

Deep Generative Models I

Variational Autoencoders – Part II

Machine Perception

Otmar Hilliges

23 April 2020

Last week(s)

Intro to generative modelling

VAE

Derivation of the ELBO

This week

More generative models

Step I: Variational autoencoders:

- Training
- Representation learning & disentanglement
- Applications

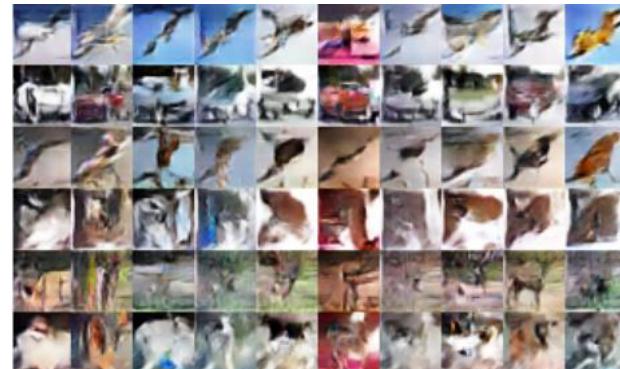
Step II: Generative adversarial networks

Generative Modelling

Given training data, generate new samples, drawn from “same” distribution



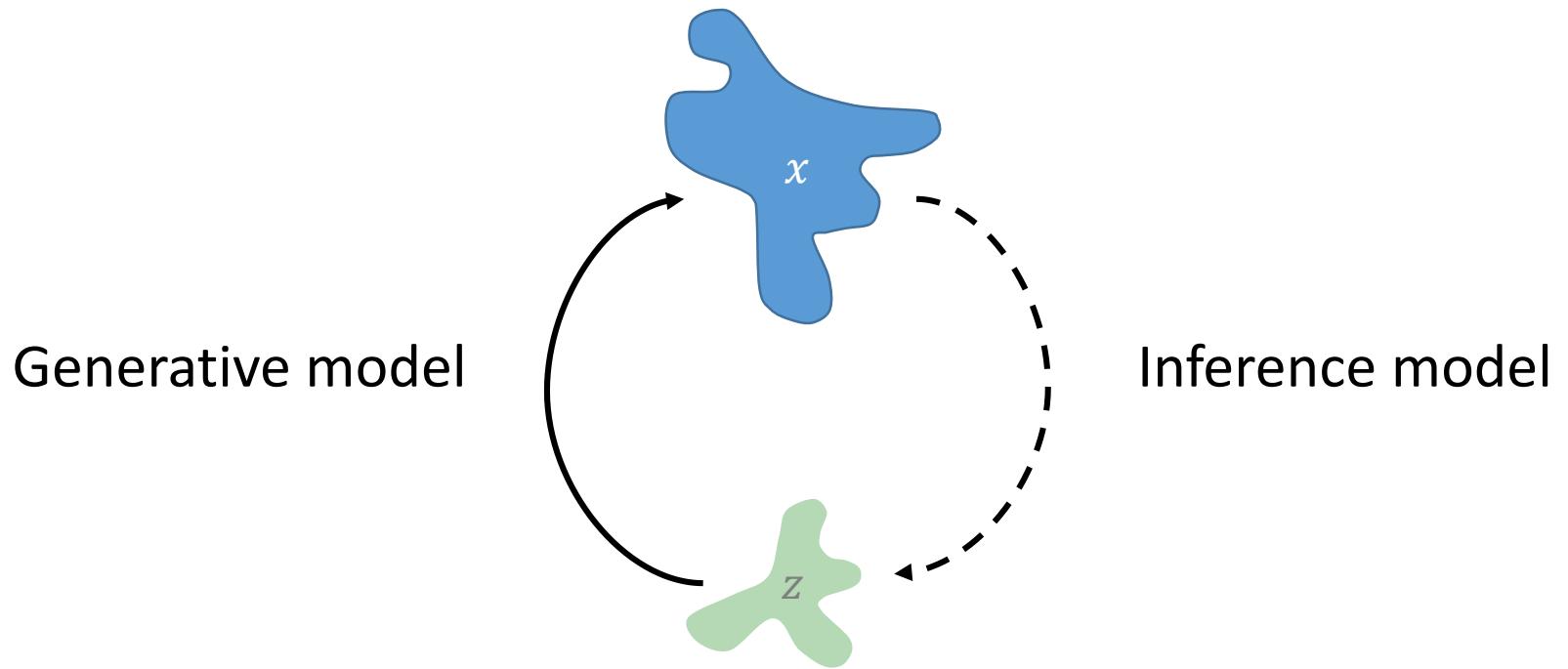
Training samples $\sim p_{data}(x)$



Generated samples $\sim p_{model}(x)$

We want $p_{model}(x)$ to be similar to $p_{data}(x)$

Variational autoencoder



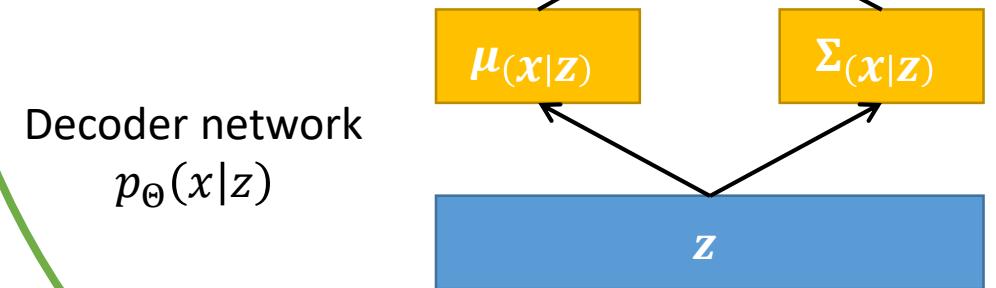
Recap: ELBO

$$\underbrace{E_z[\log p_\Theta(x^{(i)}|z)]}_{\mathcal{L}(x^{(i)}, \Theta, \phi)} - \underbrace{D_{KL}(q_\phi(z|x^{(i)}) \| p_\Theta(z^{(i)}))}_{\text{Make approximate posterior similar to prior}}$$

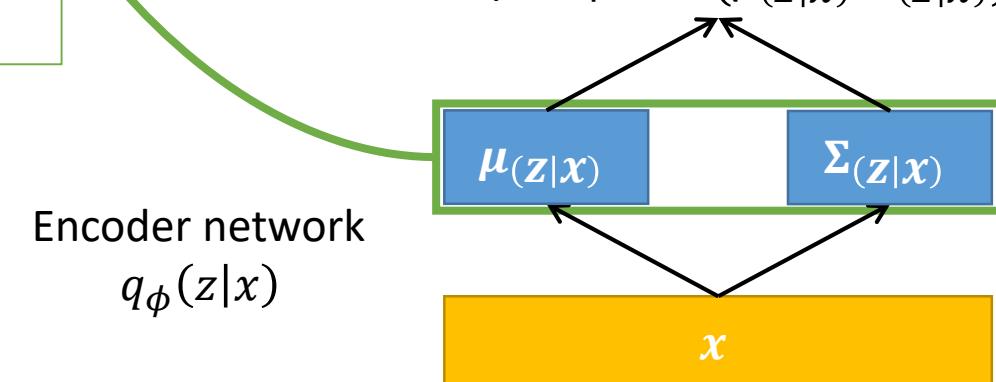
Maximize reconstruction likelihood



Sample $x|z \sim \mathcal{N}(\mu_{(x|z)}, \Sigma_{(x|z)})$



Sample $z|x \sim \mathcal{N}(\mu_{(z|x)}, \Sigma_{(z|x)})$



Encoder network
 $q_\phi(z|x)$

Decoder network
 $p_\Theta(x|z)$

Make approximate posterior
Similar to prior

Training VAEs

Training of VAEs is just backprop.

But wait, how do gradients flow through z ?

Or any other random operation (requiring sampling)?

(see blackboard)

reparametrization trick:

Given

$$z \sim N(\mu, \sigma^2) \Rightarrow \text{value changes every time we sample it}$$

Assume \exists underlying random variable

$$\epsilon \sim N(0, 1)$$

$z = \mu + \sigma \epsilon$ is still random
but does not depend on μ or σ

$$z = f(x, \epsilon, \Theta)$$

\Rightarrow can take deriv wrt μ, σ .

tells us how a small change to μ, σ
affects output if we keep ϵ fixed

Data log-likelihood

$$\log p_{\theta}(x^{(i)}) =$$

$$= \mathbb{E}_{z \sim q_{\phi}(z|x^{(i)})} [\log p_{\theta}(x^{(i)})] \quad (p_{\theta}(\cdot) \text{ doesn't depend on } z)$$

$$= \mathbb{E}_z \left[\log \frac{p_{\theta}(x^{(i)}|z) p_{\theta}(z)}{p_{\theta}(z|x^{(i)})} \right] \quad (\text{Bayes' rule})$$

$$= \mathbb{E}_z \left[\log \frac{p_{\theta}(x^{(i)}|z) p_{\theta}(z)}{p_{\theta}(z|x^{(i)})} \frac{q_{\phi}(z|x^{(i)})}{q_{\phi}(z|x^{(i)})} \right] \quad (\text{mult by constant})$$

$$= \mathbb{E}_z [\log p_{\theta}(x^{(i)}|z)] + \mathbb{E}_z [\log \frac{p_{\theta}(z)}{q_{\phi}(z|x^{(i)})}] + \mathbb{E}_z [\log \frac{q_{\phi}(z|x^{(i)})}{p_{\theta}(z|x^{(i)})}] \quad (\text{product rule logarithms})$$

(continued on next slide)

$$= \mathbb{E}_z [\log p_\theta(x^{(i)}|z)] - \mathbb{E}_z [\log \frac{q_\phi(z|x^{(i)})}{p_\theta(z)}] + \mathbb{E}_z [\log \frac{q_\phi(z|x^{(i)})}{p_\theta(z|x^{(i)})}] \quad (\log \frac{z}{y} = -\log \frac{y}{x})$$

\checkmark (with sampling)

\checkmark closed-form if Gaussian

$$= \mathbb{E}_z [\log p_\theta(x^{(i)}|z)] - D_{KL}(q_\phi(z|x^{(i)})||p_\theta(z)) + D_{KL}(q_\phi(z|x^{(i)})||p_\theta(z|x^{(i)}))$$

\uparrow
NN tractable

\uparrow
tractable

≥ 0 (see proof earlier)

"make approx. posterior
as similar as possible
to prior"

目的是: opt. only tractable

term: q_ϕ, p_θ

不需要处理 $p(x)$, $p(z|x)$

loss function $\mathcal{L}(x^{(i)}, \theta, \phi)$

$$\log(p_\theta(x^{(i)})) \geq \mathcal{L}$$

evidence lower bound "ELBO"

\Rightarrow data is at least as likely as \mathcal{L}

Training procedure:

$$\theta^*, \phi^* = \operatorname{argmax} \sum_i \mathcal{L}(x^{(i)}, \theta, \phi)$$

Gradient computation for backprop

$$\nabla_{\theta, \phi} \mathbb{E}_{z \sim q_\phi} \left[\log \frac{p_\theta(x, z)}{q_\phi(z|x^{(i)})} \right] = \mathbb{E}_z \left[\log p_\theta(x^{(i)}, z) \right] + \mathbb{E}_z \left[\log \frac{p_\theta(z)}{q_\phi(z|x^{(i)})} \right]$$

f is NN (deterministic)

$$= \nabla_{\theta, \phi} \mathbb{E}_{\epsilon \sim N(0, 1)} \left[\log \frac{p_\theta(x, f(x, \epsilon, \theta))}{q_\phi(f(x, \epsilon, \phi)|x)} \right]$$

