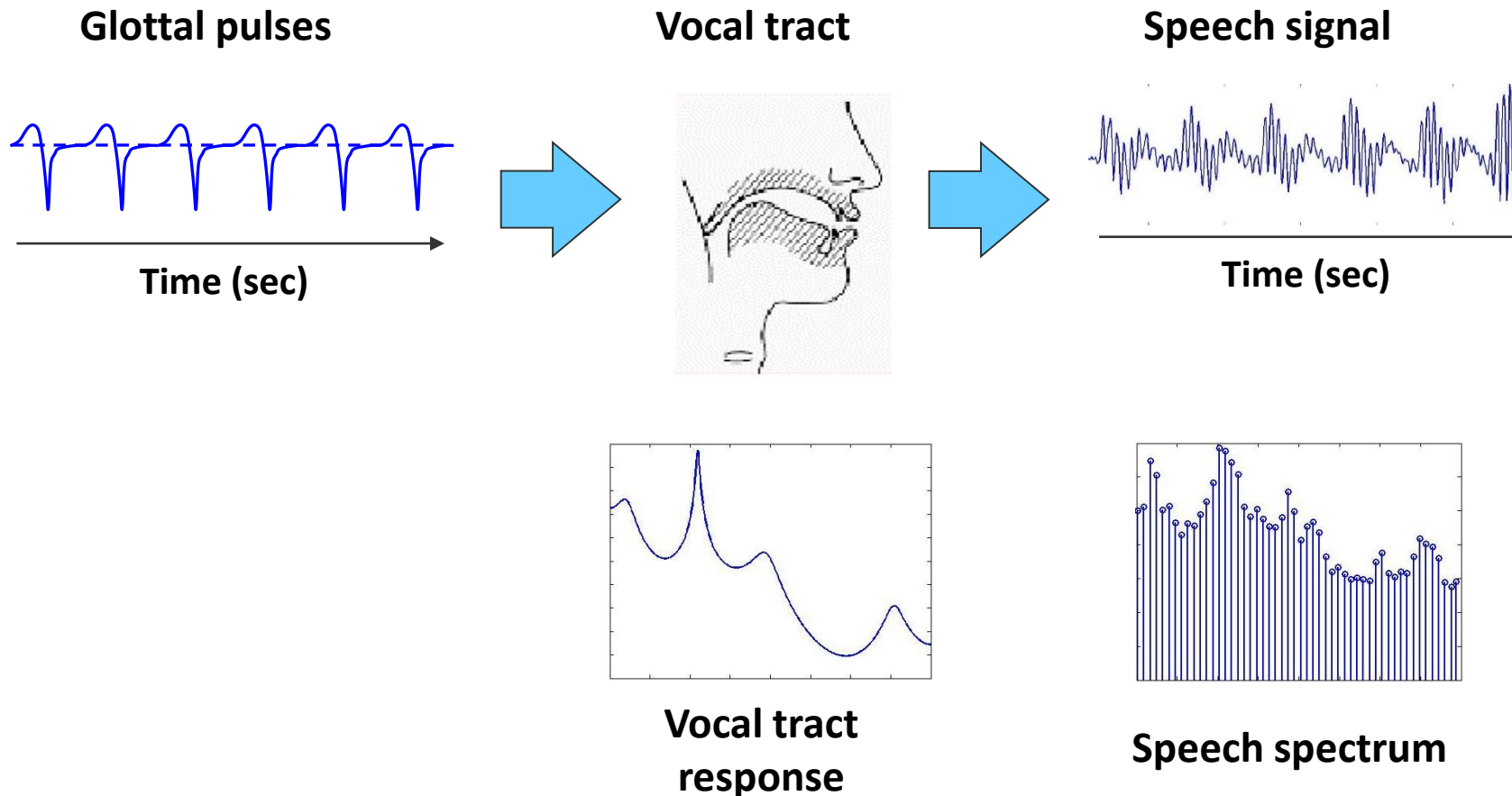
USC

SPAN



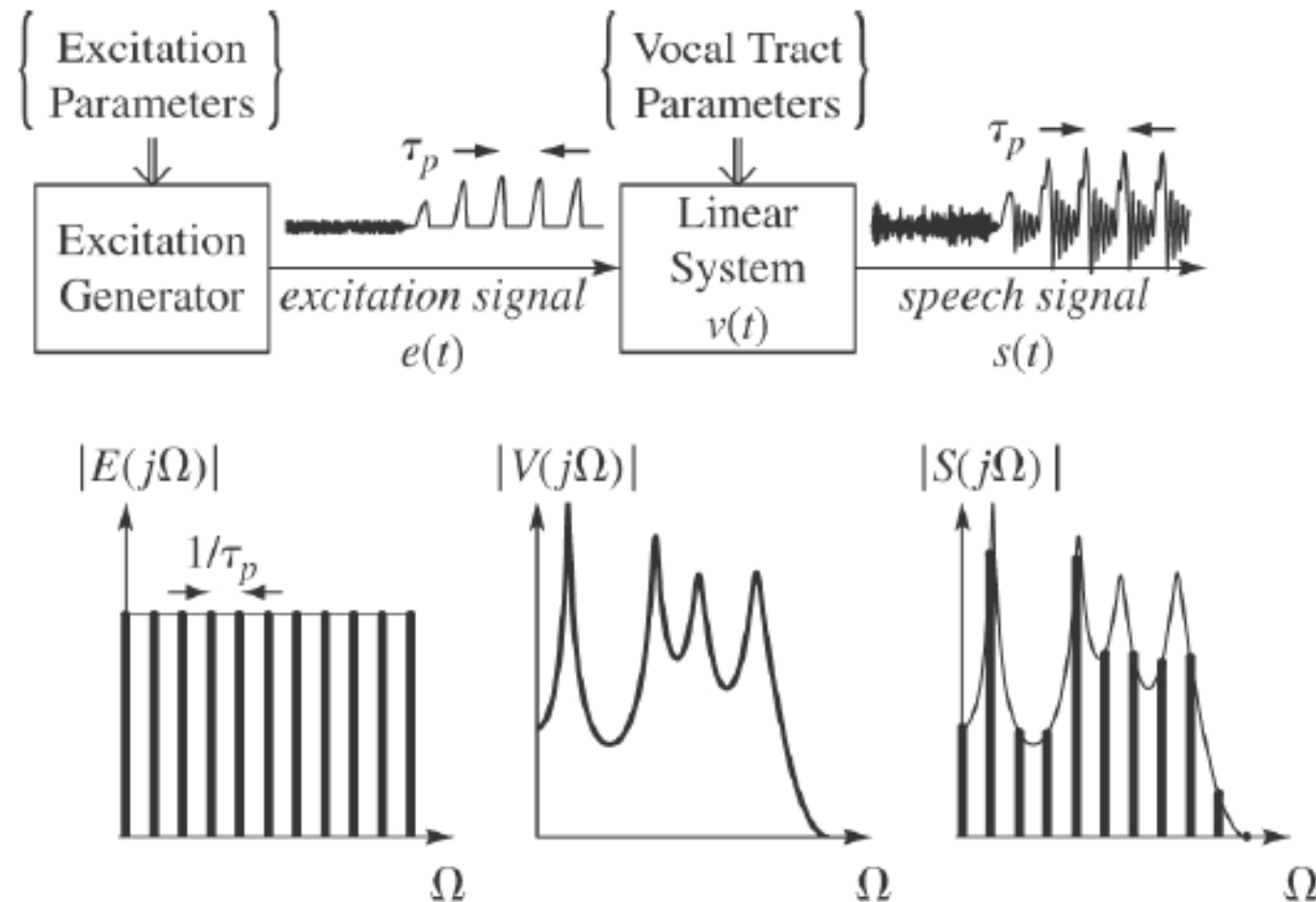
Source-Filter Model

- Features based on speech production model: Source-filter interaction
 - Anatomical structure (vocal tract / glottis) conveyed in speech spectrum

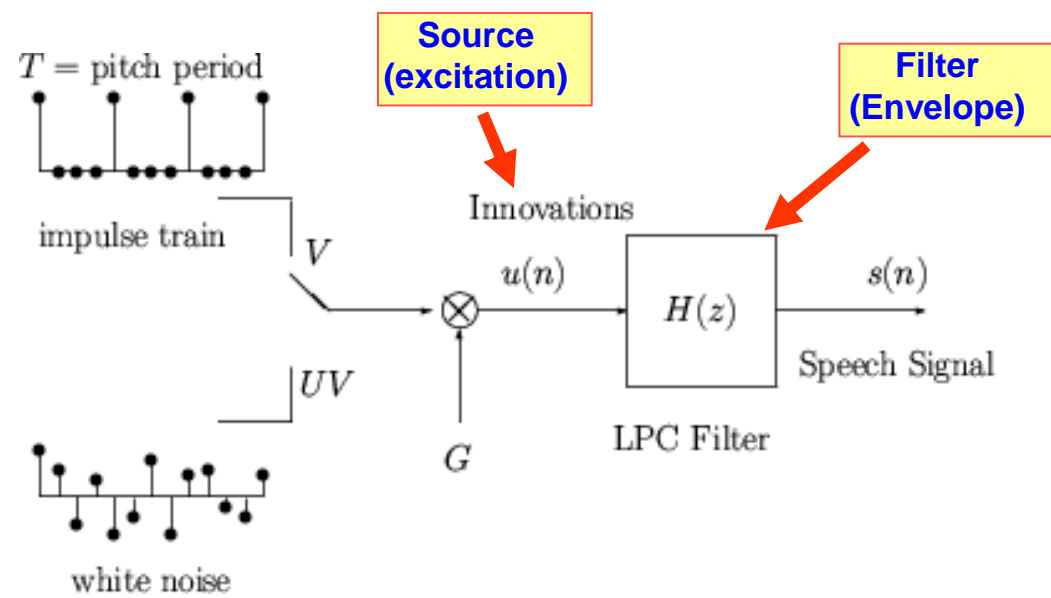


To Quatieri and Rab – Slides ➔

Source-System Model of Speech Production

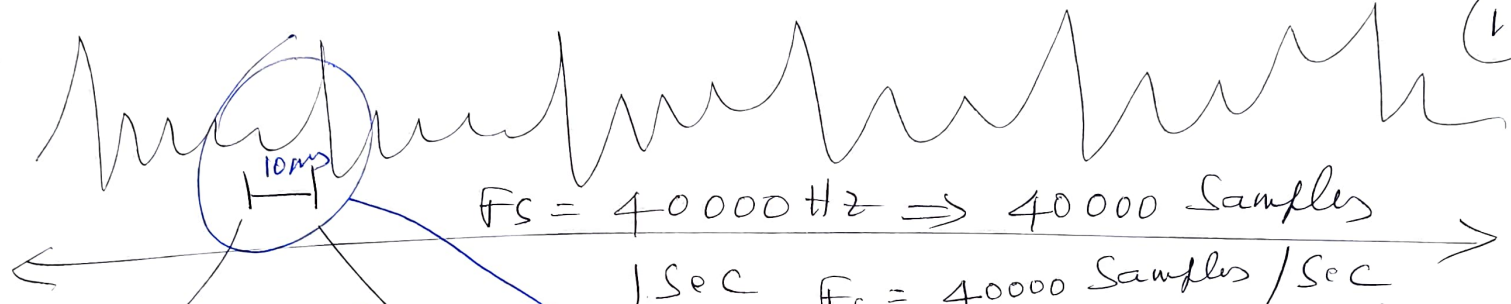


Linear Prediction based Speech Production Model



$$\frac{1}{A(z)} = \frac{1}{1 - \sum_{i=1}^p a_i z^{-i}}$$

- | | | |
|------------------------------|---------------------------|----------------------|
| •Vocal Tract | • \longleftrightarrow • | $H(z)$ (LPC Filter) |
| •Air | • \longleftrightarrow • | $u(n)$ (Innovations) |
| •Vocal Cord Vibration | • \longleftrightarrow • | V (voiced) |
| •Vocal Cord Vibration Period | • \longleftrightarrow • | T (pitch period) |
| •Fricatives and Plosives | • \longleftrightarrow • | UV (unvoiced) |
| •Air Volume | • \longleftrightarrow • | G (gain) |



$$F_s = 40000 \text{ Hz} \Rightarrow 40000 \text{ Samples}$$

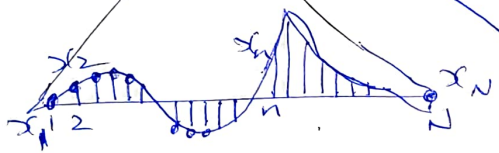
$$1 \text{ Sec } F_s = 40000 \text{ Samples / Sec}$$

