



UNIVERSITY OF  
OXFORD

# Disentangling Disentanglement in Variational Autoencoders

GenU 2019 (Originally Presented at ICML 2019)

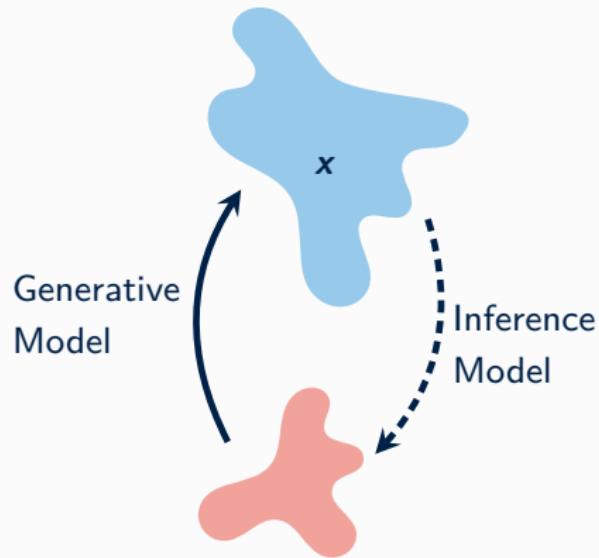
---

Emile Mathieu\*, Tom Rainforth\*, N. Siddharth\*, Yee Whye Teh

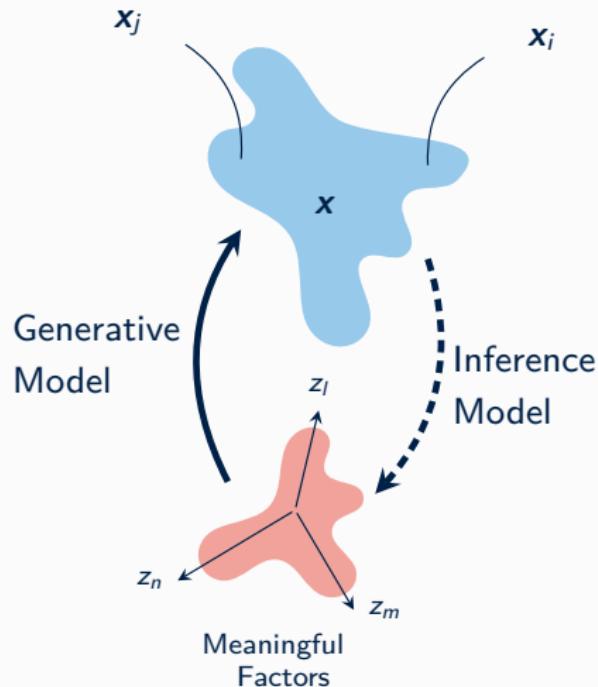
October 9th, 2019

Departments of Statistics and Engineering Science, University of Oxford

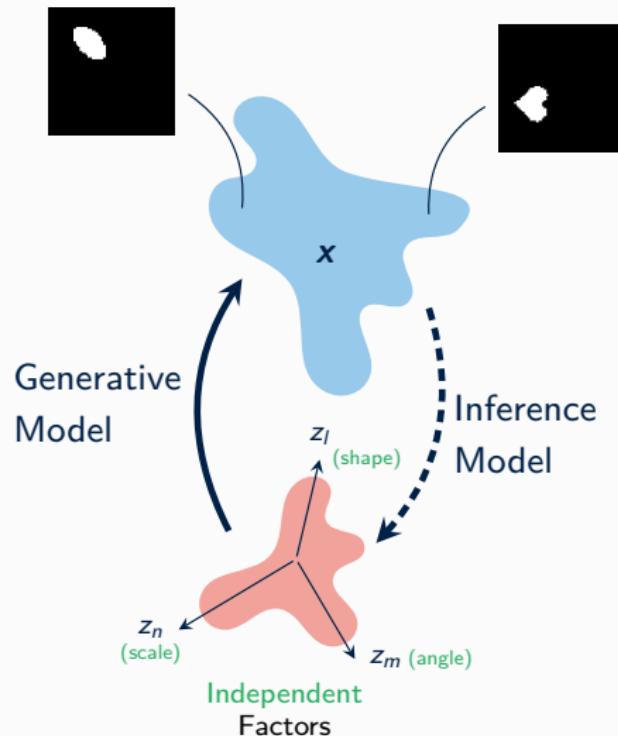
## Variational Autoencoders



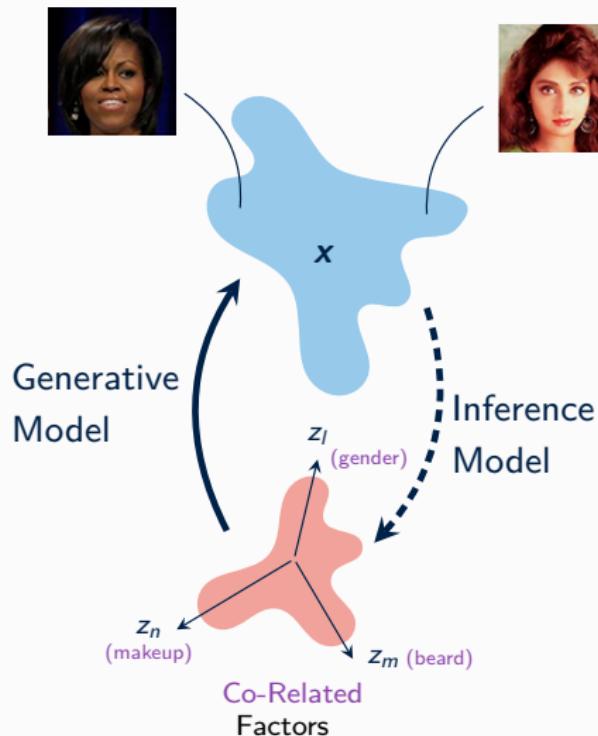
## Interpretable Representations



## Disentanglement = Independence



Decomposition  $\in \{\text{Independence, Clustering, Sparsity, ...}\}$



## Talk Outline

- Background on VAEs and disentanglement

## Talk Outline

- Background on VAEs and disentanglement
- Decomposition: a generalisation of disentanglement allowing a richer class of properties to be imposed