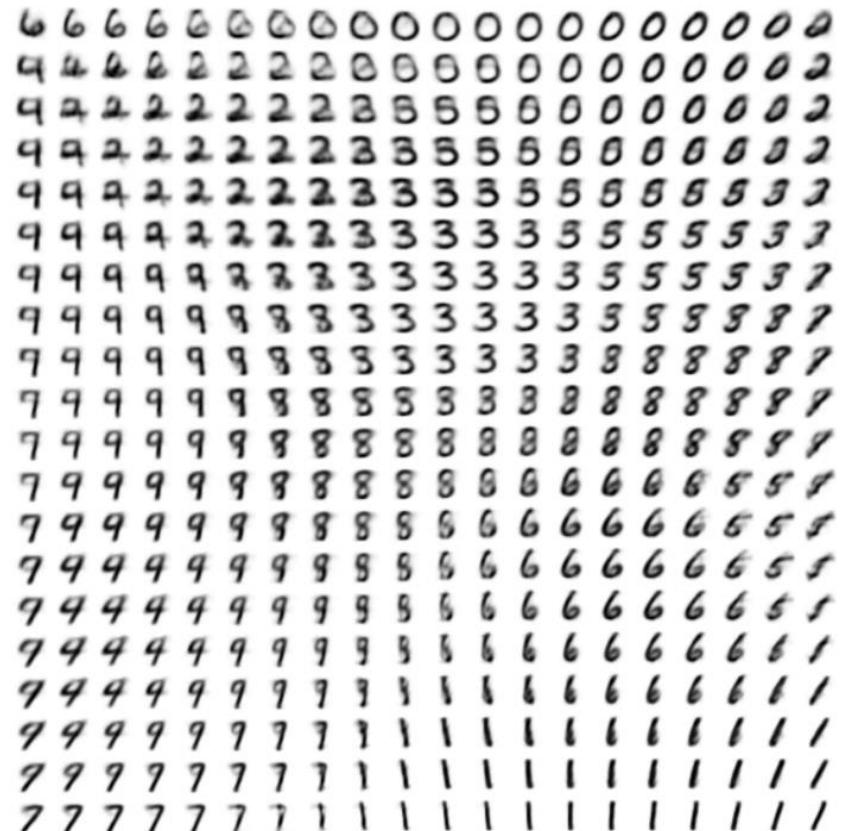
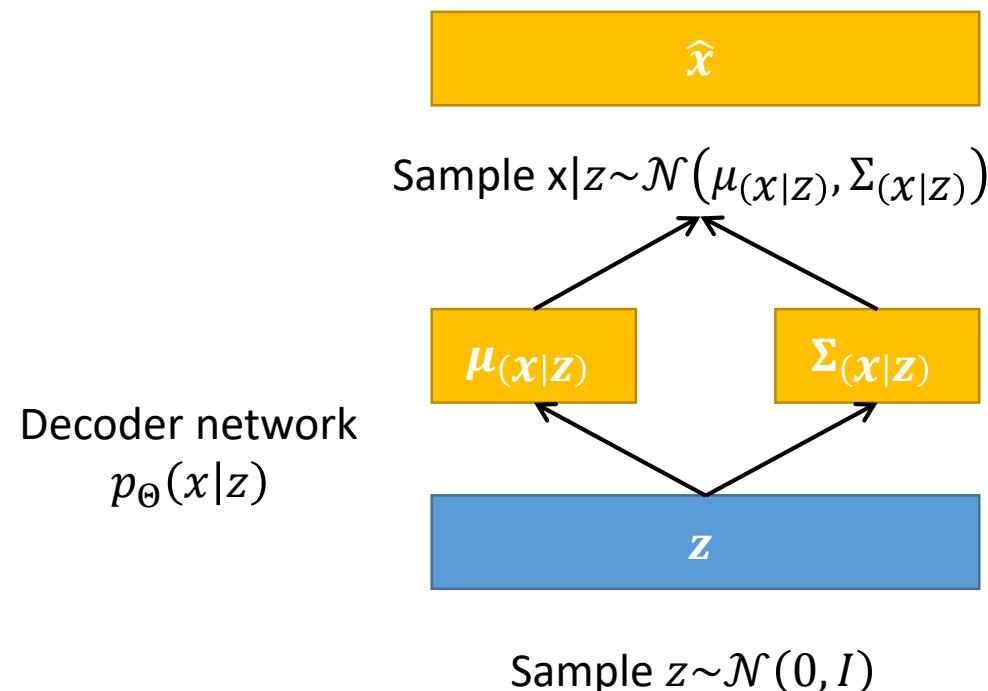
$$= \mathbb{E}_\epsilon \left[\nabla_{\theta, \phi} \log \left(\frac{p_\theta(x, f)}{q_\phi(f|x)} \right) \right]$$

$\sum = \bar{w}$
deterministic NN

$$\approx \frac{1}{k} \sum_{i=1}^k \nabla_{\theta, \phi} \log w(x_i, f_i, \theta, \phi)$$

Generating Data

At runtime only use decoder network. Sample z from prior.



