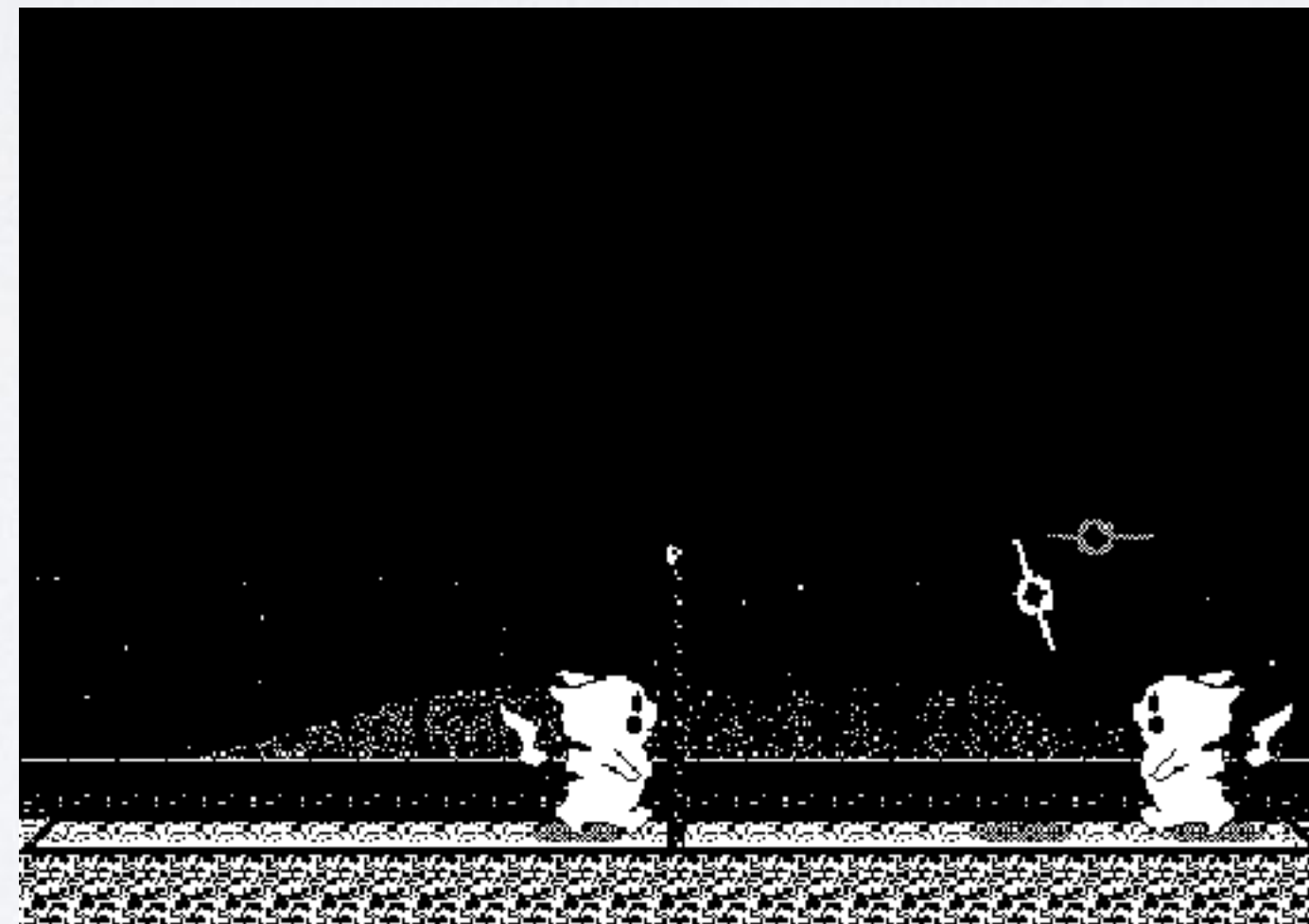
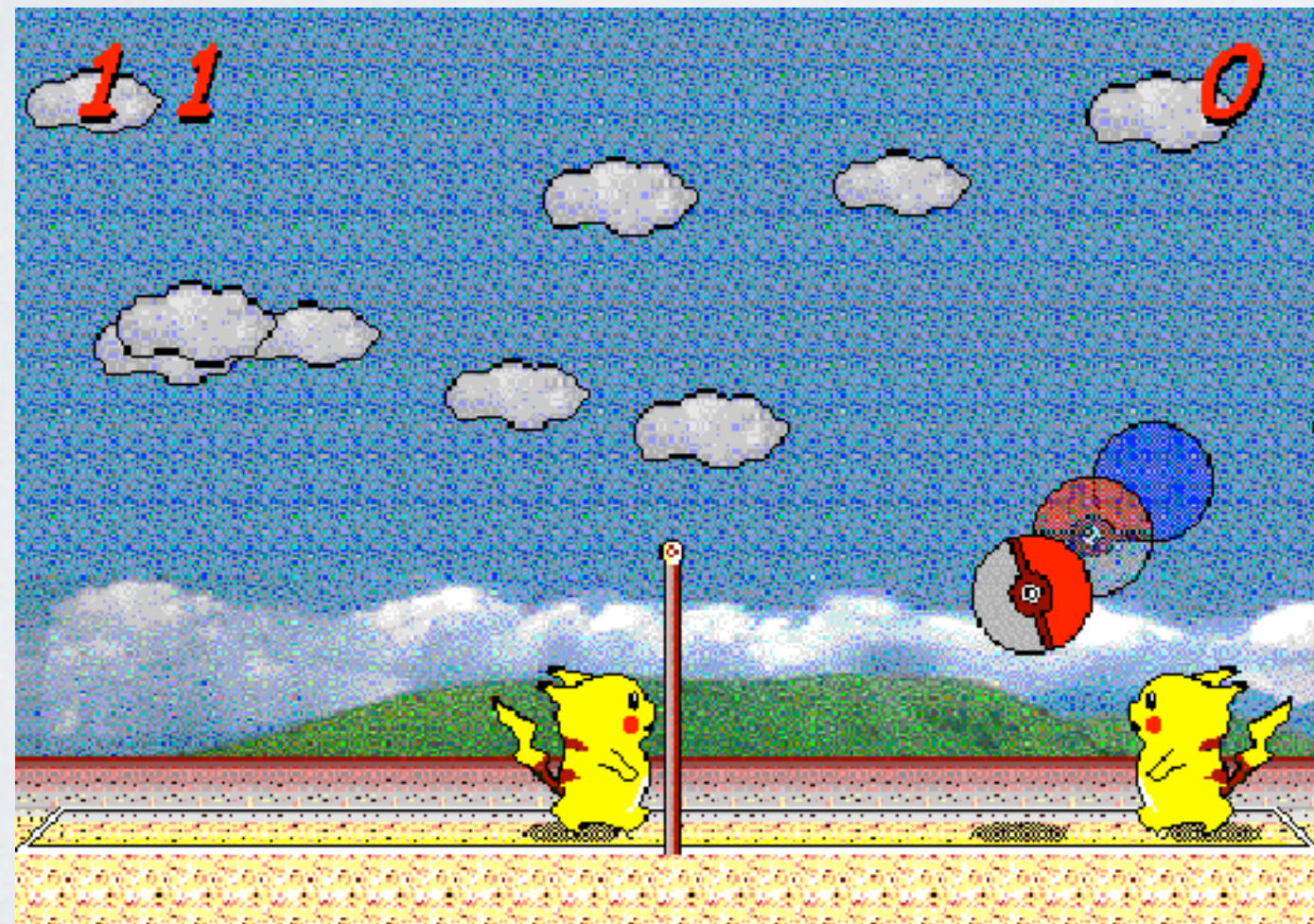


