

Disentangling Disentanglement in Variational Autoencoders

E. Mathieu*, T. Rainforth*, N. Siddharth*, Y. W. Teh.

University of Oxford, Departments of Statistics & Engineering

{emile.mathieu, rainforth, y.w.teh}@stats.ox.ac.uk nsid@robots.ox.ac.uk

Overview

We develop a generalisation of disentanglement—*decomposition* of the latent representation—characterising it as the fulfilment of:

- (a) the data encodings having an appropriate level of overlap
- (b) the aggregate encoding of the data conforming to a desired structure, represented through the prior.

Decomposition permits disentanglement, that is, explicit **independence** between latents, as a special case, but also allows for a much richer class of properties to be imposed on the learnt representation, such as **sparsity**, **clustering**, **independent subspaces**, or even intricate **hierarchical dependency** relationships.

We introduce an alternative VAE training objective that allows one to **control the level of latent overlap**, and to **enforce alternate decompositions** than independence through the **prior distribution**.

Deconstructing the β -VAE

Theorem 1. The β -VAE [1] target

$$\mathcal{L}_\beta(\mathbf{x}) = \underbrace{\mathbb{E}_{q_\phi(\mathbf{z}|\mathbf{x})}[\log p_\theta(\mathbf{x}|\mathbf{z})]}_{\text{reconstruction error}} - \underbrace{\beta \text{KL}(q_\phi(\mathbf{z}|\mathbf{x})||p_\theta(\mathbf{z}))}_{\text{per data posterior regularisation}} \quad (1)$$

can be seen as the standard ELBO, $\mathcal{L}(\mathbf{x}) (\pi_{\theta,\beta}, q_\phi)$, 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 $f_\beta(\mathbf{z}) \triangleq p_\theta(\mathbf{z})^\beta/F_\beta$

$$\mathcal{L}_\beta(\mathbf{x}) = \underbrace{\mathcal{L}(\mathbf{x}) (\pi_{\theta,\beta}, q_\phi)}_{\text{ELBO with } \beta\text{-annealed prior}} + \underbrace{(\beta-1)H_{q_\phi}}_{\text{maxent}} + \underbrace{\log F_\beta}_{\text{constant}} \quad (2)$$

where $F_\beta \triangleq \int_z p_\theta(z)^\beta dz$, and H_{q_ϕ} is the entropy of $q_\phi(\mathbf{z}|\mathbf{x})$.

Theorem 2. If $p_\theta(\mathbf{z}) = \mathcal{N}(\mathbf{z}; 0, \sigma I)$ 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 ,

$$\mathcal{L}_\beta(\mathbf{x}; \theta, \phi) = \mathcal{L}_\beta(\mathbf{x}; \theta^\dagger(R), \phi^\dagger(R)) \quad (3)$$

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

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

Implications:

- The β -VAE's disentanglement is largely down to control of overlap: it provides **no direct pressure** for independence
- Larger β are **not universally beneficial** for disentanglement
- Dispels the conjecture that the β -VAE objective encourages meaningful independent latent variables when using the standard choice of an isotropic Gaussian prior

