

Disentangled Dynamic Representations from Unordered Data

ETH zürich

Leonhard Helminger¹ Abdelaziz Djelouah Markus Gross¹ Romann M. Weber

¹ETH Zurich

Motivation

- Better interpretability through disentangled representation
- In context of video: static and dynamic information
- Pairwise comparison instead of recurrent architecture

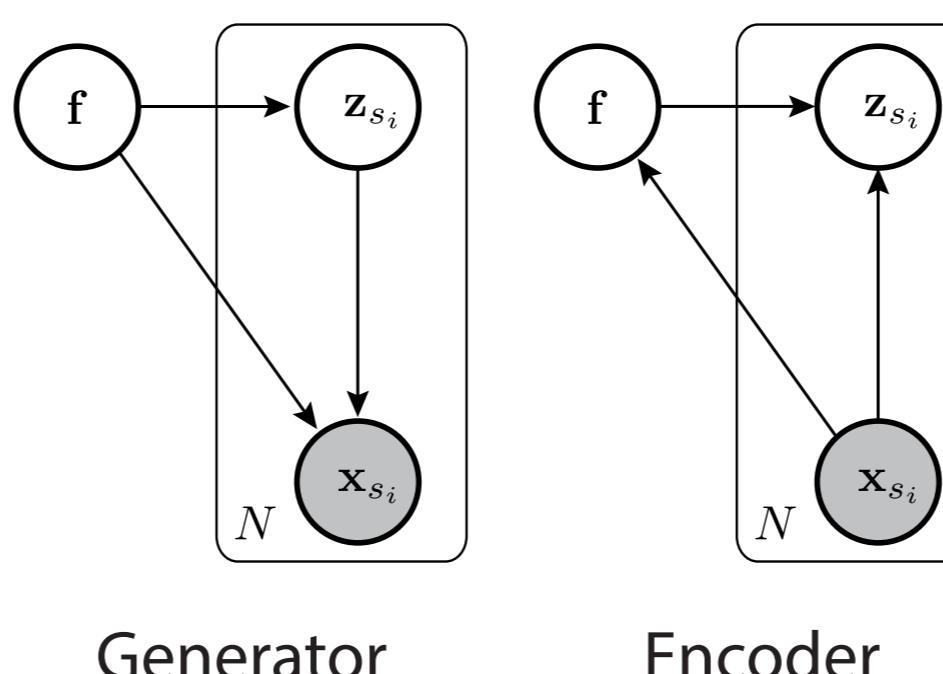
Contribution

We introduce a deep generative model that learns disentangled static and dynamic representations from data without temporal ordering

Model

Generative Model:

$$p_{\theta}(\mathbf{x}_{1:T}, \mathbf{f}, \mathbf{z}_{1:T}) = p_{\theta}(\mathbf{f}) \prod_{t=1}^T p_{\theta}(\mathbf{x}_t | \mathbf{f}, \mathbf{z}_t) p_{\theta}(\mathbf{z}_t | \mathbf{f})$$



Inference Model:

$$q_{\phi}(\mathbf{f}, \mathbf{z}_{s_1:N} | \mathbf{x}_S) = q_{\phi}(\mathbf{f} | \mathbf{x}_S) \prod_{i=1}^N q_{\phi}(\mathbf{z}_{s_i} | \mathbf{f}, \mathbf{x}_i)$$

Learning:

$$\mathcal{L} = \mathbb{E}_{q_{\phi}(\mathbf{f} | \mathbf{x}_S)} \left[\sum_{\mathbf{x} \in \mathbf{x}_S} \mathbb{E}_{q_{\phi}(\mathbf{z} | \mathbf{f}, \mathbf{x})} [\log p_{\theta}(\mathbf{x} | \mathbf{f}, \mathbf{z}) - D_{KL}(q_{\phi}(\mathbf{z} | \mathbf{f}, \mathbf{x}) || p_{\theta}(\mathbf{z} | \mathbf{f}))] \right] - D_{KL}(q_{\phi}(\mathbf{f} | \mathbf{x}_S) || p_{\theta}(\mathbf{f}))$$