$$1 \text{ sec} \approx 40000 \text{ Samples}$$

$$1000 \text{ ms} = 40000 \text{ samples}$$

$$1 \text{ ms} = 40 \text{ samples}$$

$$10 \text{ ms} = 400 \text{ samples}$$



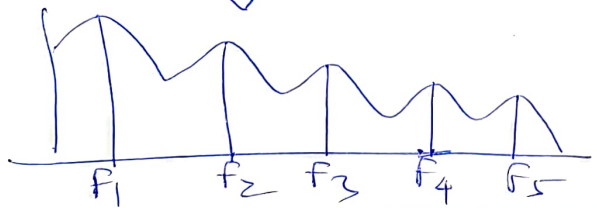
400 Samples

$$[x_1 x_2 \dots x_n \dots x_N]$$

Linear Prediction

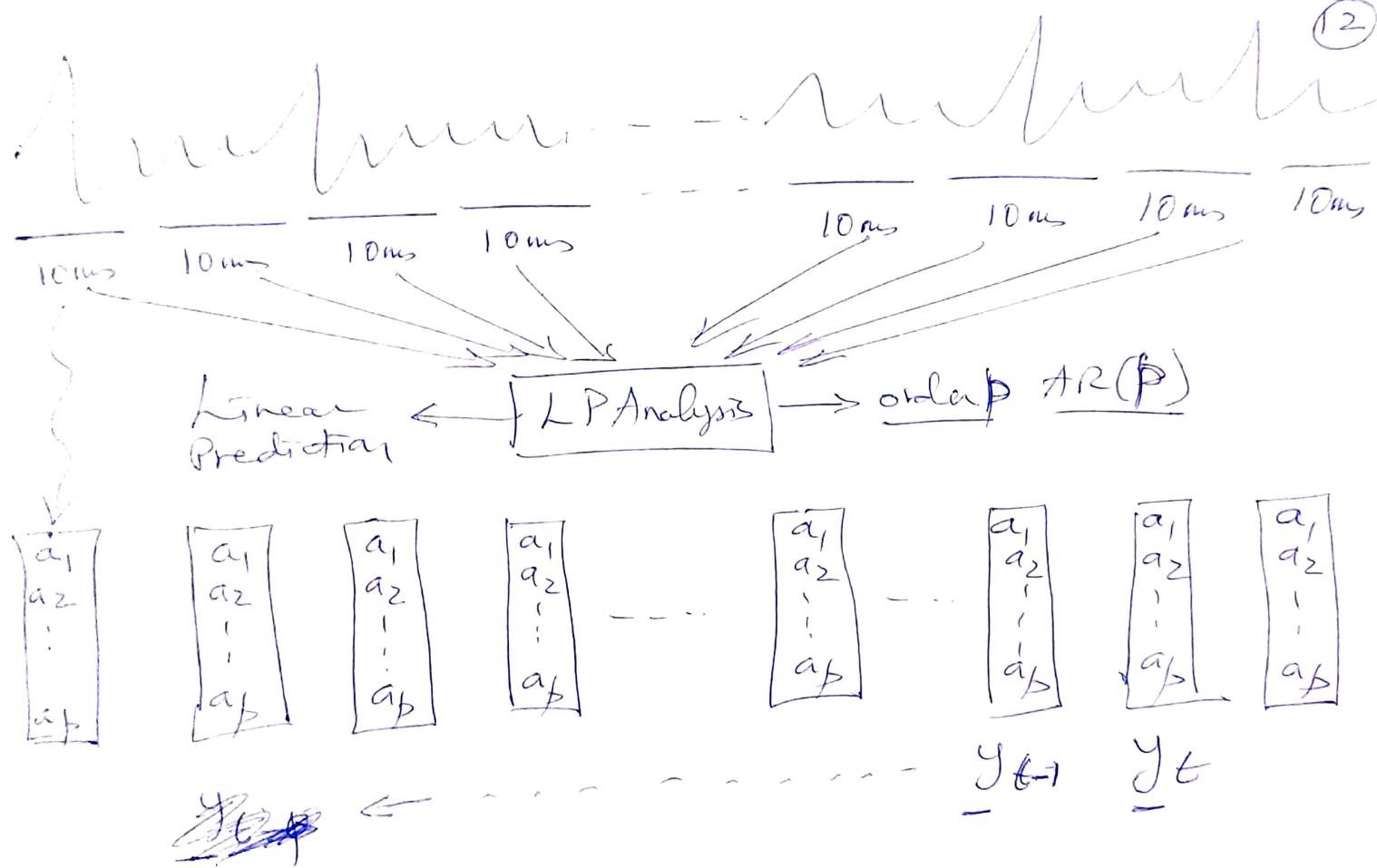
LP Analysis $p=10$

$$[a_1 a_2 \dots a_p]$$



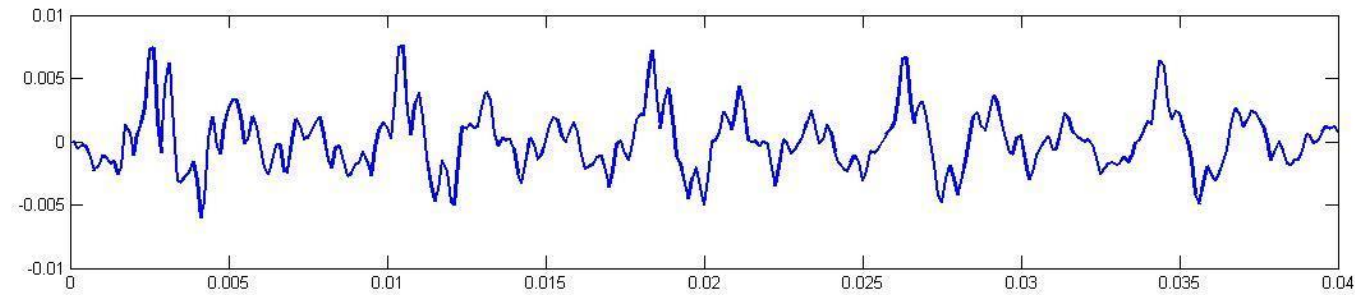
$p\text{-dim} \in \mathbb{R}^p$
feature representation
of $[x_1 x_2 \dots x_N]$

10ms short time frame
of speech.