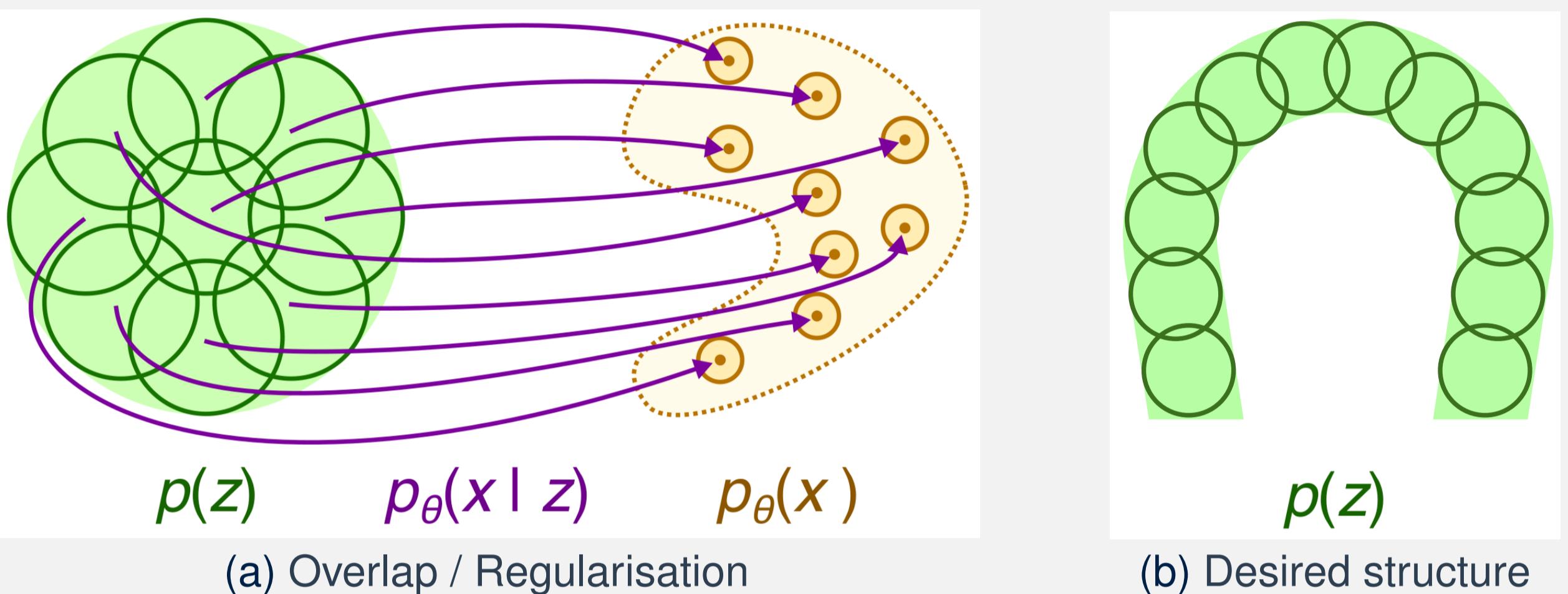
EM, TR, YWT were supported in part by the ERC under the EU's Seventh Framework Programme (FP7/2007–2013) / ERC grant agreement no. 617071. EM was also supported by Microsoft Research through its PhD Scholarship Programme. NS was funded by EPSRC/MURI grant EP/N019474/1.

[1] Irina Higgins, Loic Matthey, Arka Pal, Christopher Burgess, Xavier Glorot, Matthew Botvinick, Shakir Mohamed, and Alexander Lerchner. beta-VAE: Learning basic visual concepts with a constrained variational framework. In *Proceedings of the International Conference on Learning Representations (ICLR)*, 2016.

[2] Hyunjik Kim and Andriy Mnih. Disentangling by factorising. *CoRR*, abs/1802.05983, 2018.

[3] Loic Matthey, Irina Higgins, Demis Hassabis, and Alexander Lerchner. dSprites: Disentanglement testing sprites dataset. <https://github.com/deepmind/dsprites-dataset/>, 2017.

Decomposition allows richer structure than independence to be imposed on the latent encodings. It can be achieved with sufficient latent overlap and an appropriate prior.



