GT Label $z_1$ $z_2$ --- GT $z_t$ --- $z_T$

Loss

Predicted Label: $\hat{z}_t$

classifier

$y_1$ $y_2$ --- $y_t$ --- $y_T$

f.E.

$x_1$ $x_t$ $x_T$

Conventional
Supervised
Pipeline for ASR

$L$ ← --- $L$ : :
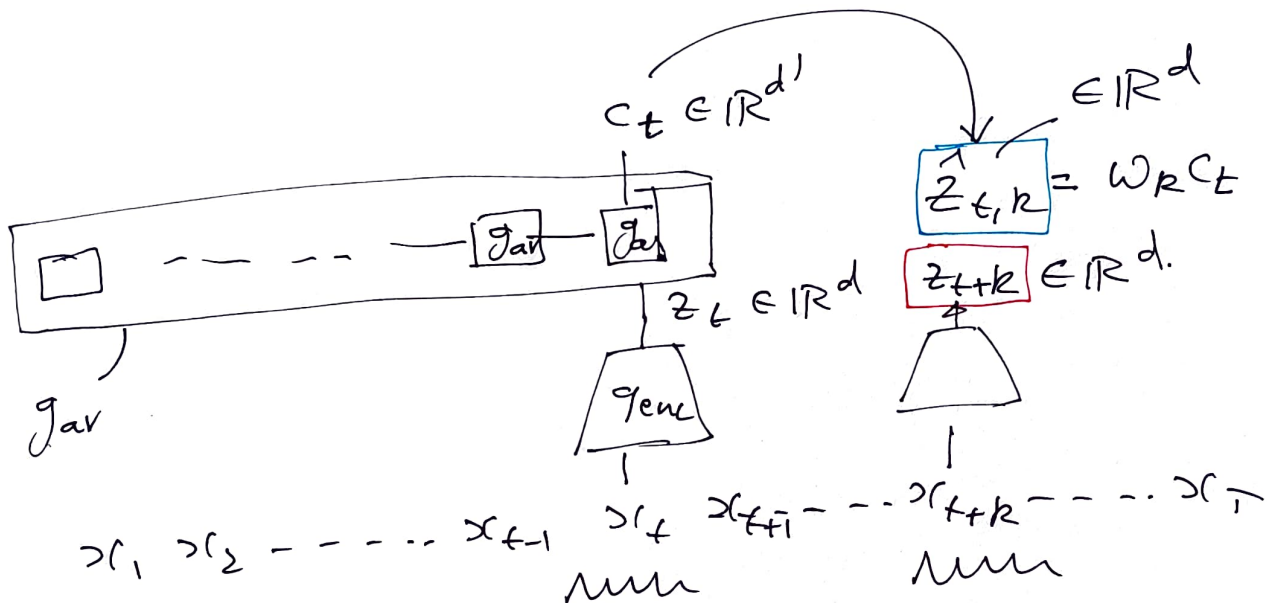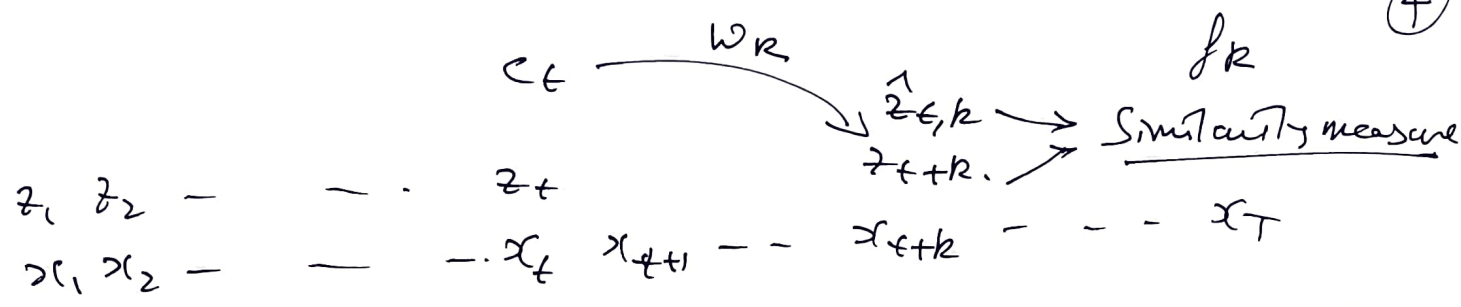
$y_t$

