UNPICKING DATA AT THE SEAMS: UNDERSTANDING DISENTANGLEMENT IN VAES

Anonymous authors

Paper under double-blind review

ABSTRACT

A generative latent variable model is said to be *disentangled*, when varying a single latent co-ordinate changes a single aspect of samples generated, e.g. object position or facial expression in an image. Related phenomena are seen in several generative paradigms, including state-of-the-art diffusion models, but disentanglement is most notably observed in Variational Autoencoders (VAEs), where oft-used *diagonal* posterior covariances are argued to be the cause. We make this picture precise. From a known exact link between optimal Gaussian posteriors and decoder derivatives, we show how diagonal posteriors "lock" a decoder's local axes so that density over the data manifold *factorises* along *independent* one-dimensional seams that map to *axis-aligned* directions in latent space. This gives a clear definition of disentanglement, explains why it emerges in VAEs and shows that, under stated assumptions, ground truth factors are *identifiable* even with a symmetric prior.

1 Introduction

Variational Autoencoders (VAEs, Kingma & Welling, 2014; Rezende et al., 2014) and variants, such as β -VAE (Higgins et al., 2017) and FactorVAE (Kim & Mnih, 2018), are often observed to *disentangle* the data, whereby changing an individual co-ordinate in latent space causes generated samples to vary in a single semantically meaningful way, such as the hair colour or facial expression in an image. This phenomenon is both of practical use, e.g. for controlled data generation, and intriguing as it is not knowingly designed into a VAE's training algorithm. Related phenomena are observed in samples of Generative Adversarial Networks (GANs) (Goodfellow et al., 2014) and diffusion models (Rombach et al., 2022; Pandey et al., 2022; Zhang et al., 2022; Yang et al., 2023).

While disentanglement lacks a formal definition, it commonly refers to identifying *generative factors* of the data (Bengio et al., 2013). Thus a better understanding of disentanglement, and how it arises "for free" in a VAE, seems relevant to many areas of machine learning, its interpretability and our understanding of the data; and may enable us to induce disentanglement reliably in domains where it cannot be readily perceived, e.g. gene sequence or protein modelling. We focus on disentanglement in VAEs, where it is well observed (Higgins et al., 2017; Burgess et al., 2018), with the expectation that a clearer understanding of it there, more as a property of the *data* than the model, may extend to other settings, such as state of the art diffusion models.

The cause of disentanglement in VAEs has been traced to *diagonal* posterior covariance matrices (Rolinek et al., 2019; Kumar & Poole, 2020), a standard choice for computational efficiency. Approximate relationships suggest that diagonal covariances promote *orthogonality* between columns in a VAE decoder's Jacobian, a property empirically associated with disentangled features (Ramesh et al., 2018; Gresele et al., 2021). Taking inspiration from this, we develop a full theoretical explanation of disentanglement and how it arises in $(\beta$ -)VAEs. Specifically, we:

- formally define disentanglement as factorising the density over a manifold into independent 1-D *seam* factors, each the push-forward of the density over an axis-aligned latent path (D1, Fig. 1);
- show that β of a β -VAE (Higgins et al., 2017) effectively controls the model variance Var[x|z], which explains why it is found to enhance disentanglement and mitigate "posterior collapse" (§B);
- identify constraints that are necessary and sufficient for disentanglement, and show that diagonal posterior covariances induce them, in aggregate and in expectation (§4); and
- prove that if a data distribution has ground truth independent factors, those factors can be *identified* by a Gaussian VAE, resolving the "unidentifiability" in non-linear ICA with Gaussian priors (§6).



Figure 1: **Disentanglement: full vs diagonal posteriors** (Σ_x) . Right singular vectors $v^i \in \mathbb{Z}$ (blue) of the decoder's Jacobian define singular vector paths (dashed blue); left singular vectors u^i define seams (dashed red). 1-D densities over seams factorise the manifold density. (left) with full posteriors, s.v. paths are not axis-aligned; the axis-traversal image in \mathcal{X} (green) does not follow the seam. (right) under C1-C2 induced by diagonal posteriors, s.v. paths axis-align and the traversal image follows the seam everywhere, and 1-D densities over seams are independent, achieving disentanglement (D1).

2 BACKGROUND

Notation: Let $x \in \mathcal{X} \doteq \mathbb{R}^m$, $z \in \mathcal{Z} \doteq \mathbb{R}^d$ denote data and latent variables $(d \leq m)$. For continuous $g: \mathcal{Z} \to \mathcal{X}$ differentiable at z, let J_z denote its Jacobian evaluated at z ($[J_z]_{ij} = \frac{\partial x_i}{\partial z_j}$) with singular value decomposition (SVD) $J_z = U_z S_z V_z^\top (U_z^\top U_z = I, V_z^\top V_z = V_z V_z^\top = I)$. Let $s^i \doteq S_{ii}$ denote the i^{th} singular value, and u^i/v^i the i^{th} left/right singular vectors (columns of U/V). We consider continuous, injective functions g differentiable g. (abbreviated g., which, e.g., admit ReLU networks. Such g define a g-dimensional manifold g and g and g are g embedded in g (see Fig. 3). Since g is injective, there exists a bijection between g and g and g and g has full-rank, where defined.

Latent Variable Model (LVM): We consider the generative model $p_{\theta}(x) = \int_{z} p_{\theta}(x|z)p(z)$ with independent z_{i} . For tractability, parameters θ are typically learned by maximising a lower bound (**ELBO**)

$$\int_{x} p(x) \log p_{\theta}(x) \geq \int_{x} p(x) \int_{z} q_{\phi}(z|x) \left(\log p_{\theta}(x|z) - \beta \log \frac{q_{\phi}(z|x)}{p(z)} \right), \tag{1}$$

where $\beta = 1$ and $q_{\phi}(z|x)$ learns to approximate the model posterior, $q_{\phi}(z|x) \to p_{\theta}(z|x) \doteq \frac{p_{\theta}(x|z)p(z)}{p_{\theta}(x)}$.

Variational Autoencoder (VAE): A VAE parameterises Eq. 1 with neural networks: a decoder network d(z) parameterises the likelihood $p_{\theta}(x|z)$; and an encoder network parameterises the typically Gaussian posteriors $q_{\phi}(z|x) = \mathcal{N}(z; e(x), \Sigma_x)$ with diagonal Σ_x . The prior p(z) is typically a standard Gaussian. We refer to a VAE with Gaussian likelihood $p_{\theta}(x|z) \doteq \mathcal{N}(x; d(z), \sigma^2 I)$ as a Gaussian VAE and to a Gaussian VAE with linear decoder d(z) = Dz, $D \in \mathbb{R}^{m \times d}$ as a linear VAE.

Disentanglement: While not well defined, disentanglement typically refers to associating distinct semantically meaningful features of the data with distinct latent co-ordinates z_i , such that data generated by varying a single z_i differ in a single semantic feature (Bengio et al., 2013; Higgins et al., 2017; Ramesh et al., 2018; Rolinek et al., 2019; Shu et al., 2019). While samples from a VAE exhibit disentanglement, setting $\beta > 1$ (a β -VAE) often enhances the effect, although at a cost to generative quality, e.g. blurrier images (Higgins et al., 2017; Burgess et al., 2018). Disentanglement relates closely to independent component analysis (**ICA**), which aims to recover statistically independent components of the data under the same LVM but with a deterministic observation model, $p_{\theta}(x|z) = \delta_{x-d(z)}$.

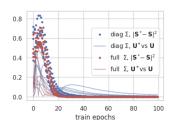
Probabilistic PCA (PPCA): (Tipping & Bishop, 1999) considers a linear Gaussian LVM

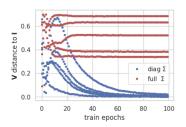
$$p(x|z) = \mathcal{N}(x; \mathbf{W}z, \sigma^2 \mathbf{I}) \qquad p(z) = \mathcal{N}(z; \mathbf{0}, \mathbf{I})$$
(2)

where $\mathbf{W} \in \mathbb{R}^{m \times d}$ and $\sigma \in \mathbb{R}$.² The exact posterior p(z|x) and MLE parameter \mathbf{W}_* are fully tractable:

$$p(z|x) = \mathcal{N}(z; \frac{1}{\sigma^2} \boldsymbol{M} \boldsymbol{W}^{\top} x, \boldsymbol{M}), \quad \boldsymbol{M} = (\boldsymbol{I} + \frac{1}{\sigma^2} \boldsymbol{W}^{\top} \boldsymbol{W})^{-1}; \quad \boldsymbol{W}_* = \boldsymbol{U}_{\boldsymbol{X}} (\boldsymbol{\Lambda}_{\boldsymbol{X}} - \sigma^2 \boldsymbol{I})^{1/2} \boldsymbol{R}$$
(3)

¹To lighten notation, explicit dependence of U, V, S, u^i, v^i, s^i on z is often suppressed where context is clear. ²We assume that data is centred, which is equivalent to including a mean parameter (Tipping & Bishop, 1999).





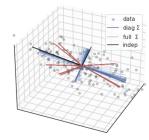


Figure 2: An LVAE breaks rotational symmetry. (*I*) both full- Σ_x and diagonal- Σ_x VAEs fit the data, i.e. learn ground truth parameters U_* , S_* ; (*c*) only in diagonal- Σ_x VAEs do right singulars v^i of the Jacobian align with standard basis vectors $z_i \in \mathcal{Z}$ (i.e. $V_* \to I$); (*r*) images of z_i : full- Σ_x VAEs map z_i to arbitrary directions (red), but diagonal- Σ_x VAEs learn (later epochs darker) to map z_i to the data's independent components (black, i.e. blue \to black).

Linear VAE: An LVAE assumes the same linear LVM as PPCA (Eq. 2) and models the likelihood $p_{\theta}(x|z) = \mathcal{N}(x; \mathbf{D}z, \sigma^2 \mathbf{I})$ as in PPCA^{EM}, differing only in *approximating* the posterior by $q_{\phi}(z|x) = \mathcal{N}(z; \mathbf{E}x, \Sigma)$, rather than computing the optimal $p_{\theta}(z|x)$. Surprisingly though, an LVAE with *diagonal* posterior covariances loses the rotational ambiguity of PPCA (Lucas et al., 2019), since

$$\Sigma_* \stackrel{(3(l))}{=} (I + \frac{1}{\sigma^2} W_*^{\top} W_*)^{-1} \stackrel{(3(r))}{=} \sigma^2 R^{\top} \Lambda_{\mathbf{X}}^{-1} R.$$
 (4)

Thus for Σ to be optimal *and* diagonal, R must belong to a finite set of signed permutations, hence the optimal decoder $D_* = U_X (\Lambda_X - \sigma^2 I)^{1/2}$ is unique up to permutation/sign (see Fig. 2). We will see that this effect, due to diagonal posterior covariances, is in fact (linear) *disentanglement* (§3).