Proposal: an objective to enforce a particular decomposition.

$$\begin{aligned} \mathcal{L}_{\alpha, \beta}(\mathbf{x}) &= \mathbb{E}_{q_\phi(\mathbf{z}|\mathbf{x})}[\log p_\theta(\mathbf{x}|\mathbf{z})] \\ &\quad - \underbrace{\beta \text{KL}(q_\phi(\mathbf{z}|\mathbf{x})||p(\mathbf{z}))}_{(a)} \\ &\quad - \underbrace{\alpha \mathbb{D}(q_\phi(\mathbf{z}), p(\mathbf{z}))}_{(b)} \end{aligned} \quad (4)$$

Reconstruct observations

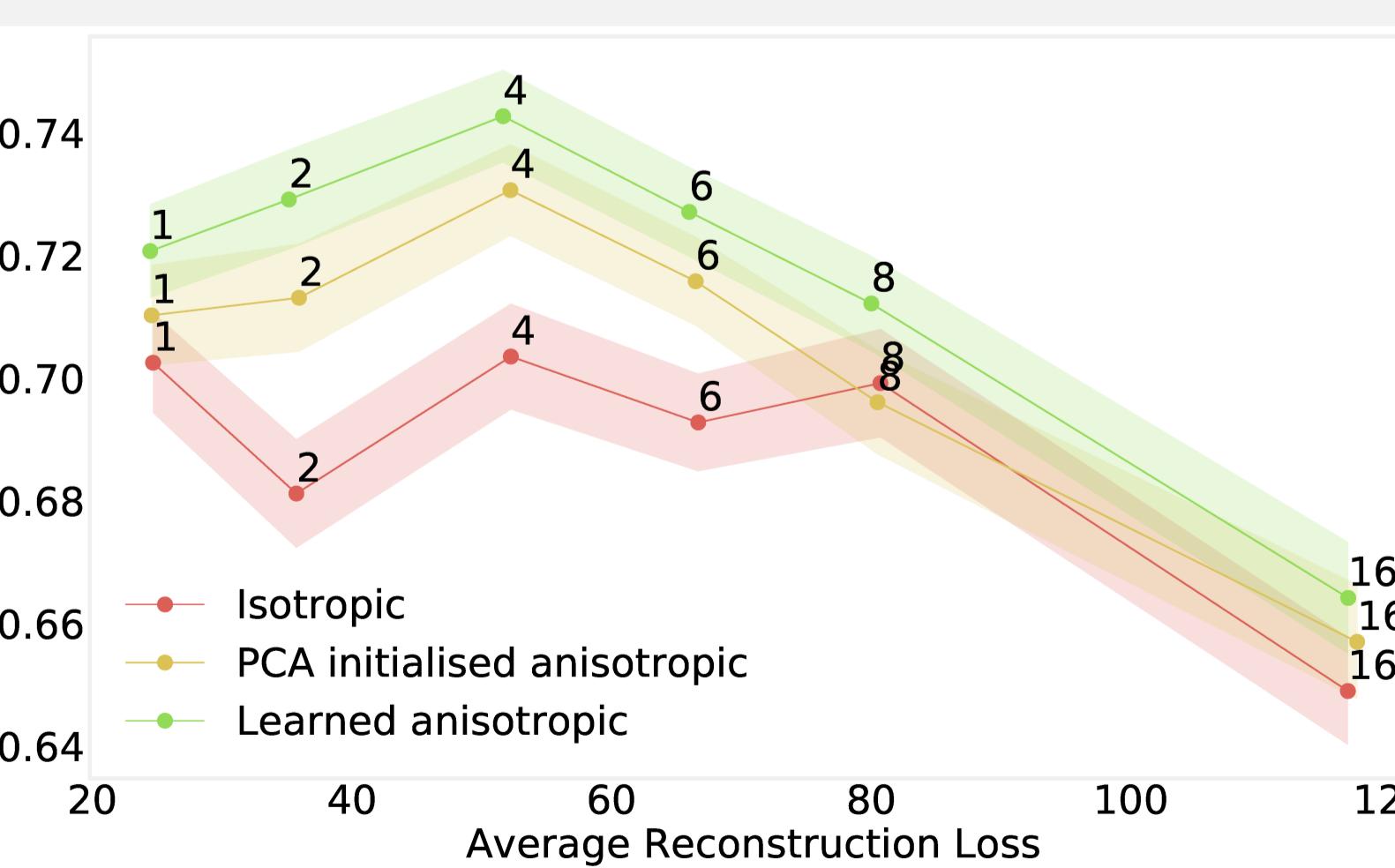
Control level of overlap