Distributions of the generative model:

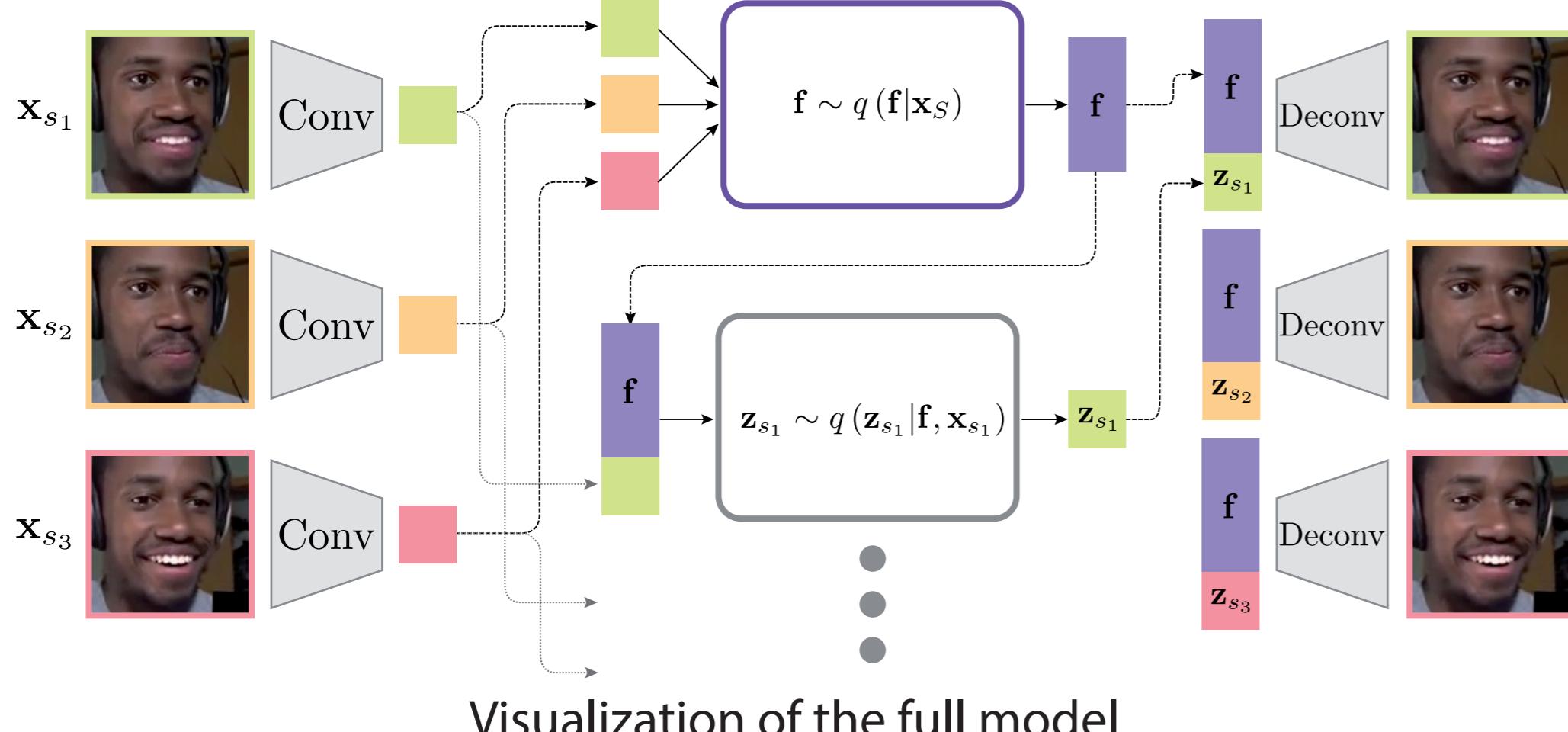
$$\begin{aligned} p_{\theta}(\mathbf{f}) &= \mathcal{N}(\mathbf{f} | \mathbf{0}, \mathbf{I}), \\ p_{\theta}(\mathbf{z} | \mathbf{f}) &= \mathcal{N}(\mathbf{z} | h_{\mu_z}(\mathbf{f}), \text{diag}(h_{\sigma_z^2}(\mathbf{f}))) \\ p_{\theta}(\mathbf{x} | \mathbf{f}, \mathbf{z}) &= \text{Ber}_p(\mathbf{x} | h_x(\mathbf{f}, \mathbf{z})) \end{aligned}$$