Further notation: Under the LVM above, we define a deterministic generative function $g: \mathbb{Z} \to \mathcal{X}$ as the map from latent variables to means $g(z) = \mathbb{E}[x|z]$, that lie on a manifold $\mathcal{M}_g = \{g(z)\} \subseteq \mathcal{X}$ (mean manifold) with push-forward density p_μ (manifold density). We will focus on Gaussian VAEs, where the data density p(x) is given by adding Gaussian noise to p_μ , i.e. convolving it with a Gaussian kernel. It is known that such data densities match if and only if their manifold densities p_μ match (e.g. Khemakhem et al., 2020), hence we focus on the manifold density p_μ .

3 DISENTANGLEMENT

We now define disentanglement; illustrate it for the linear case, justifying our disentanglement claim for LVAEs in §2; and work up to explaining how it arises in a (non-linear) Gaussian VAE. (See Fig. 1)

Definition D1 (**Disentanglement**). Let $g: \mathcal{Z} \to \mathcal{X}$ be c.i.d.a.e.. We say p_{μ} is disentangled if, for each $z \in \mathcal{Z}$, there exist 1-D densities $\{f_i\}$ such that p_{μ} factorises as

$$p_{\mu}(g(z)) = \prod_{i=1}^{d} f_i(u_i(z)), \qquad (5)$$

where each factor f_i is the 1-D push-forward of $p(z_i)$ along the <u>axis-aligned</u> line obtained by moving in the *i*-th latent coordinate while keeping all others fixed; u_i is the co-ordinate of g(z) along the image of that line; and random variables $\{u_i(z)\}$ are mutually independent under $z \sim p(z)$.

LVAE disentanglement: Consider an LVAE with diagonal posterior covariance Σ and decoder d(z) = Dz, $D \in \mathbb{R}^{m \times d}$ (§2). The mean manifold $\mathcal{M}_d \doteq \{\mu = Dz \mid z \in \mathcal{Z}\}$ is linear with Gaussian density $p_{\mu} = \mathcal{N}(\mu; \mathbf{0}, DD^{\top})$ (e.g. see Fig. 2, right), From §2, the SVD of the data matrix defines the optimal decoder $D_* = U_*S_*V_*$ (i.e. $U_* \doteq U_X$, $S_* \doteq (\Lambda_X - \sigma^2 I)^{1/2}$ and $V_* = R = I$ due to diagonal Σ). As for any Gaussian, p_{μ} factorises as a product of independent 1-D Gaussians along eigenvectors of its covariance $D_*D_*^{\top} = U_*S_*^2U_*^{\top}$, i.e. columns u^i of U_* , hence $p_{\mu} = \prod_i \mathcal{N}(u_i; 0, s^{i2})$ where $u_i \doteq u^{i\top}\mu \in \mathbb{R}$. Since $u_i = u^{i\top}D_*z = s^iz_i$, each u_i depends only on a distinct z_i , a co-ordinate in the standard basis of \mathcal{Z} (over which densities are independent 1-D Gaussian). Thus:

- p_{μ} factorises as a product of independent push-forward densities $f_i(\mu) = \mathcal{N}(\mathbf{u}^{i\top}\mu; 0, s^{i2})$; and
- the decoder maps each axis-aligned direction z_i to a distinct factor f_i ,

satisfying D1. Note that synthetic data $\mu = Dz$ generated by re-sampling z_i , holding $z_{j\neq i}$ constant, differ only in component (or "feature") u_i , agreeing with the common perception of disentanglement.

Dropping diagonality: To emphasize that disentanglement depends on diagonal posteriors, we consider *full* posterior LVAEs, where $R \neq I$ in general. The above argument follows except that columns r^i of R in Z map to independent u^i directions in X. Meanwhile, standard basis vectors in Z map in directions $u^{i\top}R$, which are arbitrary with respect to u^i directions. Hence axis-aligned traversals in latent space correspond to several *entangled* components u_i changing in generated samples. We demonstrate this empirically in Fig. 2 (see caption for details).

4 From Diagonal Posteriors to Decoder Constraints

Prior works draw a link between disentanglement in Gaussian VAEs and diagonal posteriors from an *approximate* relationship between optimal posteriors and decoder derivatives (Rolinek et al., 2019; Kumar & Poole, 2020). In fact, this relationship is *exact* by the Price/Bonnet Theorem and Opper & Archambeau (2009): the ELBO with Gaussian posteriors is optimised when their covariances satisfy

$$\Sigma_x^{-1} = \boldsymbol{I} - \frac{1}{\beta} \mathbb{E}_{q(z|x)} [\mathbf{L}_z(x)] \stackrel{*}{=} \boldsymbol{I} + \frac{1}{\beta \sigma^2} \mathbb{E}_{q(z|x)} [\boldsymbol{J}_z^{\top} \boldsymbol{J}_z + (x - d(z))^{\top} \boldsymbol{\mathsf{H}}_z], \tag{6}$$

where $\mathbf{L}_z(x) = \nabla_z^2 \log p_\theta(x|z)$ is the log likelihood Hessian; and $\mathbf{J}_z \doteq \frac{dx}{dz}$ and $\mathbf{H}_z \doteq \frac{d^2x}{dz^2}$ are the Jacobian and Hessian of the decoder (all terms evaluated at $z \in \mathcal{Z}$). Step two (*) assumes the likelihood is Gaussian. Eq. 6 immediately generalises the classical linear result in Eq. 3 and relates $\sigma^2 \doteq \operatorname{Var}[x|z]$ and $\Sigma_x \doteq \operatorname{Var}[z|x]$, showing that (un)certainty in x and z go hand in hand, as expected.

Importantly to disentanglement, Eq. 6 shows that diagonal Σ_x constrains derivatives of an optimal VAE decoder. In practice, the $J_z^{\top}J_z$ term *alone* is found to be approximately diagonal (Fig. 3, *left*) (Rolinek et al., 2019; Kumar & Poole, 2020), suggesting that *each* term diagonalises. We thus consider:

Property P1. The expectations of $J_z^{\mathsf{T}}J_z$ and $(x-d(z))^{\mathsf{T}}H_z$ in Eq. 6 are each diagonal.³

P1 has the following implications, which will prove fundamental to disentanglement:

Lemma 4.1 (Disentanglement constraints). *Under P1, in expectation for concentrated posteriors:*

- C1 Right singular vectors V_z of the decoder Jacobian J_z are standard basis vectors for all $z \in \mathcal{Z}$, i.e. after relabeling/sign flips of the latent axes, we have $V_z = I$;
- C2 The matrix of partial derivatives of singular values $(\frac{\partial s_i}{\partial z_j})_{i,j}$ is diagonal, i.e. $\frac{\partial s_i}{\partial z_j} = 0$ for all $i \neq j$. Proof. See Appendix A. C1 follows from the SVD of J_z ; C2 from observing that directions $r = x - d(z) \in \mathcal{X}$ of the directed Hessian term are, to a first approximation, tangent to the manifold.

Why β affects disentanglement: Setting $\beta > 1$ in Eq. 1 is found to increase disentanglement (Higgins et al., 2017; Burgess et al., 2018). We show that β implicitly controls the likelihood's variance in Appendix B and, as a result, $\beta > 1$ dilates posteriors (while the β -ELBO remains a valid objective). We will see that C1-C2 induce disentanglement (§5), thus Eq. 6 suggests a rationale for why $\beta > 1$ enhances disentanglement: it broadens the regions (i.e. posteriors) over which decoder derivatives are diagonalised, and so disentanglement constraints C1-C2 encouraged; and increases the overlap of posteriors where multiple constraints apply simultaneously (see Fig. 6, *right*).

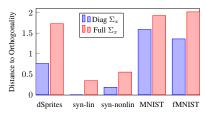
5 From Decoder Constraints to Disentanglement

5.1 THE JACOBIAN SVD: SINGULAR VECTOR PATHS AND DATA SEAMS

Constraints C1-C2, fundamental to disentanglement, are expressed in terms of SVD components of the Jacobian of a generative function g (or decoder d), we thus consider the Jacobian SVD in detail.

Local bases: For $J_z = USV^{\top}$, singular vectors (columns of V, U) respectively define (local) orthonormal bases: the V-basis, $\{v^i\}$ for Z at z; and the U-basis, $\{u^i\}$ for the tangent space to \mathcal{M}_g at x = g(z). Letting $v \doteq V^{\top}z$ and $u \doteq U^{\top}x$ denote a point z and its image x = g(z) in those bases, the chain rule gives an interpretation of the Jacobian's SVD, $J_z = USV^{\top} = \frac{\partial x}{\partial u} \frac{\partial u}{\partial v} \frac{\partial v}{\partial z}$: U and V^{\top} are simply local co-ordinate systems in each domain, and $S = \frac{du}{dv}$ is the Jacobian of a map $v \mapsto u$ expressed in those co-ordinates, under which *only respective dimensions interact* $(\frac{\partial u_i}{\partial v_i} = 0, i \neq j)$.

³Noting that these properties hold in expectation, we proceed as if they hold pointwise to simplify presentation.



217

219

220 221

222

224

225

226

227 228

229

230

231

232

233 234

235 236

237 238

239

240

241 242

243

244 245

251 252

253

254

255

256

257

258

259

260 261 262

263 264

265

266

267

268

269

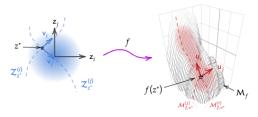


Figure 3: (*left*) Empirical support for C1: Rolinek et al. (2019) show that VAEs with diagonal Σ_x have increased orthogonality in the decoder Jacobian. (*right*) Seam factorisation: For $g: \mathbb{Z} \to \mathcal{X}$ in Lemma 5.1, Jacobian J_{z^*} and manifold $\mathcal{M}_g \subseteq \mathcal{X}$, singular-vector paths $\mathcal{V}_{z^*}^k \subseteq \mathcal{Z}$ (D3, dashed blue) following right singular vectors v^k of J_{z^*} (solid blue), map to seams $\mathcal{M}_{q,z^*}^k \subseteq \mathcal{M}_g$ (D4, dashed red) following left singular vectors u^k at $q(z^*) \in \mathcal{X}$ (solid red). $z_k \in \mathcal{Z}$ are standard basis vectors.