Impose desired structure

Prior for axis-aligned disentanglement

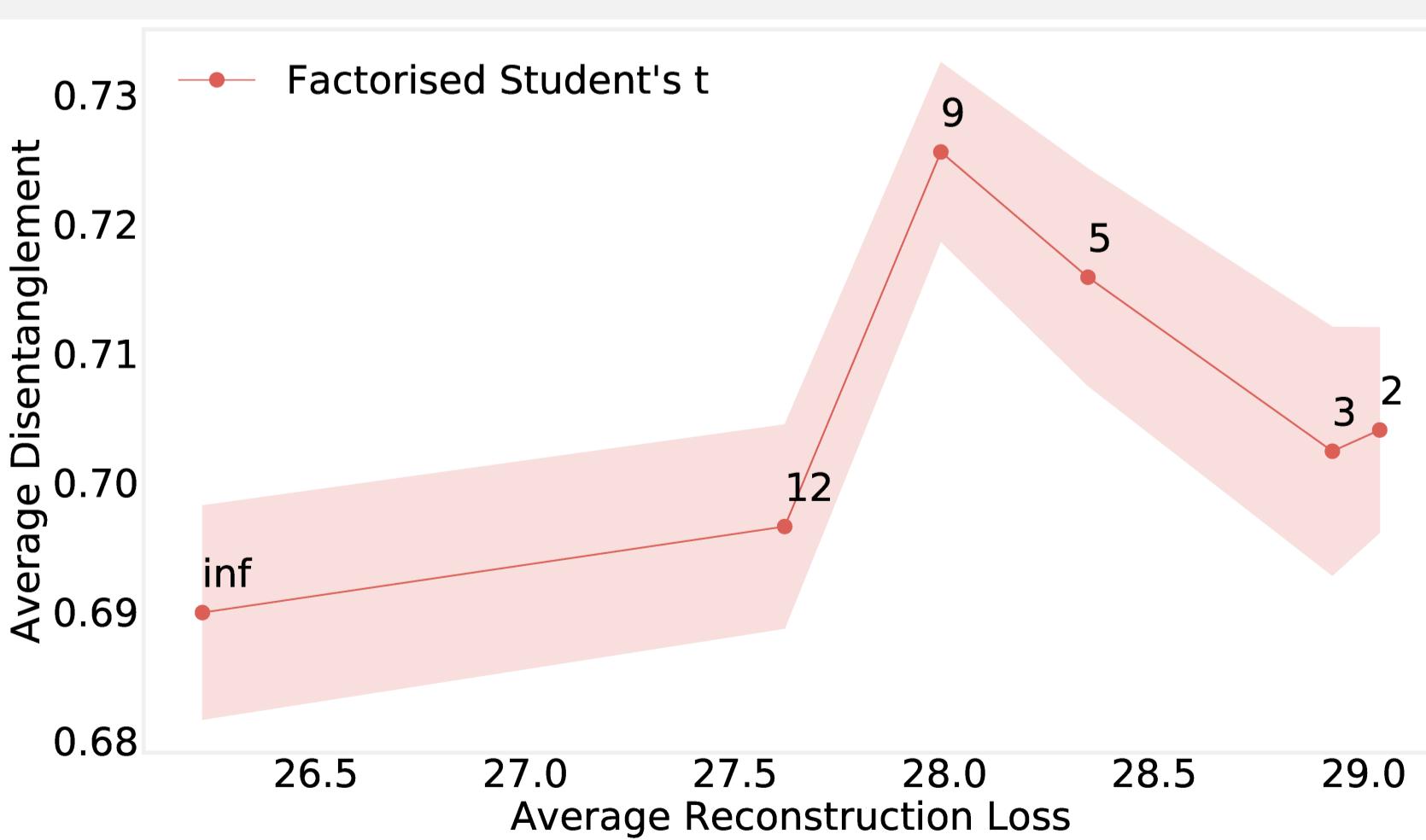
We demonstrate that substantial improvements in **disentanglement** [2] can be achieved by breaking the rotational invariance of the prior, using either

► A non-isotropic Gaussian (Fig. 2a)



(a) β -VAE (i.e. (4) with $\alpha = 0$) trained on the 2D Shapes dataset [3]. Using an anisotropic Gaussian with diagonal covariance either fixed to the principal component values or learned during training.