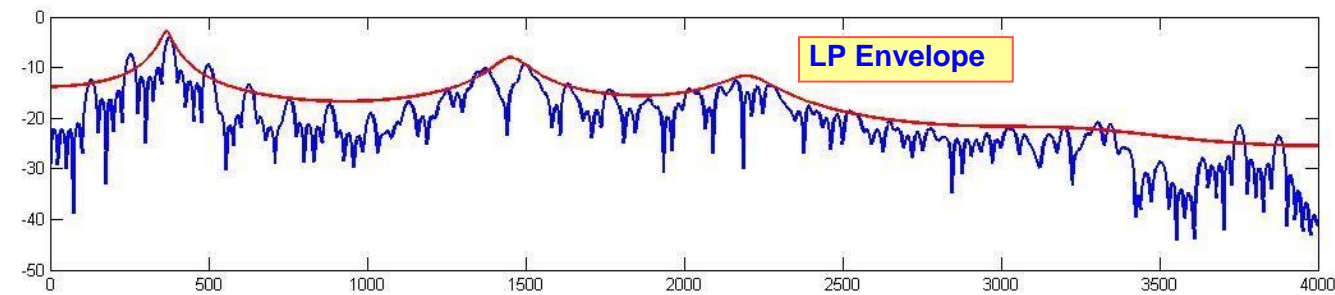


LP Analysis: Envelope (Filter) & Excitation (Source)

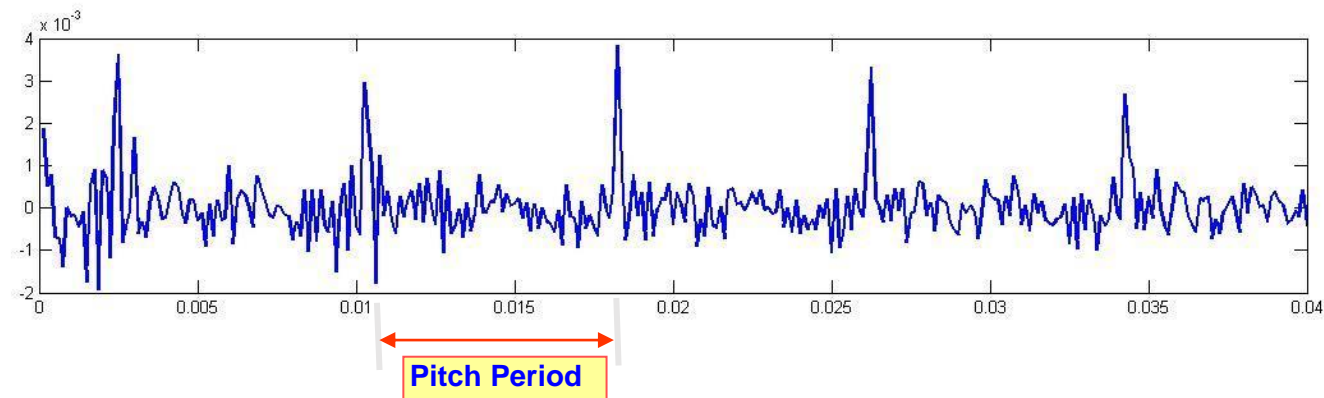
Speech Signal $S(n)$



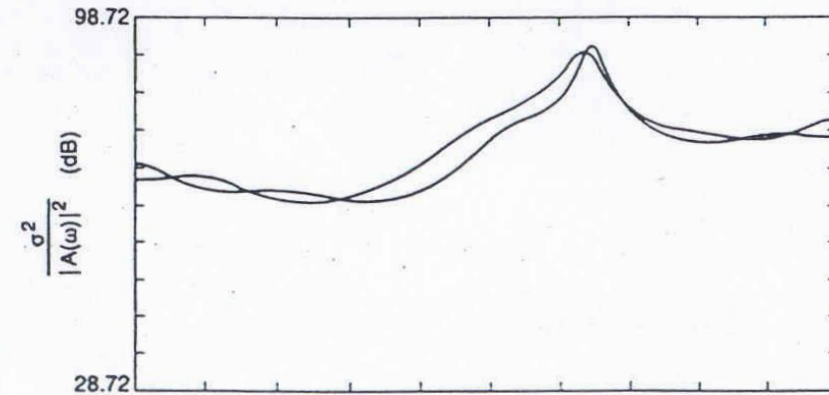
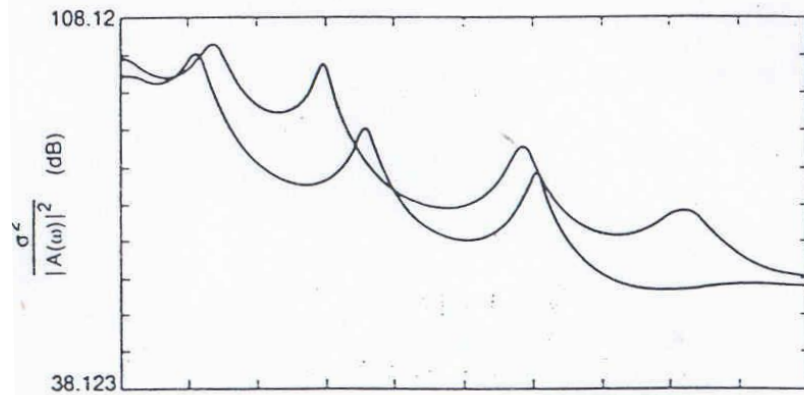
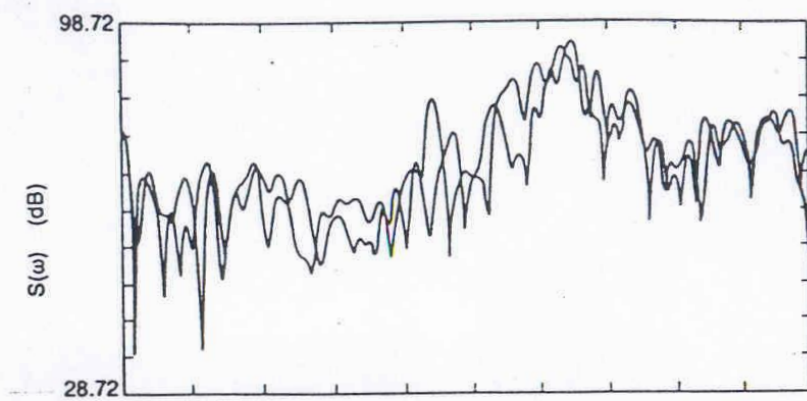
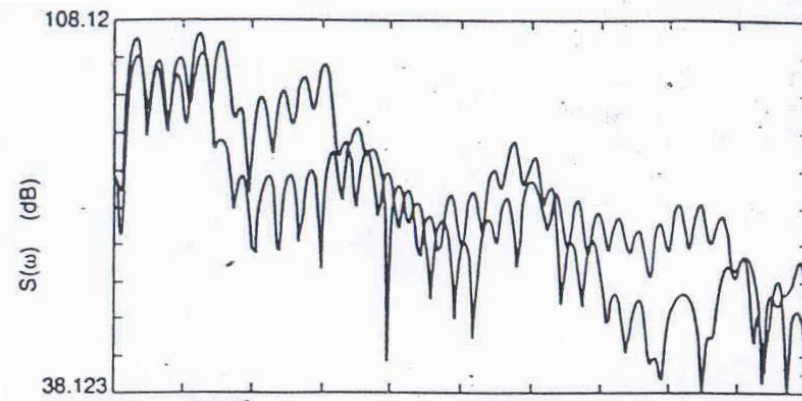
$S(w)$ with LP
spectral envelope
superimposed