DISENTANGLING FROM SEQUENTIAL DATA

SNU Datamining Laboratory
2018. 6. 18 Seminar
Sungwon, Lyu
lyusungwon@dm.snu.ac.kr

REPRESENTATION LEARNING

- Representation Learning
 - Learning representation of the data that make it easier to extract useful information when building classifiers or other predictors



(110, 5, 225, 5, 210, 10)

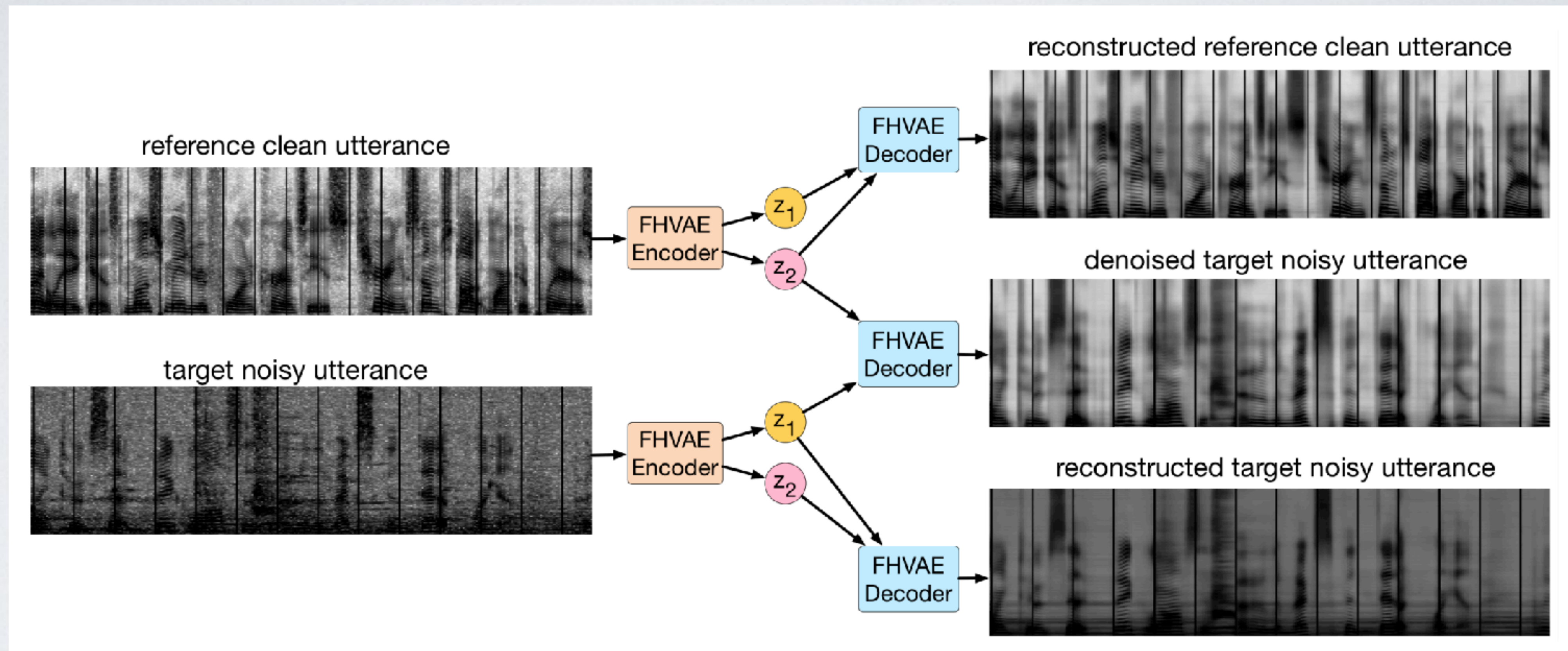
WHAT MAKES A REPRESENTATION GOOD?