Cross-modal Deep Variational Hand Pose Estimation

Adrian Spurr, Jie Song, Seonwook Park, Otmar Hilliges

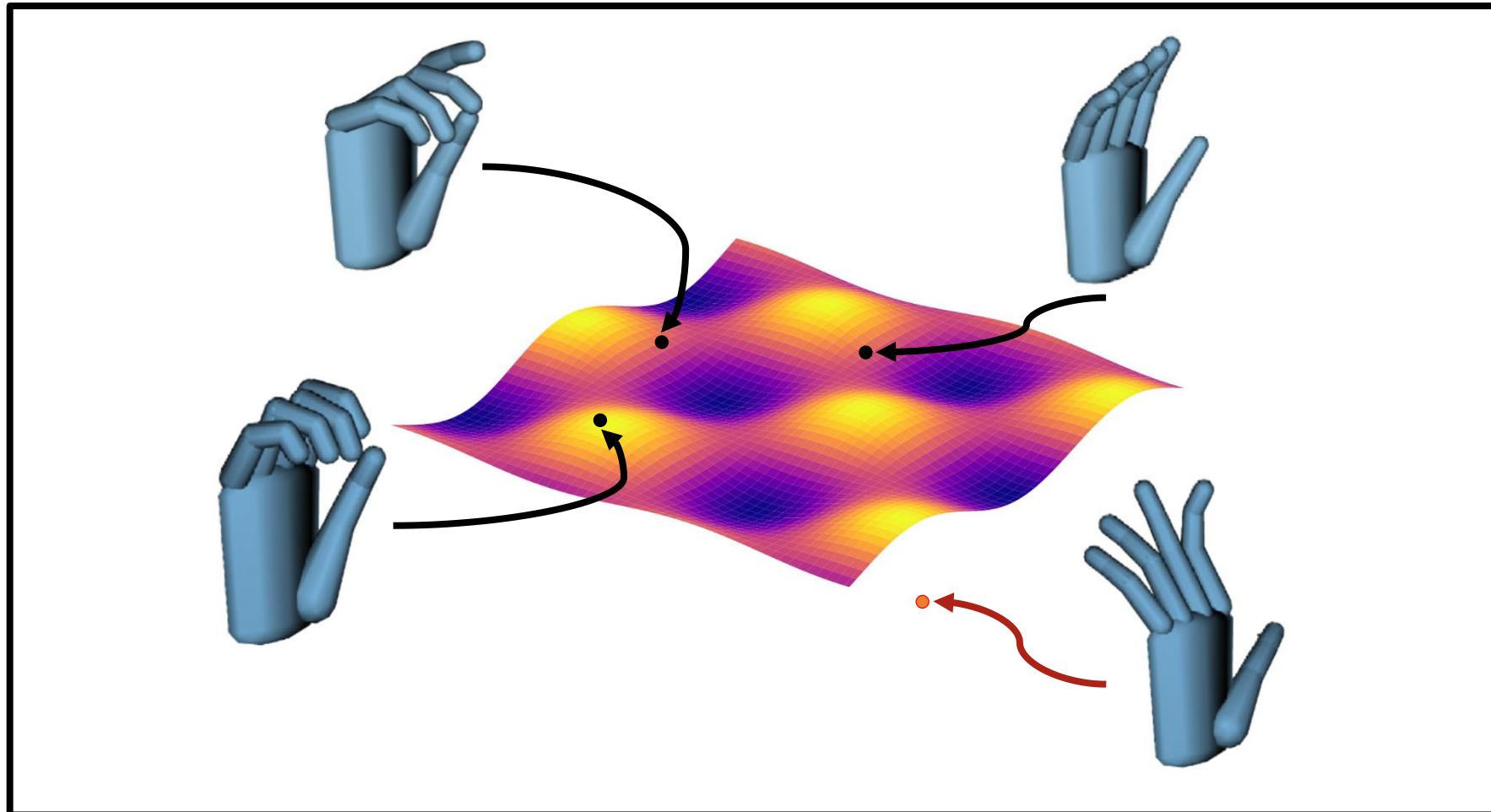
CVPR 2018

RGB hand pose estimation

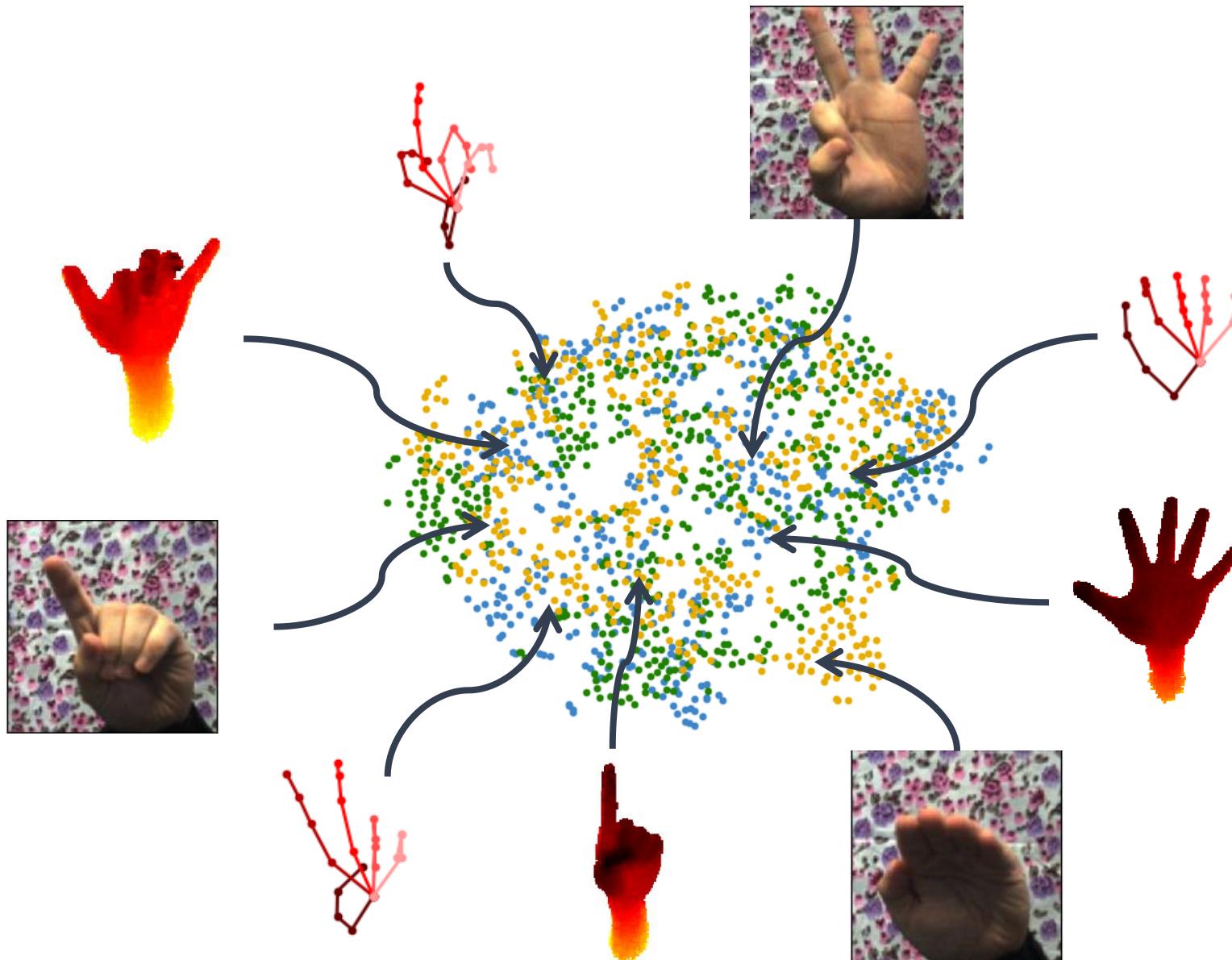


Assumption: Manifold of valid hand poses

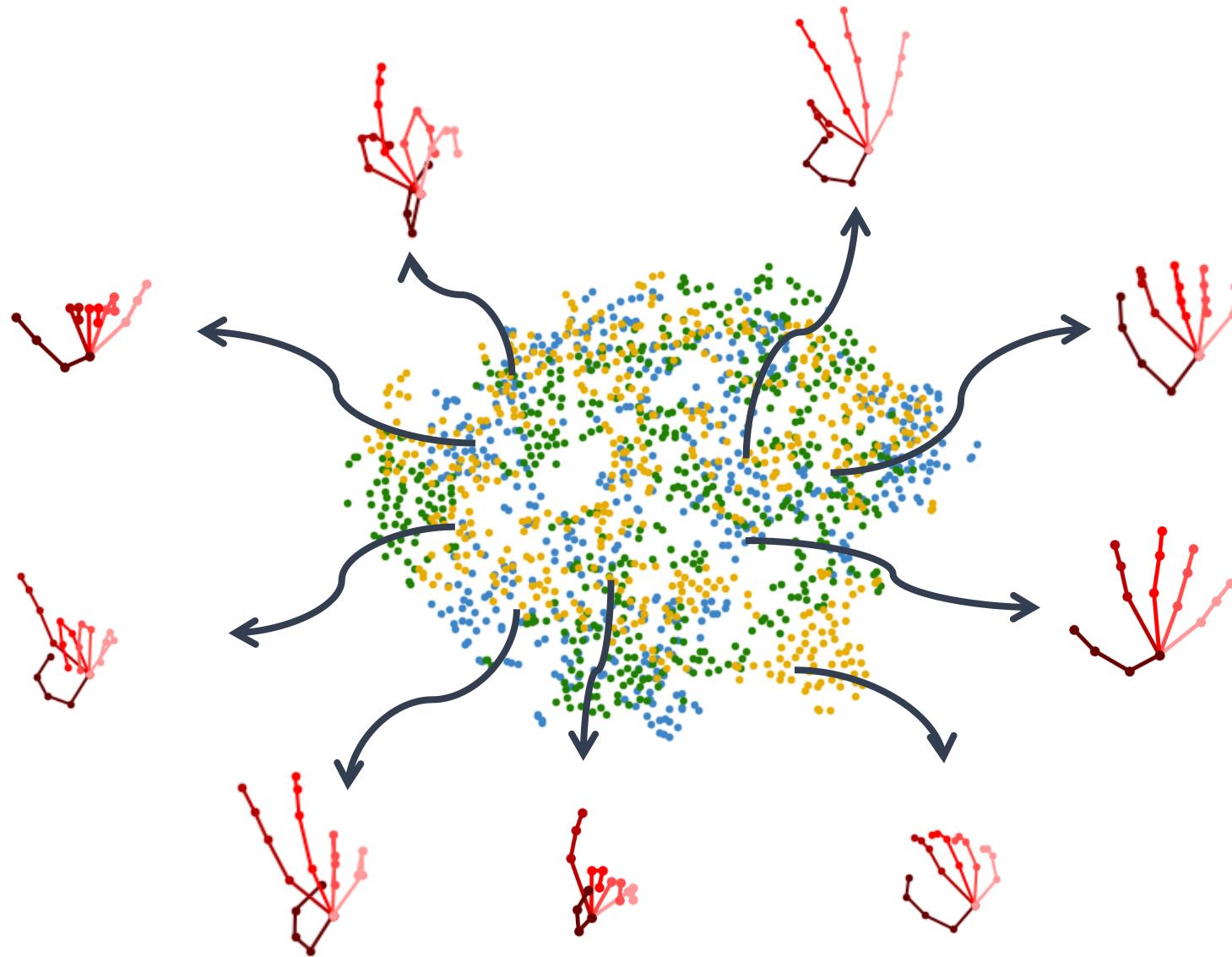
exist a low dimensional space that contains all valid hand pose

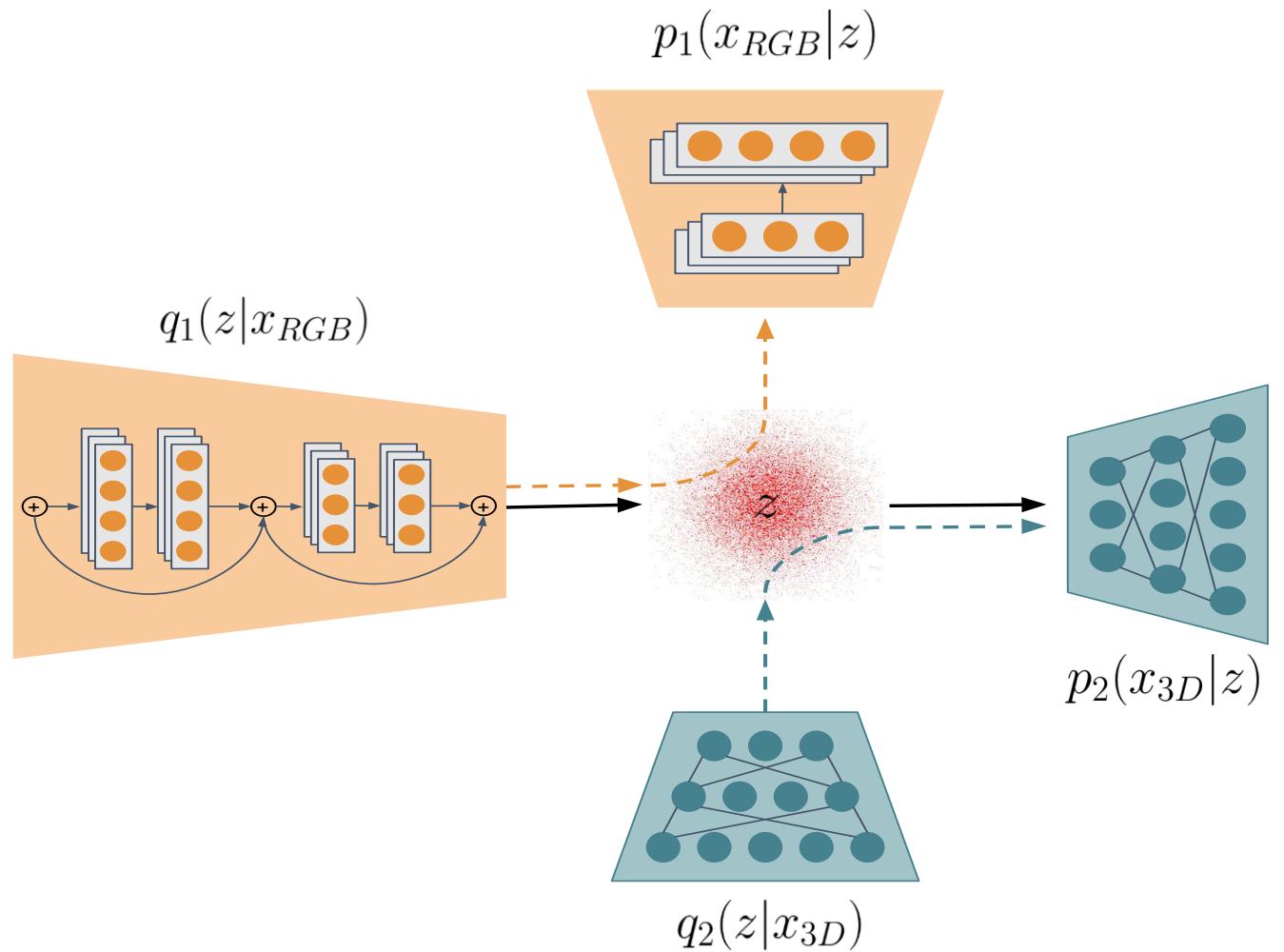


Cross-modal latent space – Inference model



Cross-modal latent space – Generative model





Cross-modal training objective

Original VAE considers only one modality → need to modify to consider multiple modalities

Consider x_i and x_t samples from different modalities (e.g RGB and 3D joint positions)

Important: Both describe inherently the same thing. In our case, the hand pose.

Task: maximize $\log p_\theta(x_t)$ under x_i

Cross-modal training objective derivation

$$\mathbb{E}_{z \sim q(z|x_i)} [\log p(x_t)]$$

$$\log p(\mathbf{x}_t) = \int_z q(z|\mathbf{x}_i) \log p(\mathbf{x}_t) dz = \int_z q(z|\mathbf{x}_i) \log \frac{p(\mathbf{x}_t)p(z|\mathbf{x}_t)q(z|\mathbf{x}_i)}{p(z|\mathbf{x}_t)q(z|\mathbf{x}_i)} dz \quad (\text{Bayes' Rule})$$

$$= \int_z q(z|\mathbf{x}_i) \log \frac{q(z|\mathbf{x}_i)}{p(z|\mathbf{x}_t)} dz + \int_z q(z|\mathbf{x}_i) \log \frac{p(\mathbf{x}_t)p(z|\mathbf{x}_t)}{q(z|\mathbf{x}_i)} dz \quad (\text{Logarithms})$$

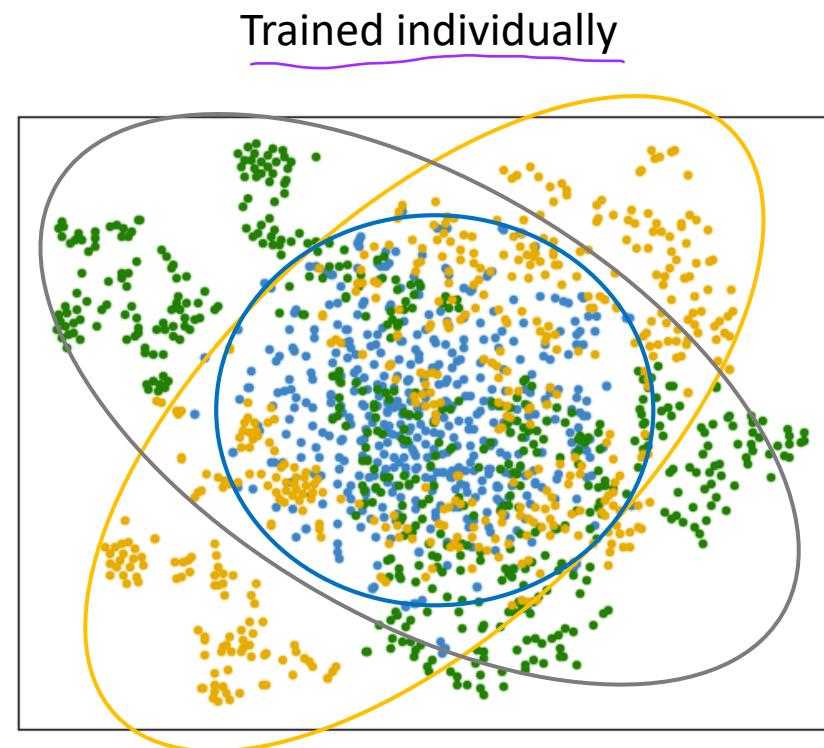
$$= D_{KL}(q(z|\mathbf{x}_i) || p(z|\mathbf{x}_t)) + \int_z q(z|\mathbf{x}_i) \log \frac{\cancel{p(\mathbf{x}_t|z)p(z)}}{q(z|\mathbf{x}_i)} dz \quad (\text{Definition KL})$$