► A product of Student-t's (Fig. 2b)



(b) Using $p_\theta^\nu(\mathbf{z}) = \prod_i \text{STUDENT-T}(z_i; \nu)$ for different degrees of freedom ν with $\beta = 1$. Note that $p_\theta^\nu(\mathbf{z}) \rightarrow \mathcal{N}(\mathbf{z}; 0, \mathbf{I})$ as $\nu \rightarrow \infty$, and reducing ν only incurs a minor increase in reconstruction loss.

Sparse Prior

We demonstrate important improvements in **sparsity** by using a relaxed **Spike and Slab** prior $p(\mathbf{z}) = \prod_d (1-\gamma) \mathcal{N}(z_d; 0, 1) + \gamma \mathcal{N}(z_d; 0, \sigma_0^2)$, and a dimension-wise MMD regulariser.

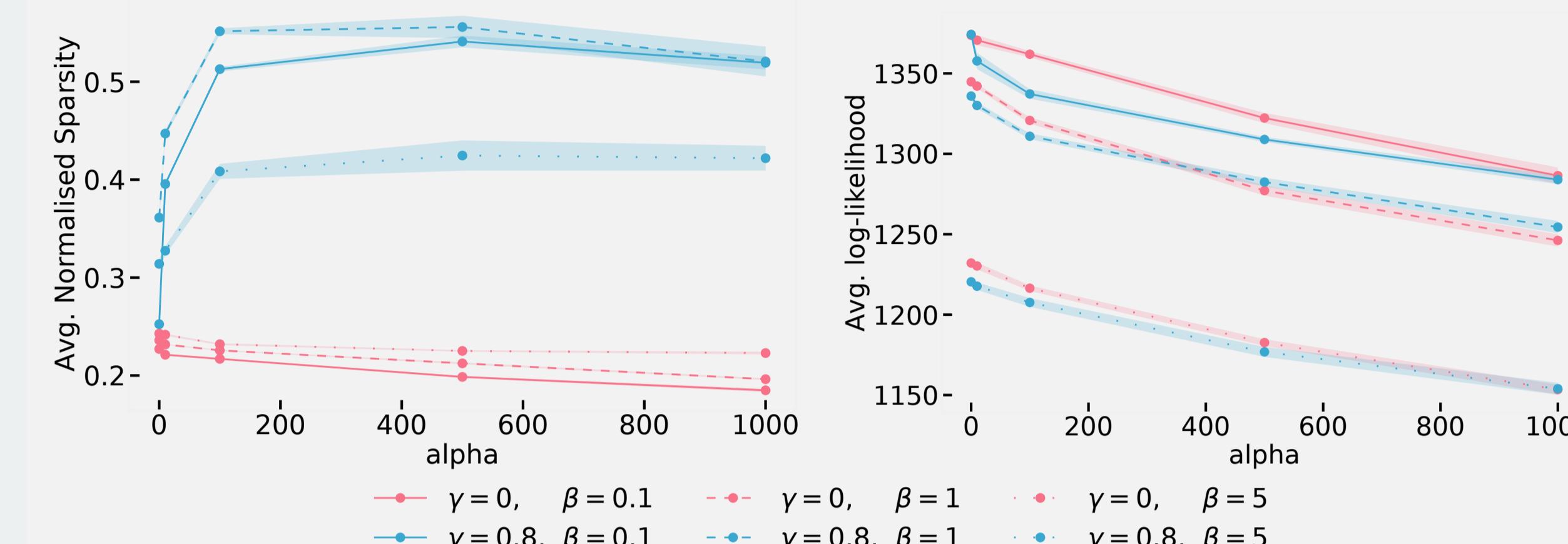


Figure 3: [Left] Sparsity vs regularisation strength α (c.f. Eq. (4), high better). [right] Avg reconstruction log-likelihood vs α (higher better).