Singular Vector Paths and Seams: The directional derivative $J_z v^i = USV^T v^i = s^i u^i$ shows that a small perturbation by right singular vector v^i at $z \in \mathcal{Z}$ translates under g to a small perturbation in direction u^i at $g(z) \in \mathcal{M}_q$. By extension, if a path in \mathcal{Z} follows v^i at each point (as a vector field) its image on \mathcal{M}_q is expected to be a path following u^i . Note that wherever the Jacobian is continuous, each of its SVD components is continuous. Since columns of an SVD can be validly permuted, their order must be fixed for paths over "i-th" singular vectors to be well defined:

Definition D2 (Regular set and continuous SVD). For c.i.d.a.e. $g: \mathcal{Z} \to \mathcal{X}$, define the regular set $\mathcal{Z}_{\mathrm{reg}} \doteq \{z \in \mathcal{Z} \mid J_z \text{ exists, has full column rank, and } s^1(z) > \cdots > s^d(z) > 0 \}.$

Vector fields $z \mapsto v^i(z)$, $z \mapsto u^i(z)$ and singular values $z \mapsto s^i(z)$ can be made continuous on each connected component of \mathcal{Z}_{reg} by fixing the SVD $J_z = U_z S_z V_z^{\top}$.⁴⁵

With this, we define paths following i-th singular vectors: singular vector paths, or s.v. paths, \mathcal{V}_z^i follow v^i in \mathcal{Z} (Def. 3, blue dashed lines in Fig. 3); and seams, $\mathcal{M}_{a,z}^i$ follow u^i over the manifold (Def. 4, red dashed lines in Fig. 3) By construction, g maps s.v. paths to seams (proved in Lemma C.1).

Singular vector paths and seams naturally extend the observation that right singular vector perturbations map to distinct left singular vector perturbations (since S is diagonal). Key to disentanglement is how 1-D densities over s.v. paths push-forward under g to 1-D densities on the manifold \mathcal{M}_g :

Lemma 5.1 (Factorisation over seams). Let $g: \mathcal{Z} \to \mathcal{X}$ be c.i.d.a.e. and the prior factorise as $p(z) = \prod_{i=1}^{d} p_i(z_i) \text{ (e.g. standard Gaussian). Then, the manifold density } p_{\mu} \text{ on } \mathcal{M}_g \text{ factorizes as}$ $p_{\mu}(g(z)) = \prod_{i=1}^{d} \frac{p_i(z_i)}{s^i(z)}, \qquad \text{for every } z \in \mathcal{Z}_{\text{reg}}. \tag{7}$ $\text{Moreover, each factor } \frac{p_i(z_i)}{s^i(z)} \text{ is the 1-D density over the } i\text{-th seam } \mathcal{M}_{g,z}^i \text{ at } x = g(z), \text{ obtained by}$

$$p_{\mu}\!\!\left(g(z)\right) = \prod_{i=1}^{n} \frac{p_{i}(z_{i})}{s^{i}(z)}, \qquad \text{for every } z \in \mathcal{Z}_{\text{reg}}.$$
 (7)

pushing forward the 1-D marginal $p_i(z_i)$ over \mathcal{V}_z^i , the i-th s.v. path though z.

Proof. See Appendix C. Follows straightforwardly from standard change–of–variables formula.

Summary. Singular paths in \mathbb{Z} are the latent curves along which q changes in the corresponding seam on \mathcal{M}_g (Lemma C.1). Lemma 5.1 states that p_μ decomposes as a product of 1-D densities along seams, each the push-forward of the marginal over the i-th latent s.v. path – precisely the factorisation condition required for disentanglement (Eq. 5). Disentanglement also requires that s.v. paths are axis-aligned (in general they may curve Fig. 3; left) and that factors are independent, i.e. only factor i changes over seam i. We now show that these extra properties follow from C1 and C2.

Theorem 5.2 (Disentanglement \Leftrightarrow C1-C2). Let $g: \mathbb{Z} \to \mathcal{X}$ be c.i.d.a.e. and the prior factorise as $p(z) = \prod_{i=1}^d p_i(z_i)$ (e.g. standard Gaussian). The push-forward density p_μ on the manifold $\mathcal{M}_g = \{g(z)\}$ is disentangled (D1) if and only if g satisfies C1 and C2 almost everywhere.

Proof. See Appendix D.1

Thm. 5.2 means that C1-C2 are precisely the constraints needed to go from the factorisation in Eq. 7 to disentanglement: C1 causes s.v. paths to axis-align and $p_i(z_i)$ to be independent; and C2 rules out s^i (of factor i) varying in latent co-ordinate z_i ($j \neq i$), ensuring that seam factors are independent.

⁴e.g. start from an arbitrary point, order singular values strictly decreasing and choose signs continuously.

⁵Restricting to \mathcal{Z}_{reg} only excludes points where J_z is undefined or has repeated singular values; these edge cases can be avoided without affecting the results. All statements are made on a fixed connected component of \mathcal{Z}_{reg} .

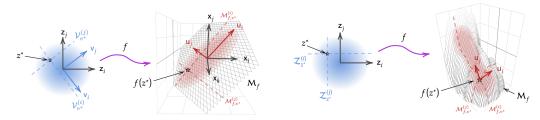


Figure 4: **Pushing forward** p(z) **from singular vector paths to seams**: s.v. path $\mathcal{V}_{z*}^i \subseteq \mathcal{Z}$ through z^* (dash blue), following right singular vectors v^i of J_z (solid blue), maps to seam $\mathcal{M}_{f,z^*}^i \subseteq \mathcal{M}_f$ through $f(z^*)$ (dash red), following left singular vectors u^i (solid red). By Lemma 5.1, 1-D marginals $p_i(z_i)$ over $\mathcal{V}_{z^*}^i$ factorising p(z) map to 1-D seam densities over \mathcal{M}_{f,z^*}^i factorising p_μ . (1) For linear f without C1, e.g. full-covariance LVAE, $\mathcal{V}_{z^*}^i$ and \mathcal{M}_{f,z^*}^i are straight lines; $p_i(z_i)$ and seam densities are independent. (r) For general c.i.d.a.e. f with C1-C2, e.g. Gaussian VAE with P1, seam densities are independent (by C2). (Essential properties of disentanglement (D1) are underlined.)

Thus, to the extent a Gaussian VAE with diagonal posteriors induces property P1 by Eq. 6, it expressly disentangles the decoder's push-forward density; and to the extent disentanglement is observed, diagonalisation constraints C1-C2 must hold. This provides a firm justification for how disentanglement emerges in VAEs, while the relationship between Eq. 6 and constraints C1-C2 also suggests a plausible rationale for why disentanglement arises inconsistently (Locatello et al., 2019).

IDENTIFIABILITY

278

279

280

281

282

283

284 285

286

287

288

289

290 291

292 293

295

296 297

298

299

300 301

302

303 304

305

306

307

308 309

310

311

312

313

314 315

321 322

323

We now investigate if a model, capable of fitting the data (i.e. data is generated under the model class), learns the *true* generative factors up to some symmetry, or could settle on a spurious factorisation.

Corollary 6.1 (LVAE Identifiability). Let data be generated under the linear Gaussian LVM Eq. 2 with ground-truth g(z) = Wz, $W = U_W S_W V_W^\top \in \mathbb{R}^{m \times d}$ of full column rank and distinct singular values. Let an LVAE with diagonal posteriors be trained on n samples, and as $n \to \infty$ its learned parameters yield $p_{\mu}^{(d)} \equiv p_{\mu}^{(g)}$ on the mean manifold. Then the LVAE achieves disentanglement (D1) and identifies ground-truth independent components on \mathcal{M}_q up to permutation and sign (**P&S**).

Proof. See Appendix D.2 (Follows from the uniqueness of the SVD).
$$\Box$$

Thus, if an LVAE learns to model the data, it learns the ground truth independent factors.

Remark 6.2 (V_W immaterial). Ground-truth right singular vectors V_W are not recoverable from p(x)under the PPCA/LVAE model; this is not a lack of identification. With a standard Gaussian prior, any orthonormal change of basis of z preserves independence and leaves p(x) unchanged. The only data-relevant object is $U_{\mathbf{w}}S_{\mathbf{w}}$; the arbitrary basis in which \mathbf{W} was written has no bearing on p(x).

The linear case hints at why independent factors may be identifiable more generally, since it depends on the Jacobian SVD, fundamental to the non-linear case. Taking this hint, we show that if a manifold density admits a seam factorisation (as in Lemma 5.1), that seam factorisation is unique (P&S) and intrinsic to p_{μ} , agnostic to any generative process or parameterisation. It follows that if the push-forward of a decoder fits p_{μ} , then its seams must align with the intrinsic seams of p_{μ} (P&S); and subsequently that a Gaussian VAE fitting p_{μ} identifies ground truth factors (P&S). (Proofs in D.3)

Lemma 6.3 (Intrinsic seams). Let $\mathcal{M} \subseteq \mathcal{X}$ carry a manifold density p_{μ} . Assume that on a regular set \mathcal{M}_{reg} there exist scalar functions $\{u_i(x)\}_{i=1}^d$ (each varying only along a 1-D curve through x, i.e. a seam) and 1-D densities $\{f_i\}$ such that $p_{\mu}(x) = \prod_{i=1}^d f_i(u_i(x)), \qquad x \in \mathcal{M}_{reg}, \tag{8}$

$$p_{\mu}(x) = \prod_{i=1}^{d} f_i(u_i(x)), \qquad x \in \mathcal{M}_{reg}, \tag{8}$$

and that the on-manifold Hessian $\mathbf{H}_x \doteq \nabla_x^2 \log p_\mu(x)$ has pairwise distinct eigenvalues a.e. on \mathcal{M}_{reg} . Then for each $x \in \mathcal{M}_{reg}$, the d seam directions (along which exactly one u_i varies) are determined intrinsically by p_{μ} , as eigenvectors of \mathbf{H}_{x} , unique up to permutation and sign (P&S).

Lemma 6.4 (A matching decoder finds seams). Under assumptions of Lemma 6.3, let $d: \mathbb{Z} \to \mathcal{X}$ be c.i.d.a.e. with factorised prior $p(z) = \prod_i p_i(z_i)$ and push-forward density $p_{\mu}^{(d)} \equiv p_{\mu}$ matching on $\mathcal{M}_d \doteq \{d(z)\} = \mathcal{M}$. If d satisfies C1–C2 a.e., then for any z and x = d(z):

- left singular vectors U_z of J_z coincide with the seam directions in Lemma 6.3 (up to P&S);
- the images under d of singular-vector paths (D3) are exactly the seams through x;
- along the *i*-th seam, the factor f_i is the 1-D push-forward of $p_i(z_i)$ (as in Lemma 5.1).

Theorem 6.5 (Gaussian VAE Identifiability). Let data be generated by c.i.d.a.e. $g: \mathcal{Z} \to \mathcal{X}$ with factorised prior $p(z) = \prod_{i=1}^d p_i(z_i)$. Let a Gaussian VAE with diagonal posteriors learn a decoder $d: \mathcal{Z} \to \mathcal{X}$. Suppose both g and d satisfy C1–C2 and manifold densities match: $p_{\mu}^{(d)} \equiv p_{\mu}^{(g)}$ on $\mathcal{M} = \{g(z)\} = \{d(z)\}$. If, eigenvalues of the tangent Hessian (see proof) are pairwise distinct a.e., then d identifies ground-truth independent components on \mathcal{M}_g , up to permutation and sign (P&S).