~~≥ 0~~

$$\geq \int_z q(z|\mathbf{x}_i) \log \underline{p(\mathbf{x}_t|z)} dz - \int_z q(z|\mathbf{x}_i) \log \frac{q(z|\mathbf{x}_i)}{p(z)} dz \quad (\log \frac{x}{y} = -\log \frac{y}{x})$$

$$= \mathbb{E}_{z \sim q(z|\mathbf{x}_i)} [\log p(\mathbf{x}_t|z)] - D_{KL}(q(z|\mathbf{x}_i) || p(z)) \quad (\mathbb{E} \text{ over } z, \text{KL})$$

t-SNE embedding comparisons

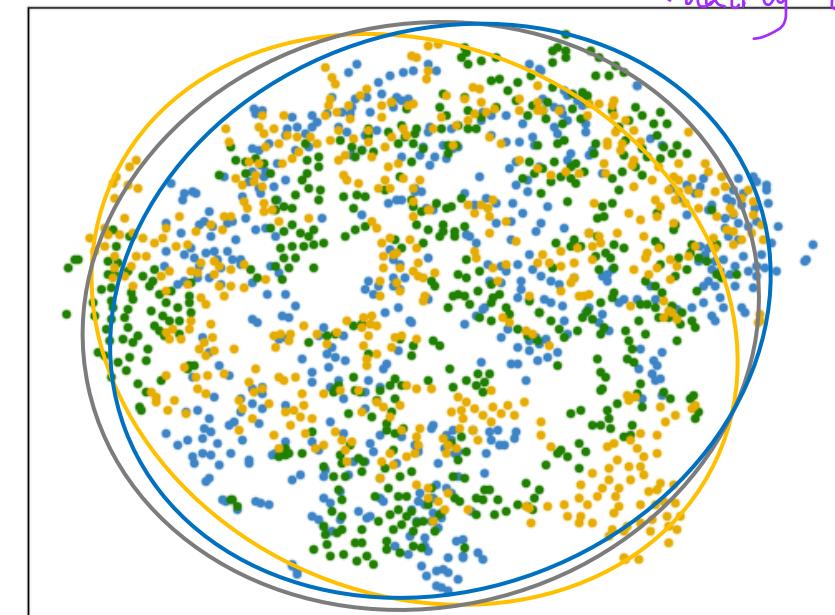


Non-overlapping partially overlapping

Latent space

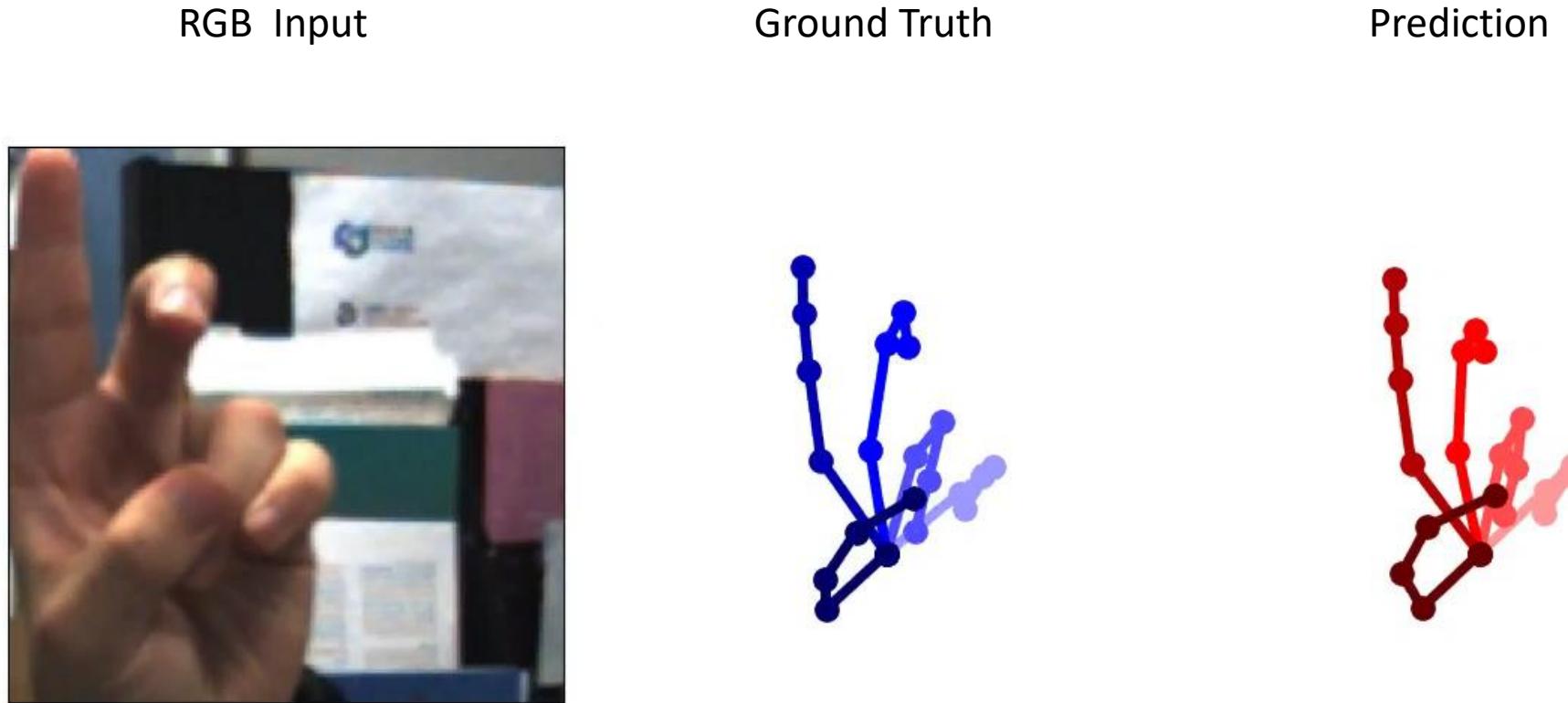
overlapping

Trained cross-modal



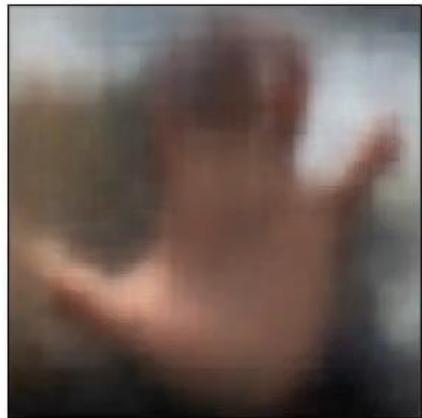
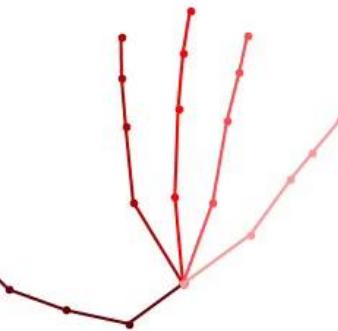
one latent space similar configurations, indep. of their modality end up being close together

Posterior estimates: RGB (real) to 3D joints

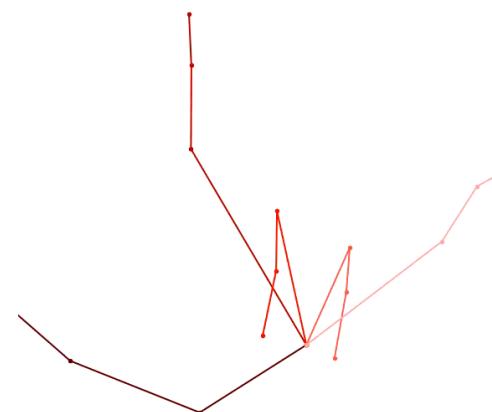
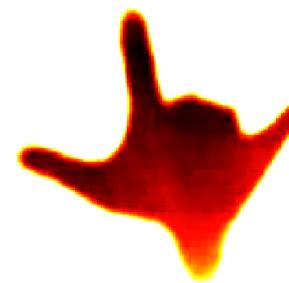


Synthetic hands: Latent Space Walk

generated images



[Synthetic samples RGB]



[Synthetic samples depth]

Learning useful representations



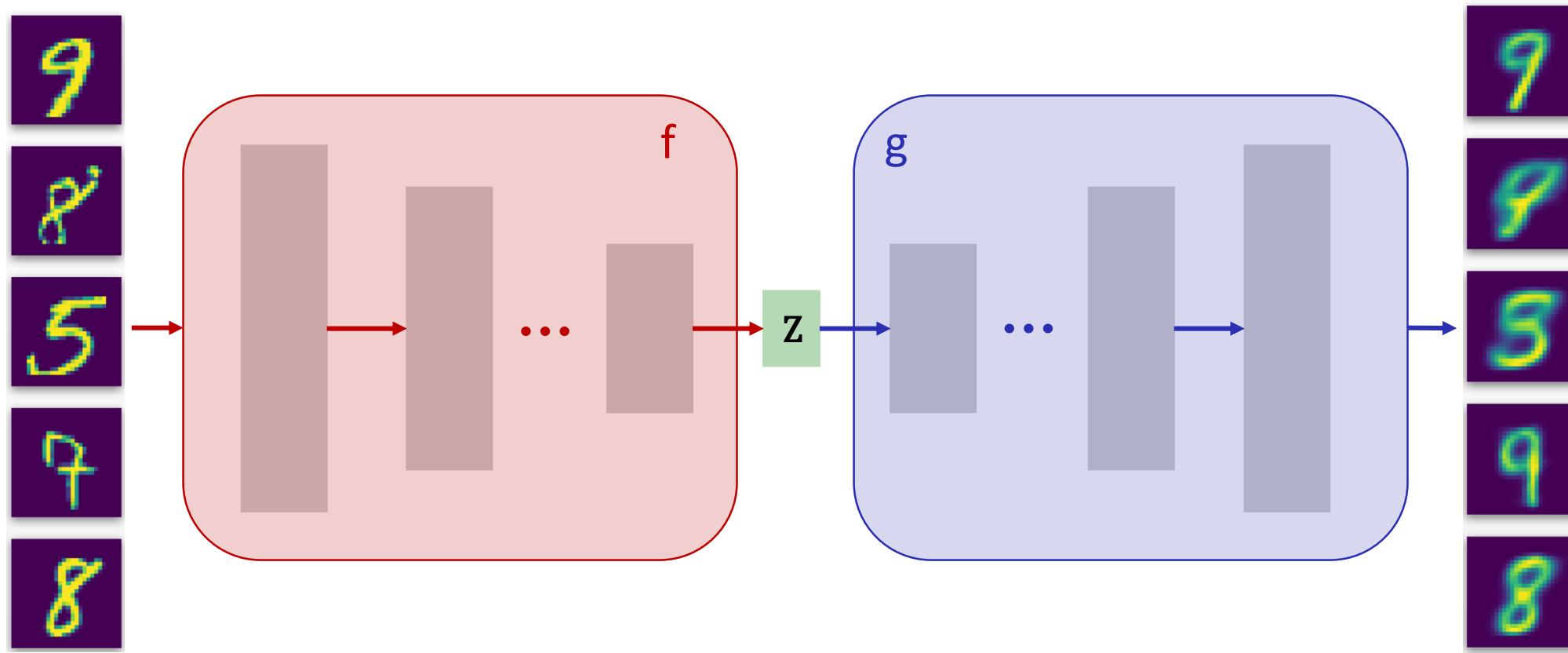
A 10x10 grid of handwritten digits from 0 to 9, arranged in a single row. The digits are written in a cursive style. Some digits are highlighted in bold or outlined in black.