- Smoothness: $x \approx y$ then $f(x) \approx f(y)$
- Multiple Explanatory factors: Generalize many configurations of factors
- A hierarchical organization of explanatory factors
- Semi-supervised learning
- Shared factors across tasks
- Manifolds: Low meaningful dimensions
- Natural clustering: Named, categorized
- Temporal and Spatial coherence
- Sparsity: insensitive to small variations of x

DISENTANGLING

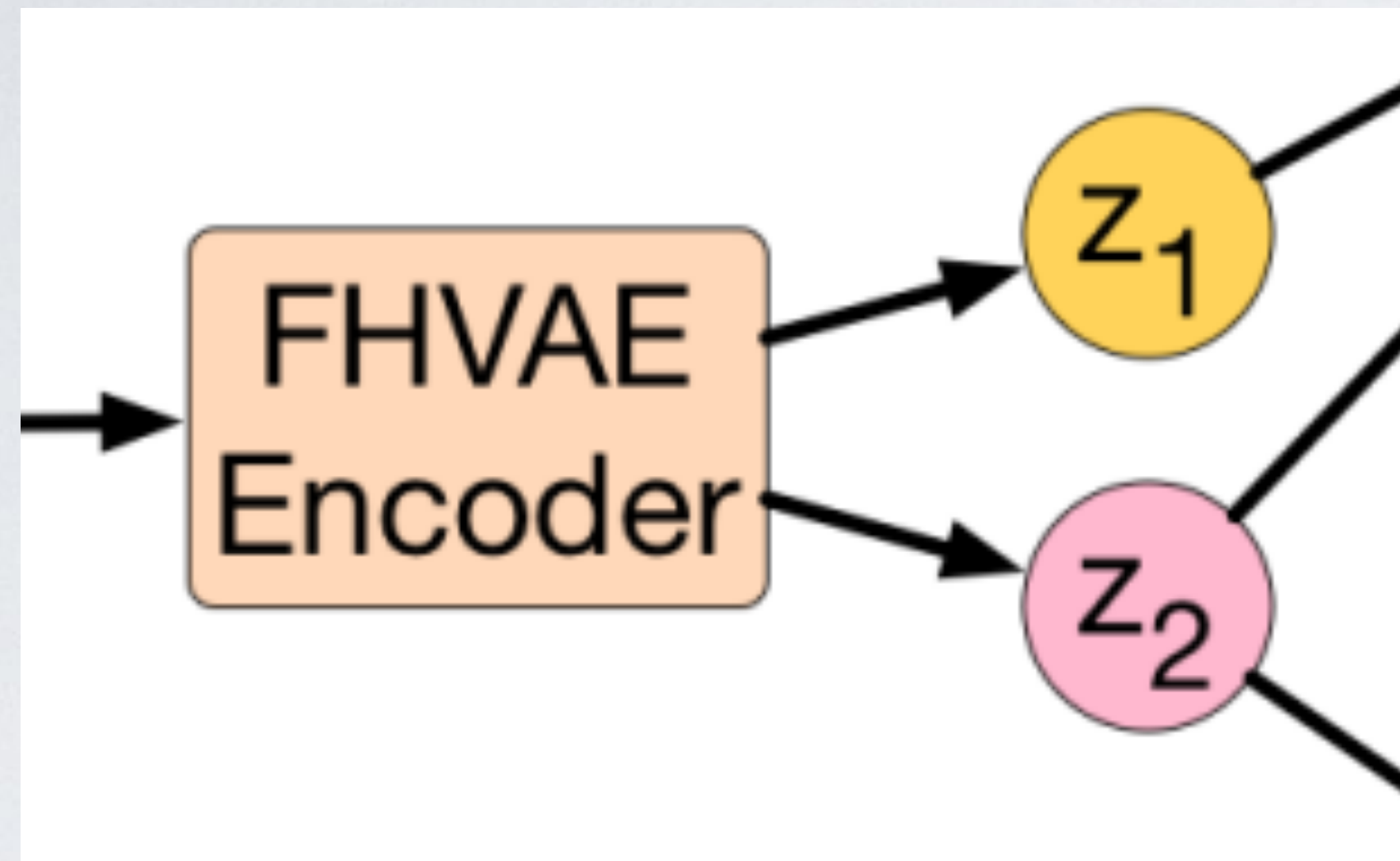
- DC-IGN
 - Clamping a part of the hidden units for a pair of data points that are known to match in all but one factors of variation
- InfoGAN
 - Maximize the mutual information between a latent code c and x through the use of an auxiliary distribution $Q(c|x)$
- Beta-Vae
 - Encourages the latent representation to be factorised by adding beta to VAE objective