Thus, if a Gaussian VAE fits the push-forward of a Gaussian distribution under the conditions of 5.2, then the VAE identifies and disentangles the ground truth generative factors (up to permutation/sign).

Remark 6.6. P&S symmetry is optimal since seams follow u_i , with no inherent order or orientation.

We can also consider fitting a Gaussian VAE to data sampled from the push-forward of other priors.

Corollary 6.7 (BSS). In Theorem 6.5, if priors $p^{(g)}(z)$ and $p^{(d)}(z)$ factorise and $p^{(g)}_{\mu} \equiv p^{(d)}_{\mu}$ with C1-C2 holding a.e., then the seam decomposition p_{μ} on \mathcal{M} is unique up to permutation and sign, and g and $p^{(g)}(z)$ are recoverable up to an axis-aligned diffeomorphism $\phi \doteq g^{-1} \circ d : \mathcal{Z} \to \mathcal{Z}$.

Proof. Immediate from proof of Theorem 6.5, which does not depend on the form of p(z). The diffeomorphism follows since $\mathcal{M}_q = \mathcal{M}_d$ and by injectivity of d, g.

Remark 6.8 (Gaussian p(z) unidentifiability). Classical non-linear ICA aims to identify ground truth factors of the model in Theorem 6.5 with deterministic $p_{\theta}(x|z) = \delta_{x-d(z)}$, which is impossible if p(z) is Gaussian (Khemakhem et al., 2020; Locatello et al., 2019). Corollary 6.1 and Theorem 6.5 show that under a probabilistic formulation together with constraints C1–C2 induced by diagonal posteriors, we obtain identifiability without requiring extra side information, e.g. auxiliary variables.⁶

In summary, we have defined disentanglement, shown that it holds if and only if C1-C2 hold, and that seam factors are unique and therefore identifiable, up to expected symmetry.

7 EMPIRICAL SUPPORT

We include empirical results to illustrate disentanglement and support our claims. From our analysis, we expect (i) diagonal posteriors to promote diagonalised derivative terms; and (ii) diagonalised derivatives to correlate with disentanglement.

Both (i) and (ii) are well illustrated in the linear case where ground truth factors are known analytically. Fig. 2 shows results for diagonal and full covariance LVAEs learning Gaussian parameters (see caption for details). All models learn optimal parameters as expected (left); but, as (i) predicts, only diagonal covariances cause right singular vectors of J_z to converge to standard basis vectors, $V \to I$ (C1), (centre); hence latent traversals map to independent components along left singular vectors u^i (right), yielding disentanglement (D1), predicted by (ii). This evidences diagonal covariances "breaking the rotational symmetry" of a Gaussian prior. While the linear case seems trivial, it is a fundamental demonstration since our analysis shows that disentanglement in general follows the same rationale. Interestingly Fig. 3 shows that learning parameters a disentangled model is notably slower (left), due to the rate at which $V \to I$ (centre). A diagonal covariance model must find one of a finite set of solutions among the infinite solutions of a full covariance model.

Various studies show further empirical support. Supporting (i) and (ii), Rolinek et al. (2019) show that columns of a decoder's Jacobian are more orthogonal (i.e. $V \rightarrow I$, C1) in diagonal covariance VAEs than those with full covariance (Fig. 3, *left*), and that diagonality correlates with disentanglement. Supporting (ii), Kumar & Poole (2020) show that directly inducing column-orthogonality in the decoder Jacobian promotes disentanglement.

⁶Unidentifiability proofs typically make use of an arbitrary rotation applied to the Gaussian prior (cf R in Eq. 3), but we have seen that C1 removes such symmetry, even in the linear case (Fig. 2). (See also Remark 6.2.)

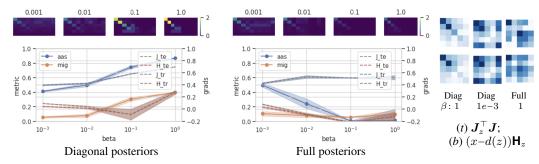


Figure 5: **Diagonal vs Full Posteriors**: (*left*) (*bottom*) disentanglement metrics and estimated diagonality of Eq. 6 terms (see Appendix E). With diagonal posteriors, disentanglement and diagonality are correlated (supporting P1), relative to full posteriors; (*top*) heatmaps of mutual information between model latents and ground truth factors. (*right*) derivatives in Eq. 6. terms are less diagonal for lower beta or full posteriors. (All results averaged over multiple runs)

Complementing prior work, we train diagonal and full covariance Gaussian VAEs (d=10) on the dSprites dataset, for which the 5 ground truth generative factors are known (all results averaged over 5 runs). How well each latent co-ordinate identifies a ground truth factor can be estimated explicitly from their *mutual information*, and a function of mutual information often captures a model's overall disentanglement, e.g. the mutual information gap (MIG, Chen et al. (2018)). Fig. 5 (main plots) reports MIG, axis alignment score (AAS) (a novel metric based on the entropy of the mutual information distribution) and derivative diagonality estimates against β (see Appendix E for details). For diagonal covariance VAEs (*left*), disentanglement metrics and diagonality broadly increase together with β . For full posteriors (*right*) no clear trend is observed. The heatmaps above (aligned with plots by β) show mutual information between each latent co-ordinate and ground truth factor (ordered greedily to put highest mutual information scores along the diagonal). We see that for diagonal covariance VAEs, disentanglement increases with β and that individual latent co-ordinates (horizontal) correlate with distinct ground truth factors (vertical), whereas that is not clearly observed for full covariances. To illustrate, Fig. 5 (right) shows heatmaps of the $d \times d$ derivative terms in Eq. 6 (e.g. $J_{\perp}^{-1}J_{z}$. Jacobian term top, Hessian term bottom), each for diagonal covariances, $\beta = 1$ (*l*); diagonal covariances, $\beta = 0.001$ (c); and full covariances, $\beta = 1$ (r) (each averaged over a batch).

Lastly, we make a prediction based on our understanding of the interplay between β and disentanglement (§4). Our analysis suggests that higher β increases expected noise, which weakens reconstructions, but promotes disentanglement via widened posteriors; whereas lower β assumes less noise, which tightens reconstructions, but disentanglement is restricted under concentrated posteriors to smaller, potentially disconnected, regions. For both clear samples *and* disentanglement, this suggests starting with high β , so that disentanglement is induced broadly; and reducing β over training so that the model manifold is pulled tighter to the data and images become sharper. While not the focus of this work, we run a preliminary experiment on the dSprites dataset as a proof of concept. Results in Appendix F include baselines for constant high and low β values (1 and 0.001); and the effect of exponentially reducing β (1 \rightarrow 0.001) over training. The baselines illustrate the expected results described above and, promisingly in line with our prediction, reducing β shows both sharp reconstructions and good disentanglement. We note a resemblance to the "de-noising" process in denoising autoencoders and diffusion models and see this as an interesting future direction.

8 Related Work

Many works study aspects of VAEs, or disentanglement in other modelling paradigms. Here, we focus on those at the nexus, investigating disentanglement in VAEs.

Higgins et al. (2017) showed that disentanglement in VAEs is enhanced by setting $\beta>1$ in the (β -)ELBO (Eq. 1). Burgess et al. (2018) conjectured that diagonal posterior covariances may cause disentanglement. Rolinek et al. (2019) showed supporting empirical evidence (3) and derived an approximate relationship between diagonal posteriors and Jacobian orthogonality, conjectured to cause disentanglement. Kumar & Poole (2020) simplified and generalised their argument, arriving at an (independent) approximation of Eq. 6. By comparison, our work makes precise the link between posterior covariances and decoder derivatives, formally defines disentanglement and proves that it arises if and only if derivative constraints hold.

Lucas et al. (2019); Bao et al. (2020) and Koehler et al. (2022) properties of LVAEs, notably Lucas et al. (2019) show the equivalence of a Gaussian β -VAE and a Gaussian VAE with a different variance assumption, which we generalise in Appendix B. Zietlow et al. (2021) suggest disentanglement is sensitive to perturbations to the data distribution. Reizinger et al. (2022) seek to relate the ELBO to *independent mechanism analysis* (Gresele et al., 2021), which encourages column-orthogonality in the mixing function of ICA. We note that orthogonality of the Jacobian (C1) is not sufficient to guarantee disentanglement/identification of independent components, as that also requires (C2).

Ramesh et al. (2018) trace independent factors by following leading left singular vectors (our *seams*) of the Jacobian of a GAN generator, whereas Chadebec & Allassonnière (2022) and Arvanitidis et al. (2018) consider paths in latent space defined by the inverse image of paths over the data manifold (our *s.v. paths*). Pan et al. (2023) claim that the data manifold is identifiable from a geometric perspective assuming Jacobian-orthogonality, differing to our probabilistic factorisation approach. Bhowal et al. (2024) consider the encoder/decoder dissected into linear and non-linear components, loosely resembling our view of the Jacobian from its SVD. However, the decoder function is quite different to its Jacobian, and dissecting a function into linear/non-linear components is not well defined, whereas the SVD is unique (up to permutation and sign).

Buchholz et al. (2022) analyse function classes that are identifiable by ICA, proving that *conformal maps* are identifiable, but *orthogonal coordinate transformations* (OCTs), defined to satisfy C1, are not. This is relatively close to our work in spirit, differing in that our identifiability proof (Theorem 6.5) relies on a stochastic model and an additional constraint, C2.

9 Conclusion

 Unsupervised disentanglement of generative factors of the data is of fundamental interest in machine learning. Thus, irrespective of the popularity of VAEs as a generative model class, understanding how they disentangle the data *for free* may have useful implications for other paradigms, and we take significant strides in this respect, in particular giving disentanglement a formal definition (D1): as factorising the manifold density into independent components, each the image of an axis-aligned traversal in latent space. We also give a simple interpretation of β in a β -VAE, as adjusting the assumed variance or entropy of the likelihood, which justifies why $\beta > 1$ promotes disentanglement while degrading generative quality; and also why $\beta < 1$ mitigates *posterior collapse*.