1	4	7	8	7	0	1	2	1	2
9	1	9	0	6	3	7	9	1	7
7	6	7	9	0	4	2	0	8	7
2	0	9	3	6	8	7	6	3	1
4	2	1	6	3	1	5	4	4	3
2	5	4	4	5	4	6	3	3	7
0	4	4	6	9	3	8	8	3	7
1	0	0	9	1	1	2	1	1	9
1	2	1	0	0	9	3	4	9	2
8	5	8	5	7	3	6	8	7	1

Goal: Learn features that capture similarities and dissimilarities

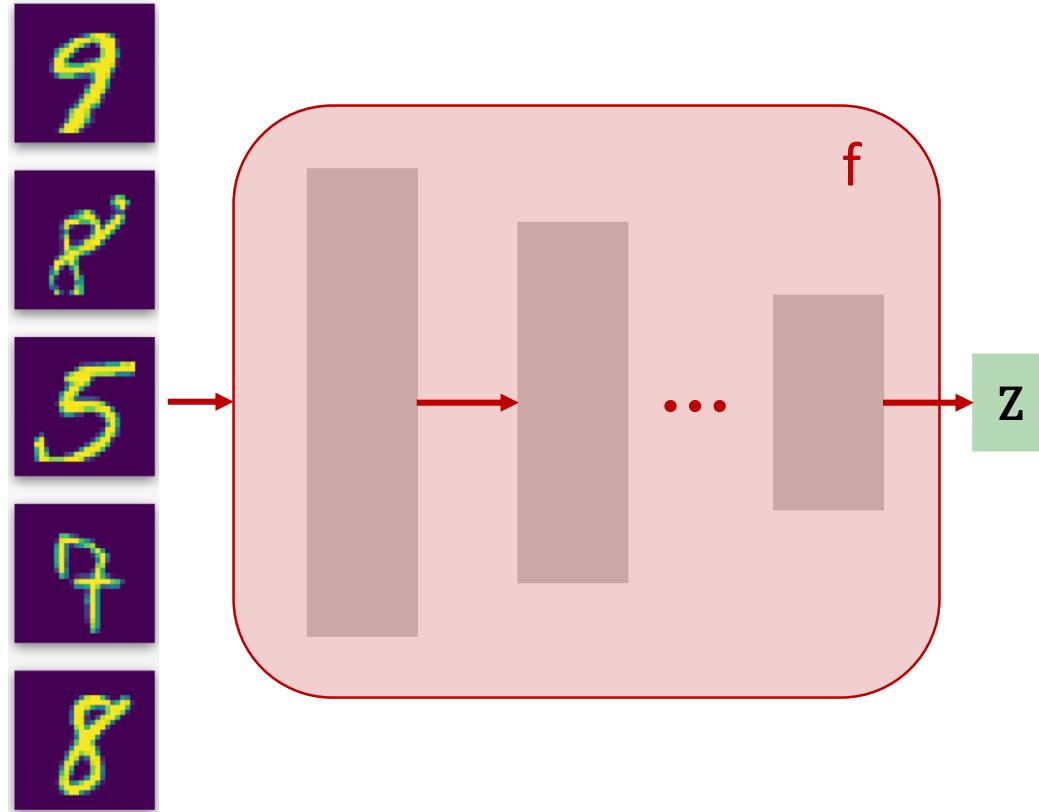
Requirement: Objective that defines notion of utility (task-dependent)

Autoencoders



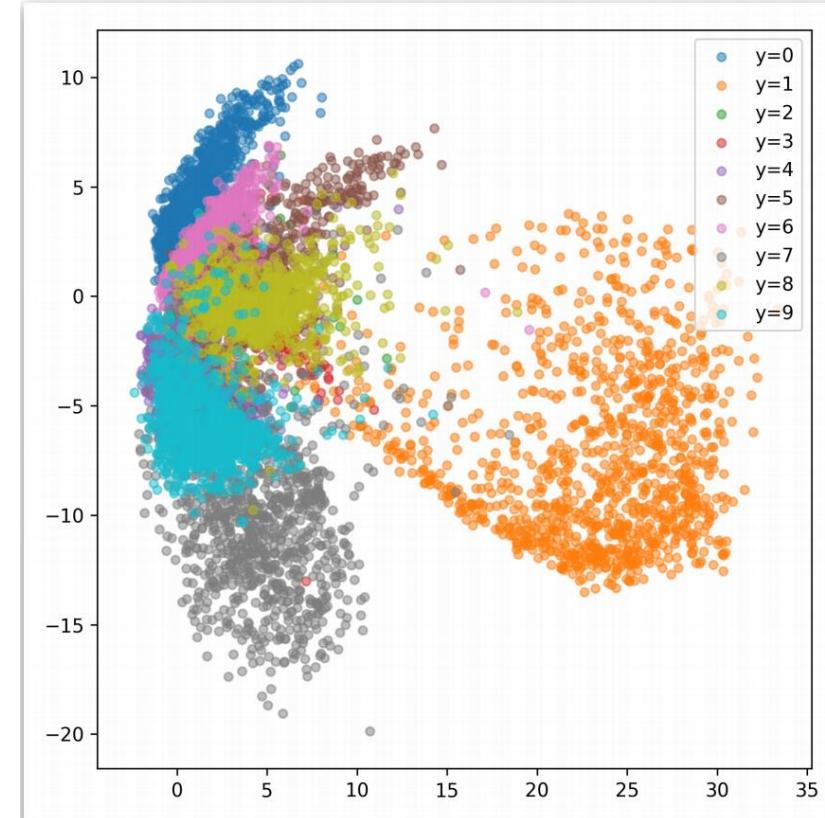
Notion of utility: Ability to reconstruct pixels from code
(features are a compressed representation of original data)

Autoencoders



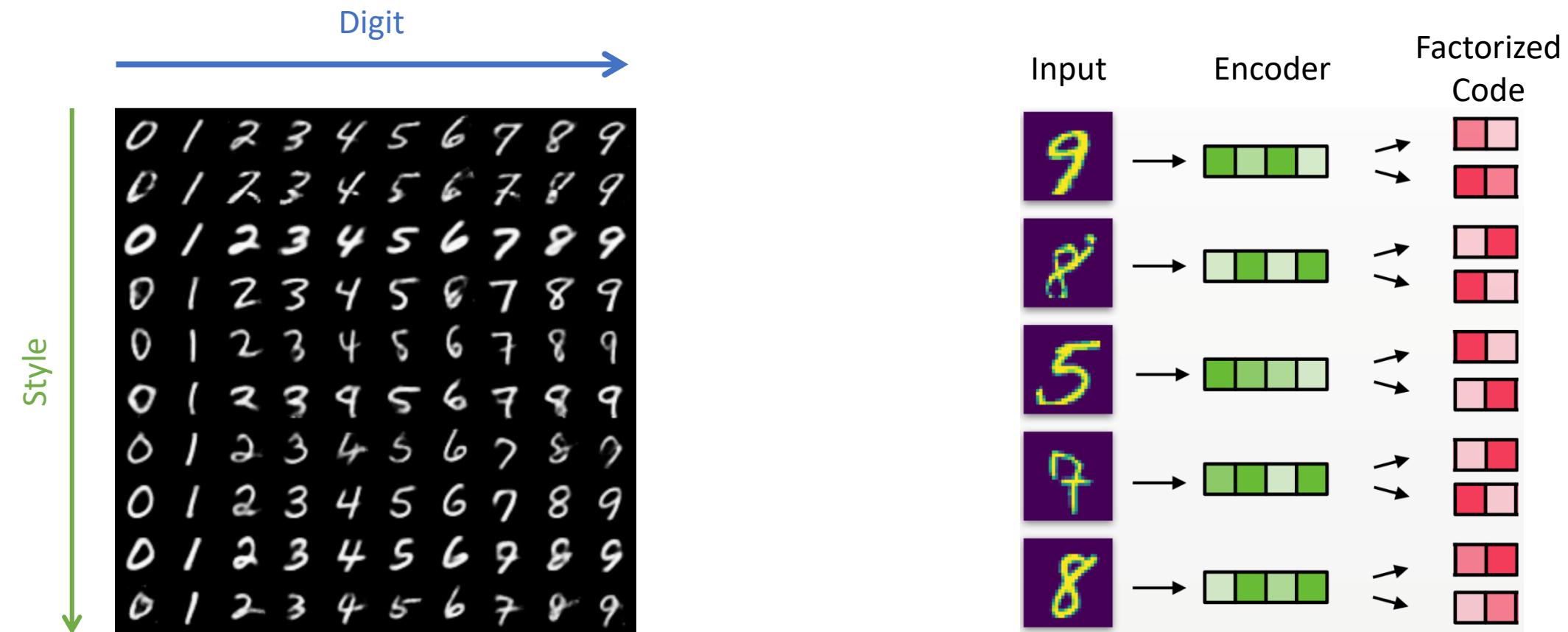
纠缠表示

Learned representation



Entangled Representation: Individual dimensions in code encode some unknown combination of features in the data.