MOTIVATION



<https://www.youtube.com/watch?v=naJZITvCfI4&feature=youtu.be>
<https://www.youtube.com/watch?v=VMX3IZYWYdg&feature=youtu.be>

FHVAE



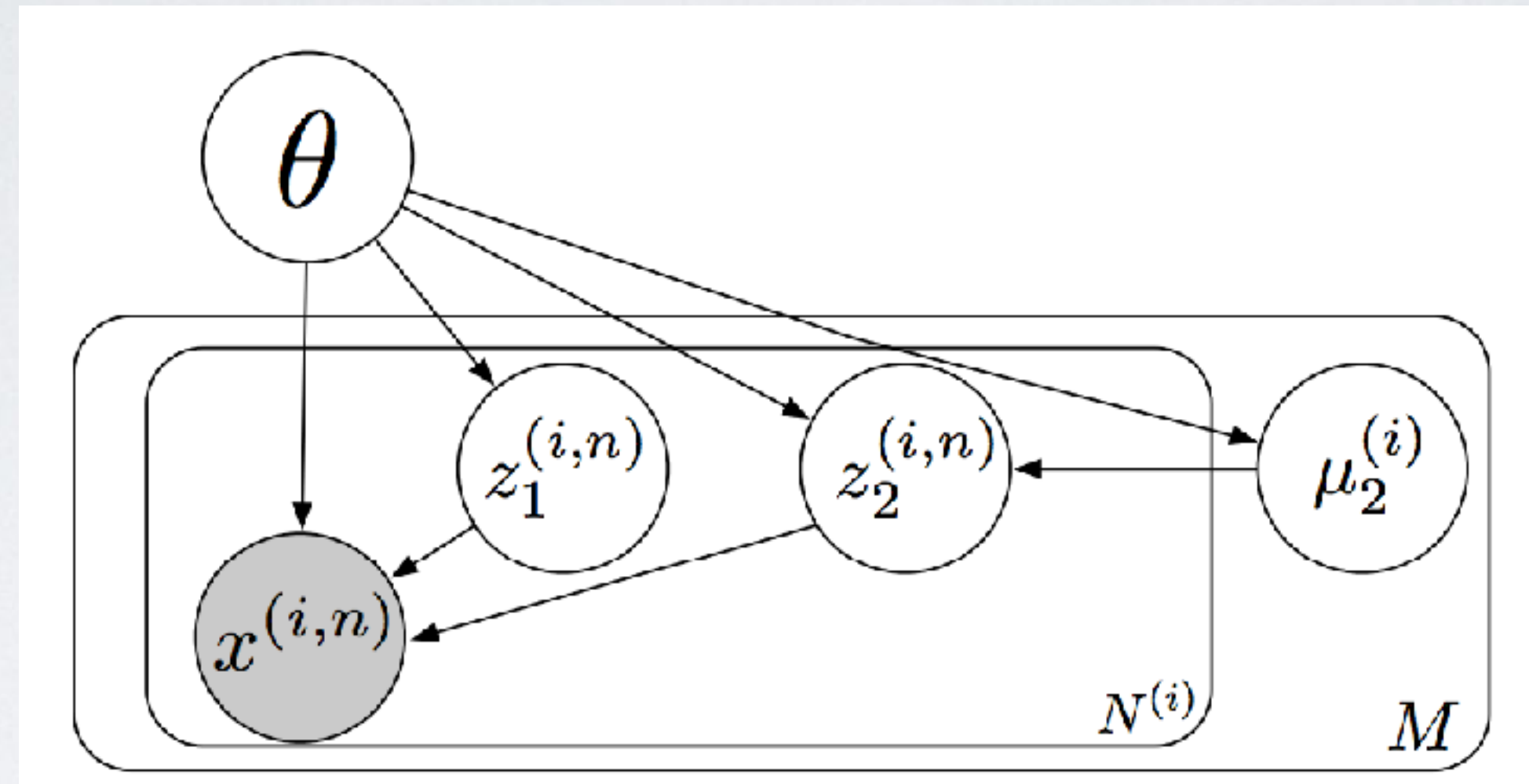
- Factorized Hierarchical VAE
- Sequence-level attributes
 - Speaker identity, Character style
 - Latent sequence variable (Z_1)
- Segment-level attributes
 - Phonetic content, Action
 - Latent segment variable (Z_2)
- Notation

$$\mathcal{D} = \{\mathbf{X}^{(i)}\}_{i=1}^M, \mathbf{X}^{(i)} = \{\mathbf{x}^{(i,n)}\}_{n=1}^{N^{(i)}}$$