$g_\gamma$
$z_t$

$h_\theta$

$x_t$ $x_{t+k}$ $x_{t+k}$ $x_T$

$x_1$ $x_2$ --- /iʌ/

shall we go for coffee

/w/ /ig/

Transfer to
Downstream
Task

②

$L$

$z_{t,R}$

$R_\gamma$

$c_t$

$x_{f+R}$
or
$z_{t+R}$.

$h_{t+1}$

$\eta_{t-1}$

$\eta_{t}$

Label $z_t$

$q_\phi$

To Downstream

$h_\theta *$

$z_1$

$z_{t-2}$   $z_{t-1}$   $z_t$

$g_{enc}$

$z_{t+R}$

$x_1, x_2 - - - x_f - - - x_T$

$x_1$   $x_2$ — — — .   $x_{t-1}$   $x_t$   $x_{f+1} - - -$      $x_{t+R}$

$R$

$$c_t \in \mathbb{R}^{d'}$$

$$\in \mathbb{R}^d$$

$$\hat{z}_{t,k} = W_k c_t$$

$$z_{t+k} \in \mathbb{R}^d.$$

$$z_t \in \mathbb{R}^d$$

$$g_{enc}$$

$$x_1, x_2 \; \text{------} \; x_{t-1} \quad x_t \quad x_{t+1} \text{---} \; x_{t+k} \text{----} \; x_T$$

$$g_{ar}$$

$$z_t = g_{enc}(x_t) \quad \rightarrow \quad \text{10 ms latent variable}$$

$$c_t = g_{ar}(z_{\leq t})$$

$$c_t \xrightarrow{\quad W_R \quad} \hat{z}_{t,k} \xrightarrow{\quad f_R \quad} \underline{\text{Similarity measure}}$$

$$z_{t+k}.$$

$$z_1, z_2 \quad -- \quad -\cdot \quad z_t$$

$$x_1, x_2 \quad - \quad - \quad -\cdot x_t \quad x_{t+1} \quad -- \quad x_{t+k} \quad - \quad - \quad - \quad x_T$$

□ Ideally, if we have a generative model $p\left(x_{t+k} \mid c_t\right)$ it is easy to predict future observations (as in a LM)

— probabilistic "goodness" of $x_{t+k} \mid c_t$

□ But, use a metric that preserves Mutual Information (MI) between $x_{t+k}$ & $c_t$

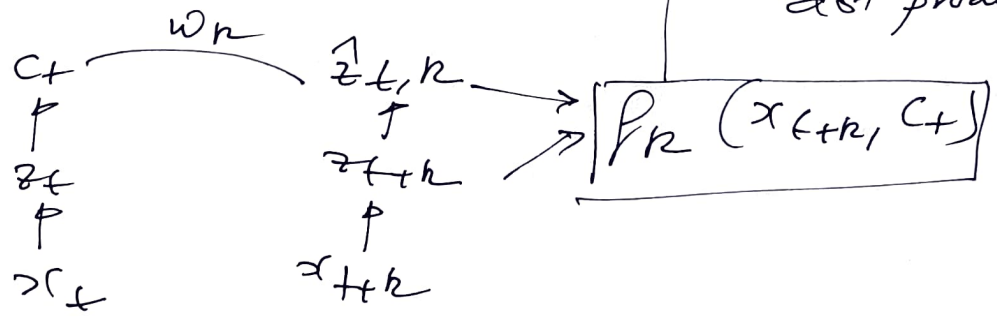$$\text{as} \quad f_R(x_{t+k}, c_t) \propto \frac{p(x_{t+k} \mid c_t)}{p(x_{t+k})}$$

where $f_R(x_{t+k}, c_t) = \exp\left(z_{t+k}^T W_R c_t\right)$

ie

$$f_k(x_{t+k}, c_t) \propto \frac{p(x_{t+k} \mid c_t)}{p(x_{t+k})}$$

realized as a scalar "similarity" score
between $x_{t+k}$ $[$ie $z_{t+k}]$ & $c_t$ $[$ie $W_k c_t = \hat{z}_{t,k}]$

$$f_k(x_{t+k}, c_t) = \exp\left( z_{t+k}^T \overbrace{W_k c_t}^{\hat{z}_{t,k}} \right)$$