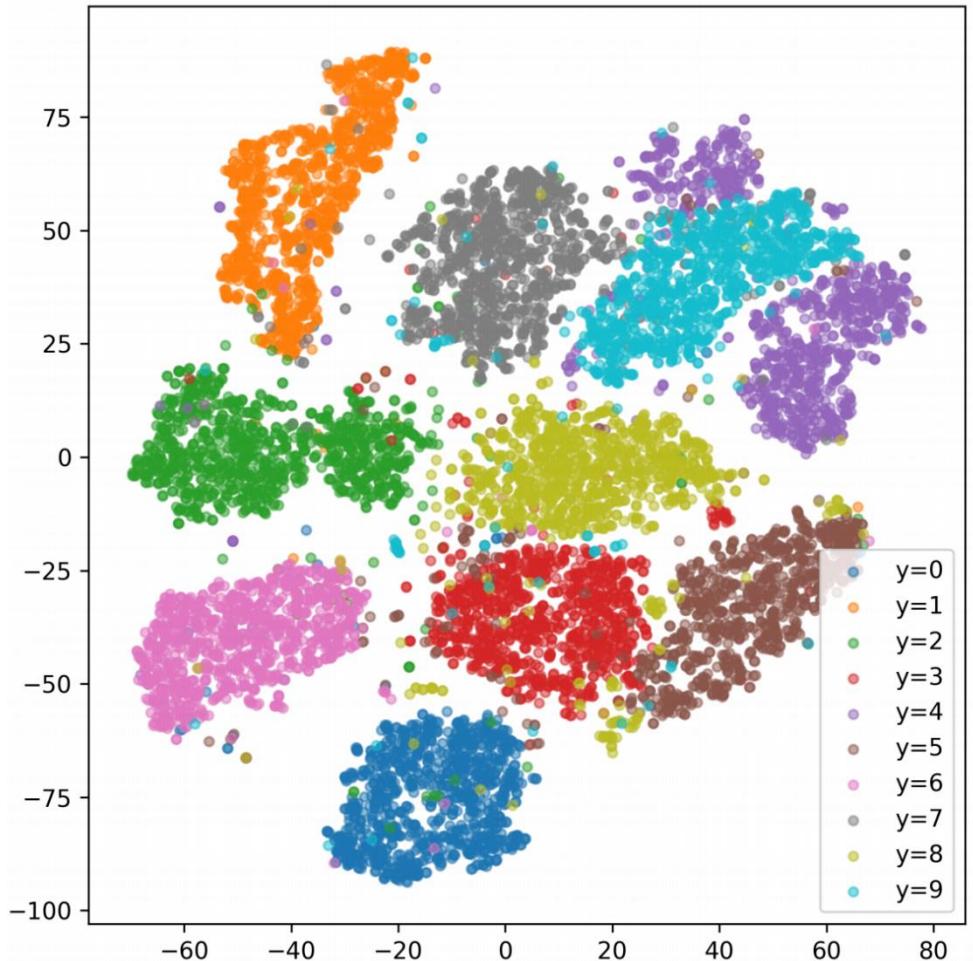
Goal: Disentangled Representations



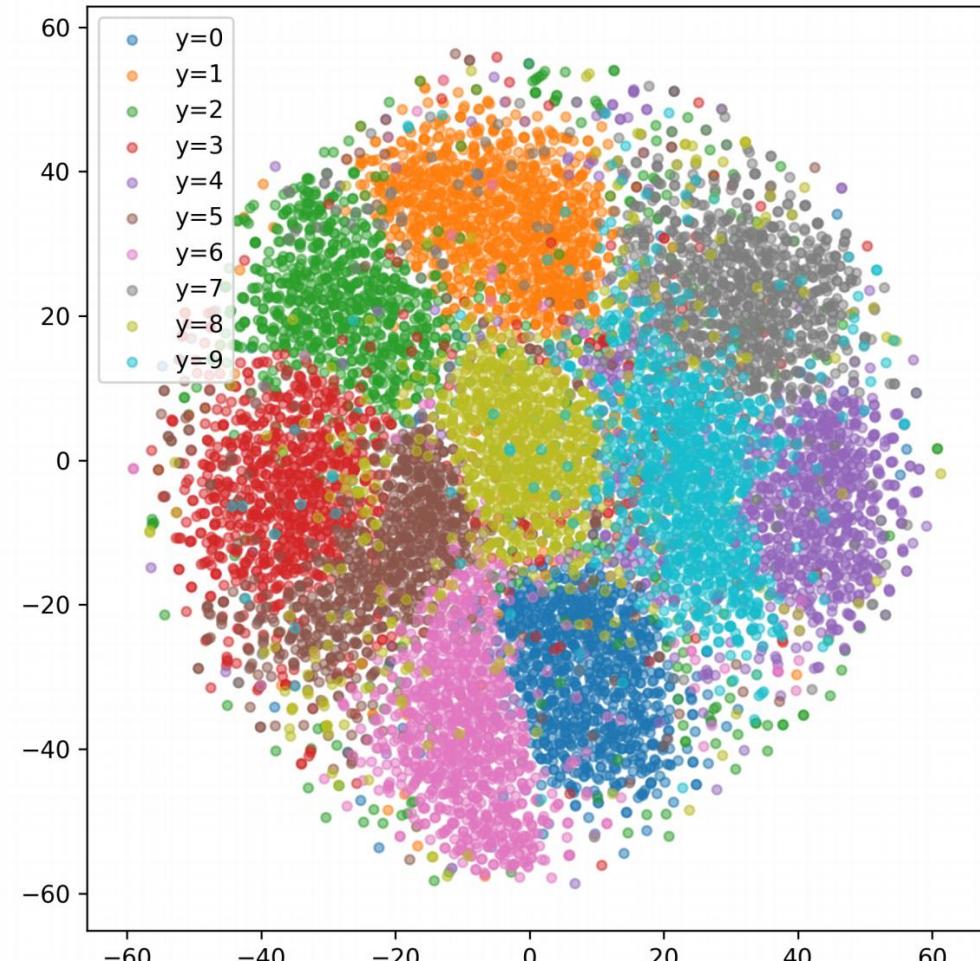
Goal: Learn features that correspond to distinct factors of variation
One Notion of Utility: statistical independence

Regular vs Variational Autoencoders

AE (z-dim 50, TSNE)



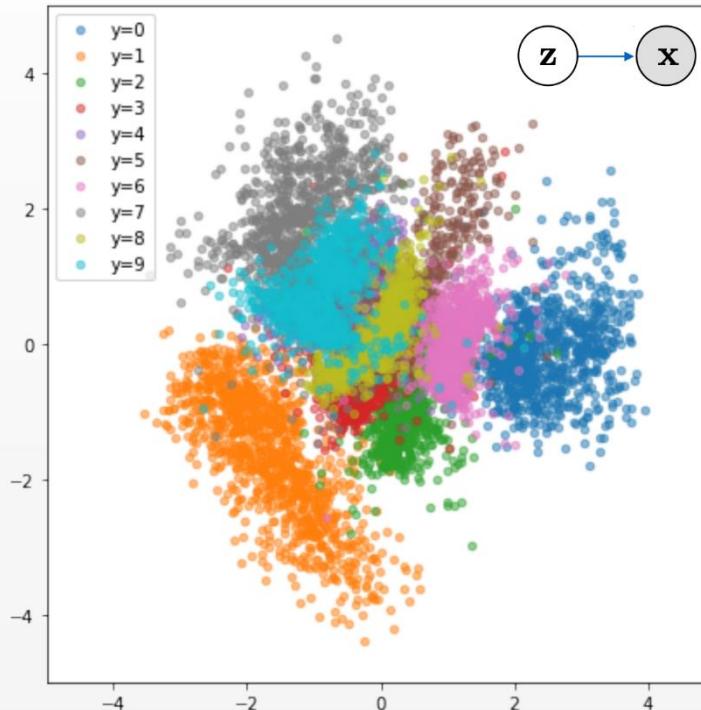
VAE (z-dim 50, TSNE)



Issue: KL-Divergence regularizes values of z but
Representations are still entangled

Unsupervised vs Semi-supervised

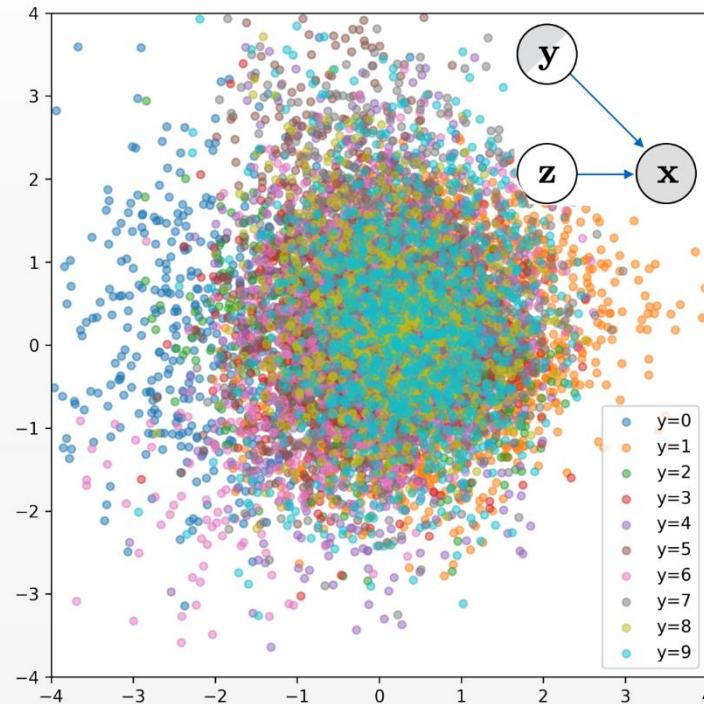
Unsupervised, Entangled



Latent code z represents both style and digit

$$p_{\theta}(x|y, z)$$

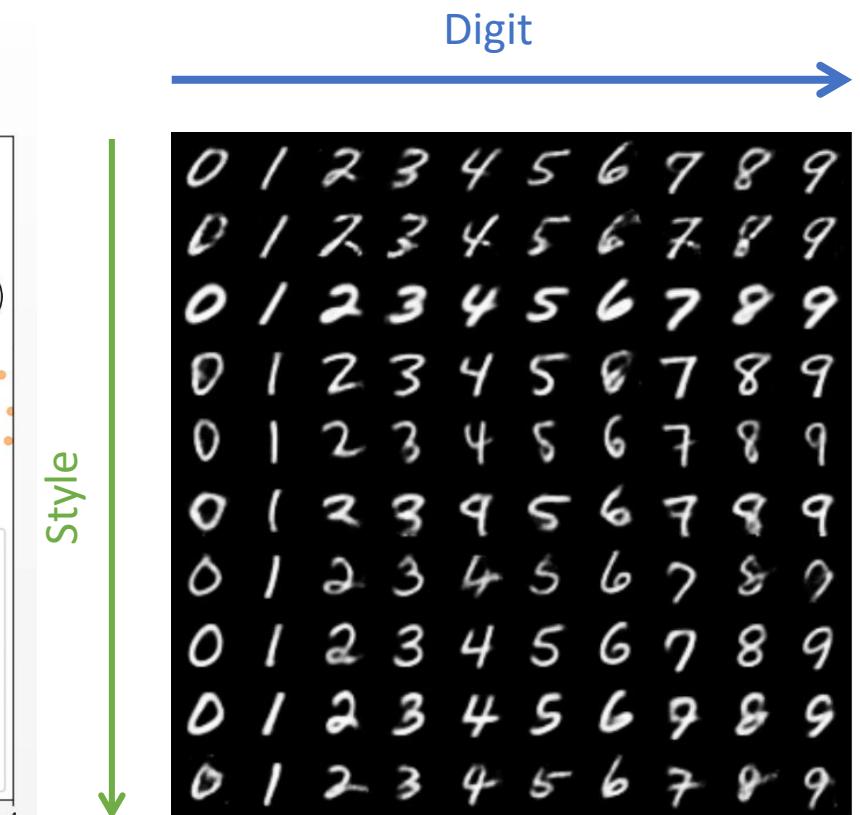
Semi-supervised, Disentangled



Style variable z is conditionally independent from digit y (*)

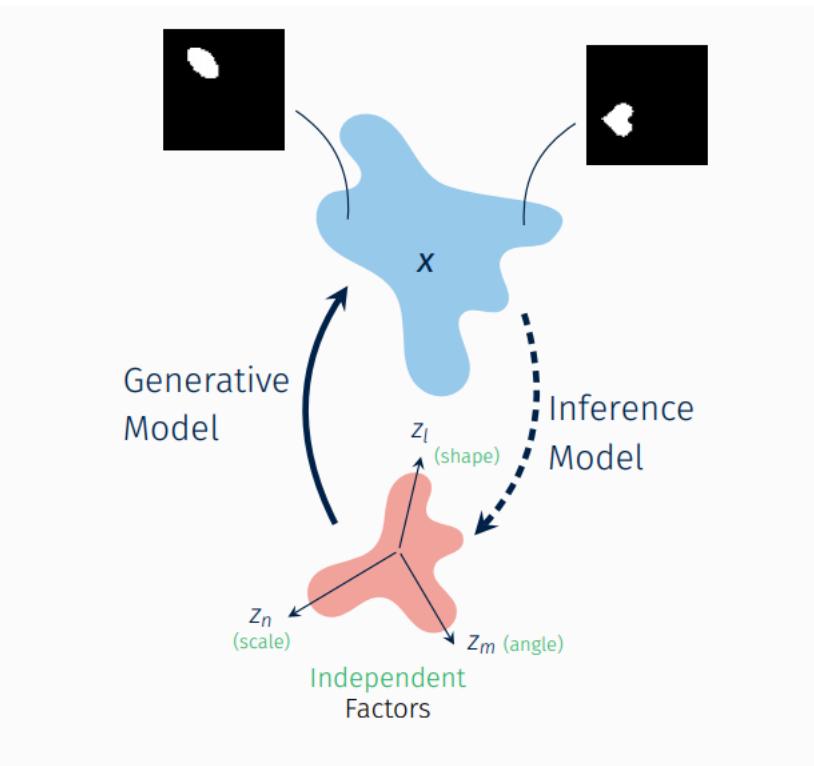
$$q_{\phi}(y, z|x)$$

[image source: Babak Esmaeli]



Separate label y from “nuisance” variable z

Learning factorized representations (unsupervised)



[image source: Emile Mathieu]