$$\mathbf{Z}_1 = \{\mathbf{z}_1^{(n)}\}, \mathbf{Z}_2 = \{\mathbf{z}_2^{(n)}\}$$

FHVAE

- Generative model



- Equation

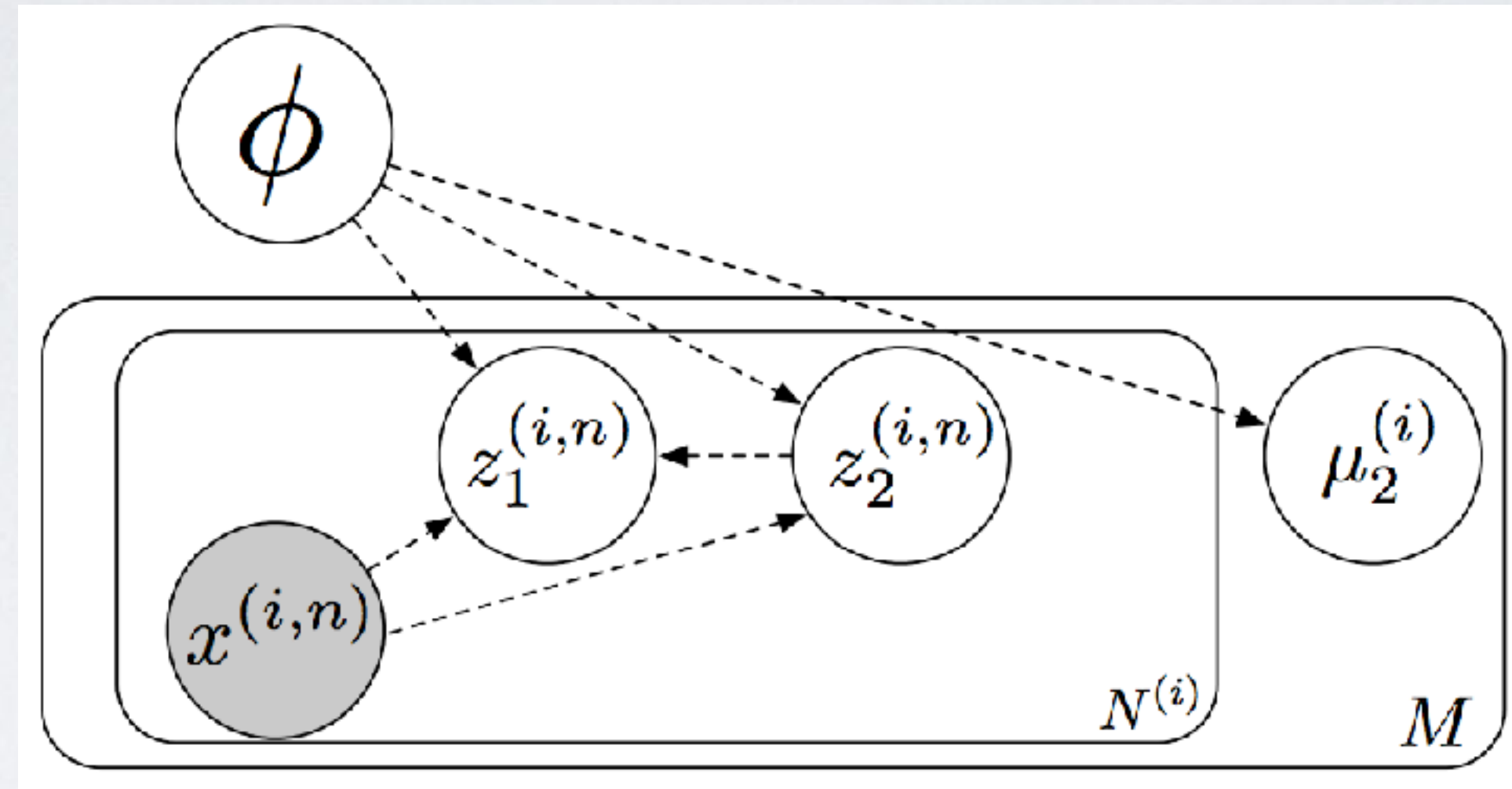
$$p_{\theta}(\mathbf{X}, \mathbf{Z}_1, \mathbf{Z}_2, \boldsymbol{\mu}_2) = p_{\theta}(\boldsymbol{\mu}_2) \prod_{n=1}^N p_{\theta}(\mathbf{x}^{(n)} | \mathbf{z}_1^{(n)}, \mathbf{z}_2^{(n)}) p_{\theta}(\mathbf{z}_1^{(n)}) p_{\theta}(\mathbf{z}_2^{(n)} | \boldsymbol{\mu}_2).$$

$$p_{\theta}(\mathbf{x} | \mathbf{z}_1, \mathbf{z}_2) = \mathcal{N}(\mathbf{x} | f_{\mu_x}(\mathbf{z}_1, \mathbf{z}_2), \text{diag}(f_{\sigma_x^2}(\mathbf{z}_1, \mathbf{z}_2)))$$

$$p_{\theta}(\mathbf{z}_1) = \mathcal{N}(\mathbf{z}_1 | \mathbf{0}, \sigma_{\mathbf{z}_1}^2 \mathbf{I}), \quad p_{\theta}(\mathbf{z}_2 | \boldsymbol{\mu}_2) = \mathcal{N}(\mathbf{z}_2 | \boldsymbol{\mu}_2, \sigma_{\mathbf{z}_2}^2 \mathbf{I}), \quad p_{\theta}(\boldsymbol{\mu}_2) = \mathcal{N}(\boldsymbol{\mu}_2 | \mathbf{0}, \sigma_{\boldsymbol{\mu}_2}^2 \mathbf{I})$$