Much of our work is not restricted to VAEs, but is general to a class of push-forward distributions, i.e. a factorised prior combined with conditinuous, injective, differentiable almost everywhere function, as arise also in GANs and FLOWs. We show that the combination of two relatively simple concepts, the first derivative or Jacobian and the SVD, combine in a surprisingly elegant way by which, under suitable conditions, a push-forward density factorises. Indeed the entire factorisation in the latent space essentially projects onto the manifold and independent densities over *singular vector paths* in latent space push-forward to independent densities over *seams* on the manifold. We show that the constraints needed to achieve disentanglement are precisely those imposed, albeit in aggregate and in expectation, by a VAE with diagonal covariances, justifying both why disentanglement arises and perhaps also why it can be elusive. Significantly, given the proven *impossibility* of ICA under a Gaussian prior (precisely the setting of a Gaussian VAE), we prove that, under suitable assumptions, a Gaussian VAE does identify ground truth independent factors, up to their inherent symmetry group.

Neural networks models are often considered too complex to explain, despite their increasingly widespread deployment in everyday applications. An improved theoretical understanding seems essential to safely take advantage of machine learning progress in potentially critical systems. Our work aims to be a useful step in that direction, providing novel insight into how a density decomposes over generative factors. Interestingly, our work shows that, regardless of the non-linear complexity of a VAE, how it pushes forward the prior density can be considered relatively simply.

VAEs (and variants e.g. AEs, SAEs) form part of a pipeline in many state-of-the-art models, e.g. latent diffusion (e.g. Rombach et al., 2022; Pandey et al., 2022; Yang et al., 2023; Zhang et al., 2022) and LLMs; other recent works show that supervised learning (Dhuliawala et al., 2024) and self-supervised learning (Bizeul et al., 2024) can be viewed as latent models trained under ELBO variants. In future work we will look to see if our results transfer to such other learning paradigms.

⁷We report apparent discrepancies in Reizinger et al. (2022) in Appendix G.

REFERENCES

- Georgios Arvanitidis, Lars Kai Hansen, and Søren Hauberg. Latent space oddity: On the curvature of deep generative models. In *ICLR*, 2018.
- Xuchan Bao, James Lucas, Sushant Sachdeva, and Roger B Grosse. Regularized linear autoencoders recover the principal components, eventually. In *NeurIPS*, 2020.
 - Yoshua Bengio, Aaron Courville, and Pascal Vincent. Representation learning: A review and new perspectives. In *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2013.
 - Pratik Bhowal, Achint Soni, and Sirisha Rambhatla. Why do variational autoencoders really promote disentanglement? In *ICML*, 2024.
 - Alice Bizeul, Bernhard Schölkopf, and Carl Allen. A Probabilistic Model to explain Self-Supervised Representation Learning. In *TMLR*, 2024.
 - Samuel R Bowman, Luke Vilnis, Oriol Vinyals, Andrew M Dai, Rafal Jozefowicz, and Samy Bengio. Generating sentences from a continuous space. In *Conference on Computational Natural Language Learning*, 2015.
 - Simon Buchholz, Michel Besserve, and Bernhard Schölkopf. Function classes for identifiable nonlinear independent component analysis. In *NeurIPS*, 2022.
 - Christopher P Burgess, Irina Higgins, Arka Pal, Loic Matthey, Nick Watters, Guillaume Desjardins, and Alexander Lerchner. Understanding disentangling in β -vae. *arXiv preprint arXiv:1804.03599*, 2018.
 - Clément Chadebec and Stéphanie Allassonnière. A geometric perspective on variational autoencoders. In *NeurIPS*, 2022.
 - Ricky TQ Chen, Xuechen Li, Roger B Grosse, and David K Duvenaud. Isolating sources of disentanglement in variational autoencoders. *Advances in neural information processing systems*, 31, 2018.
 - Shehzaad Dhuliawala, Mrinmaya Sachan, and Carl Allen. Variational Classification. TMLR, 2024.
 - Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. Generative adversarial nets. *NeurIPS*, 2014.
 - Luigi Gresele, Julius Von Kügelgen, Vincent Stimper, Bernhard Schölkopf, and Michel Besserve. Independent mechanism analysis, a new concept? In *NeurIPS*, 2021.
 - Irina Higgins, Loic Matthey, Arka Pal, Christopher Burgess, Xavier Glorot, Matthew Botvinick, Shakir Mohamed, and Alexander Lerchner. β -VAE: Learning Basic Visual Concepts with a Constrained Variational Framework. In *ICLR*, 2017.
 - Ilyes Khemakhem, Diederik Kingma, Ricardo Monti, and Aapo Hyvarinen. Variational autoencoders and nonlinear ica: A unifying framework. In *AIStats*, 2020.
 - Hyunjik Kim and Andriy Mnih. Disentangling by factorising. In *ICML*, 2018.
- Diederik P Kingma and Max Welling. Auto-encoding variational bayes. In *ICLR*, 2014.
 - Frederic Koehler, Viraj Mehta, Chenghui Zhou, and Andrej Risteski. Variational autoencoders in the presence of low-dimensional data: landscape and implicit bias. In *ICLR*, 2022.
- Abhishek Kumar and Ben Poole. On Implicit Regularization in β -VAEs. In *ICML*, 2020.
- Francesco Locatello, Stefan Bauer, Mario Lucic, Gunnar Raetsch, Sylvain Gelly, Bernhard Schölkopf, and Olivier Bachem. Challenging common assumptions in the unsupervised learning of disentangled representations. In *ICML*, 2019.
 - James Lucas, George Tucker, Roger B Grosse, and Mohammad Norouzi. Don't Blame the ELBO! a Linear VAE Perspective on Posterior Collapse. In *NeurIPS*, 2019.

Manfred Opper and Cédric Archambeau. The variational gaussian approximation revisited. Neural computation, 21(3):786–792, 2009. Ziqi Pan, Li Niu, and Liqing Zhang. Geometric inductive biases for identifiable unsupervised learning of disentangled representations. In AAAI, 2023. Kushagra Pandey, Avideep Mukherjee, Piyush Rai, and Abhishek Kumar. Diffusevae: Efficient, controllable and high-fidelity generation from low-dimensional latents. In TMLR, 2022. Aditya Ramesh, Youngduck Choi, and Yann LeCun. A spectral regularizer for unsupervised disentanglement. arXiv preprint arXiv:1812.01161, 2018. Patrik Reizinger, Luigi Gresele, Jack Brady, Julius Von Kügelgen, Dominik Zietlow, Bernhard Schölkopf, Georg Martius, Wieland Brendel, and Michel Besserve. Embrace the gap: Vaes perform independent mechanism analysis. In NeurIPS, 2022. Danilo Jimenez Rezende and Fabio Viola. Taming vaes. arXiv preprint arXiv:1810.00597, 2018. Danilo Jimenez Rezende, Shakir Mohamed, and Daan Wierstra. Stochastic backpropagation and approximate inference in deep generative models. In *ICML*, 2014. Michal Rolinek, Dominik Zietlow, and Georg Martius. Variational Autoencoders Pursue PCA Directions (by Accident). In CVPR, 2019. Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models. In CVPR, 2022. Rui Shu, Yining Chen, Abhishek Kumar, Stefano Ermon, and Ben Poole. Weakly supervised disentanglement with guarantees. In *ICLR*, 2019. Michael E Tipping and Christopher M Bishop. Probabilistic principal component analysis. *Journal* of the Royal Statistical Society Series B: Statistical Methodology, 61(3):611–622, 1999. Tao Yang, Yuwang Wang, Yan Lu, and Nanning Zheng. Disdiff: unsupervised disentanglement of diffusion probabilistic models. In NeurIPS, 2023. Zijian Zhang, Zhou Zhao, and Zhijie Lin. Unsupervised representation learning from pre-trained diffusion probabilistic models. In *NeurIPS*, 2022. Dominik Zietlow, Michal Rolinek, and Georg Martius. Demystifying inductive biases for (beta-) vae based architectures. In ICML, 2021.

A PROOF OF DECODER DERIVATIVE CONSTRAINTS

Property P1. The expectations of $J_z^{\top}J_z$ and $(x-d(z))^{\top}H_z$ in Eq. 6 are each diagonal.⁸

Lemma 4.1 (**Disentanglement constraints**). *Under P1*, in expectation for concentrated posteriors:

- C1 Right singular vectors V_z of the decoder Jacobian J_z are standard basis vectors for all $z \in \mathcal{Z}$, i.e. after relabeling/sign flips of the latent axes, we have $V_z = I$;
- **C2** The matrix of partial derivatives of singular values $(\frac{\partial s_i}{\partial z_j})_{i,j}$ is diagonal, i.e. $\frac{\partial s_i}{\partial z_j} = 0$ for all $i \neq j$.

Proof. (Preliminaries): Recall $p(z) = \mathcal{N}(0, I)$ and $p(x \mid z) = \mathcal{N}(x; d(z), \sigma^2 I)$. Let $q(z \mid x) = \mathcal{N}(z; e(x), \Sigma_x)$ be the trained posterior with Σ_x diagonal (by assumption). Denote the SVD of the decoder Jacobian by

$$J_z = U_z S_z V_z^{\top}, \qquad U_z \in \mathbb{R}^{n \times d}, \quad S_z = \text{Diag}(s_1(z), \dots, s_k(z)), \quad s_i(z) > 0.$$

Assume full column rank on the manifold $(s_i(z)>0)$. For Gaussian likelihood with variance σ^2 , the Hessian of the log-likelihood w.r.t. z can be written $\nabla_z^2 \log p(x\,|\,z) = -\frac{1}{\sigma^2} \big(\boldsymbol{J}_z^\top \boldsymbol{J}_z - \sum_{\ell=1}^n r_\ell \, \boldsymbol{H}_\ell(z) \big)$, where r=x-d(z) and $\boldsymbol{H}_\ell(z)\in\mathbb{R}^{k\times k}$ is the Hessian of the ℓ -th decoder coordinate, $[\boldsymbol{H}_\ell]_{pq}=\partial^2 d_\ell/\partial z_p\,\partial z_q$. Combined with the Opper–Archambeau fixed-point yields

$$\Sigma_x^{-1} = \boldsymbol{I} - \mathbb{E}_q \left[\nabla_z^2 \log p(x|z) \right] = \boldsymbol{I} + \frac{1}{\sigma^2} \mathbb{E}_q \left[\boldsymbol{J}_z^\top \boldsymbol{J}_z \right] - \frac{1}{\sigma^2} \mathbb{E}_q \left[\sum_{\ell=1}^n r_\ell \, \boldsymbol{H}_\ell(z) \right]. \tag{9}$$

C1: For diagonal Σ_x in Eq. 9 under the "no cancellation" P1, we must have

$$\mathbb{E}_q[\boldsymbol{J}_z^{\top}\boldsymbol{J}_z]$$
 is diagonal, $\mathbb{E}_q[\sum_{\ell=1}^n r_{\ell} \boldsymbol{H}_{\ell}(z)]$ is diagonal. (10)

If $q(z \mid x)$ is concentrated at e(x), the first expectation is well approximated by its value at z = e(x) up to o(1) errors, hence $\boldsymbol{J}_{e(x)}^{\top}\boldsymbol{J}_{e(x)} = \mathrm{Diag}(s_1^2,\ldots,s_k^2)$; and right singular vectors at e(x) are the standard basis up to signed permutations. By relabelling latent axes and absorbing signs, $\boldsymbol{V}_{e(x)} = \boldsymbol{I}$. Continuity and data coverage then give $\boldsymbol{V}_z = \boldsymbol{I}$ for all z visited by the encoder, establishing $\boldsymbol{C1}$.