## Talk Outline

- Background on VAEs and disentanglement
- Decomposition: a generalisation of disentanglement allowing a richer class of properties to be imposed
- Theoretical analysis of  $\beta$ -VAE

## Talk Outline

- Background on VAEs and disentanglement
- Decomposition: a generalisation of disentanglement allowing a richer class of properties to be imposed
- Theoretical analysis of  $\beta$ -VAE
- An objective for encouraging decomposition

## Talk Outline

- Background on VAEs and disentanglement
- Decomposition: a generalisation of disentanglement allowing a richer class of properties to be imposed
- Theoretical analysis of  $\beta$ -VAE
- An objective for encouraging decomposition
- Experiments that showcase examples decompositions:
  - independence
  - clustering
  - sparsity

## Talk Outline

- Background on VAEs and disentanglement
- Decomposition: a generalisation of disentanglement allowing a richer class of properties to be imposed
- Theoretical analysis of  $\beta$ -VAE
- An objective for encouraging decomposition
- Experiments that showcase examples decompositions:
  - independence
  - clustering
  - sparsity
- Points of discussion

## Variational Autoencoders (VAEs)

- Have observations  $\{\mathbf{x}_1, \dots, \mathbf{x}_n\}$  from data generating process  $p_{\mathcal{D}}(\mathbf{x})$

## Variational Autoencoders (VAEs)

- Have observations  $\{\mathbf{x}_1, \dots, \mathbf{x}_n\}$  from data generating process  $p_{\mathcal{D}}(\mathbf{x})$
- Want to learn a corresponding latent variable model  $p(z)p_{\theta}(\mathbf{x}|z)$

$$\arg \min_{\theta} \text{KL}(p_{\mathcal{D}}(\mathbf{x}) \parallel p_{\theta}(\mathbf{x})) = \arg \max_{\theta} \mathbb{E}_{p_{\mathcal{D}}(\mathbf{x})}[\log p_{\theta}(\mathbf{x})]$$

## Variational Autoencoders (VAEs)

- Have observations  $\{\mathbf{x}_1, \dots, \mathbf{x}_n\}$  from data generating process  $p_{\mathcal{D}}(\mathbf{x})$
- Want to learn a corresponding latent variable model  $p(\mathbf{z})p_{\theta}(\mathbf{x}|\mathbf{z})$

$$\arg \min_{\theta} \text{KL}(p_{\mathcal{D}}(\mathbf{x}) \parallel p_{\theta}(\mathbf{x})) = \arg \max_{\theta} \mathbb{E}_{p_{\mathcal{D}}(\mathbf{x})} [\log p_{\theta}(\mathbf{x})]$$

- As this intractable, instead introduce inference model  $q_{\phi}(\mathbf{z}|\mathbf{x})$  and replace  $\log p_{\theta}(\mathbf{x})$  with the lower bound

$$\begin{aligned}\mathcal{L}(\mathbf{x}; \theta, \phi) &\triangleq \log p_{\theta}(\mathbf{x}) - \text{KL}(q_{\phi}(\mathbf{z}|\mathbf{x}) \parallel p_{\theta}(\mathbf{z}|\mathbf{x})) \\ &= \mathbb{E}_{q_{\phi}(\mathbf{z}|\mathbf{x})} [\log p_{\theta}(\mathbf{x}|\mathbf{z})] - \text{KL}(q_{\phi}(\mathbf{z}|\mathbf{x}) \parallel p(\mathbf{z})).\end{aligned}\tag{1}$$

## Variational Autoencoders (VAEs)

- Have observations  $\{\mathbf{x}_1, \dots, \mathbf{x}_n\}$  from data generating process  $p_{\mathcal{D}}(\mathbf{x})$
- Want to learn a corresponding latent variable model  $p(\mathbf{z})p_{\theta}(\mathbf{x}|\mathbf{z})$

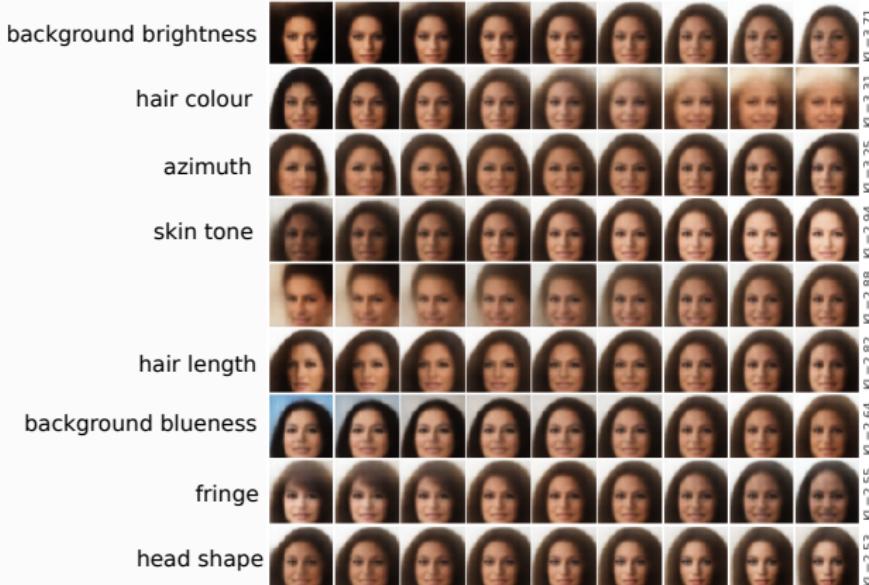
$$\arg \min_{\theta} \text{KL}(p_{\mathcal{D}}(\mathbf{x}) \parallel p_{\theta}(\mathbf{x})) = \arg \max_{\theta} \mathbb{E}_{p_{\mathcal{D}}(\mathbf{x})} [\log p_{\theta}(\mathbf{x})]$$