FHVAE

- Inference model



- Equation

$$q_{\phi}(\mathbf{Z}_1^{(i)}, \mathbf{Z}_2^{(i)}, \mu_2^{(i)} | \mathbf{X}^{(i)}) = q_{\phi}(\mu_2^{(i)}) \prod_{n=1}^{N^{(i)}} q_{\phi}(z_1^{(i,n)} | \mathbf{x}^{(i,n)}, z_2^{(i,n)}) q_{\phi}(z_2^{(i,n)} | \mathbf{x}^{(i,n)})$$

$$q_{\phi}(\mu_2^{(i)}) = \mathcal{N}(\mu_2^{(i)} | g_{\mu_{\mu_2}}(i), \sigma_{\tilde{\mu}_2}^2 \mathbf{I}), \quad q_{\phi}(z_2 | \mathbf{x}) = \mathcal{N}(z_2 | g_{\mu_{z_2}}(\mathbf{x}), \text{diag}(g_{\sigma_{z_2}^2}(\mathbf{x})))$$

$$q_{\phi}(z_1 | \mathbf{x}, z_2) = \mathcal{N}(z_1 | g_{\mu_{z_1}}(\mathbf{x}, z_2), \text{diag}(g_{\sigma_{z_1}^2}(\mathbf{x}, z_2))),$$

FHVAE

- Factorising variational lower bound

$$\begin{aligned}\mathcal{L}(\theta, \phi; \mathbf{X}) &= \sum_{n=1}^N \mathcal{L}(\theta, \phi; \mathbf{x}^{(n)} | \tilde{\boldsymbol{\mu}}_2) + \log p_{\theta}(\tilde{\boldsymbol{\mu}}_2) + \text{const} \\ \mathcal{L}(\theta, \phi; \mathbf{x}^{(n)} | \tilde{\boldsymbol{\mu}}_2) &= \mathbb{E}_{q_{\phi}(\mathbf{z}_1^{(n)}, \mathbf{z}_2^{(n)} | \mathbf{x}^{(n)})} [\log p_{\theta}(\mathbf{x}^{(n)} | \mathbf{z}_1^{(n)}, \mathbf{z}_2^{(n)})] \\ &\quad - \mathbb{E}_{q_{\phi}(\mathbf{z}_2^{(n)} | \mathbf{x}^{(n)})} [D_{KL}(q_{\phi}(\mathbf{z}_1^{(n)} | \mathbf{x}^{(n)}, \mathbf{z}_2^{(n)}) || p_{\theta}(\mathbf{z}_1^{(n)}))] \\ &\quad - D_{KL}(q_{\phi}(\mathbf{z}_2^{(n)} | \mathbf{x}^{(n)}) || p_{\theta}(\mathbf{z}_2^{(n)} | \tilde{\boldsymbol{\mu}}_2)).\end{aligned}$$

- Segment Variational lower bound

$$\mathcal{L}(\theta, \phi; \mathbf{x}^{(n)}) = \mathcal{L}(\theta, \phi; \mathbf{x}^{(n)} | \tilde{\boldsymbol{\mu}}_2) + \frac{1}{N} \log p_{\theta}(\tilde{\boldsymbol{\mu}}_2) + \text{const.}$$

- Discriminative objective

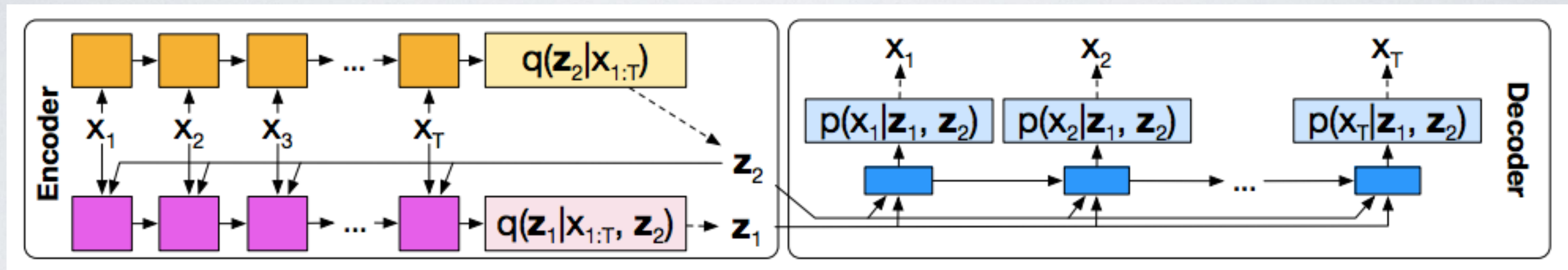
$$\begin{aligned}\log p(i | \mathbf{z}_2^{(i,n)}) &= \log p(\mathbf{z}_2^{(i,n)} | i) - \log \sum_{j=1}^M p(\mathbf{z}_2^{(i,n)} | j) \quad (p(i) \text{ is assumed uniform}) \\ &:= \log p_{\theta}(\mathbf{z}_2^{(i,n)} | \tilde{\boldsymbol{\mu}}_2^{(i)}) - \log \left(\sum_{j=1}^M p_{\theta}(\mathbf{z}_2^{(i,n)} | \tilde{\boldsymbol{\mu}}_2^{(j)}) \right),\end{aligned}$$