Excitation Signal
 $E(n)$

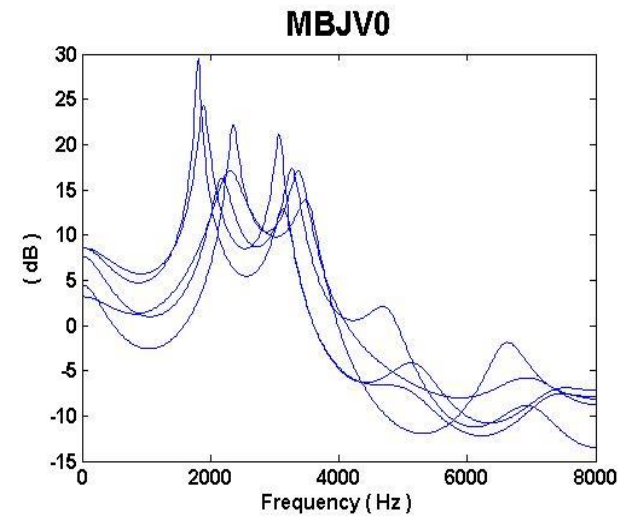
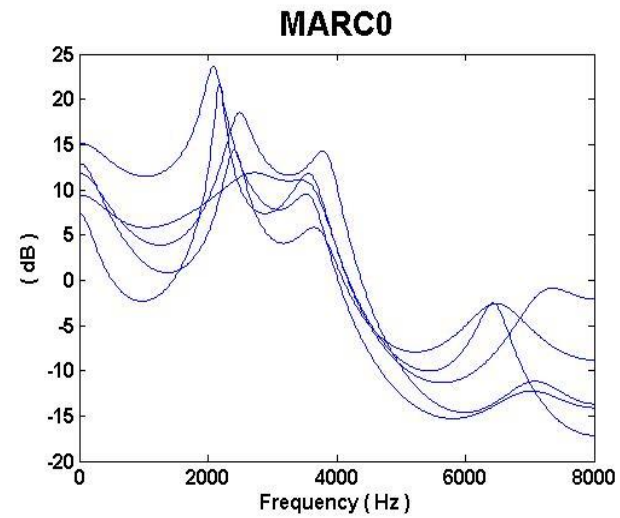
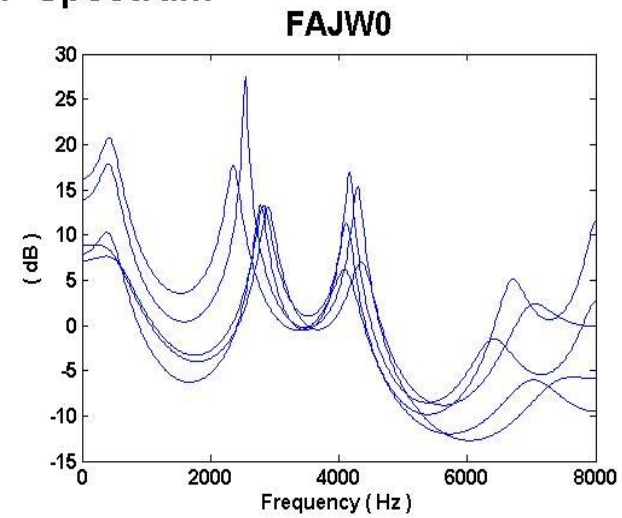
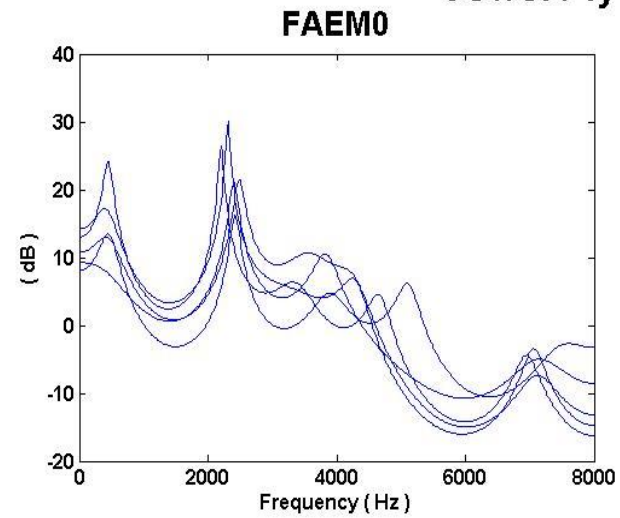


