

# Mathematical Analysis

Transmitter (Speech Recognition):

→ Semantic Encoder: input spectrum → text-related semantic features.

• Input: speech sample  $m = [m_1, m_2, \dots, m_Q]$

• It is divided into  $N$  frames.

→ Applying Hanning window to all  $n$ , as

$$w(n) = 0.54 - 0.46 \cos\left(\frac{2\pi n}{N-1}\right), 0 \leq n < N$$

→ Let  $X$  be the vector of length  $N$  of FFT,

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

→ After applying Hanning window, FFT, logarithmic operations and normalization, we get a spectrogram that is the input to our semantic encoder.

→ The ultimate goal of the ASR task is to recover the final text transcription  $\hat{t}$ , as close to  $t$ .

→ Semantic Encoder is made up of two CNN and GRU (Bidirectional RNN).

→ The input spectra  $S$ , are first converted to intermediate features using several convolutional layers.

→ The no. of filters, in each CNN module is  $E_p$ ,  $p \in [1, 2, \dots, P]$ .

→ The output of the last CNN module is

$$b \in \mathbb{R}^{B \times C_p \times D_p \times E_p}$$

→ Then  $b$  is fed into  $\odot$  BRNN modules, that gives  $d \in \mathbb{R}^{B \times G \times H_2}$ , where the no. of GRU units in BRNN module is  $H_2, q \in [1, 2, \dots, Q]$ .

→ Finally the text related semantic features  $p$  are obtained from  $d$  by passing through multiple cascaded dense layers and a softmax layer.

Channel Encoder:

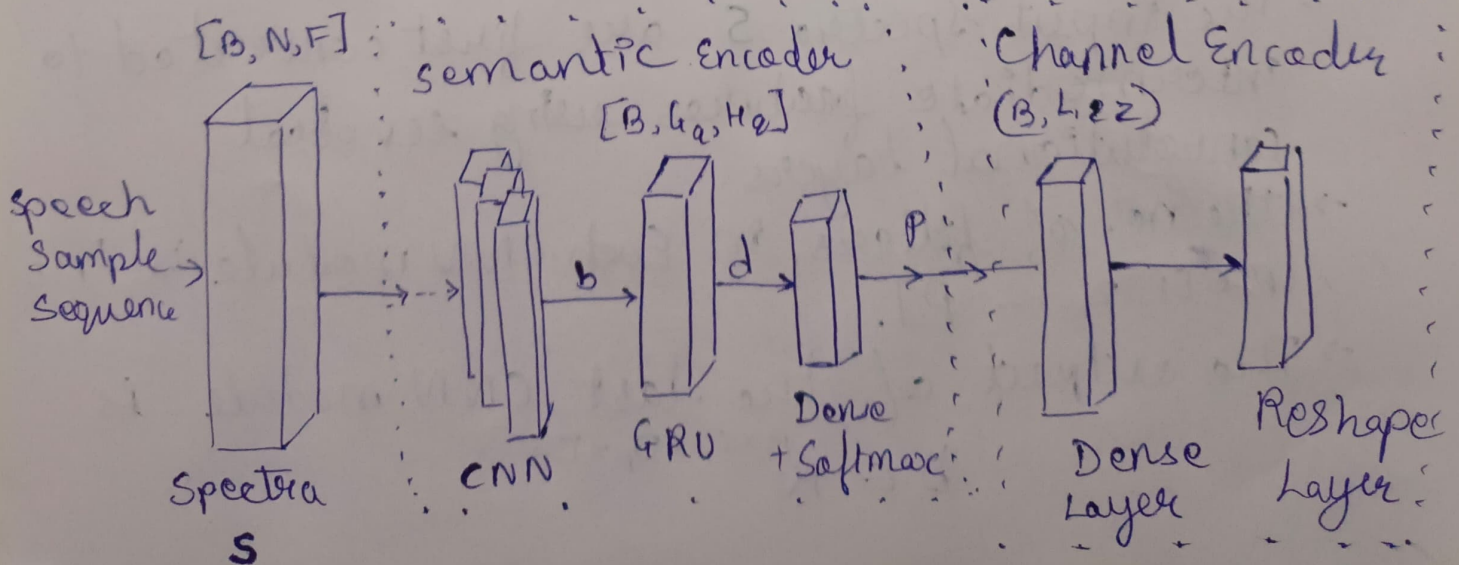
→ These features are mapped into symbols  $\vec{x}$  by the channel encoder to be transmitted as

$$\vec{x} = \underbrace{\vec{T}_\beta^c}_{\text{channel Encoder}} \left( \underbrace{\vec{T}_\alpha^s}_{\text{Semantic Encoder}} (\vec{s}) \right)$$

• Assumption:  $E \|\vec{x}\|^2 = 1$

• System Model assuming perfect CSI:

$$\vec{y} = \vec{h} * \vec{x} + \vec{w}; \quad \vec{w} \sim \text{CN}(0, \sigma^2 \mathbf{I})$$





