

Hidden Markov Model Variational Autoencoder for Acoustic Unit Discovery

Janek Ebberts¹, Jahn Heymann¹, Lukas Drude¹, Thomas Glarner¹, Reinhold Haeb-Umbach¹, Bhiksha Raj²,
¹ Paderborn University, Germany ²Carnegie Mellon University, United States

Introduction

- Acoustic unit discovery (AUD)
 - ▶ Learning AUs (phonetic inventory) from raw speech
 - ▶ Unsupervised training of generative model
 - ▶ SOTA: GMM/HMM
- Known from ASR: Superiority of DNNs over GMMs
 - ▶ But: Discriminative DNNs not transferable to AUD
- Variational Autoencoder (VAE) [1]
 - ▶ Deep generative model
 - ▶ Sophisticated data distribution modeling by DNN
 - ▶ Efficient variational inference by DNN
- **Here:** Marrying VAE with HMM for AUD with sophisticated emission distribution modeling

VAE

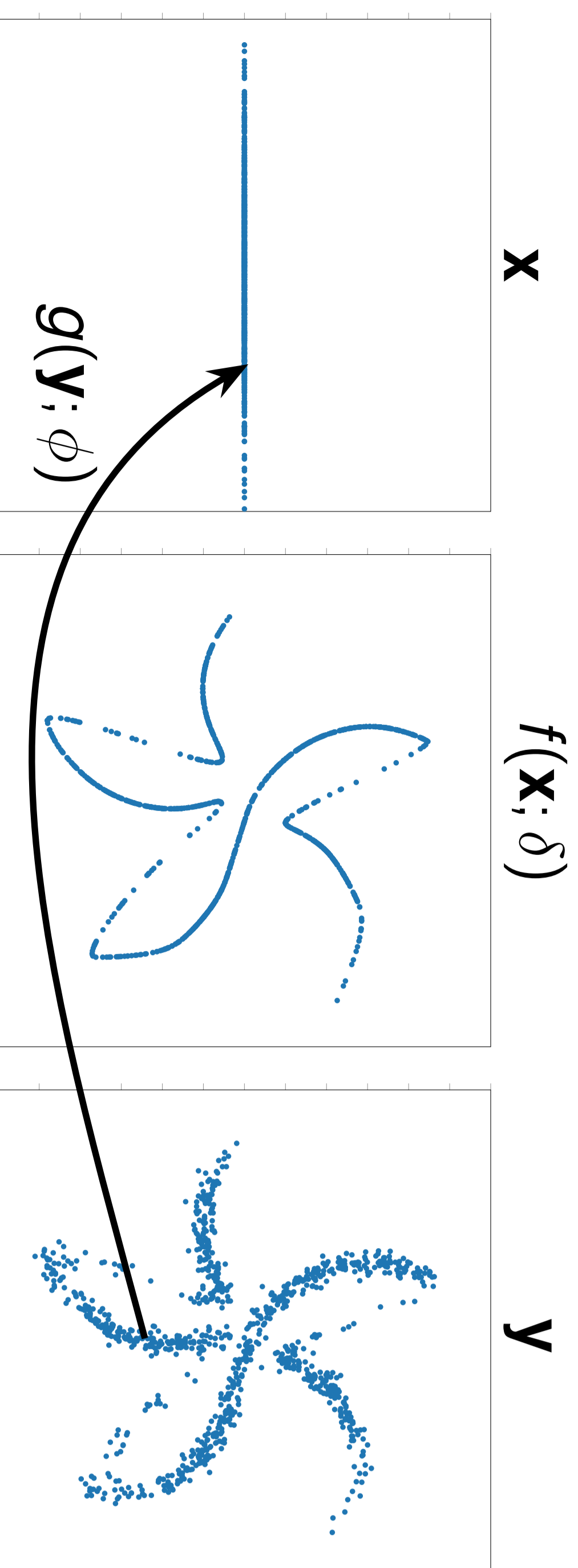
- Given: arbitrary distributed observations \mathbf{y}
- Assuming: simply structured latent codes $\mathbf{x} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$
- Modeling observations by decoder $f(\mathbf{x}; \delta)$ and Gaussian observation noise \mathbf{v} :

$$\mathbf{y} = f(\mathbf{x}; \delta) + \mathbf{v}; \quad \mathbf{v} \sim \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I})$$

$$\Rightarrow p(\mathbf{y}|\mathbf{x}; \delta) = \mathcal{N}(\mathbf{y}; f(\mathbf{x}; \delta), \sigma^2 \mathbf{I})$$
- Variational inference by encoder $g(\mathbf{y}; \phi)$:

$$q(\mathbf{x}|\mathbf{y}; \phi) = \mathcal{N}(\mathbf{x}; \mu_{\mathbf{x}|\mathbf{y}}, \Sigma_{\mathbf{x}|\mathbf{y}})$$