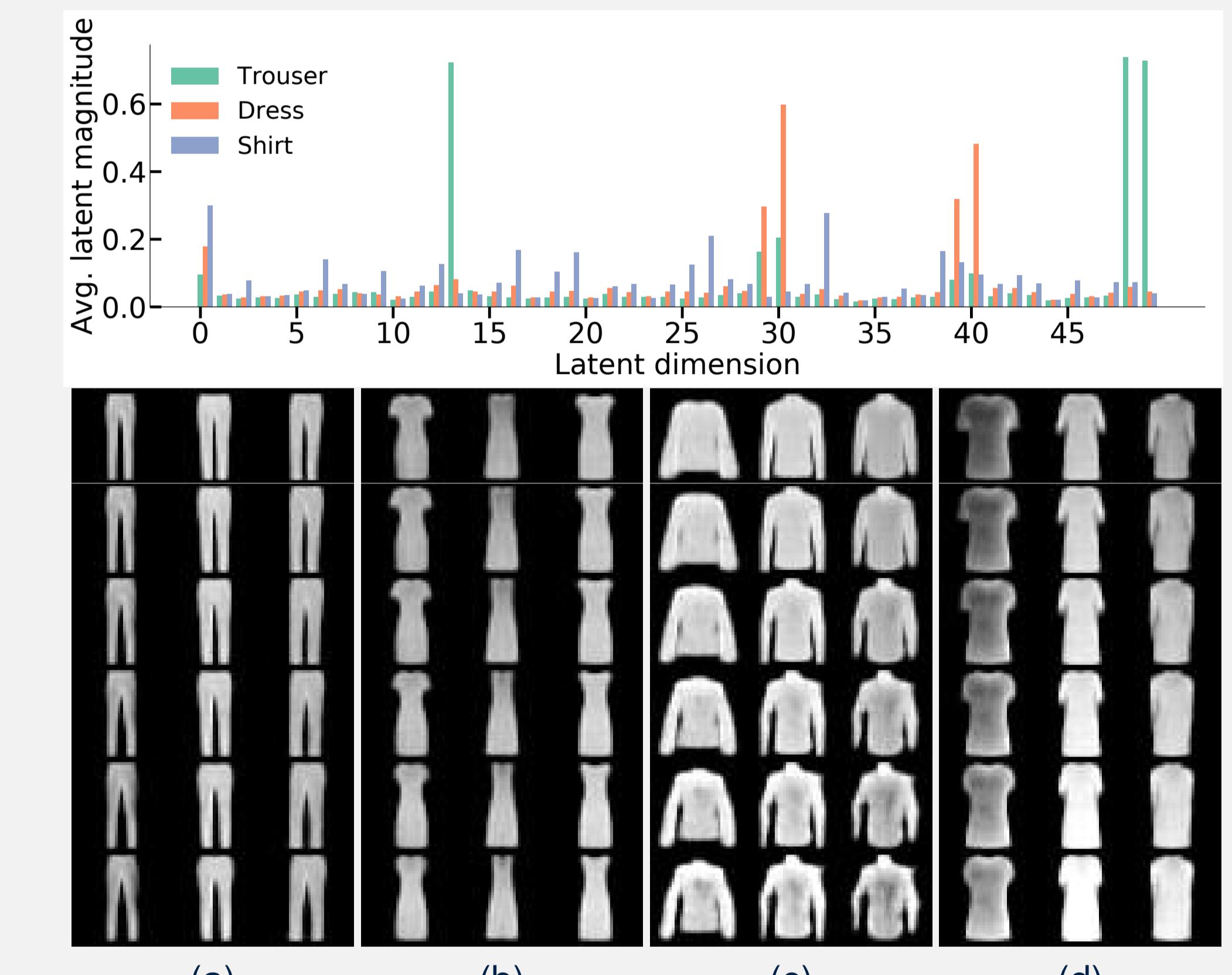


Figure 4: Qualitative evaluation of sparsity.