$\underbrace{\phantom{z_{t+k}^T W_k c_t}}$ dot product

$$
\begin{array}{ccc}
c_t & \xrightarrow{W_k} & \hat{z}_{t,k} \\
\uparrow p & & \uparrow \\
z_t & & z_{t+k} \\
\uparrow p & & \uparrow p \\
x_t & & x_{t+k}
\end{array}
\longrightarrow \boxed{f_k(x_{t+k}, c_t)}
$$

# Mutual Information [Information Gain]

- Between 2 random variables $X$ & $Y$

- Measure of the Mutual Dependence between the 2 variables

- Quantifies "amount of information" obtained about one r-v by observing the other r-v.

- Uncertainty about $X$ (or $Y$) reduced once $Y$ (or $X$) is observed

→ Information gain.

MI

$$I(x;y) = D_{KL}\left(P_{(x,y)} \,||\, P_x \otimes P_y\right) \quad \text{(7)}$$
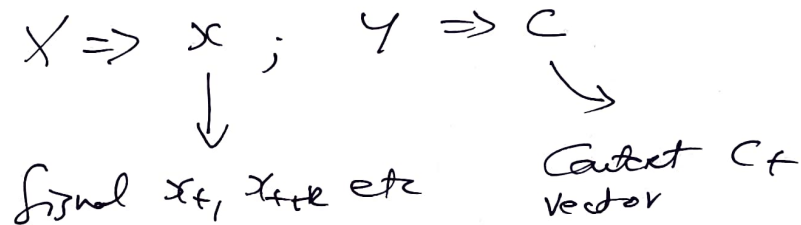
$\Rightarrow$ for PMFs

$$I(X;Y) = \sum_{y \in Y} \sum_{x \in X} P_{x,y}(x,y) \log\left[\frac{P_{x,y}(x,y)}{P_x(x) \cdot P_y(y)}\right]$$

$I(x,y) = 0$    if and only if $x$ & $y$ are independent variables

ie $P_{x,y}(x,y) = P_x(x) \cdot P_y(y)$

In our setting

$$X \Rightarrow x \quad ; \quad Y \Rightarrow c$$

$X \Rightarrow x \downarrow$ signal $x_t, x_{t+k}$ etc

$Y \Rightarrow c \searrow$ Context $c_t$ vector

$$I(x; c) = \sum_{x,c} p(x,c) \lg \frac{p(x,c)}{p(x) \cdot p(c)}$$

$$= \sum_{x,c} p(x,c) \lg \frac{p(x/c)}{p(x)}$$

$$p(x,c) = p(x/c) \cdot p(c)$$

Maximise MI between encoded representations $c_t$ & $z_{t+k}$ in the place of $x_{t+k}$

in place of

of interest

$$I\left(x_{t+k}, c_t\right)$$

$$= \sum_{x,c} p\left(x_{t+k}, c_t\right) \log \frac{p\left(x_{t+k}/c_t\right)}{p\left(x_{t+k}\right)}$$

If Loss between $c_t$ & $x_{t+k}$ ( or $z_{t+k}$ ) is of the form

$$\mathcal{L}_N = - \underset{X}{\mathbb{E}}\left[\log \frac{f_R\left(x_{t+k}, c_t\right)}{\sum_{x_j \in X} f_R\left(x_j, c_t\right)}\right]$$