- Objective function

$$\mathcal{L}^{dis}(\theta, \phi; \mathbf{x}^{(i,n)}) = \mathcal{L}(\theta, \phi; \mathbf{x}^{(i,n)}) + \alpha \log p(i | \mathbf{z}_2^{(i,n)})$$

IMPLEMENTATION

- Sequence to sequence autoencoder model



$$\begin{aligned}
 (\mathbf{h}_{\mathbf{z}_2,t}, \mathbf{c}_{\mathbf{z}_2,t}) &= \text{LSTM}(x_{t-1}, \mathbf{h}_{\mathbf{z}_2,t-1}, \mathbf{c}_{\mathbf{z}_2,t-1}; \phi_{\text{LSTM}, \mathbf{z}_2}) \\
 q_{\phi}(\mathbf{z}_2 | x_{1:T}) &= \mathcal{N}(\mathbf{z}_2 | \text{MLP}(\mathbf{h}_{\mathbf{z}_2,T}; \phi_{\text{MLP}_{\mu}, \mathbf{z}_2}), \text{diag}(\exp(\text{MLP}(\mathbf{h}_{\mathbf{z}_2,T}; \phi_{\text{MLP}_{\sigma^2}, \mathbf{z}_2})))) \\
 (\mathbf{h}_{\mathbf{z}_1,t}, \mathbf{c}_{\mathbf{z}_1,t}) &= \text{LSTM}([x_{t-1}; \mathbf{z}_2], \mathbf{h}_{\mathbf{z}_1,t-1}, \mathbf{c}_{\mathbf{z}_1,t-1}; \phi_{\mathbf{z}_1}) \\
 q_{\phi}(\mathbf{z}_1 | x_{1:T}, \mathbf{z}_2) &= \mathcal{N}(\mathbf{z}_1 | \text{MLP}(\mathbf{h}_{\mathbf{z}_1,T}; \phi_{\text{MLP}_{\mu}, \mathbf{z}_1}), \text{diag}(\exp(\text{MLP}(\mathbf{h}_{\mathbf{z}_1,T}; \phi_{\text{MLP}_{\sigma^2}, \mathbf{z}_1})))) \\
 (\mathbf{h}_{\mathbf{x},t}, \mathbf{c}_{\mathbf{x},t}) &= \text{LSTM}([\mathbf{z}_1; \mathbf{z}_2], \mathbf{h}_{\mathbf{x},t-1}, \mathbf{c}_{\mathbf{x},t-1}; \phi_{\mathbf{x}}) \\
 p_{\theta}(x_t | \mathbf{z}_1, \mathbf{z}_2) &= \mathcal{N}(x_t | \text{MLP}(\mathbf{h}_{\mathbf{x},t}; \phi_{\text{MLP}_{\mu}, \mathbf{x}}), \text{diag}(\exp(\text{MLP}(\mathbf{h}_{\mathbf{x},t}; \phi_{\text{MLP}_{\sigma^2}, \mathbf{x}}))))
 \end{aligned}$$

EXPERIMENT

- Speaker Verification
- Automatic Speech Recognition

Table 1: Comparison of speaker verification equal error rate (EER) on the TIMIT test set

Features	Dimension	α	Raw	LDA (12 dim)	LDA (24 dim)
i-vector	48	-	10.12%	6.25%	5.95%
	100	-	9.52%	6.10%	5.50%
	200	-	9.82%	6.54%	6.10%
μ_2	16	0	5.06%	4.02%	-
	16	10^{-1}	4.91%	4.61%	-
	16	10^0	3.87%	3.86%	-
	16	10^1	2.38%	2.08%	-
	32	10^1	2.38%	2.08%	1.34%
μ_1	16	10^0	22.77%	15.62%	-
	16	10^1	27.68%	22.17%	-
	32	10^1	22.47%	16.82%	17.26%

Table 2: TIMIT test phone error rate of acoustic models trained on different features and sets

Train Set and Configuration			Test PER by Set		
ASR	FHVAE	Features	Male	Female	All
Train All	-	FBank	20.1%	16.7%	19.1%
Train Male	-	FBank	21.0%	32.8%	25.2%
	Train All, $\alpha = 10$	z_1	22.0%	26.2%	23.5%

Table 3: Aurora-4 test word error rate of acoustic models trained on different features and sets

Train Set and Configuration			Test WER by Set				
ASR	{FH-, β -}VAE	Features	Clean	Noisy	Channel	NC	All
Train All	-	FBank	3.60%	7.06%	8.24%	18.49%	11.80%
Train Clean	-	FBank	3.47%	50.97%	36.99%	71.80%	55.51%
	Dev, $\beta = 1$	z (β -VAE)	4.95%	23.54%	31.12%	46.21%	32.47%
	Dev, $\beta = 2$	z (β -VAE)	3.57%	27.24%	30.56%	48.17%	34.75%
	Dev, $\beta = 4$	z (β -VAE)	3.89%	24.40%	29.80%	47.87%	33.38%
	Dev, $\beta = 8$	z (β -VAE)	5.32%	34.84%	36.13%	58.02%	42.76%
	Dev, $\alpha = 10$	z_1 (FHVAE)	5.01%	16.42%	20.29%	36.33%	24.41%
	Dev, $\alpha = 10$	z_2 (FHVAE)	41.08%	68.73%	61.89%	86.36%	72.53%
	Dev\NC, $\alpha = 10$	z_1 (FHVAE)	5.25%	16.52%	19.30%	40.59%	26.23%