# System Model Assuming Imperfect CSI.

→ Imperfect CSI at  $R_x$  is more realistic assumption

System model:

$$y = \hat{h} x + n$$

Received signal at  $R_x$  :  $(-\sqrt{E_s}, \sqrt{E_s})$   $n \sim \mathcal{CN}(0, N_0)$

Let us assume the imperfect CSI to  $\hat{h}$ , and the correlation between  $h$  and  $\hat{h}$  be  $\rho$ .

$$\text{i.e. } E[h^* \hat{h}] = \rho, \quad 0 \leq \rho < 1.$$

Assumptions:

① Imperfect CSI is denoted as  $\hat{h}$ , is available at the receiver for detection ( $R_x$ )

②  $h \sim \mathcal{CN}(0, 1)$ ,  $\rho$  = correlation between  $h$  &  $\hat{h}$ .

$$E[h^* \hat{h}] = \rho, \quad 0 \leq \rho < 1$$

By Maximum Likelihood detection rule:

$$\begin{aligned} d &= \operatorname{Re}\{\hat{h}^* y\} \\ &= \operatorname{Re}\{\hat{h}^* (h x + n)\} \end{aligned}$$

If  $d \geq 0$ , the detected symbol is 1.  
If  $d < 0$ , the detected symbol is 0.

$h$  can be expressed in terms of  $\hat{h}$  and Random error  $s$  as follow:

$$h = \rho \hat{h} + \sqrt{1-\rho^2} \cdot s,$$

$$s \sim \mathcal{CN}(0, 1).$$

$$\begin{aligned} d &= \text{Re} \{ \hat{h}^* (\rho \hat{h} + \sqrt{1-\rho^2} s) x + n \} \\ &= \text{Re} \{ (\rho |\hat{h}|^2 + \sqrt{1-\rho^2} \hat{h}^* s) x + \hat{h}^* n \} \\ &= \text{Re} \{ \rho |\hat{h}|^2 x + \sqrt{1-\rho^2} \hat{h}^* s x + \hat{h}^* n \} \\ &= \rho |\hat{h}|^2 x + \text{Re} \{ \sqrt{1-\rho^2} \hat{h}^* s x + \hat{h}^* n \} \end{aligned}$$

$$z = \frac{d}{|\hat{h}|} = \underbrace{\rho |\hat{h}| x}_{\text{signal part}} + \underbrace{\text{Re} \left\{ \sqrt{1-\rho^2} \frac{\hat{h}^* s x}{|\hat{h}|} + \frac{\hat{h}^* n}{|\hat{h}|} \right\}}_{\substack{\text{Not a noise part} \\ \text{but due to } s \text{ (Random error)}}}$$

↑  
noise part

$\hat{h}$  and  $s$  both should be independent.

Nakagami-m Fading channel:

Gray coded discrete (square and rectangular)  $M$ -ary OAM ~~modules~~ modems with  $M = 2^n$  ( $n = 2, 3, \dots, N$ ) are used for the channel adaptive scheme.

→ The transmissions are assumed to be over a slowly flat-fading channel model assumed to follow a Nakagami- $m$  distribution

Received signal  $y = at + gn$

where  $a \rightarrow$  Nakagami- $m$ -fading coefficient  
with pdf as

$$p_a(a) = \frac{2}{\Gamma(m)} \left(\frac{m}{\Omega}\right)^m a^{2m-1} \exp\left\{-\frac{ma^2}{\Omega}\right\}$$

where  $m$  is the Nakagami fading parameter and  $\Gamma(m)$  is the Gamma function defined by  $\Gamma(m) = \int_0^{\infty} y^{m-1} e^{-y} dy$ .

and  $\Omega = E[|a|^2]$   
Expectation

By simplifying the instantaneous SNR will be

$$\gamma = \frac{|a|^2 E_s}{N_0}$$

for Nakagami- $m$  fading channel,  
pdf of  $\gamma$  is given by

$$p_\gamma(\gamma) = \left(\frac{m}{\gamma}\right)^m \frac{\gamma^{m-1} \exp\left\{-\frac{m\gamma}{\Omega}\right\}}{\Gamma(m)}$$

$$\bar{\gamma} = E[\gamma] = \frac{\Omega E_s}{N_0}$$

$$\underline{\gamma \geq 0.}$$



The received SNR is split into  $(N+1)$  fading regions (bins), with region  $n$  having a corresponding mode  $m_n$ .

$\{r_n^m\}_n^N = 2$  contains lower threshold for  $N$  fading regions

$r_1$  is set to 0 dB and  $r_{N+1}$  to  $\infty$ .

Thus Received SNR,  $r$ , falls within region  $n$  ( $r_n \leq r < r_{n+1}$ )

$r_1 \leq r < r_2$  (outage)