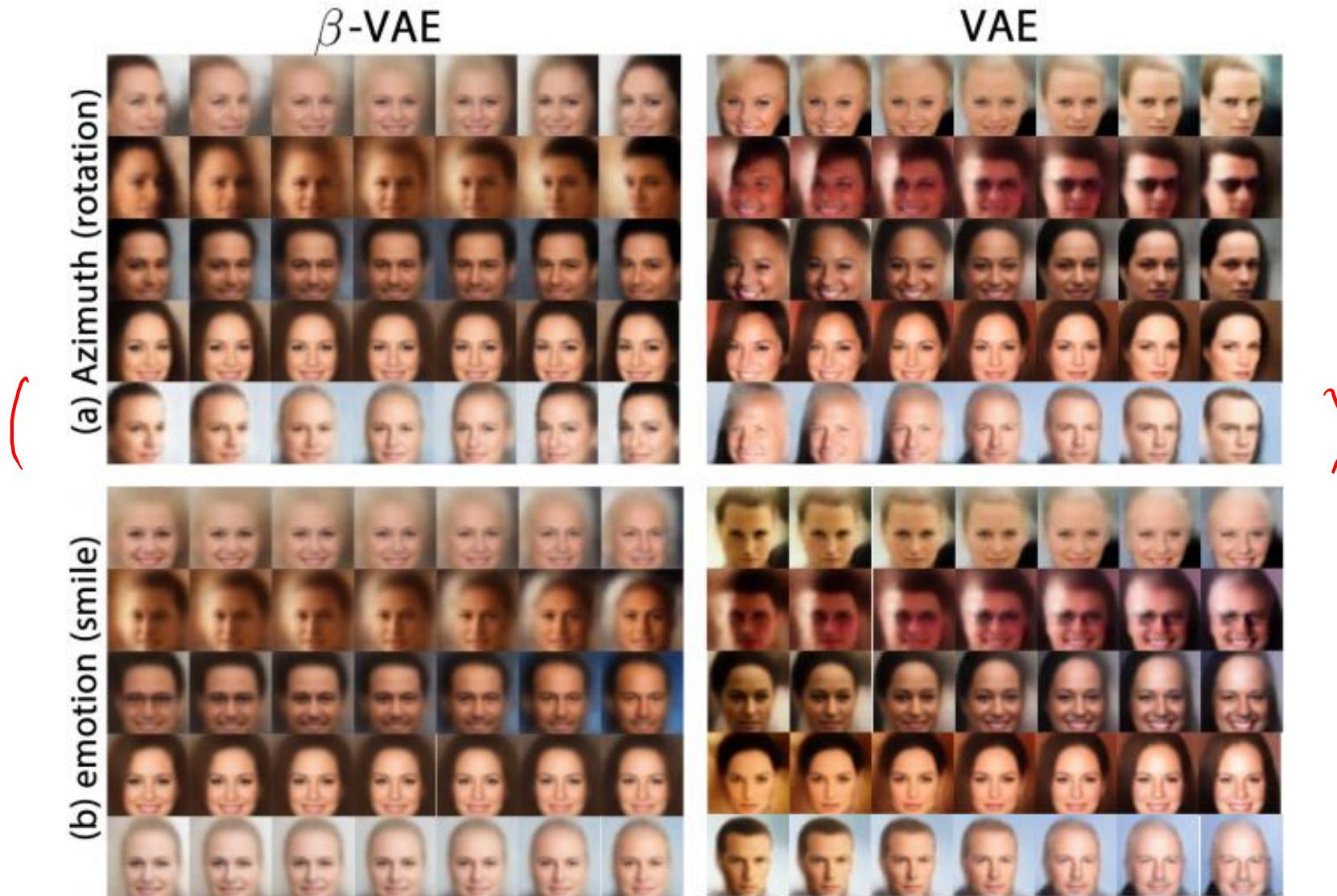
β -VAE

Goal: Learn disentangled representation *without* supervision

Idea: Provide framework for automated discovery of interpretable factorised latent representations

Approach: modification of the VAE, introducing an adjustable hyperparameter beta that balances latent channel capacity and independence constraints with reconstruction accuracy.

β -VAE



β -VAE

$$p(\mathbf{x}|\mathbf{z}) \approx p(\mathbf{x}|\mathbf{v}, \mathbf{w})$$

Conditionally independent factors $\mathbf{v} \in \mathbb{R}^K$, where $\log p(\mathbf{v}|\mathbf{x}) = \sum_k \log p(v_k|\mathbf{x})$

Conditionally dependent factors $\mathbf{w} \in \mathbb{R}^H$

$\mathbf{z} \in \mathbb{R}^M$, where $M \geq K$

model : which factor $\in \mathcal{V}$, which $\in \mathcal{W}$

Assume that samples \mathbf{x} are generated by controlling the generative factors \mathbf{v} and \mathbf{w} .

β -VAE – Training Objective

$$\max_{\phi, \theta} \mathbb{E}_{x \sim \mathcal{D}} \left[\mathbb{E}_{z \sim q_\phi(z|x)} \log p_\theta(x|z) \right]$$

subject to $D_{KL}(q_\phi(z|x) \| p_\theta(z)) < \delta$

$$\begin{aligned}
 &= \mathbb{E}_z [\log p_\theta(x^{(i)}|z)] - \mathbb{E}_z [\log \frac{q_\phi(z|x^{(i)})}{p_\theta(z)}] + \mathbb{E}_z [\log \frac{q_\phi(z|x^{(i)})}{p_\theta(z|x^{(i)})}] \quad (\log \frac{x}{y} = -\log \frac{y}{x}) \\
 &\quad \checkmark \text{(with sampling)} \quad \checkmark \text{closed-form if Gaussian} \quad \times \\
 &= \underbrace{\mathbb{E}_z [\log p_\theta(x^{(i)}|z)]}_{\parallel \text{reconstruction} \parallel} - D_{KL} \underbrace{\frac{q_\phi(z|x^{(i)})}{p_\theta(z)}}_{\text{tractable}} + \underbrace{D_{KL} \frac{q_\phi(z|x^{(i)})}{p_\theta(z|x^{(i)})} \| p_\theta(z|x^{(i)})}_{\geq 0 \text{ (see proof earlier)}} \\
 &\quad \downarrow \quad \downarrow \quad \downarrow \\
 &\quad \text{goes high if our sample is similar to true distribution} \quad \text{"make approx. posterior as similar as possible to prior"} \\
 &\quad \text{loss function} \quad \mathcal{L}(x^{(i)}, \theta, \phi) \quad \text{Training procedure:} \\
 &\quad \log(p_\theta(x^{(i)})) \geq \mathcal{L} \quad \theta^*, \phi^* = \arg \max \sum_i \mathcal{L}(x^{(i)}, \theta, \phi) \\
 &\quad \text{evidence lower bound "ELBO"} \\
 &\quad \Rightarrow \text{data is at least as likely as } \mathcal{L}
 \end{aligned}$$

β -VAE – Training Objective

$\beta=1$, \Rightarrow ELBO

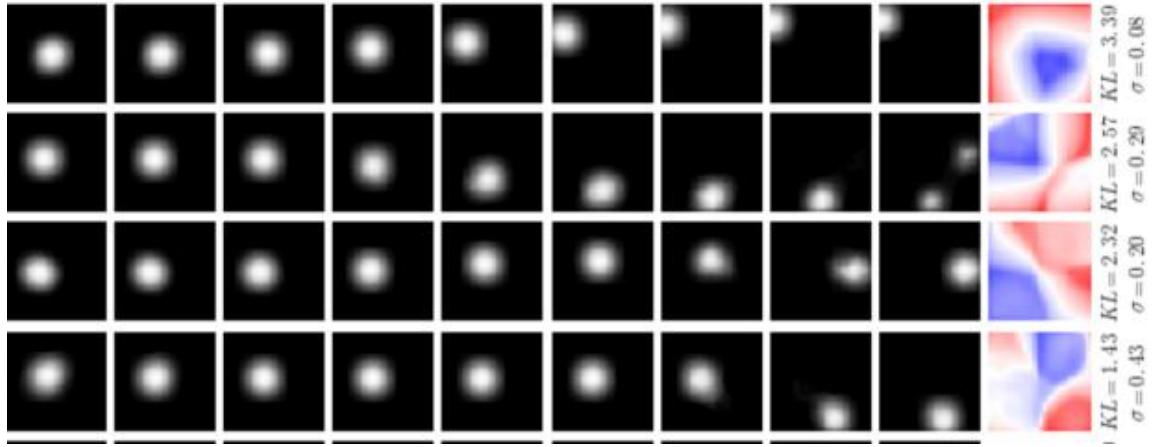
$\beta > 1$, \Rightarrow put more pressure on KL

$$\begin{aligned}
 \mathcal{F}(\theta, \phi, \beta) &= \mathbb{E}_{\mathbf{z} \sim q_\phi(\mathbf{z}|\mathbf{x})} \log p_\theta(\mathbf{x}|\mathbf{z}) - \beta(D_{\text{KL}}(q_\phi(\mathbf{z}|\mathbf{x})||p_\theta(\mathbf{z})) - \delta) \\
 &= \mathbb{E}_{\mathbf{z} \sim q_\phi(\mathbf{z}|\mathbf{x})} \log p_\theta(\mathbf{x}|\mathbf{z}) - \beta D_{\text{KL}}(q_\phi(\mathbf{z}|\mathbf{x})||p_\theta(\mathbf{z})) + \beta\delta \\
 &\geq \mathbb{E}_{\mathbf{z} \sim q_\phi(\mathbf{z}|\mathbf{x})} \log p_\theta(\mathbf{x}|\mathbf{z}) - \beta D_{\text{KL}}(q_\phi(\mathbf{z}|\mathbf{x})||p_\theta(\mathbf{z})) \quad ; \text{ Because } \beta, \delta \geq 0
 \end{aligned}$$

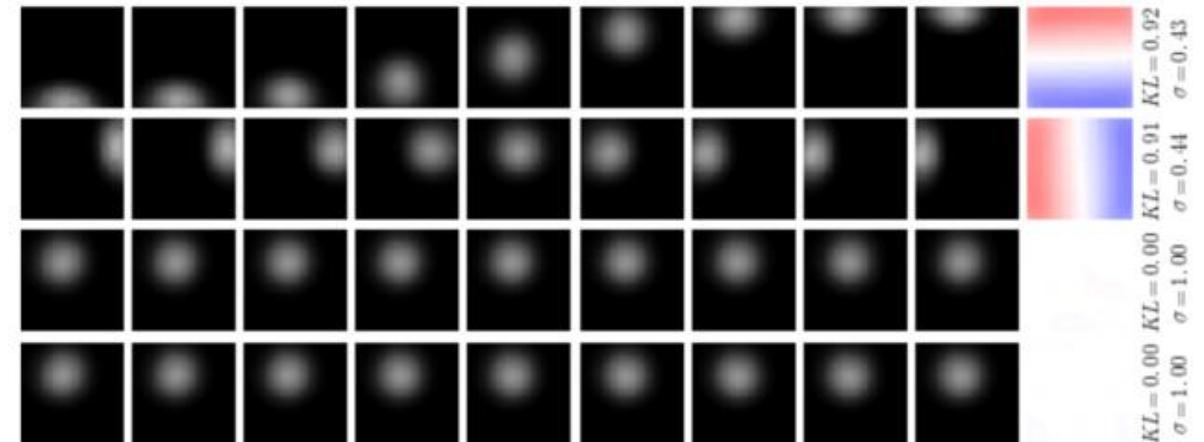
$$L_{\text{BETA}}(\phi, \beta) = -\mathbb{E}_{\mathbf{z} \sim q_\phi(\mathbf{z}|\mathbf{x})} \log p_\theta(\mathbf{x}|\mathbf{z}) + \beta D_{\text{KL}}(q_\phi(\mathbf{z}|\mathbf{x})||p_\theta(\mathbf{z}))$$

Entangled versus Disentangled Representations

vanilla VAE



β -VAE

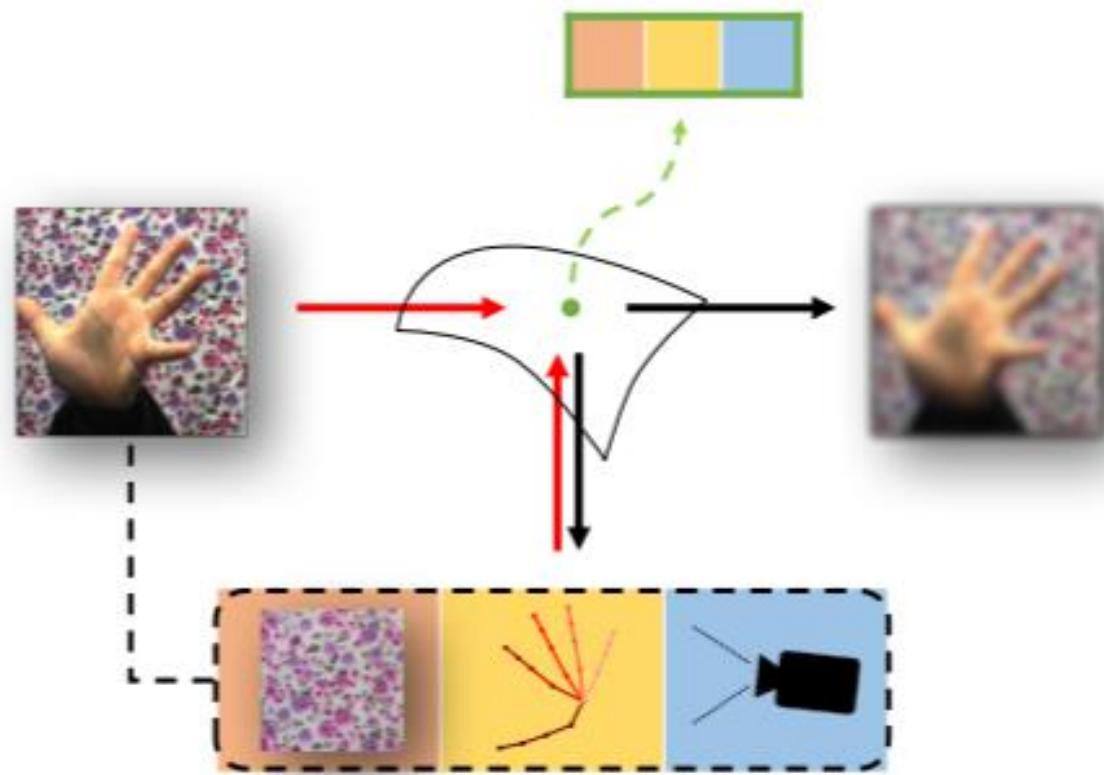


Disentangling Latent Hands for Image Synthesis and Pose Estimation

Linlin Yang, Angela Yao

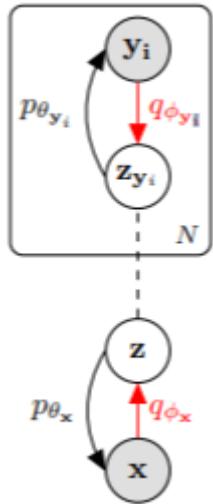
CVPR 2019

Disentangling latent factors



Disentangling latent factors

$$L(\phi_x, \phi_{y_1}, \phi_{y_2}, \theta_x, \theta_{y_1}, \theta_{y_2}) = ELBO_{dis}(x, y_1, y_2, \phi_{y_1}, \phi_{y_2}, \theta_x, \theta_{y_1}, \theta_{y_2}) + ELBO_{emb}(x, y_1, y_2, \phi_x).$$