X: N Samples ⟨ 1 +ve $x_{t+k}$ / N-1 -ve $x_j$ ⟩

Then it can be shown

$$I\left(x_{t+k}, c_t\right) \geq \log(N) - \mathcal{L}_N$$

Ideally, we need to optimise $g_{enc}$, $g_{ar}$ & $W_k$

to maximize <u>mutual information</u> between $c_t$ & $x_{t+k}$

— Instead, work with $MI$ between $c_t$ & $z_{t+k}$

— As a proxy to <u>max $MI$</u> $\Rightarrow$ <u>min Info NCE loss $L_N$</u>

$$c_t \xleftarrow{\phantom{MI}} \quad \xrightarrow{\phantom{xx}} MI \text{ or } L_N$$

$$\uparrow$$
$$z_t \qquad \underline{MI} \qquad z_{t+k}$$
$$\uparrow \qquad \uparrow$$
$$x_{t-1}, x_{t}, x_{t+1} \;-\;-\;--\; x_{t+k} \;-\;-\;-\; x_T$$
$$x_1, x_2 \;-\;-\;-\;-\;$$

. $\quad I(x_{t+k}, c_t) \geq \log N - L_N \longrightarrow$ <u>Info NCE loss</u>

$\hat{L}_N \downarrow \Rightarrow$ Lower Bound of $I(x_{t+k}, c_t) \uparrow$

$\quad\quad\quad\;\Rightarrow$ Good for $MI(x_{t+k}, c_t)$

$N \uparrow \Rightarrow$ Good for  "

# Implication of Baird-Result

$$\text{Max} \quad I(x_{t+k}, c_t)$$

$$\Rightarrow \quad \min_{g_{enc}, g_{ar}, W_k} L_N$$

$$= \quad \min_{g_{enc}, g_{ar}, W_k} \left[ -\mathbb{E}_X \left( \log \frac{f_R(x_{t+k}, c_t)}{\sum\limits_{x_j \in X} f_R(x_j, c_t)} \right) \right]$$

$$= \quad \min_{g_{enc}, g_{ar}, W_k} \left[ -\mathbb{E}_X \left\{ \log \frac{\exp\left(z_{t+k}^T W_k c_t\right)}{\sum\limits_{x_j \in X} \exp\left(z_j^T W_k c_t\right)} \right\} \right]$$

Self-Supervision by CPC: Find optimal Model Parameters

$$\theta^* = g_{enc}^*, g_{ar}^*, w_R^*.$$

$$= \underset{g_{enc}, g_{ar}, W_R}{\arg\min} -\mathbb{E}_x \left[ \lg \frac{\exp\left(z_{t+k}^T W_R c_t\right)}{\sum\limits_{z_j \in X} \exp\left(z_j^T W_R c_t\right)} \right]$$

$$\Rightarrow \quad \begin{aligned} z_{t+k} &= g_{enc}\left(x_{t+k}\right) \\ c_t &= g_{ar}\left(z_{\le t}\right) \\ z_j &= g_{enc}\left(x_j\right) \end{aligned}$$