Distributions of the inference model:

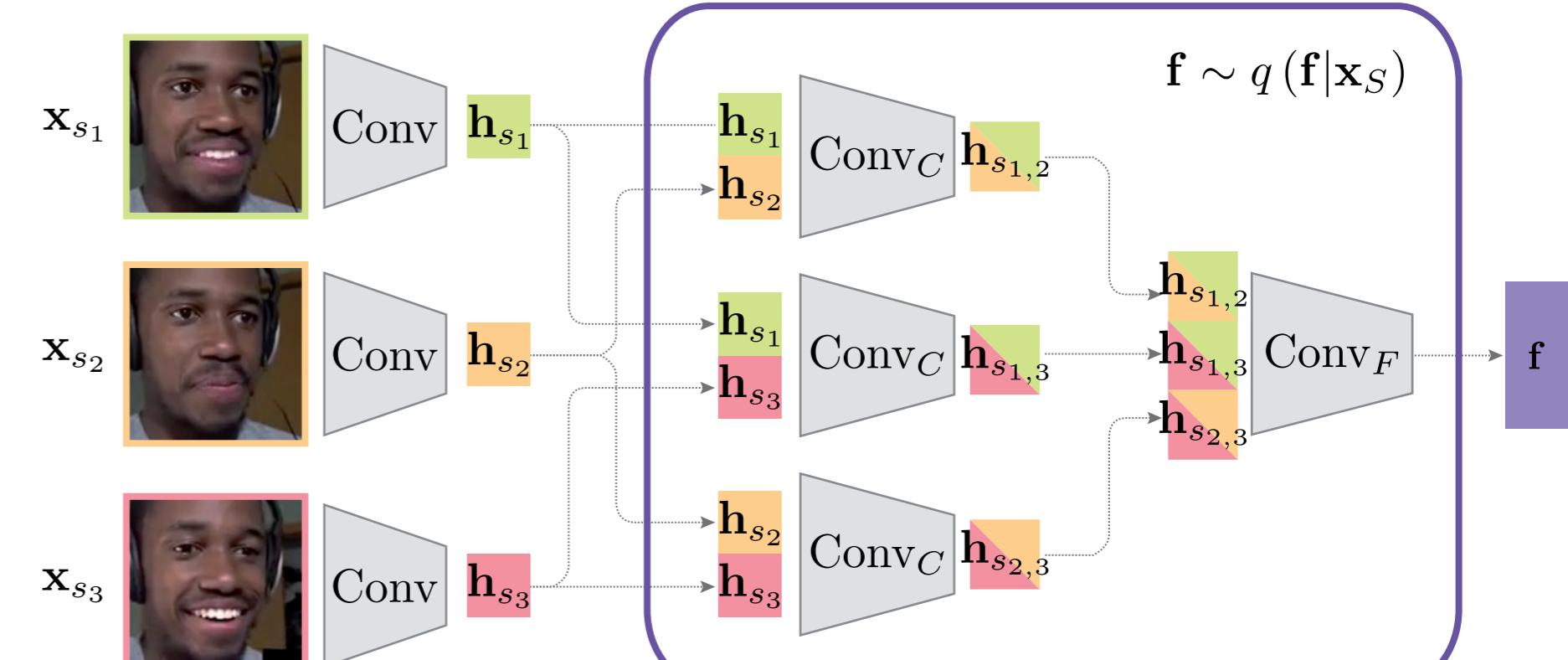
$$\begin{aligned} q_{\phi}(\mathbf{f} | \mathbf{x}_S) &= \mathcal{N}(\mathbf{f} | g_{\mu_f}(\mathbf{x}_S), \text{diag}(g_{\sigma_f^2}(\mathbf{x}_S))), \\ q_{\phi}(\mathbf{z} | \mathbf{f}, \mathbf{x}) &= \mathcal{N}(\mathbf{z} | g_{\mu_z}(\mathbf{f}, \mathbf{x}), \text{diag}(g_{\sigma_z^2}(\mathbf{f}, \mathbf{x}))) \end{aligned}$$