- As this intractable, instead introduce inference model  $q_{\phi}(\mathbf{z}|\mathbf{x})$  and replace  $\log p_{\theta}(\mathbf{x})$  with the lower bound

$$\begin{aligned}\mathcal{L}(\mathbf{x}; \theta, \phi) &\triangleq \log p_{\theta}(\mathbf{x}) - \text{KL}(q_{\phi}(\mathbf{z}|\mathbf{x}) \parallel p_{\theta}(\mathbf{z}|\mathbf{x})) \\ &= \mathbb{E}_{q_{\phi}(\mathbf{z}|\mathbf{x})} [\log p_{\theta}(\mathbf{x}|\mathbf{z})] - \text{KL}(q_{\phi}(\mathbf{z}|\mathbf{x}) \parallel p(\mathbf{z})).\end{aligned}\tag{1}$$

- VAE views this as a deep stochastic autoencoder, where  $q_{\phi}(\mathbf{z}|\mathbf{x})$  is the encoder and  $p_{\theta}(\mathbf{x}|\mathbf{z})$  the decoder

# Disentanglement



<sup>1</sup>Hyunjik Kim and Andriy Mnih. "Disentangling by Factorising". In: *International Conference on Machine Learning*. 2018.

# Disentanglement

- Care about the pushforward distribution of the data in latent space:

$$q_\phi(\mathbf{z}) = \mathbb{E}_{p_{\mathcal{D}}(\mathbf{x})} [q_\phi(\mathbf{z}|\mathbf{x})] \quad (2)$$

# Disentanglement

- Care about the pushforward distribution of the data in latent space:

$$q_\phi(\mathbf{z}) = \mathbb{E}_{p_{\mathcal{D}}(\mathbf{x})} [q_\phi(\mathbf{z}|\mathbf{x})] \quad (2)$$

- Approaches try to learn representations where the dimensions of  $\mathbf{z}$  are independent, i.e.

$$q_\phi(\mathbf{z}) = \prod_d q_\phi(z_d) \quad (3)$$

# Disentanglement

- Care about the pushforward distribution of the data in latent space:

$$q_\phi(\mathbf{z}) = \mathbb{E}_{p_{\mathcal{D}}(\mathbf{x})} [q_\phi(\mathbf{z}|\mathbf{x})] \quad (2)$$

- Approaches try to learn representations where the dimensions of  $\mathbf{z}$  are independent, i.e.

$$q_\phi(\mathbf{z}) = \prod_d q_\phi(z_d) \quad (3)$$

- Some papers require  $z_d$  to correspond to “true generative factors”

## Decomposition: A Generalization of Disentanglement

More general framework for imposing *meaningful* structure on learned representations

## Decomposition: A Generalization of Disentanglement

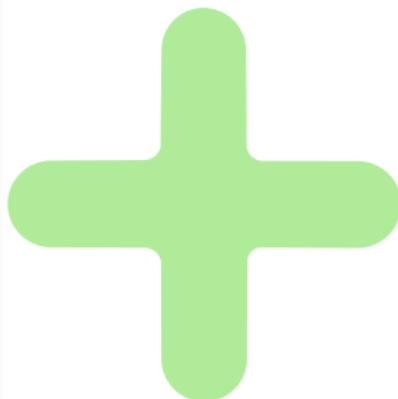
More general framework for imposing *meaningful* structure on learned representations

Characterise decomposition as the fulfilment of two factors:

- (a) matching between the marginal posterior  $q_\phi(\mathbf{z})$  and structured prior  $p(\mathbf{z})$  to constrain with the required decomposition.
- (b) appropriate level of overlap between encodings in the latent space,