Clustered prior: Gaussian mixture

As shown in Figure 5, our framework allows to control

- The level of overlap
- The form of the marginal posterior

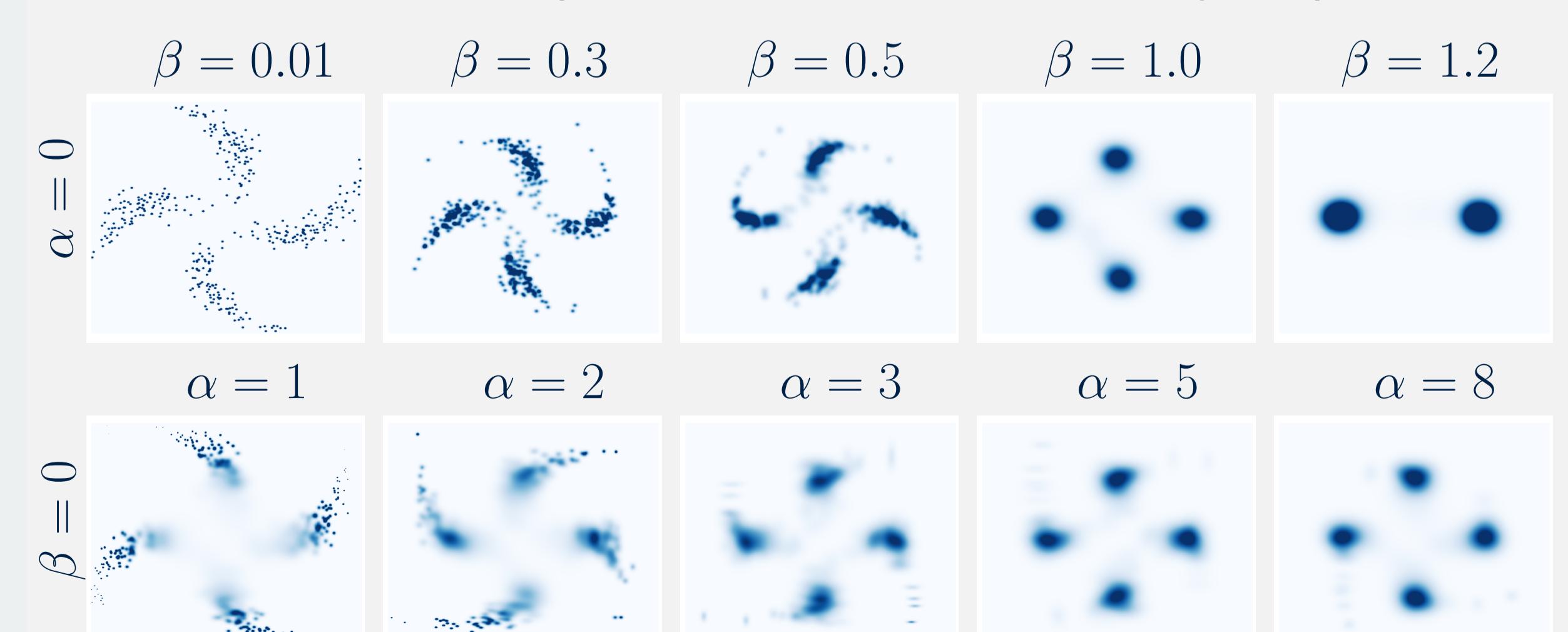


Figure 5: Density of aggregate posterior $q_\phi(\mathbf{z})$ for different values of α and β . Models are trained on the "pinwheels" dataset, minimising objective (4).