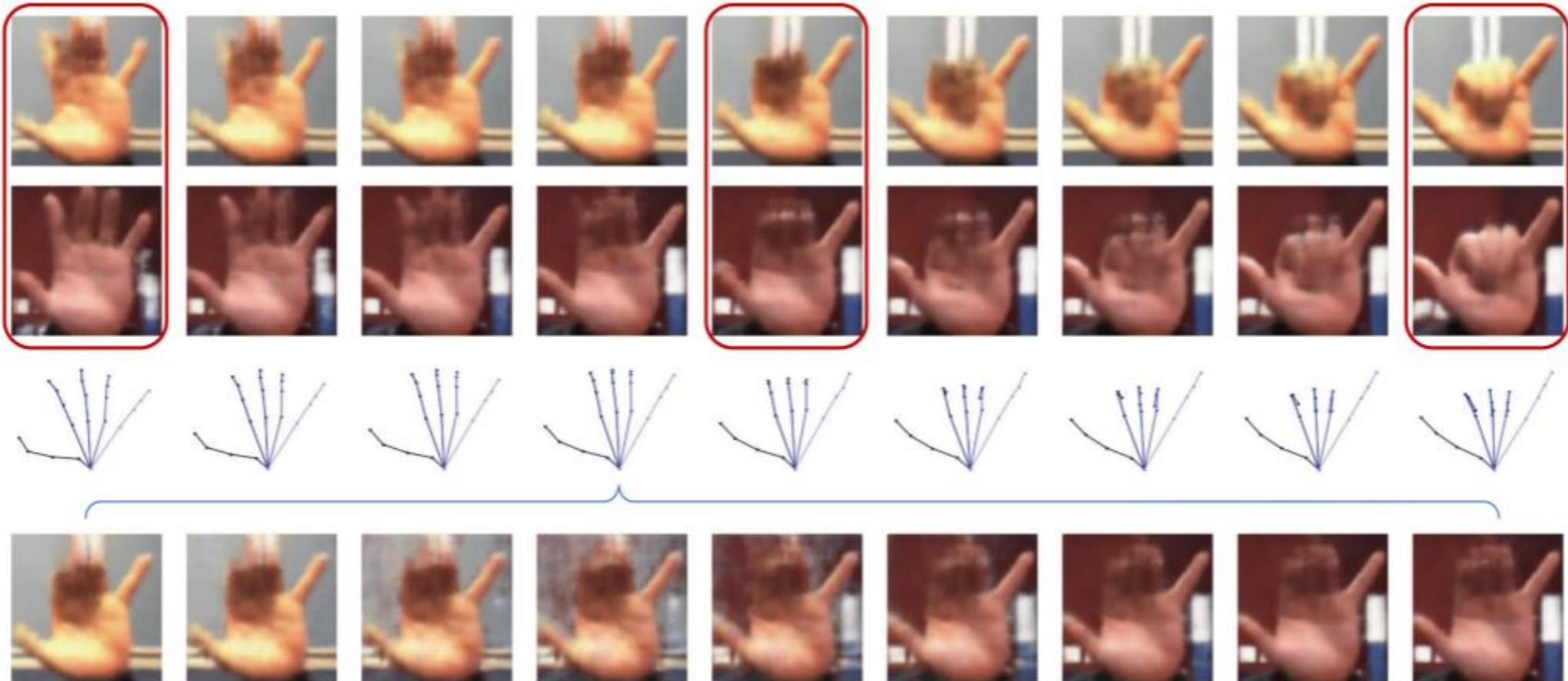
where

$$\begin{aligned} \log p(x, y_1, y_2) &\geq ELBO_{dis}(x, y_1, y_2, \phi_{y_1}, \phi_{y_2}, \theta_x, \theta_{y_1}, \theta_{y_2}) \\ &= \lambda_x \mathbb{E}_{z \sim q_{\phi_{y_1}, \phi_{y_2}}} [\log p_{\theta_x}(x|z)] \\ &\quad + \lambda_{y_1} \mathbb{E}_{z_{y_1} \sim q_{\phi_{y_1}}} [\log p_{\theta_{y_1}}(y_1|z_{y_1})] \\ &\quad + \lambda_{y_2} \mathbb{E}_{z_{y_2} \sim q_{\phi_{y_2}}} [\log p_{\theta_{y_2}}(y_2|z_{y_2})] \\ &\quad - \beta D_{KL} (q_{\phi_{y_1}, \phi_{y_2}}(z|y_1, y_2) || p(z)) \end{aligned}$$

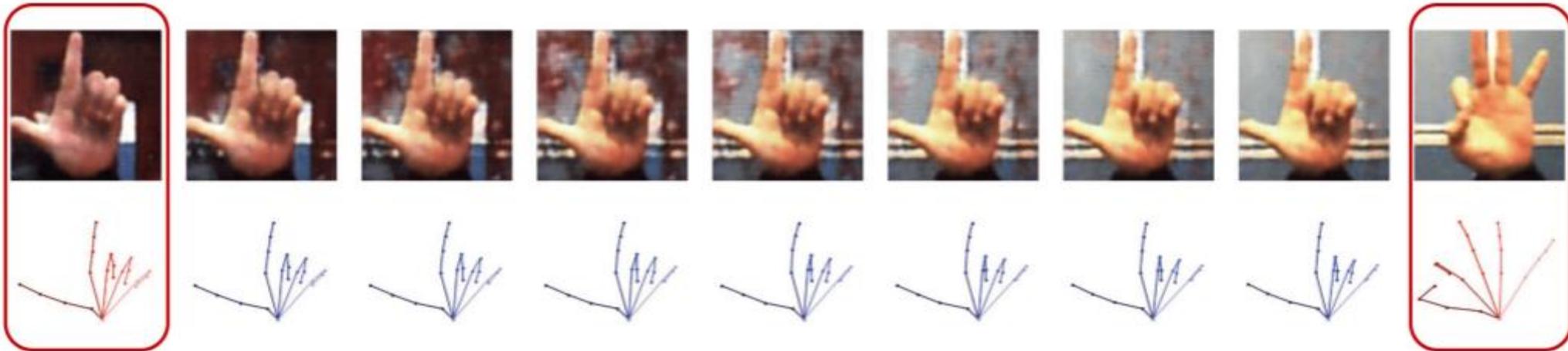
and

$$\begin{aligned} ELBO_{emb}(x, y_1, y_2, \phi_x) &= \lambda'_x \mathbb{E}_{z \sim q_{\phi_x}} [\log p_{\theta_x}(x|z)] \\ &\quad + \lambda'_{y_1} \mathbb{E}_{z_{y_1} \sim q_{\phi_x}} [\log p_{\theta_{y_1}}(y_1|z_{y_1})] \\ &\quad + \lambda'_{y_2} \mathbb{E}_{z_{y_2} \sim q_{\phi_x}} [\log p_{\theta_{y_2}}(y_2|z_{y_2})] \\ &\quad - \beta' D_{KL} (q_{\phi_x}(z|x) || p(z)) \end{aligned}$$

Latentspace walk with disentangled z



Latentspace walk with disentangled z

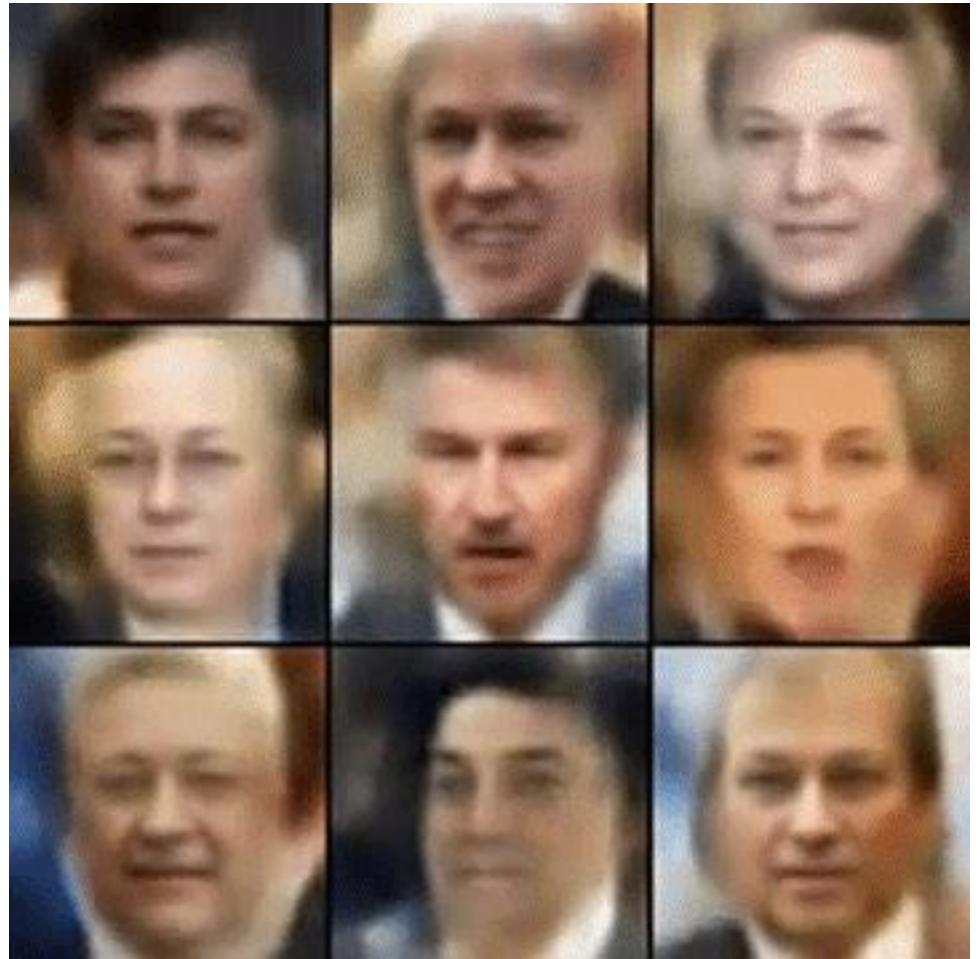


Limitations of VAEs

Tendency to generate blurry images.

Believed to be due to injected noise and weak inference models (Gaussian assumption of latent samples to simplistic, model capacity to weak)

More expressive model → substantially better results (e.g Kingma et al. 2016, *Improving variational inference with inverse autoregressive flow*)



Next

Deep Generative Modelling II: Generative Adversarial Networks (GANs)