Spectral slices

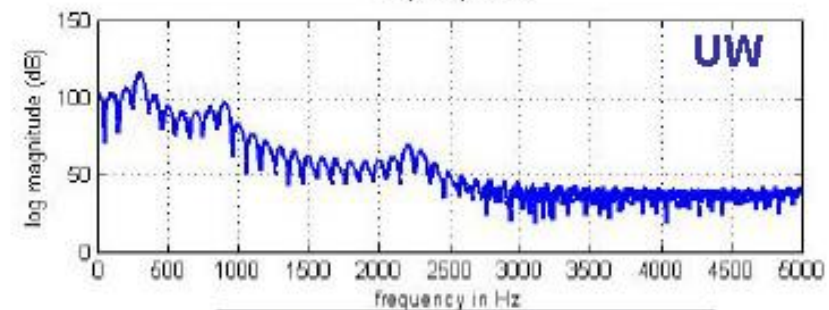
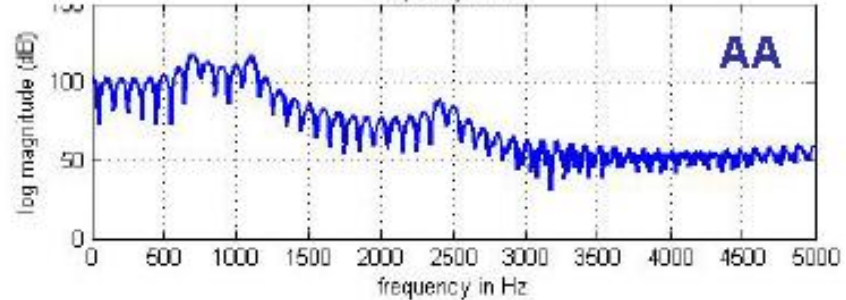
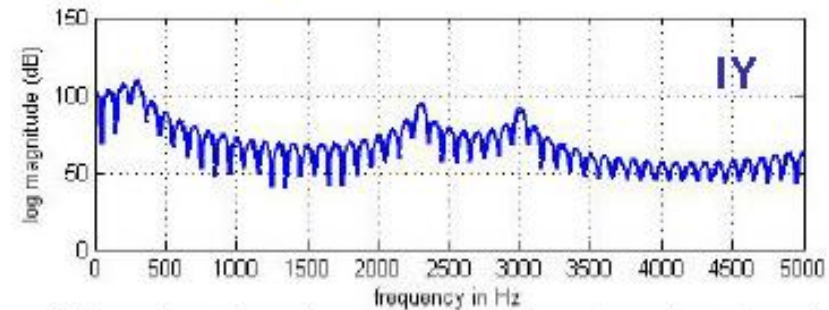
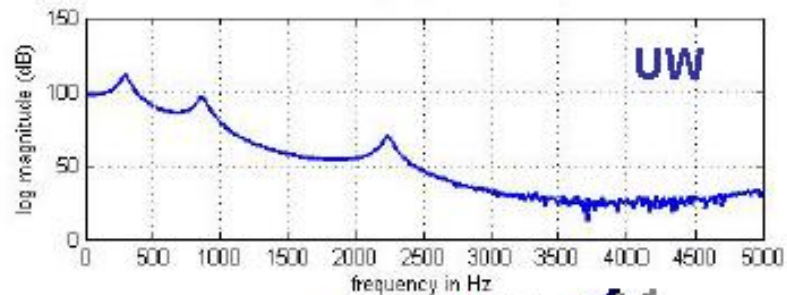
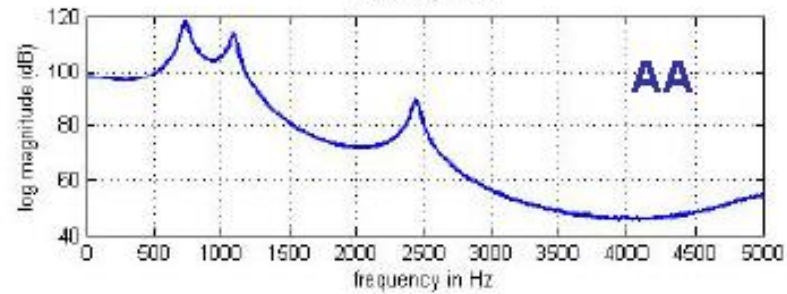
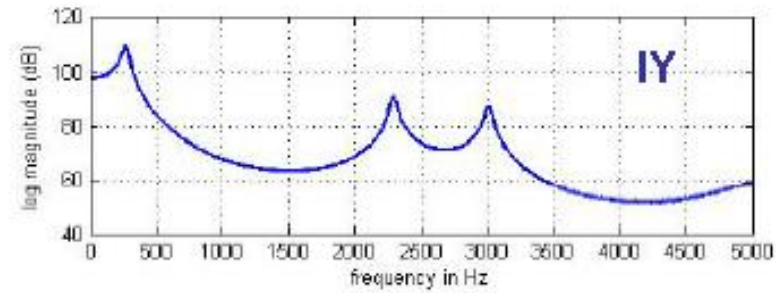