# Decomposition: An Analysis

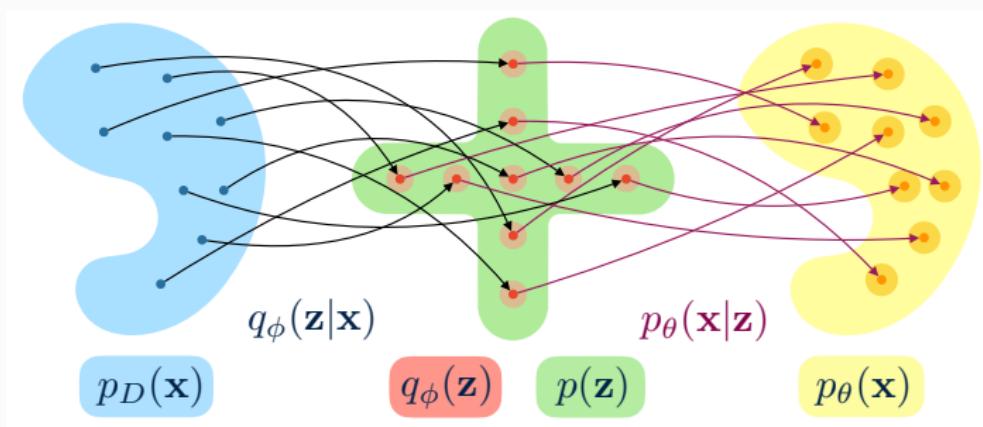
Desired Structure



$p(\mathbf{z})$

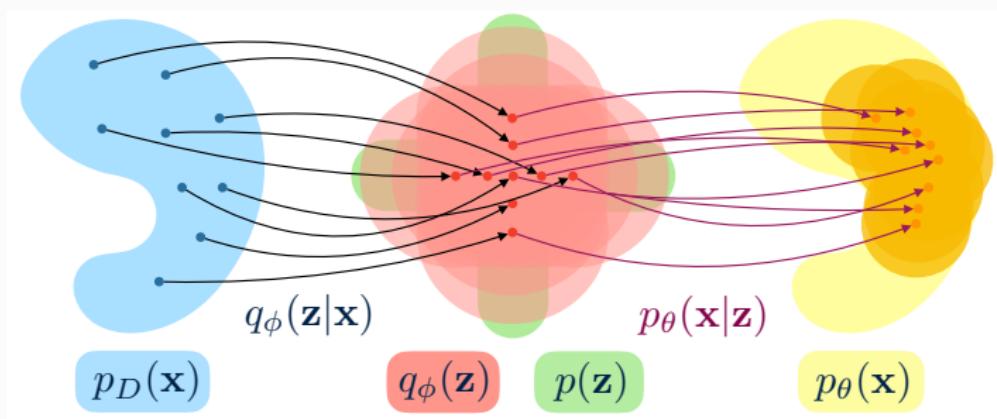
# Decomposition: An Analysis

Insufficient Overlap



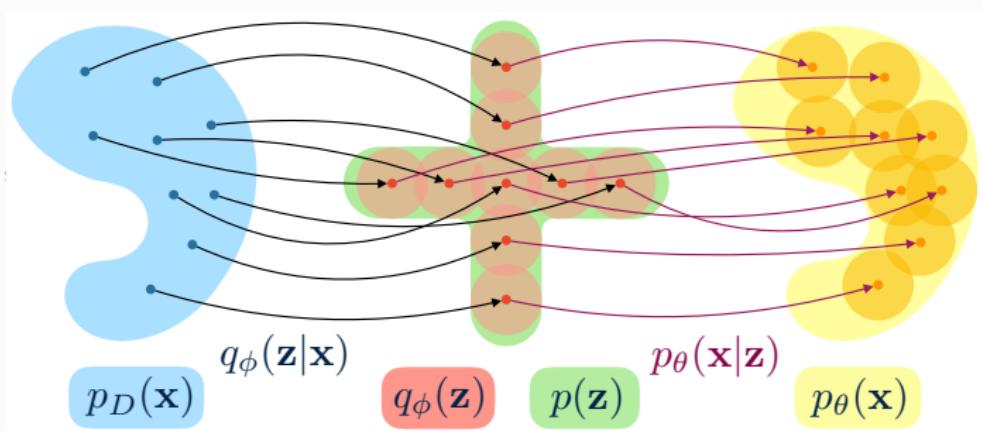
# Decomposition: An Analysis

Too Much Overlap



# Decomposition: An Analysis

Appropriate Overlap



# Deconstructing the $\beta$ -VAE

The  $\beta$ -VAE<sup>2</sup> adjusts the standard ELBO by scaling the KL term:

$$\mathcal{L}_\beta(\mathbf{x}) = \mathbb{E}_{q_\phi(\mathbf{z}|\mathbf{x})}[\log p_\theta(\mathbf{x}|\mathbf{z})] - \beta \cdot \text{KL}(q_\phi(\mathbf{z}|\mathbf{x})||p(\mathbf{z})) \quad (4)$$

---

<sup>2</sup>Irina Higgins et al. "beta-VAE: Learning basic visual concepts with a constrained variational framework". In: *Proceedings of the International Conference on Learning Representations*. 2016.

# Deconstructing the $\beta$ -VAE — Entropy Regulariser

## Theorem 1

$\mathcal{L}_\beta(\mathbf{x})$  can be interpreted in terms of the standard ELBO for an adjusted target  $\pi_{\theta,\beta}(\mathbf{x}, \mathbf{z}) \triangleq p_\theta(\mathbf{x} | \mathbf{z}) f_\beta(\mathbf{z})$  with annealed prior  $f_\beta(\mathbf{z}) \triangleq p(\mathbf{z})^\beta / F_\beta$

$$\begin{aligned}\mathcal{L}_\beta(\mathbf{x}) &= \mathbb{E}_{q_\phi(\mathbf{z}|\mathbf{x})}[\log p_\theta(\mathbf{x}|\mathbf{z})] - \beta \cdot \text{KL}(q_\phi(\mathbf{z}|\mathbf{x}) || p(\mathbf{z})) \\ &= \underbrace{\mathcal{L}(\mathbf{x}) (\pi_{\theta,\beta}, q_\phi)}_{\text{ELBO with } \beta\text{-annealed prior}} + \underbrace{(\beta - 1) \cdot H_{q_\phi}}_{\text{maxent}} + \underbrace{\log F_\beta}_{\text{constant}}\end{aligned}$$

# Deconstructing the $\beta$ -VAE — Entropy Regulariser

## Theorem 1

$\mathcal{L}_\beta(\mathbf{x})$  can be interpreted in terms of the standard ELBO for an adjusted target  $\pi_{\theta,\beta}(\mathbf{x}, \mathbf{z}) \triangleq p_\theta(\mathbf{x} | \mathbf{z})f_\beta(\mathbf{z})$  with annealed prior  $f_\beta(\mathbf{z}) \triangleq p(\mathbf{z})^\beta / F_\beta$

$$\begin{aligned}\mathcal{L}_\beta(\mathbf{x}) &= \mathbb{E}_{q_\phi(\mathbf{z}|\mathbf{x})}[\log p_\theta(\mathbf{x}|\mathbf{z})] - \beta \cdot \text{KL}(q_\phi(\mathbf{z}|\mathbf{x})||p(\mathbf{z})) \\ &= \underbrace{\mathcal{L}(\mathbf{x})(\pi_{\theta,\beta}, q_\phi)}_{\text{ELBO with } \beta\text{-annealed prior}} + \underbrace{(\beta - 1) \cdot H_{q_\phi}}_{\text{maxent}} + \underbrace{\log F_\beta}_{\text{constant}}\end{aligned}$$