## $L_N$ is a NCE loss — towards maximizing MI.

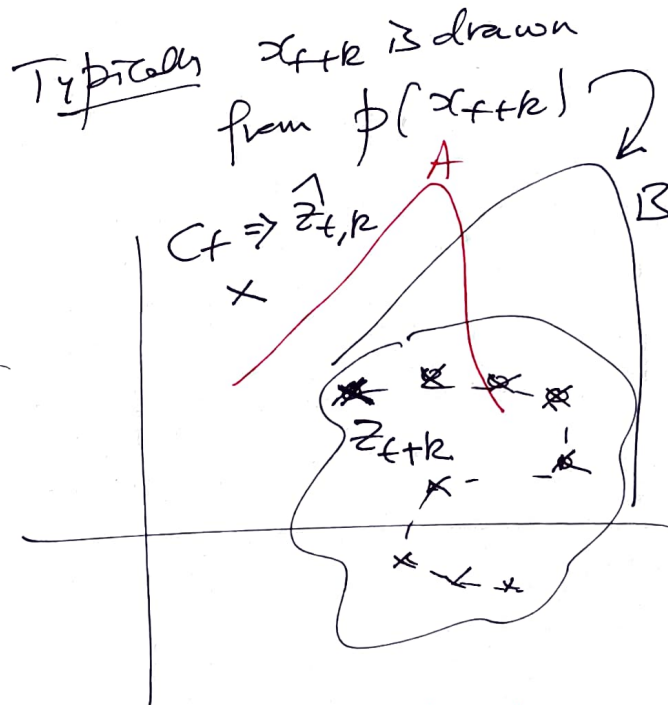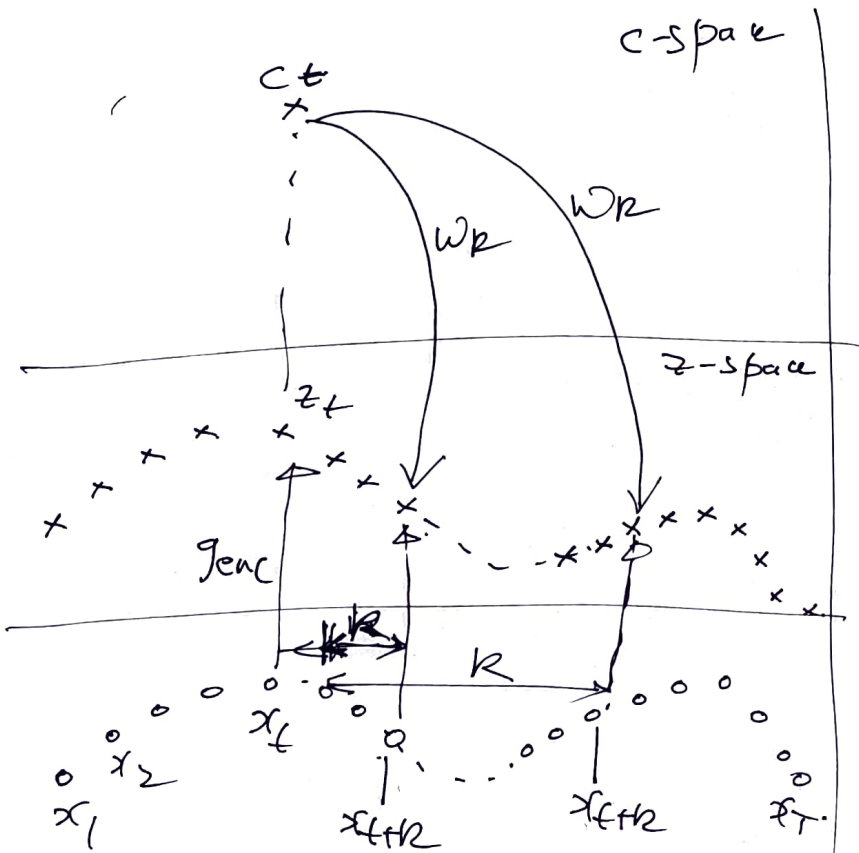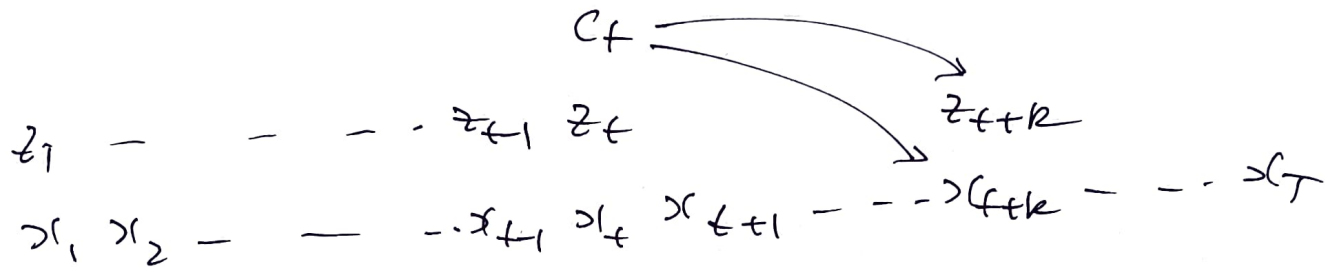$L_N$ : Given a set of $X = \{x_1, x_2, \dots x_N\}$

of $N$ random samples

- has 1 +ve sample from $p(x_{t+k}/c_t)$

- $N-1$ -ve samples from $p(x_{t+k})$

Data Distn.