Spectral Envelopes

Vowel / iy / LP Spectrum



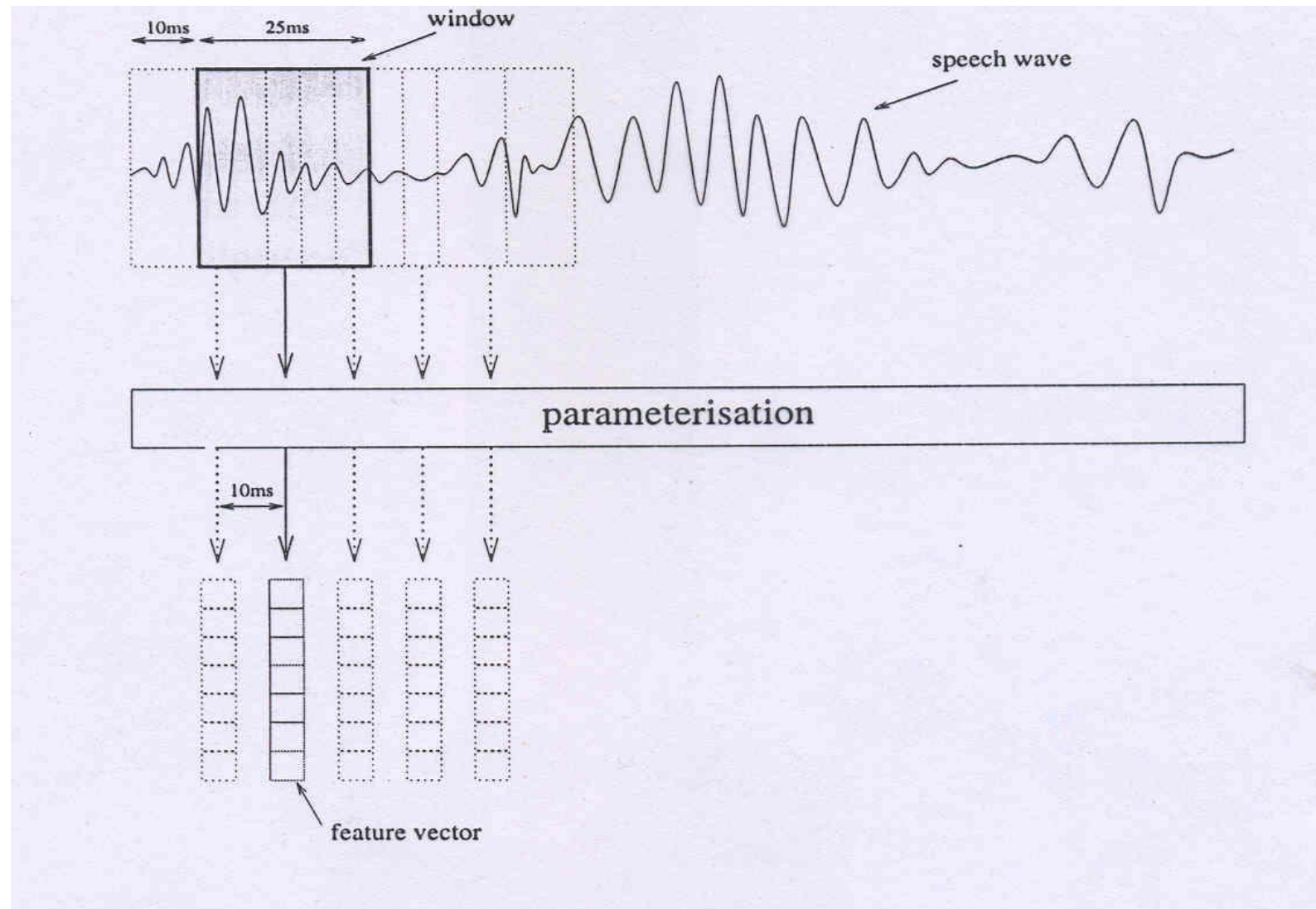
Canonic Vowel Spectra



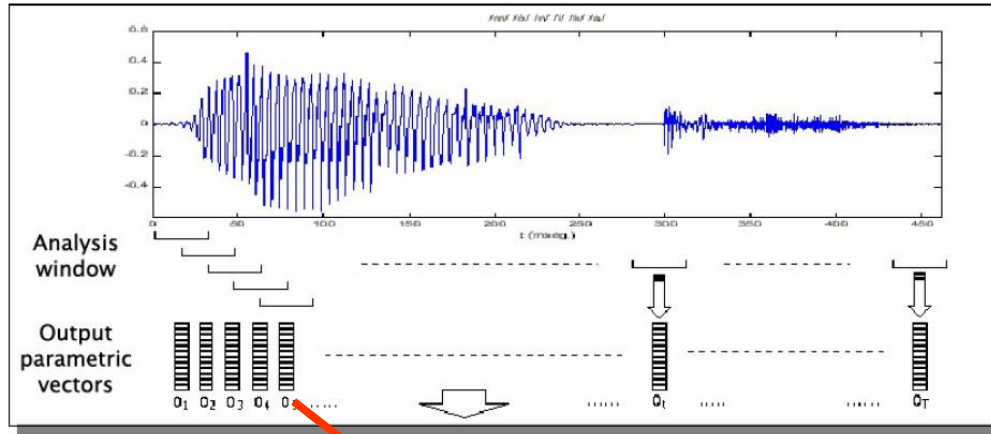
100 Hz Fundamental

54

Short-time Analysis and Parameterization



Feature Space



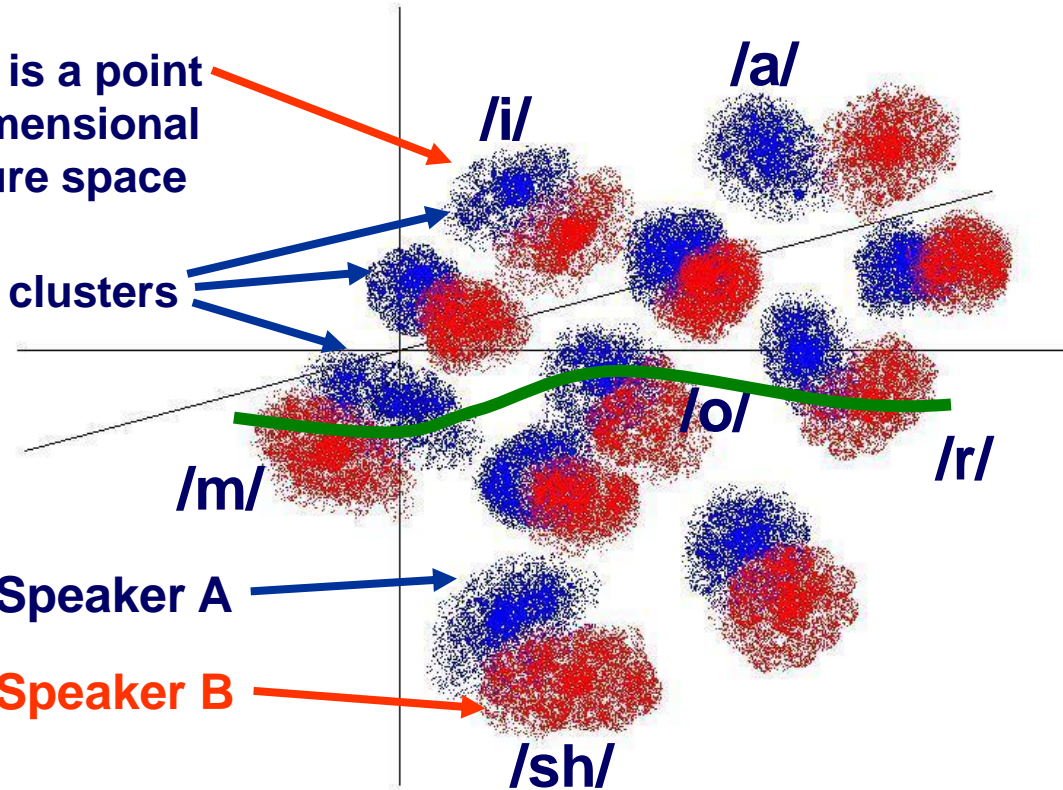
Bag of vectors representation
of speaker's acoustic space