(Directed Hessian is Tangent to Manifold): Assume the model is well trained so that $d(e(x)) \approx x$. For $z = e(x) + \delta$ with δ small, a first-order Taylor expansion gives

$$d(z) = d(e(x)) + \mathbf{J}_{e(x)} \delta + O(\|\delta\|^2), \Rightarrow r = x - d(z) = -\mathbf{J}_{e(x)} \delta + O(\|\delta\|^2).$$

Thus, to first order, r lies in the column space of $J_{e(x)}$, i.e. in the span of the left singular vectors $\{u_i(e(x))\}_{i=1}^k$ (columns of U_z):

$$r = \mathbf{U}_z a, \qquad a \in \mathbb{R}^k$$
 (11)

To simplify, we drop the expectation and consider $\sum_{\ell=1}^n r_\ell \, \boldsymbol{H}_\ell(z)$ to be diagonal for all $r \in \operatorname{span}(\boldsymbol{U}_z)$. Thus $\sum_{\ell=1}^n u_{i\ell}^\top \boldsymbol{H}_\ell(z)$ is diagonal for all rows u_i of \boldsymbol{U}_z and, by definition of slices \boldsymbol{H}_ℓ and \boldsymbol{J}_z ,

$$\left[\sum_{\ell=1}^{n} u_{i\ell}^{\top} \boldsymbol{H}_{\ell}(z)\right]_{pq} = \sum_{\ell=1}^{n} u_{i\ell}^{\top} \frac{\partial^{2} d_{\ell}}{\partial z_{p} z_{q}} = (\boldsymbol{U}_{z}^{\top} \frac{\partial \boldsymbol{J}_{z}}{\partial z_{p}})_{iq}$$
(12)

(Diagonality of $(\frac{\partial s_i}{\partial z_j})_{i,j}$):

Differentiating $U_z^{\top}U_z = I_k$ to get $\frac{\partial U_z^{\top}}{\partial z_j}U_z + U_z^{\top}\frac{\partial U_z}{\partial z_j} = \mathbf{0}$ shows $\Omega_j(z) := U_z^{\top}\frac{\partial U_z}{\partial z_j} \in \mathbb{R}^{d\times d}$ is skew-symmetric. Differentiating $J_z = U_z S_z$ w.r.t. z_j and premultiplying by U_z^{\top} then gives

$$U_z^{\top} \frac{\partial J_z}{\partial z_i} = \Omega_j(z) S_z + \frac{\partial S_z}{\partial z_i}.$$
 (13)

⁸Noting that these properties hold in expectation, we proceed as if they hold pointwise to simplify presentation. ⁹we use: $\Sigma_x^{-1} = -\mathbb{E}_q[\nabla_z^2 \log p(x,z)] = -\mathbb{E}_q[\nabla_z^2 \log p(z) + \nabla_z^2 \log p(x|z)]; \quad \nabla_z^2 \log p(z) = -\mathbf{I}.$

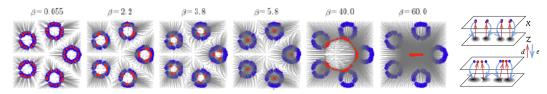


Figure 6: Illustrating $\beta \propto \mathrm{Var}[x|z]$ (blue = data, red = reconstruction): (1) For low β (β = 0.55), $\mathrm{Var}[x|z]$ is low (by Eq. 6), and data must be well reconstructed (right, top). As β increases, $\mathrm{Var}[x|z]$ and so $\mathrm{Var}[z|x]$ increase, and posteriors of nearby samples $\{x_i\}_i$ increasingly overlap (right, bottom). For z in overlapping $\{q(z|x_i)\}_i$, the decoder $\mathbb{E}[x|z]$ maps to a weighted average of $\{x_i\}_i$. Initially, close neighbours reconstruct to their mean (β = 2.2, 3.8), then small circles "become neighbours" and map to their centres. Finally (β = 60), all samples reconstruct to the global centroid. (reproduced with permission from Rezende & Viola, 2018) (r) illustrating posterior overlap, (t) low β , (t) higher β .

Since $\Omega_j(z)$ is skew-symmetric, all diagonal entries are zero; since $\frac{\partial S_z}{\partial z_j}$ is diagonal, all non-diagonal entries are zero. Thus, of respective entries, only one is non-zero and can be considered separately.

From Eq. 12, only $(\Omega_j(z) S_z)_{:j}$ elements can be non-zero, thus all elements of $\Omega_j(z)$ must be zero (by skew-sym.), ruling out rotation in the tangent plane¹⁰

For $\frac{\partial S_z}{\partial z_j}$, diagonality (only $(\frac{\partial S_z}{\partial z_j})_{kk}$ elements non-zero) and Eq. 12 (only $(\frac{\partial S_z}{\partial z_j})_{:j}$ elements non-zero) imply that only elements $(\frac{\partial S_z}{\partial z_j})_{jj} = \frac{\partial s_j}{\partial z_j}$ can be non-zero, eliminating mixed partials

$$\frac{\partial s_i}{\partial z_i}(z) = 0$$
 for all $i \neq j$,

i.e. the Jacobian of the singular-value map $s(z) = (s_1(z), \dots, s_k(z))$ is diagonal, proving C2. \square

B β CONTROLS NOISE VARIANCE

Choosing $\beta > 1$ in Eq. 1 can enhance disentanglement (Higgins et al., 2017; Burgess et al., 2018) and has been viewed as re-weighting ELBO components or as a Lagrange multiplier. We show that β implicitly controls the likelihood's variance and that the " β -ELBO" remains a valid objective.

Dividing the ELBO by a constant and suitably adjusting the learning rate leaves the VAE training algorithm unchanged, hence consider Eq. 1 divided through by β with the log likelihood scaled by β^{-1} . For a Gaussian VAE with $\mathrm{Var}[x|z] = \sigma^2$, this exactly equates to a standard VAE with variance $\beta\sigma^2$ (Lucas et al., 2019). More generally, scaling the log likelihood by β^{-1} is equivalent to an *implicit likelihood* $p_{\theta}(x|z)^{1/\beta}$, where β acts as a *temperature* parameter: $\beta \to \infty$ increases the effective entropy towards uniform (the model assumes more noise in the data, fitting more loosely), and $\beta \to 0$ reduces it to a delta (reconstructions should be tight). Optimal posteriors fit to the implicit likelihood, $q_{\phi}(z|x) \propto p_{\theta}(x|z)^{1/\beta}p(z)$, which thus dilate $(\beta > 1)$ or concentrate $(\beta < 1)$. This generalises the Gaussian result (Lucas et al., 2019), showing that the β -ELBO is simply the ELBO for a different likelihood model. β

Empirical support: Our claim, in effect that $\mathrm{Var}[x|z] \propto \beta$, is well illustrated on synthetic data in Fig. 6 (Rezende & Viola, 2018, see caption for details). It also immediately explains blur in β -VAEs since $\beta > 1$ simply assumes more noise. It also explains why $\beta < 1$ helps mitigate posterior collapse (Bowman et al., 2015), i.e. when a VAE's likelihood is sufficiently expressive that it can directly model the data distribution, p(x|z) = p(x), leaving latent variables redundant (posterior "collapses" to prior). As $\beta \to 0$, the effective variance of $p_{\theta}(x|z)$, and the distributions it can describe, reduces. Thus for some $\beta < 1$ the effective variance falls below $\mathrm{Var}[x]$, rendering posterior collapse impossible as some variance in x can only be explained by z. Thus our claim that β controls effective variance explains well-known empirical observations, which in turn provide empirical support for the claim.

¹⁰we have: $\Omega_j(z)_{kj} = -\Omega_j(z)_{jk} = 0$, if $j \neq k$; and $\Omega_j(z)_{jj} = 0$ (skew-sym.).

¹¹Technically, the β-ELBO's value is incorrect without renormalising the implicit likelihood, but that is typically irrelevant, e.g. for commonly used Gaussian likelihoods, only the quadratic "MSE" term appears in the loss.

C APPENDIX: SINGULAR VECTOR PATHS AND SEAMS

Definition D3 (*i*-th singular-vector path). Let $g: \mathbb{Z} \to \mathcal{X}$ be c.i.d.a.e.. For $z^* \in \mathcal{Z}_{reg}$, $i \in \{1, \ldots, d\}$, the *i*-th singular-vector path (s.v. path) through z^* is any C^1 curve $t \mapsto z_t^i$ with $z_0^i = z^*$ satisfying

$$\frac{d}{dt}z_t^i = \mathbf{v}^i(z_t^i)$$
 for t in its maximal interval $I_{z^*,i} \subseteq \mathbb{R}$.

We denote the path set by $\mathcal{V}_{z^*}^i \doteq \{z_t^i : t \in I_{z^*,i}\} \subseteq \mathcal{Z}_{reg}$. (See Fig. 3, left, dash blue lines).

Definition D4 (i-th seam). Let $g: \mathbb{Z} \to \mathcal{X}$ be c.i.d.a.e. with manifold $\mathcal{M}_g = \{g(z)\}$. For $z^* \in \mathcal{Z}_{reg}$, $i \in \{1, \ldots, d\}$, the i-th seam through $g(z^*)$ is any C^1 curve $t \mapsto x_t^i$ in \mathcal{M}_g with $x_0^i = g(z^*)$ satisfying

$$\frac{d}{dt}x_t^i = s^i(g^{-1}(x_t^i)) u^i(g^{-1}(x_t^i))$$
 for t in $I_{z^*,i}$.

We denote the path set $\mathcal{M}_{q,z^*}^i \doteq \{x_t^i : t \in I_{z^*,i}\} \subseteq \mathcal{M}_q$ and define seam coordinate

$$u_i(t) \doteq \int_0^t s^i(g^{-1}(x_\tau^i)) d\tau, \quad so \quad \frac{d}{dt}u_i(t) = s^i(g^{-1}(x_t^i)), \quad u_i(0) = 0.$$

 u_i measures position along the seam in units of s^i (strictly monotone as $s^i > 0$). (See Fig. 3, right).