$p(x_{t+k}/c_t)$ : ~~Each~~ Conditional pdf of $x_{t+k}$ given $c_t$

MI Concept $\rightarrow$ $\begin{cases} \text{Once } c_t \text{ is observed available,} \\ \text{uncertainty about } x_{t+k} \text{ reduces} \end{cases}$

Noise Distn

$p(x_{t+k}) \rightsquigarrow$ pdf when no knowledge of $c_t$ is available.

$$c_t \longrightarrow z_{t+k}$$

$$z_1 \quad - \quad - \quad - \quad - \cdot z_{t-1} \; z_t \qquad z_{t+k}$$

$$x_1, x_2 \quad - \quad - \quad - \cdots x_{t+1} \; x_t \; x_{t+1} \; - \; - \; - \rightarrow c_{t+k} \; - \; - \; - \cdot \rightarrow x_T$$



c-space

$c_t$

$W_R$  $W_R$

z-space

$z_t$

$g_{enc}$

$\frac{1}{k}R$

$R$

$x_t$ $q$

$x_{t+k}$  $x_{t+k}$  $x_T$

$x_1$ $x_2$

Typically $x_{t+k}$ is drawn from $p(x_{t+k})$

$c_t \Rightarrow \hat{z}_{t,k}$

A

B

$z_{t+k}$

A: $p(x_{t+k} \mid c_t)$

B: $p(x_{t+k})$

$$x_t \longrightarrow z_t \longrightarrow c_t$$

On knowing $c_t$, pick $x_{t+k}$ for some 'k'

then given / an observing $c_t$

$$p(x_{t+k}) \longrightarrow p(x_{t+k}/c_t)$$

Noise distribution

Data Distribution

$\left[\begin{array}{l} \text{-ve samples are} \\ \text{drawn from this} \\ \text{distr} \end{array}\right]$

$\left[\begin{array}{l} \text{+ve sample } x_{t+k} \\ \text{is supposed to come} \\ \text{from this distribution} \end{array}\right]$

<u>Verify</u> this $\longrightarrow$ <u>N C E</u>

## InfoNCE loss

$$\mathcal{L}_N = - \mathop{E}_{X}\left[ \log \frac{f_k(x_{t+k}, c_t)}{\sum\limits_{x_j \in X} f_k(x_j, c_t)} \right]$$

Is a categorical cross-entropy loss

$-$ of classifying the $+$ve sample 'correctly'

ie as coming from Data Distn. &

Not the Noise Distn

$$P\left[\begin{array}{l} x_{(t+k} \text{ was drawn from The Conditional distn } p\left(x_{t+k}/c_t\right) \\ \text{rather than the proposal or 'noise' distn } p(x_{t+k}) \end{array}\right]$$

$$= p\left(x_{t+k} \text{ is a } +ve \text{ sample} \mid x, c_t \right)$$

$$P\left(x_{t+k} \text{ is from Data Distn} \,\middle|\, X, c_t\right) = \frac{P\left(x_{t+k}, c_t\right)}{P(c_t)}$$

a +ve Sample

$$= \frac{\dfrac{P\left(x_{t+k} \mid c_t\right)}{P\left(x_{t+k}\right)}}{\displaystyle\sum_{j=1}^{N} \dfrac{P\left(x_j \mid c_t\right)}{P\left(x_j\right)}}$$

(jump) ——→

Eqn (5)

$$\frac{P\left(x_{t+k}^{\text{+ve}} \,\middle|\, c_t\right) \cdot P(c_t)}{= P\left(x_{t+k}, c_t\right)} \Bigg) \times$$

$$\frac{\text{This is the } \text{Prob associated with Correct class}}{@ \text{ optimal (minimization of InfoNCE Loss } L_N)}$$