Each vector is a point
in the 13-dimensional
MFCC feature space

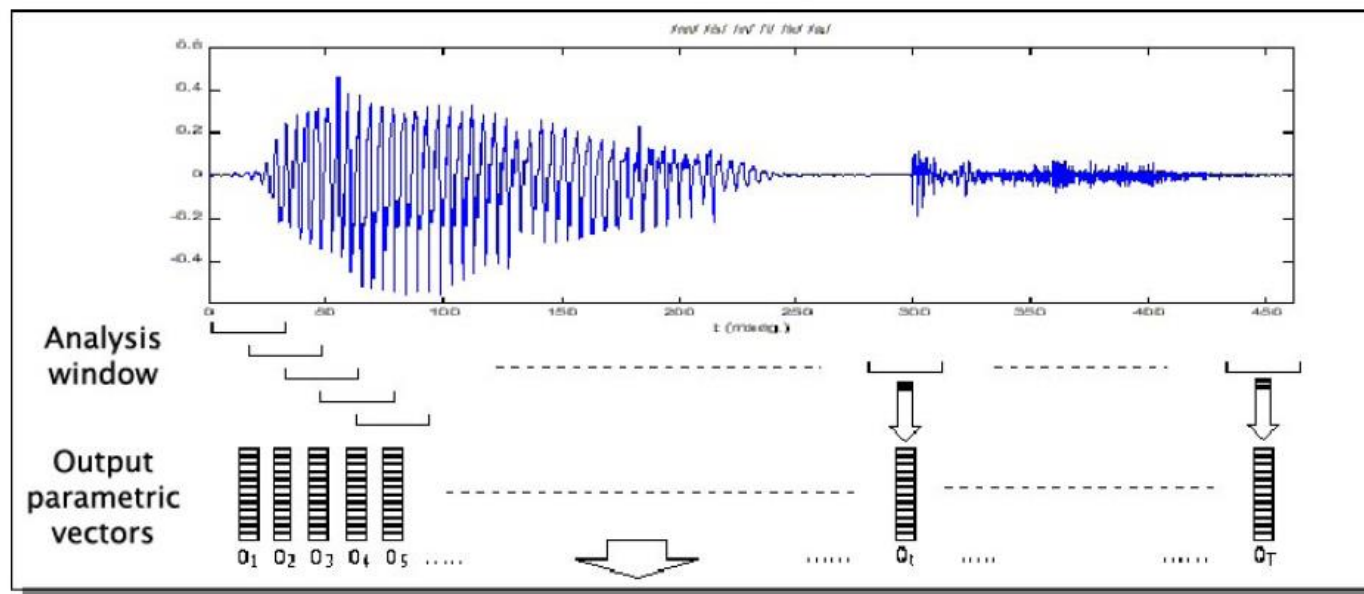
Phone clusters

■ Speaker A

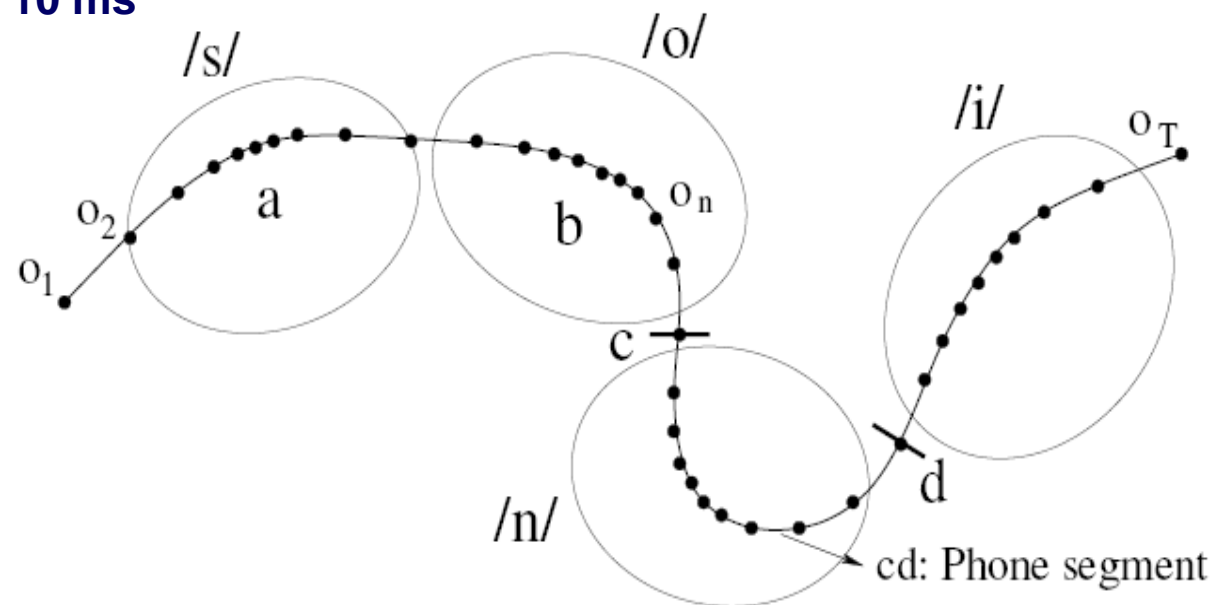
■ Speaker B

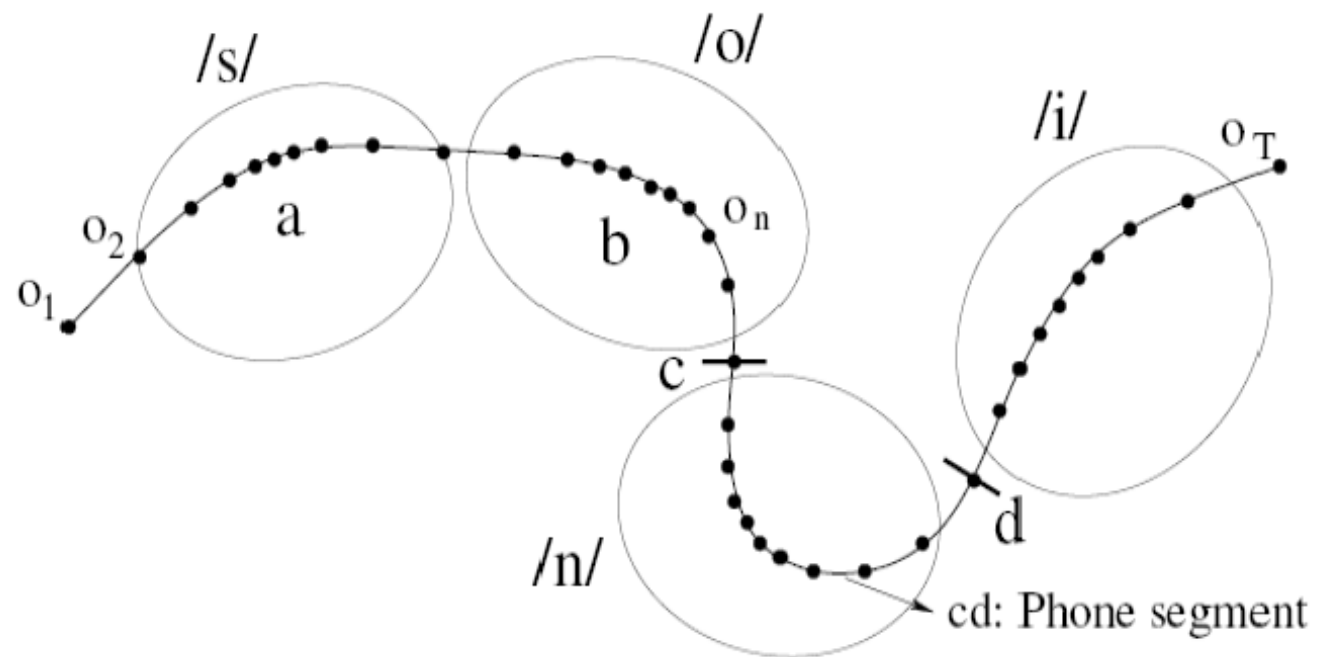


Feature Space



One feature vector every 10 ms





SPEECH RECOGNITION ALGORITHMS

- ❑ TAKE THIS FEATURE VECTOR SEQUENCE
- ❑ AS INPUT AND DETERMINE “WHAT HAS BEEN SAID”
- ❑ e.g. SEQUENCE OF PHONES / SEQUENCE OF WORDS etc.

Thank you !!

Spectrogram

- Speech is a continuous evolution of the vocal tract
- Spectrogram shows time-frequency evolution
- Represented as a time-series of short-time spectra