Lemma C.1 (Paths \mapsto seams). Let $g: \mathcal{Z} \to \mathcal{X}$ be c.i.d.a.e., $z^* \in \mathcal{Z}_{reg}$, $i \in \{1, ..., d\}$, and let $\mathcal{V}^i_{z^*}$ be the i-th s.v. path through z^* . Then the image of $\mathcal{V}^i_{z^*}$ under g is the i-th seam through $g(z^*)$: $\mathcal{M}^i_{g,z^*} = \{g(z) : z \in \mathcal{V}^i_{z^*}\}$.

Proof. For $x_t^i \doteq g(z_t^i)$, by the chain rule and SVD: $\frac{dx_t^i}{dt} = J_{z_t^i} \frac{dz_t^i}{dt} = J_{z_t^i} \boldsymbol{v}^i(z_t^i) = s^i(z_t^i) \, \boldsymbol{u}^i(z_t^i)$, so x_t^i satisfies Def. 4.

Lemma 5.1 (Factorisation over seams). Let $g: \mathbb{Z} \to \mathcal{X}$ be c.i.d.a.e. and the prior factorise as $p(z) = \prod_{i=1}^d p_i(z_i)$ (e.g. standard Gaussian). Then, the manifold density p_μ on \mathcal{M}_g factorizes as

$$p_{\mu}(g(z)) = \prod_{i=1}^{a} \frac{p_{i}(z_{i})}{s^{i}(z)}, \qquad \text{for every } z \in \mathcal{Z}_{\text{reg}}. \tag{7}$$

Moreover, each factor $\frac{p_i(z_i)}{s^i(z)}$ is the 1-D density over the i-th seam $\mathcal{M}_{g,z}^i$ at x = g(z), obtained by pushing forward the 1-D marginal $p_i(z_i)$ over \mathcal{V}_z^i , the i-th s.v. path though z.

Proof. By a standard change–of–variables (on embedded manifolds) and $|J_z^\top J_z| = \prod_{i=1}^d s^i(z)^2$,

$$p_{\mu}(g(z)) = \det(\boldsymbol{J}_{z}^{\top} \boldsymbol{J}_{z})^{-1/2} p(z) = \frac{\prod_{i=1}^{d} p_{i}(z_{i})}{\prod_{i=1}^{d} s^{i}(z)},$$
 yielding Eq. 7.

For each i and x=g(z), the change–of–variables formula along $\mathcal{M}_{g,z}^i$ (i-th seam through x with coordinate u_i , Def. 4) at t=0 gives the local pushed 1-D seam-density

$$f_i^{(z)}(u_i(t)) \doteq \frac{p_i([z_i^i]_i)}{s^i(z_i^i)}, \qquad \frac{d}{dt}u_i(t) = s^i(z_t^i), \ u_i(0) = 0,$$
 (14)

where $t\mapsto z^i_t$ is the *i*-th singular-vector path through z (Def. 3) with local co-ordinate $z^i_i(t)$. Evaluating at t=0: $f^{(z)}_i(u_i(0))=\frac{p_i(z_i)}{s^i(z)}$, shows that seam-densities are the factors of Eq. 7. \square

¹²Different choices of sign for v^i reverse the time direction $(t \mapsto -t)$ but generate the same path set.

D APPENDIX: DISENTANGLEMENT AND IDENTIFIABILITY PROOFS

D.1 PROOF OF NON-LINEAR DISENTANGLEMENT

Theorem 5.2 (Disentanglement \Leftrightarrow **C1-C2).** *Let* $g: \mathcal{Z} \to \mathcal{X}$ *be* c.i.d.a.e. *and the prior factorise* as $p(z) = \prod_{i=1}^d p_i(z_i)$ (e.g. standard Gaussian). The push-forward density p_μ on the manifold $\mathcal{M}_q = \{g(z)\}$ is disentangled (D1) if and only if g satisfies C1 and C2 almost everywhere.

Proof. (C1/2 \Rightarrow D1) By Thm. 5.1, the manifold density factorises pointwise as $p_{\mu}(d(z)) = \prod_{i=1}^{d} \frac{p_{i}(z_{i})}{s^{i}(z)}$. By **C1**, $V_{z} = I$ for all z, so the i-th singular-vector path through z is exactly the axis-aligned line $\{z': [z']_{i} \text{ varies}, [z']_{\neg i} = [z]_{\neg i}\}$; by Lemma C.1 its image is the i-th seam through x = g(z) following u^{i} . By **C2**, $s^{i}(z)$ depends only on z_{i} . Define the seam coordinate u_{i} along the i-th seam as in D4; then u_{i} is a strictly monotone function of z_{i} , hence the 1-D push-forward of p_{i} along that seam is $f_{i}(u_{i}) = \left|\frac{du_{i}}{dz_{i}}\right|^{-1} p_{i}(z_{i}) = \frac{p_{i}(z_{i})}{s^{i}(z_{i})}$, Thus $p_{\mu}(g(z)) = \prod_{i} f_{i}(u_{i}(z))$ with each f_{i} evaluated on the i-th seam. Finally, since u_{i} is monotone in z_{i} and $\{z_{i}\}$ are independent, the random variables $\{u_{i}\}$ are independent; hence factors $\{f_{i}(u_{i})\}$ are statistically independent as required by D1.

(D1 \Leftarrow C1/2) Assume p_{μ} is disentangled under g. By D1, each factor f_i is obtained by pushing forward $p(z_i)$ along an axis-aligned line indirection z^i , and the i-th seam follows $J_z z^i = J_z^i$ at z (column i of J_z). By factor independence (D1), only f_i can change along the i-th seam, hence all other factors must be orthogonal to J_i , i.e. $J_z^{\top} J_z$ is diagonal or $J_z = U_z S_z$ (C1). Since f_i depends on $s^i \doteq [S_z]_{ii}$ and only s^i can change along seam i, then $\frac{\partial s_i}{\partial z_i} = 0$, $i \neq j$ (C2).

D.2 PROOF OF LINEAR VAE IDENTIFIABILITY

Corollary 6.1 (**LVAE Identifiability**). Let data be generated under the linear Gaussian LVM Eq. 2 with ground-truth $g(z) = \mathbf{W} z$, $\mathbf{W} = \mathbf{U}_{\mathbf{W}} \mathbf{S}_{\mathbf{W}} \mathbf{V}_{\mathbf{W}}^{\top} \in \mathbb{R}^{m \times d}$ of full column rank and distinct singular values. Let an LVAE with diagonal posteriors be trained on n samples, and as $n \to \infty$ its learned parameters yield $p_{\mu}^{(d)} \equiv p_{\mu}^{(g)}$ on the mean manifold. Then the LVAE achieves disentanglement (D1) and identifies ground-truth independent components on \mathcal{M}_g up to permutation and sign (**P&S**).

Proof. (Ground truth) With $\mu = \mathbf{W}z$ and $u \doteq \mathbf{U}_{\mathbf{W}}^{\top} \mu = \mathbf{S}_{\mathbf{W}}(\mathbf{V}_{\mathbf{W}}^{\top}z)$, then $z \sim \mathcal{N}(0, \mathbf{I})$ and orthonormal $\mathbf{V}_{\mathbf{W}}$ imply $\mathbf{V}_{\mathbf{W}}^{\top}z \sim \mathcal{N}(0, \mathbf{I})$, hence $\{u_i \doteq u_{\mathbf{W},i}\}$ are independent, $u_i \sim \mathcal{N}(0, s_{\mathbf{W},i}^2)$ and $p_{\mu}^{(g)} = \prod_{i=1}^d \mathcal{N}(u_i; 0, s_{\mathbf{W},i}^2)$.

(Model) Let d(z) = Dz with SVD $D = U_D S_D V_D^{\top}$. For an LVAE, the Hessian term in Eq. 6 is zero and Assumption 1 is trivially satisfied. Thus, by Lemma 4.1, right singular paths are axis—aligned (C1, note C2 is vacuous). Therefore, by the disentanglement theorem Theorem 5.2, $p_{\mu}^{(d)}$ is disentangled and factorizes into statistically independent components along the decoder's seams (columns of U_D). Since $u_{D,i} = s_{D,i}z_i$ and $z_i \sim \mathcal{N}(0,1)$, each seam factor is Gaussian with variance $s_{D,i}^2$, i.e. $p_{\mu}^{(d)} = \prod_{i=1}^d \mathcal{N}(u_{D,i};0,s_{D,i}^2)$.

(*Matching*) Equality $p_{\mu}^{(d)} \equiv p_{\mu}^{(g)}$ and *distinct* $\{s_{W,i}\}$ imply uniqueness of the Gaussian product decomposition, up to permutation. Thus the LVAE's independent components (seam factors) match ground-truth components up to permutation/sign, i.e. identifiability and disentanglement on \mathcal{M}_q . \square

D.3 Proof of Gaussian VAE Identifiability

Lemma 6.3 (Intrinsic seams). Let $\mathcal{M} \subseteq \mathcal{X}$ carry a manifold density p_{μ} . Assume that on a regular set \mathcal{M}_{reg} there exist scalar functions $\{u_i(x)\}_{i=1}^d$ (each varying only along a 1-D curve through x, i.e. a seam) and 1-D densities $\{f_i\}$ such that

$$p_{\mu}(x) = \prod_{i=1}^{n} f_i(u_i(x)), \qquad x \in \mathcal{M}_{reg}, \tag{8}$$

 $p_{\mu}(x) = \prod_{i=1}^{a} f_{i}(u_{i}(x)), \quad x \in \mathcal{M}_{reg},$ and that the on-manifold Hessian $\mathbf{H}_{x} \doteq \nabla_{x}^{2} \log p_{\mu}(x)$ has pairwise distinct eigenvalues a.e. on \mathcal{M}_{reg} . Then for each $x \in \mathcal{M}_{reg}$, the d seam directions (along which exactly one u_{i} varies) are determined intrinsically by p_{μ} , as eigenvectors of \mathbf{H}_{x} , unique up to permutation and sign (P&S).

Proof. For $x \in \mathcal{M}_{reg}$, let u^i, \ldots, u^d be unit tangent directions at x such that, along u^i , only u_i varies locally while $u_{i\neq i}$ remain constant, thus $\{u^i\}$ are orthonormal. Stack them as columns of $U_x \in \mathbb{R}^{m \times d}$ and define seam coordinates $u(x) \doteq (u_1(x), \dots, u_d(x))$. Thus

$$\nabla_u \log p_\mu(x) = \boldsymbol{U}_x^\top \nabla_x \log p_\mu(x), \qquad \qquad \nabla_u^2 \log p_\mu(x) = \boldsymbol{U}_x^\top \boldsymbol{\mathsf{H}}_x \boldsymbol{U}_x.$$

By Eq. 8, $\log p_{\mu}(x) = \sum_{i=1}^{d} \log f_i(u_i(x))$, hence

$$\left[\nabla_u^2 \log p_{\mu}(x)\right]_{ij} = \begin{cases} \frac{\partial^2}{\partial u_i^2} \log f_i(u_i(x)) & (i=j), \\ 0 & (i\neq j), \end{cases}$$

i.e. $\nabla_u^2 \log p_\mu(x)$ is diagonal, hence $\mathbf{H}_x = \mathbf{U}_x^\top [\nabla_u^2 \log p_\mu(x)] \mathbf{U}_x$ is an eigendecomposition with distinct eigenvalues (by assumption). Thus, u^i are eigenvectors of $\nabla^2_u \log p_\mu(x)$ and are unique up to P&S, as are seam directions following them.