RELATED WORKS

- A Deep Generative Model for Disentangled Representations of Sequential Data
 - Applied almost the same architecture to video



(a) random test data sequences



(b) reconstruction



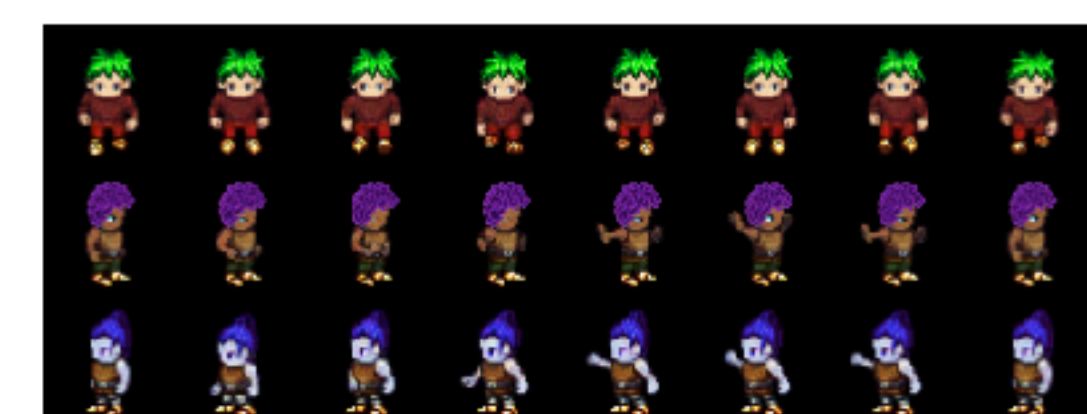
(c) reconstruction with randomly sampled f



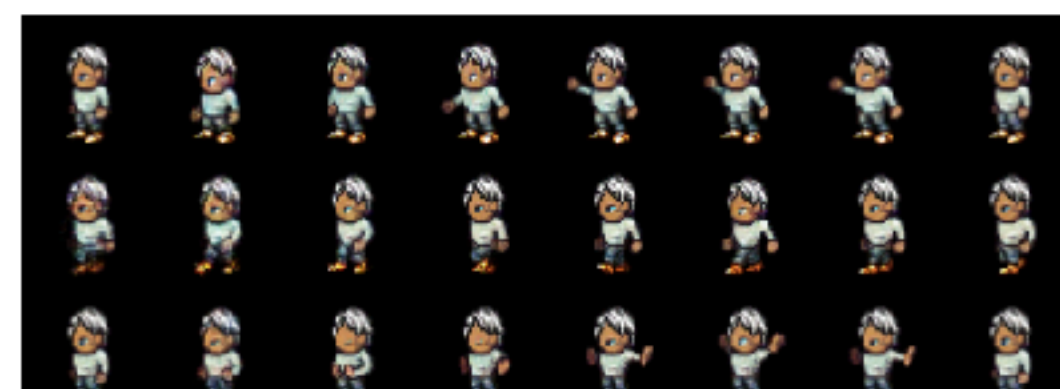
(d) reconstruction with randomly sampled $z_{1:T}$



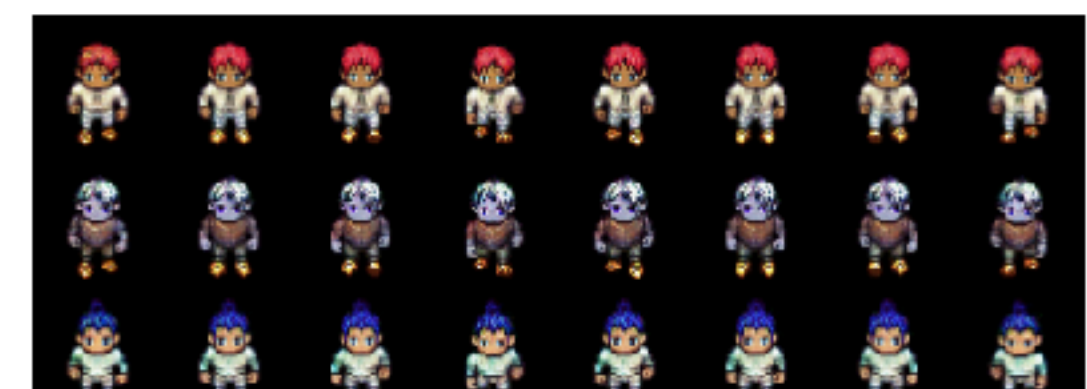
(e) reconstruction with swapped encoding f



(f) reconstruction with swapped encoding $z_{1:T}$



(g) generated sequences with fixed f



(h) generated sequences with fixed $z_{1:T}$

RELATED WORKS

- Multimodal Unsupervised Image-to-Image Translation
 - Used Adversarial training to disentangle style and content feature