## Implications

$\beta$ -VAE disentangles largely by controlling the level of overlap

# Deconstructing the $\beta$ -VAE — Entropy Regulariser

## Theorem 1

$\mathcal{L}_\beta(\mathbf{x})$  can be interpreted in terms of the standard ELBO for an adjusted target  $\pi_{\theta,\beta}(\mathbf{x}, \mathbf{z}) \triangleq p_\theta(\mathbf{x} | \mathbf{z}) f_\beta(\mathbf{z})$  with annealed prior  $f_\beta(\mathbf{z}) \triangleq p(\mathbf{z})^\beta / F_\beta$

$$\begin{aligned}\mathcal{L}_\beta(\mathbf{x}) &= \mathbb{E}_{q_\phi(\mathbf{z}|\mathbf{x})}[\log p_\theta(\mathbf{x}|\mathbf{z})] - \beta \cdot \text{KL}(q_\phi(\mathbf{z}|\mathbf{x}) || p(\mathbf{z})) \\ &= \underbrace{\mathcal{L}(\mathbf{x})(\pi_{\theta,\beta}, q_\phi)}_{\text{ELBO with } \beta\text{-annealed prior}} + \underbrace{(\beta - 1) \cdot H_{q_\phi}}_{\text{maxent}} + \underbrace{\log F_\beta}_{\text{constant}}\end{aligned}$$

## Implications

$\beta$ -VAE disentangles largely by controlling the level of overlap  
It places no direct pressure on the latents to be independent!

# Deconstructing the $\beta$ -VAE — Rotational Invariance

## Theorem 2

If  $p(\mathbf{z}) = \mathcal{N}(\mathbf{z}; 0, \sigma I)$  and  $q_\phi(\mathbf{z} | \mathbf{x}) = \mathcal{N}(\mathbf{z}; \mu_\phi(\mathbf{x}), S_\phi(\mathbf{x}))$ , then for all rotation matrices  $R$ ,

$$\mathcal{L}_\beta(\mathbf{x}; \theta, \phi) = \mathcal{L}_\beta(\mathbf{x}; \theta^\dagger(R), \phi^\dagger(R))$$

where  $\theta^\dagger(R)$  and  $\phi^\dagger(R)$  are transformed networks such that

$$p_{\theta^\dagger}(\mathbf{x} | \mathbf{z}) = p_\theta(\mathbf{x} | R^T \mathbf{z}),$$
$$q_{\phi^\dagger}(\mathbf{z} | \mathbf{x}) = \mathcal{N}(\mathbf{z}; R\mu_\phi(\mathbf{x}), RS_\phi(\mathbf{x})R^T).$$

# Deconstructing the $\beta$ -VAE — Rotational Invariance

## Theorem 2

If  $p(\mathbf{z}) = \mathcal{N}(\mathbf{z}; 0, \sigma I)$  and  $q_\phi(\mathbf{z} | \mathbf{x}) = \mathcal{N}(\mathbf{z}; \mu_\phi(\mathbf{x}), S_\phi(\mathbf{x}))$ , then for all rotation matrices  $R$ ,

$$\mathcal{L}_\beta(\mathbf{x}; \theta, \phi) = \mathcal{L}_\beta(\mathbf{x}; \theta^\dagger(R), \phi^\dagger(R))$$

where  $\theta^\dagger(R)$  and  $\phi^\dagger(R)$  are transformed networks such that

$$p_{\theta^\dagger}(\mathbf{x} | \mathbf{z}) = p_\theta(\mathbf{x} | R^T \mathbf{z}),$$
$$q_{\phi^\dagger}(\mathbf{z} | \mathbf{x}) = \mathcal{N}(\mathbf{z}; R\mu_\phi(\mathbf{x}), RS_\phi(\mathbf{x})R^T).$$

## Implications

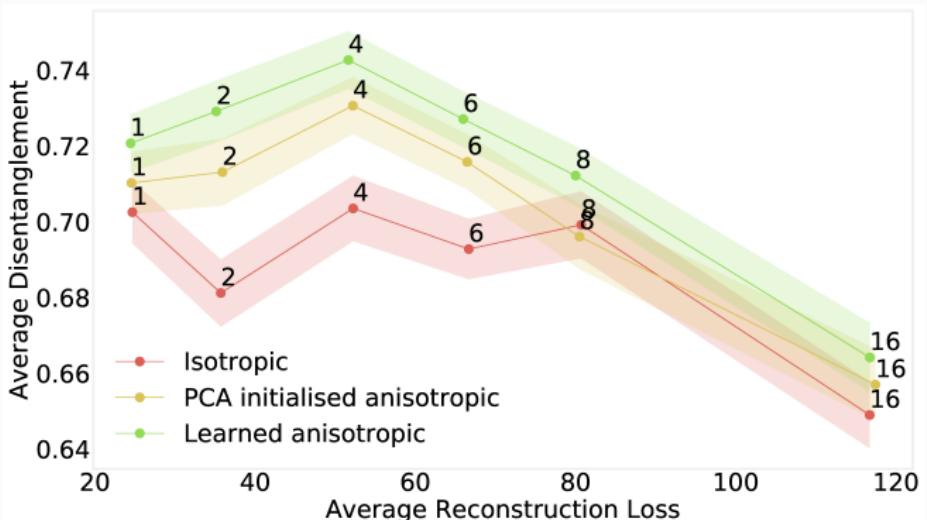
The  $\beta$ -VAE does not directly encourage meaningful representations when  $p(\mathbf{z}) = \mathcal{N}(\mathbf{z}; 0, \sigma I)$