where $\mathbf{x}_S = (x_{s_1}, \dots, x_{s_N})$, $N = 3$

Details



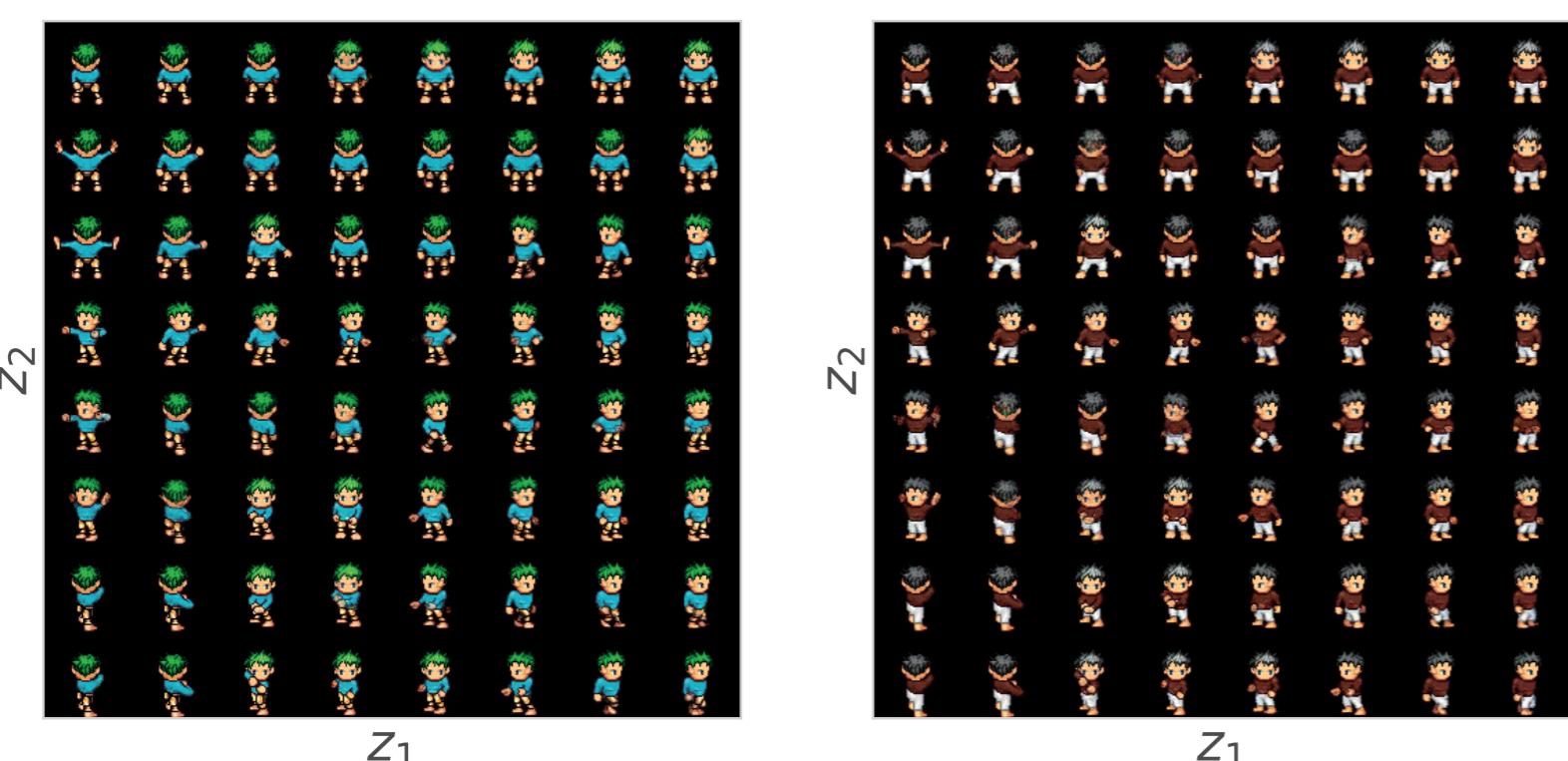