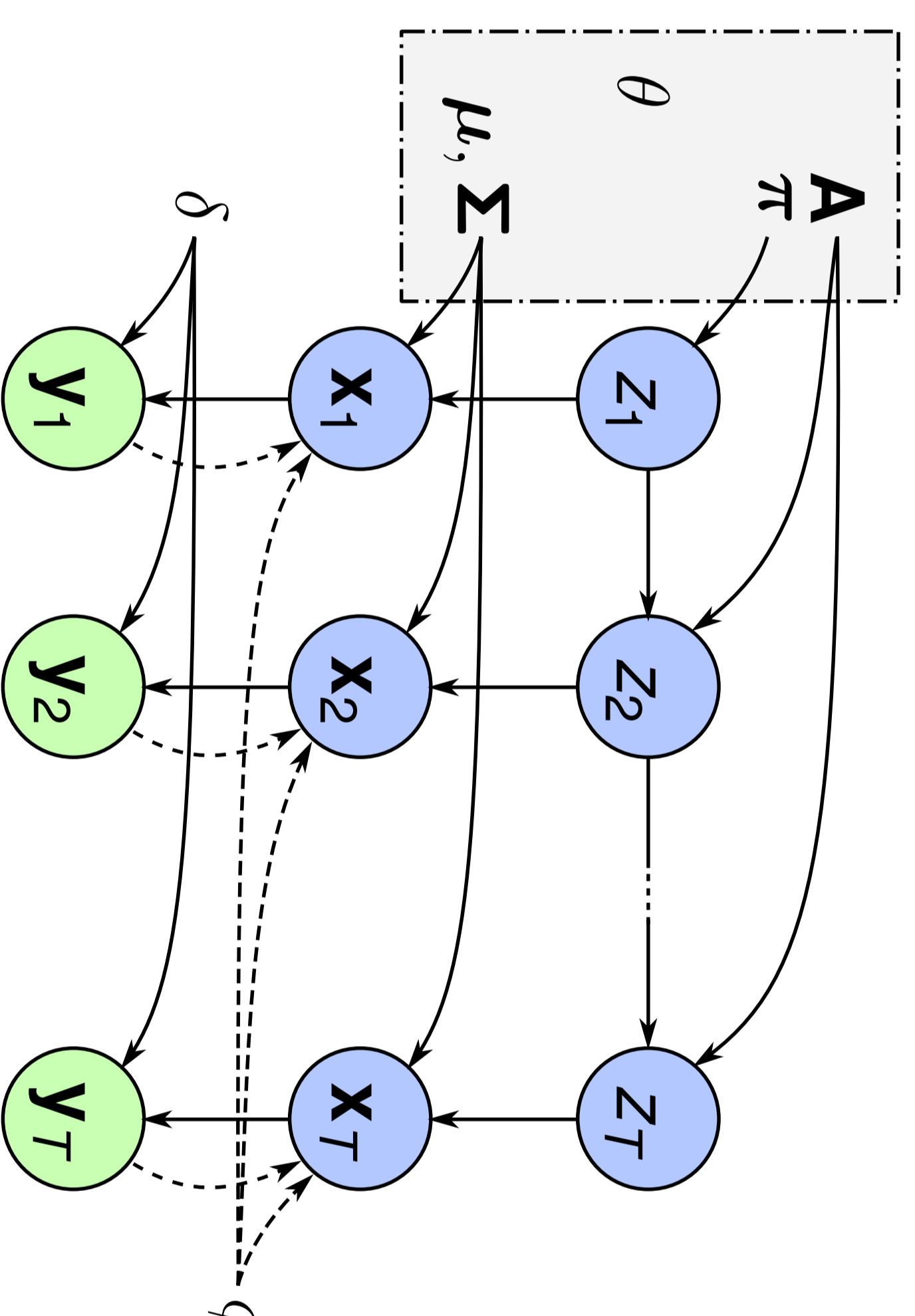
$$(\mu_{\mathbf{x}|\mathbf{y}}, \ln \sigma_{\mathbf{x}|\mathbf{y}}) = g(\mathbf{y}; \phi)$$



HMMVAE

- Introducing latent classes (states) z
 - Class specific codes: $p(\mathbf{x}|z; \theta) = \mathcal{N}(\mathbf{x}; \mu_z, \Sigma_z)$
-

- Modeling temporal correlations by HMM:



Inference & Training

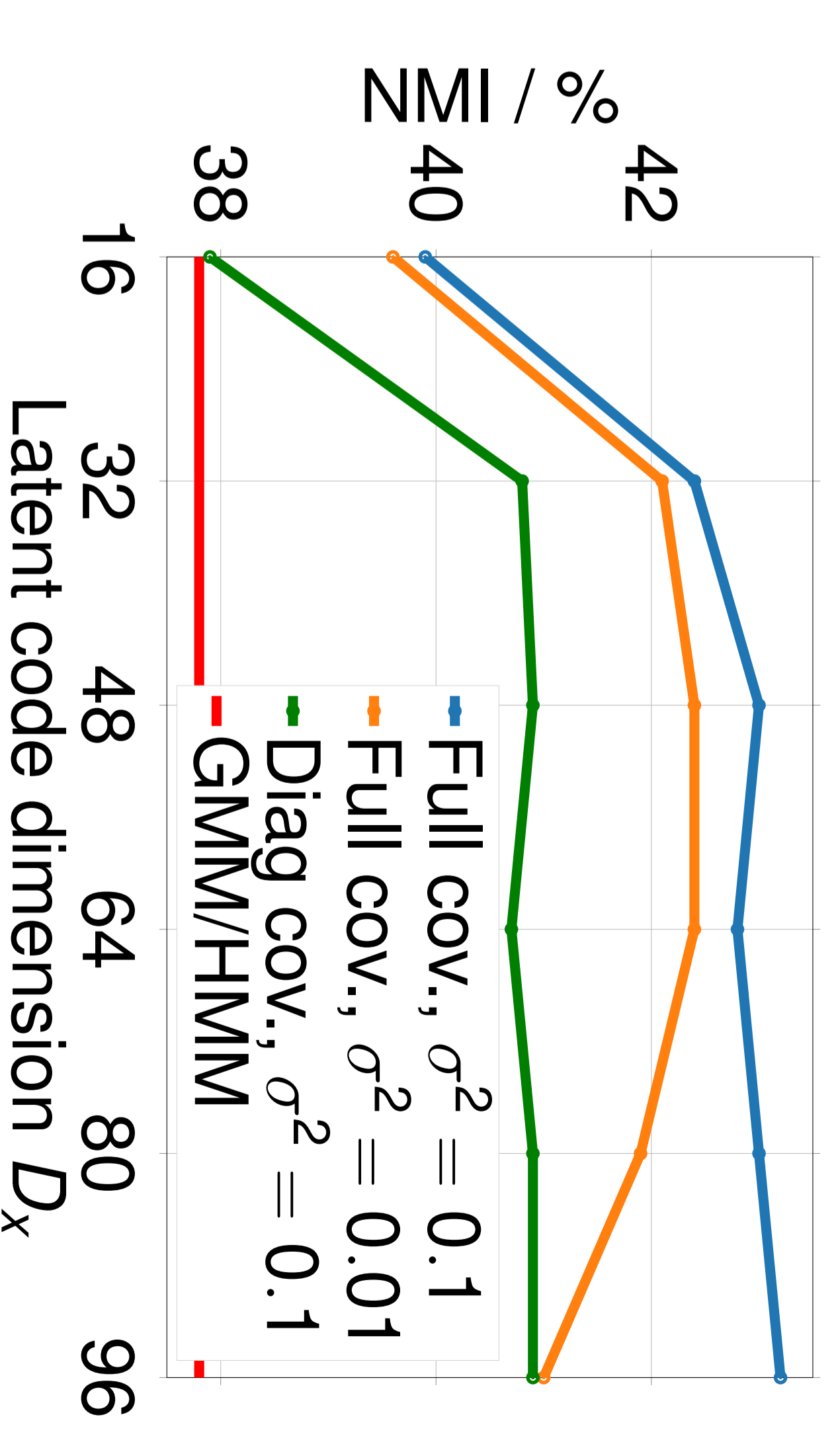
- Acoustic score: $\ln b_z(\mathbf{y}) = -H(q(\mathbf{x}|\mathbf{y}), p(\mathbf{x}|z))$ enables forward-backward and Viterbi:
- Joint training of all parameters θ, δ, ϕ
- Objective: Maximization of variational lower bound

$$\mathcal{L}(\mathbf{Y}; \theta, \delta, \phi) = \mathbb{E}_{q(\mathbf{x}|\mathbf{y}; \phi)} [\ln p(\mathbf{Y}|\mathbf{X}; \delta)] - \text{KL}(q(\mathbf{X}, \mathbf{Z}|\mathbf{Y}; \phi) \| p(\mathbf{X}, \mathbf{Z}; \theta))$$
- Optimizer: Adam with stepsize $\alpha=10^{-3}$

Experiments

- Goal: discovering latent acoustic units
- Database: Timit
 - ▶ 4620 train sentences (3.14h), 1680 test sentences (0.81h)
- Features: 13 element MFCCs with Δ and $\Delta\Delta$
- Initialization:
 - ▶ Using segmentations found by unsupervised GMM/HMM [2]
 - ▶ 72 AUs, each modeled by three states (left-right topology)

Model	Train.	NMI	eq. PER	PER
GMM/HMM	FB	37.8%	65.4%	Performance measures
HMMVAE	Viterbi	42.8%	58.9%	Normalized mutual information (NMI)
HMMVAE	FB	42.6%	59.0%	Equivalent phone error rate (eq. PER)



Conclusions

- Extended VAE by an HMM in latent code space to capture temporal correlations
- Iterative EM-like algorithm for inference and optimization
- Applied HMMVAE to unsupervised AUD task
- Significantly improved AUD performance over variational GMM/HMM in terms of NMI and eq. PER
- Outlook: Bayesian parameter estimation

References

- [1] "Auto-Encoding Variational Bayes", D. P. Kingma and M. Welling
[2] "Variational Inference for Acoustic Unit Discovery", L. Ondel, L. Burget, and J. Cernocky