## Decomposition: Objective

$$\mathcal{L}_{\alpha, \beta}(\mathbf{x}) = \mathbb{E}_{q_{\phi}(\mathbf{z}|\mathbf{x})}[\log p_{\theta}(\mathbf{x} | \mathbf{z})] \quad \text{Reconstruct observations}$$
$$- \beta \cdot \text{KL}(q_{\phi}(\mathbf{z} | \mathbf{x}) \| p(\mathbf{z})) \quad \text{Control level of overlap}$$
$$- \alpha \cdot \mathbb{D}(q_{\phi}(\mathbf{z}), p(\mathbf{z})) \quad \text{Impose desired structure}$$

# Decomposition: Generalising Disentanglement

**Independence:**  $p(z) = \mathcal{N}(\mathbf{0}, \sigma^*)$



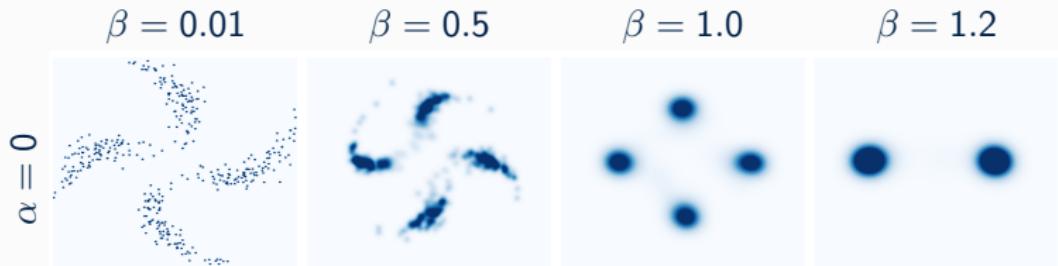
**Figure 1:**  $\beta$ -VAE trained on 2D Shapes<sup>3</sup> computing disentanglement<sup>4</sup>.

<sup>3</sup>Loic Matthey et al. *dSprites: Disentanglement testing Sprites dataset*. <https://github.com/deepmind/dsprites-dataset/>. 2017.

<sup>4</sup>Kim and Mnih, "Disentangling by Factorising".

# Decomposition: Generalising Disentanglement

**Clustering:**  $p(z) = \sum_{k=1}^4 \rho_k \cdot \mathcal{N}(\mu_k, \sigma_k)$



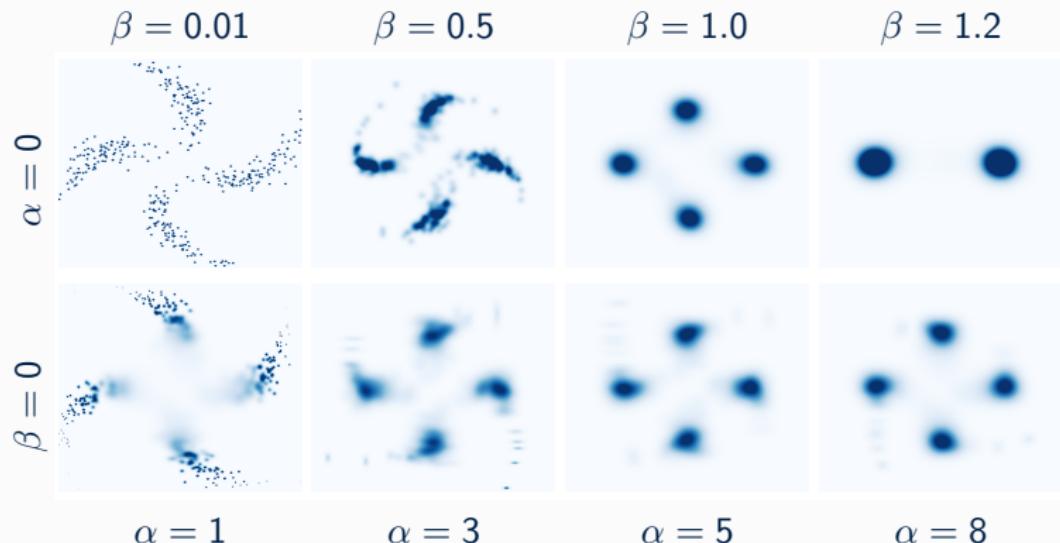
**Figure 2:** Density of aggregate posterior  $q_\phi(z)$  with different  $\alpha, \beta$  for the pinwheel dataset.<sup>5</sup>

---

<sup>5</sup><http://hips.seas.harvard.edu/content/synthetic-pinwheel-data-matlab>.

# Decomposition: Generalising Disentanglement

**Clustering:**  $p(z) = \sum_{k=1}^4 \rho_k \cdot \mathcal{N}(\mu_k, \sigma_k)$

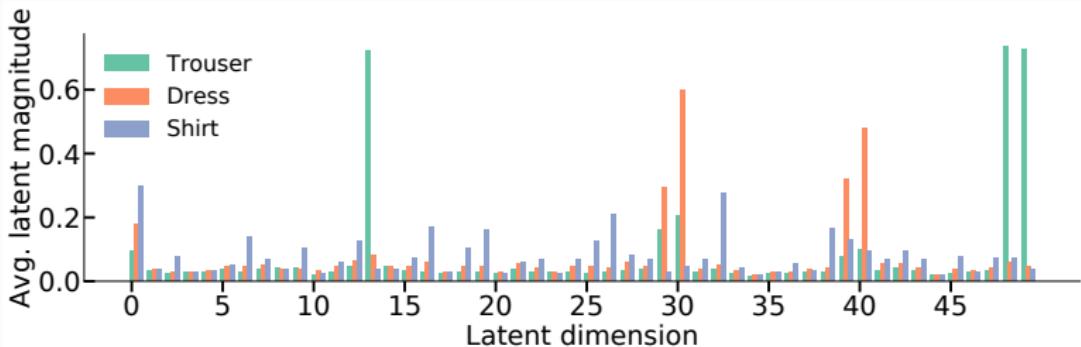


**Figure 2:** Density of aggregate posterior  $q_\phi(z)$  with different  $\alpha, \beta$  for the pinwheel dataset.<sup>5</sup>

<sup>5</sup><http://hips.seas.harvard.edu/content/synthetic-pinwheel-data-matlab>.