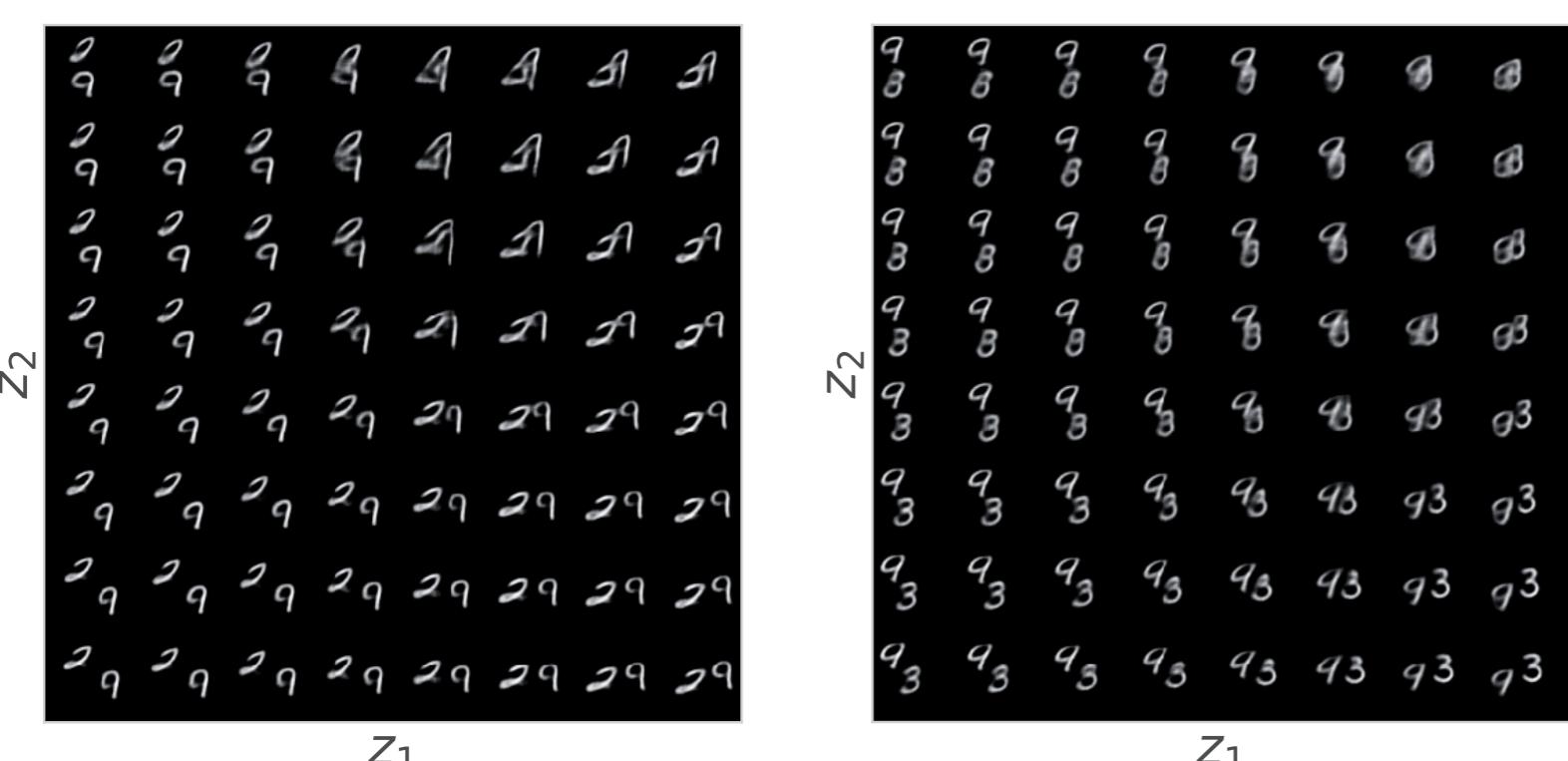
Visualization of the full model



Visualization of the encoder for the static latent variable model

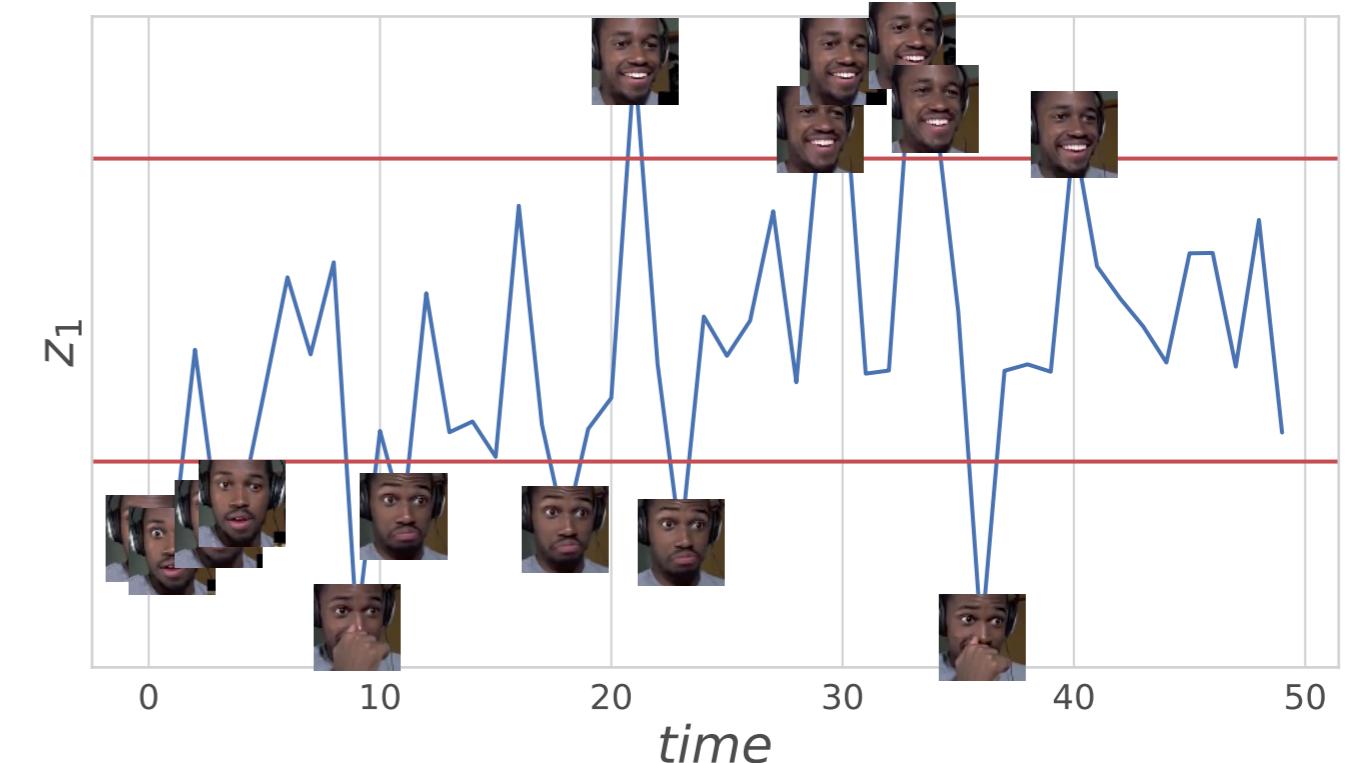
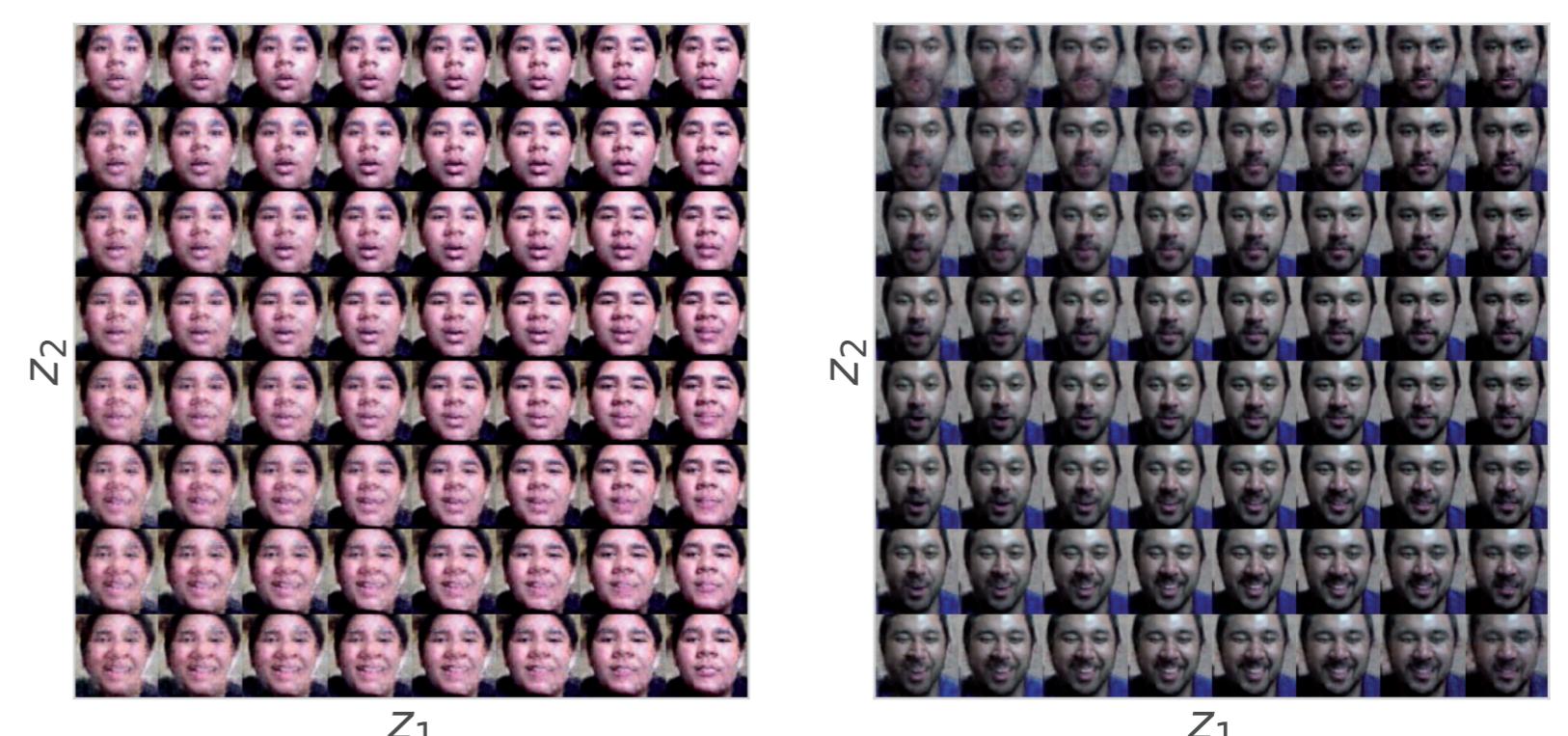
Results

Learned dynamic data manifold of the 2d latent space



- Digits and style of the handwritten numbers are consistent over the spanned latent space
- Coherence of the dynamic space between different identities

Naive emotion detection in a 2d latent space



- Two-dimensional latent space captures dynamics of faces
- Naive analysis of the dynamic plot can already extract some meaningful facial expressions (smiling / astonished)