# Decomposition: Generalising Disentanglement

**Sparsity:**  $p(z) = \prod_d (1 - \gamma) \cdot \mathcal{N}(z_d; 0, 1) + \gamma \cdot \mathcal{N}(z_d; 0, \sigma_0^2)$

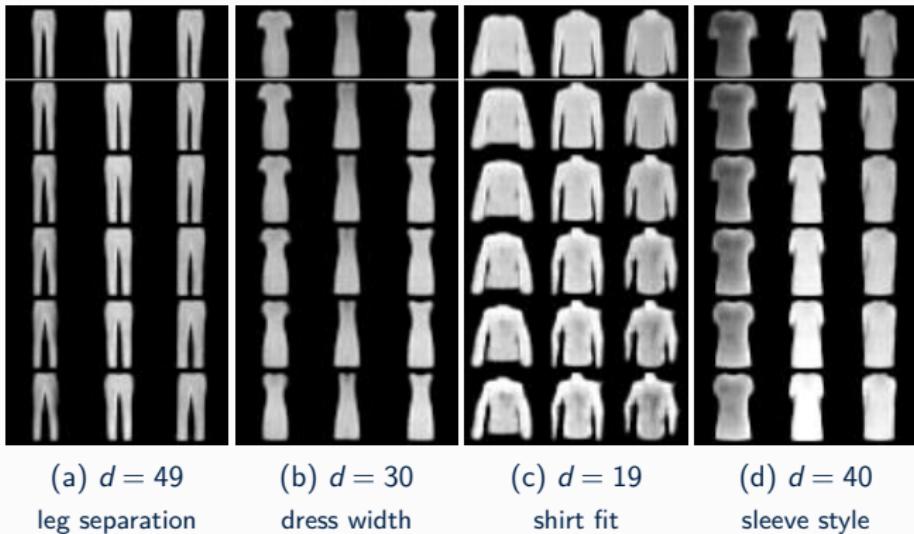


**Figure 3:** Sparsity of learnt representations for the *Fashion-MNIST*<sup>6</sup> dataset.

<sup>6</sup>Han Xiao, Kashif Rasul, and Roland Vollgraf. *Fashion-MNIST: a Novel Image Dataset for Benchmarking Machine Learning Algorithms*. Aug. 28, 2017. arXiv: cs.LG/1708.07747 [cs.LG].

# Decomposition: Generalising Disentanglement

**Sparsity:**  $p(z) = \prod_d (1 - \gamma) \cdot \mathcal{N}(z_d; 0, 1) + \gamma \cdot \mathcal{N}(z_d; 0, \sigma_0^2)$

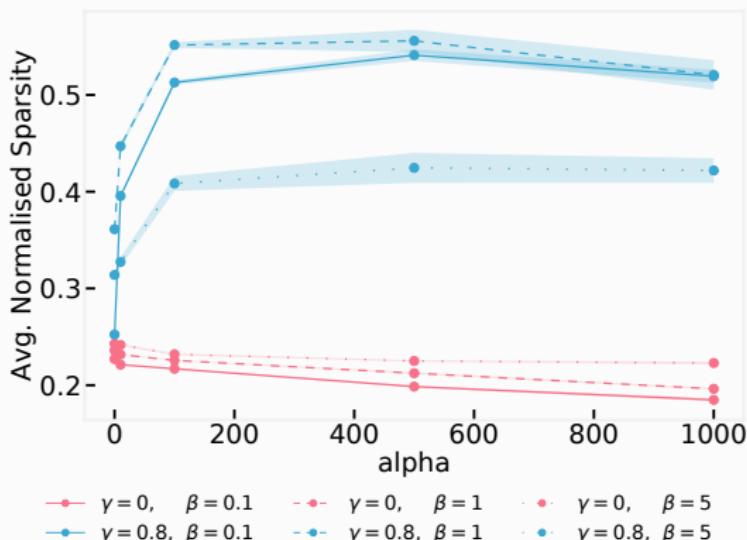


**Figure 3:** Latent space traversals for “active” dimensions<sup>6</sup>.

<sup>6</sup>Xiao, Rasul, and Vollgraf, *Fashion-MNIST: a Novel Image Dataset for Benchmarking Machine Learning Algorithms*.

# Decomposition: Generalising Disentanglement

**Sparsity:**  $p(z) = \prod_d (1 - \gamma) \cdot \mathcal{N}(z_d; 0, 1) + \gamma \cdot \mathcal{N}(z_d; 0, \sigma_0^2)$



**Figure 3:** Sparsity vs regularisation strength  $\alpha$  (higher better)<sup>6</sup>.

<sup>6</sup>Xiao, Rasul, and Vollgraf, *Fashion-MNIST: a Novel Image Dataset for Benchmarking Machine Learning Algorithms*.

# Recap

We propose and develop:

# Recap

We propose and develop:

- Decomposition: a generalisation of disentanglement involving:

# Recap

We propose and develop:

- Decomposition: a generalisation of disentanglement involving:
  - (a) match between  $q_\phi(\mathbf{z})$  and  $p(\mathbf{z})$

# Recap

We propose and develop:

- Decomposition: a generalisation of disentanglement involving:
  - (a) match between  $q_\phi(\mathbf{z})$  and  $p(\mathbf{z})$
  - (b) appropriate overlap of latent encodings

# Recap

We propose and develop:

- Decomposition: a generalisation of disentanglement involving:
  - (a) match between  $q_\phi(\mathbf{z})$  and  $p(\mathbf{z})$
  - (b) appropriate overlap of latent encodings
- A theoretical analysis of the  $\beta$ -VAE objective showing it primarily only contributes to overlap.

# Recap

We propose and develop:

- Decomposition: a generalisation of disentanglement involving:
  - (a) match between  $q_\phi(\mathbf{z})$  and  $p(\mathbf{z})$
  - (b) appropriate overlap of latent encodings
- A theoretical analysis of the  $\beta$ -VAE objective showing it primarily only contributes to overlap.
- An objective that incorporates both factors (a) and (b).

# Recap

We propose and develop:

- Decomposition: a generalisation of disentanglement involving:
  - (a) match between  $q_\phi(\mathbf{z})$  and  $p(\mathbf{z})$
  - (b) appropriate overlap of latent encodings
- A theoretical analysis of the  $\beta$ -VAE objective showing it primarily only contributes to overlap.
- An objective that incorporates both factors (a) and (b).
- Experiments that showcase efficacy at different decompositions:
  - independence   • clustering   • sparsity

## Discussion Points

How Do We Characterize Overlap?

## Discussion Points

### How Do We Characterize Overlap?

- Closely linked to the mutual information  $I(x; z)$

## Discussion Points

### How Do We Characterize Overlap?

- Closely linked to the mutual information  $I(x; z)$
- But two can noticeably vary, particularly for complicated encoders

## Discussion Points

### How Do We Characterize Overlap?

- Closely linked to the mutual information  $I(\mathbf{x}; \mathbf{z})$
- But two can noticeably vary, particularly for complicated encoders
- Can achieve any  $I(\mathbf{x}; \mathbf{z})$  with negligible overlap if

$$q_{\phi}(\mathbf{z}|\mathbf{x}) = \lambda \cdot \text{Uniform}(\|\mu_{\phi}(\mathbf{x}) - \mathbf{z}\|_2 < \epsilon) + (1 - \lambda) \cdot p(\mathbf{z}) \quad (5)$$

# Discussion Points

## How Do We Characterize Overlap?

- Closely linked to the mutual information  $I(\mathbf{x}; \mathbf{z})$
- But two can noticeably vary, particularly for complicated encoders
- Can achieve any  $I(\mathbf{x}; \mathbf{z})$  with negligible overlap if

$$q_{\phi}(\mathbf{z}|\mathbf{x}) = \lambda \cdot \text{Uniform}(\|\mu_{\phi}(\mathbf{x}) - \mathbf{z}\|_2 < \epsilon) + (1 - \lambda) \cdot p(\mathbf{z}) \quad (5)$$

- $I(\mathbf{x}; \mathbf{z})$  does not distinguish between large overlap with a small number of datapoints and small overlap with a large number

## Discussion Points

### How Do We Characterize Overlap?

- Closely linked to the mutual information  $I(\mathbf{x}; \mathbf{z})$
- But two can noticeably vary, particularly for complicated encoders
- Can achieve any  $I(\mathbf{x}; \mathbf{z})$  with negligible overlap if

$$q_{\phi}(\mathbf{z}|\mathbf{x}) = \lambda \cdot \text{Uniform}(\|\mu_{\phi}(\mathbf{x}) - \mathbf{z}\|_2 < \epsilon) + (1 - \lambda) \cdot p(\mathbf{z}) \quad (5)$$

- $I(\mathbf{x}; \mathbf{z})$  does not distinguish between large overlap with a small number of datapoints and small overlap with a large number
- It can also fail to account for *locality* in the latent space

### Can VAEs Uncover True Generative Factors?

---

<sup>7</sup> Francesco Locatello et al. "Challenging Common Assumptions in the Unsupervised Learning of Disentangled Representations". In: *International Conference on Machine Learning* (2019).

## Discussion Points

### Can VAEs Uncover True Generative Factors?

- Empirically this is well known to be very difficult

---

<sup>7</sup>Locatello et al., "Challenging Common Assumptions in the Unsupervised Learning of Disentangled Representations".

## Discussion Points

### Can VAEs Uncover True Generative Factors?

- Empirically this is well known to be very difficult
- Locatello et al. (2019)<sup>7</sup> suggest unsupervised disentanglement is impossible without inductive biases due to equivalence classes in the generative model

---

<sup>7</sup>Locatello et al., "Challenging Common Assumptions in the Unsupervised Learning of Disentangled Representations".

## Discussion Points

### Can VAEs Uncover True Generative Factors?

- Empirically this is well known to be very difficult
- Locatello et al. (2019)<sup>7</sup> suggest unsupervised disentanglement is impossible without inductive biases due to equivalence classes in the generative model
- This fails to account for the stochasticity of the encoder so do not directly apply to VAEs

---

<sup>7</sup>Locatello et al., "Challenging Common Assumptions in the Unsupervised Learning of Disentangled Representations".

## Discussion Points

### Can VAEs Uncover True Generative Factors?

- Empirically this is well known to be very difficult
- Locatello et al. (2019)<sup>7</sup> suggest unsupervised disentanglement is impossible without inductive biases due to equivalence classes in the generative model
- This fails to account for the stochasticity of the encoder so do not directly apply to VAEs
- When, and to what extent, can we expect learned representations to mimic the true generative process?

---

<sup>7</sup>Locatello et al., "Challenging Common Assumptions in the Unsupervised Learning of Disentangled Representations".

## Discussion Points

### Can VAEs Uncover True Generative Factors?

- Empirically this is well known to be very difficult
- Locatello et al. (2019)<sup>7</sup> suggest unsupervised disentanglement is impossible without inductive biases due to equivalence classes in the generative model
- This fails to account for the stochasticity of the encoder so do not directly apply to VAEs
- When, and to what extent, can we expect learned representations to mimic the true generative process?
- Should disentanglement/decomposition metrics include any notion of a true generative process?

---

<sup>7</sup>Locatello et al., "Challenging Common Assumptions in the Unsupervised Learning of Disentangled Representations".



Emile Mathieu



Tom Rainforth



N. Siddharth



Yee Whye Teh

Code



Paper



[iffsid/disentangling-disentanglement](https://github.com/iffsid/disentangling-disentanglement)

ICML 2019 arXiv:1812.02833