Implication for identifiability. Lemma 6.3 isolates the intrinsic geometry of p_{μ} : once p_{μ} is fixed, the seams and their directions are fixed (P&S). Any Gaussian VAE decoder d matching p_{μ} and satisfying C1-C2 must therefore align its singular paths with those seams and inherit the same seam

Lemma 6.4 (A matching decoder finds seams). *Under assumptions of Lemma 6.3, let* $d: \mathbb{Z} \to \mathcal{X}$ be c.i.d.a.e. with factorised prior $p(z) = \prod_i p_i(z_i)$ and push-forward density $p_{\mu}^{(d)} \equiv p_{\mu}$ matching on $\mathcal{M}_d \doteq \{d(z)\} = \mathcal{M}$. If d satisfies CI-C2 a.e., then for any z and x = d(z):

- left singular vectors U_z of J_z coincide with the seam directions in Lemma 6.3 (up to P&S);
- the images under d of singular-vector paths (D3) are exactly the seams through x;
- along the i-th seam, the factor f_i is the 1-D push-forward of $p_i(z_i)$ (as in Lemma 5.1).

Proof. Since $p_{\mu}^{(d)} \equiv p_{\mu}$, both induce the same $\mathbf{H}_{x} = \nabla_{x}^{2} \log p_{\mu}(x)$. By Lemma 5.1 and C1, $\log p_{\mu}(x) = \sum_{i=1}^{d} \left(\log p_{i}(z_{i}) - \log s_{i}(z)\right)$, with s_{i} depending only on z_{i} by C2. Letting $u = \mathbf{U}_{z}^{\top}x$, the gradient along the manifold (*on-manifold score*) is given by

$$\nabla_x \log p_{\mu}(x) = U_z \nabla_u \log p_{\mu}(x), \qquad \left[\nabla_u \log p_{\mu}(x)\right]_i = \frac{1}{s_i(z)} \frac{\partial}{\partial z_i} \left(\log p_i(z_i) - \log s_i(z)\right).$$

Differentiating again along the manifold gives the *on-manifold Hessian* \mathbf{H}_x

$$\nabla_x^2 \log p_{\mu}(x) = \boldsymbol{U}_z \left[\nabla_u^2 \log p_{\mu}(x) \right] \boldsymbol{U}_z^{\gamma \top}, \qquad \left[\nabla_u^2 \log p_{\mu}(x) \right]_{ij} = \begin{cases} \frac{1}{s_i} \frac{\partial}{\partial z_i} \left[\nabla_u \log p_{\mu}(x) \right]_i & (i = j) \\ 0 & (i \neq j) \end{cases}$$

Hence,

$$\mathbf{H}_{x} = \mathbf{U}_{z} \operatorname{Diag}\left(\frac{1}{s_{i}(z)} \frac{\partial}{\partial z_{i}} \left[\nabla_{u} \log p_{\mu}(x)\right]_{i}\right) \mathbf{U}_{z}^{\top}$$

is an eigendecomposition. With a simple spectrum, eigenvectors are unique up to P&S and coincide with the intrinsic seam directions in Lemma 6.3. By Lemma C.1, singular-vector paths map to seams; and by Lemma 5.1, the factor along seam i equals the 1-D push-forward of $p_i(z_i)$.

Note that both proofs adopt the same technique, but Lemma 6.3 is entirely intrinsic to the manifold (hence no mention of a Jacobian), whereas Lemma 6.4 is with reference to a parameterisation of the manifold by a function d.

D.4 PROOF OF GAUSSIAN VAE IDENTIFIABILITY

Corollary D.1 (Gaussian VAE Identifiability). Let data be generated by c.i.d.a.e. $g: \mathbb{Z} \to \mathcal{X}$ with factorised prior $p(z) = \prod_{i=1}^d p_i(z_i)$; let a Gaussian VAE with diagonal posteriors learn a decoder $d: \mathbb{Z} \to \mathcal{X}$. Suppose both g and d satisfy C1–C2 a.e. and manifold densities match, $p_{\mu}^{(d)} \equiv p_{\mu}^{(g)}$, on the common manifold $\mathcal{M} = \{g(z)\} = \{d(z)\}$. If the tangent Hessian \mathbf{H}_x has a simple spectrum a.e., then d identifies the ground-truth seam decomposition (independent components) of $p_{\mu}^{(g)}$, up to P&S.

Proof. (Matching U) Equality $p_{\mu}^{(g)} \equiv p_{\mu}^{(d)}$ implies the same \mathbf{H}_x . By Lemma 6.3, its eigenvectors are intrinsic and unique (P&S), so $U_z^{(d)} = U_z^{(g)}$ (P&S, herein assume indices are relabelled to match).

(Matching S) In this common basis, under C1–C2, on-manifold scores are equal:

$$\left[\nabla_{u} \log p_{\mu}^{(g)}(x)\right]_{i} = \frac{1}{s_{i}^{(g)}(z)} \frac{\partial}{\partial z_{i}} \left(\log p_{i}(z_{i}) - \log s_{i}^{(g)}(z)\right) = \left[\nabla_{u} \log p_{\mu}^{(d)}(x)\right]_{i},$$

hence integrating along seam i (with $z_{\neg i}$ fixed, using equality of p_{μ} at a reference point to fix the integration constant), 1-D seam factors $\frac{p_i(z_i)}{s_i^{\gamma}}$ match. Since p_i is fixed, $s_i^{(d)} = s_i^{(g)}$, i.e. $S_z^{(d)} = S_z^{(g)}$.

With $V_z^{(d)} = V_z^{(g)} = I$ by C1, it follows that $J_z^{(d)} = J_z^{(g)}$ and the seam decomposition is identified up to P&S.

E DISENTANGLEMENT METRICS

Axis alignment score (AAS): Given a matrix of mutual information values, between each latent co-ordinate and each ground truth factor, one can normalise over rows or columns to compute a "distribution" of mutual information.

The entropy of each distribution gives a measure of how narrowly or sparsely information about a ground truth factor is captured across latents or the spread of information about each factor captured by a single latent. In either case, a "high entropy" distribution means information is widely spread, while low entropy means information about a factor is concentrated in a single latent, i.e. disentangled.

Entropy of the mutual information distribution can be computed row-wise or column-wise. AAS is a holistic metric combining the intuitions of both options into a single, robust score that evaluates how close matrix M is to a permuted diagonal form (zero entropy, perfect disentanglement).

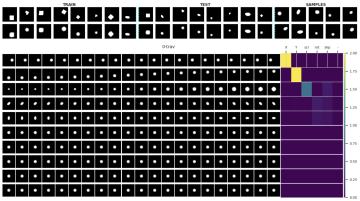
In a perfectly disentangled MI matrix, the sum of peak values per row equals the sum of peak values per column, and both equal the total sum of the matrix. AAS measures the ratio of the "sum of peaks" to the "total sum":

```
sum_col_max = sum(max(mut_info, dim=0))
sum_row_max = sum(max(mut_info, dim=1)
aas = 0.5 * (sum_row_max + sum_col_max) / sum(mut_info)
```

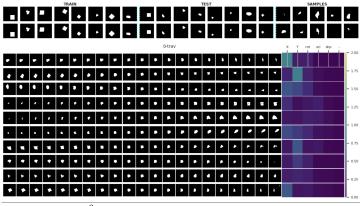
Normalised off diagonal: For gradient terms (here, a $d \times d$ matrix M) we compute a measure of diagonality by computing the ratios of normalised off-diagonal absolute values to on-diagonal values.

```
 d = M.shape[1] \\ num_off_diag = m * (m - 1) \\ M = abs(M) \\ M = diag(M)^(-0.5) * M * diag(M)^(-0.5) # normalise \\ mean_off_diag = (sum(M) - sum(diag(M))) / num_off_diag
```

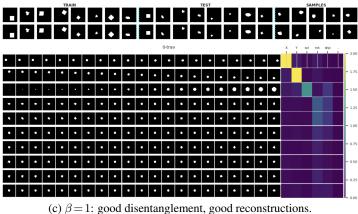
REDUCING β OVER TRAINING

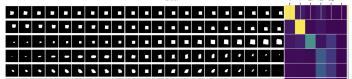


(a) $\beta = 1$: good disentanglement, blurry reconstructions.



(b) $\beta = 10^{-3}$: no clear disentanglement, good reconstructions.





(d) Traversals from a random test sample

Figure 7: **Testing the** β **-hypothesis:** (top) high β (1) gives best disentanglement (see heatmap) but blurry images (see top rows); (mid) low β (0.001) gives poor disentanglement but good reconstructions; (bottom) lowering β over training (1 \rightarrow 0.001) gives good disentanglement (see heatmap) and good reconstructions.

G MATERIAL ERRORS IN REIZINGER ET AL. (2022)

We note what appear to be several fundamental mathematical errors in the proof of Theorem 1 in Reizinger et al. (2022) rendering it invalid. Theorem 1 claims an approximation to the exact relationship given in Eq. 6

- 1. p.33, after "triangle inequality": $\left| \mathbb{E} \left[\|a\|^2 \|b\|^2 \right] \right| \leq \mathbb{E} \left[\|a b\|^2 \right]$, where a = x f, $b = -\sum \frac{\partial f}{\partial z_k} \dots$
 - (dropping expectations for clarity) this has the form $|||a||^2 ||b||^2| \le ||a b||^2$ (*)
 - true triangle inequality: $|||a|| ||b||| \le ||a b|| \implies |||a|| ||b|||^2 \le ||a b||^2$ (by squaring) this differs to (*) since norms are squared inside the absolute operator on the L.H.S.
 - counter-example to (*): $b=x>0, a=x+1 \implies |||a||^2-||b||^2|=|2x+1|>1=||a-b||^2$
- 2. next step, p.33: $\mathbb{E} \big[\| (c-e) (d-e) \|^2 \big] \le \mathbb{E} \big[\| c-e \|^2 + \| d-e \|^2 \big]$ where $c=x, \ d=f(z) \sum \frac{\partial f}{\partial z_k} ..., \ e=f(\mu)$
 - this has the form of the standard triangle inequality $||a-b|| \le ||a|| + ||b||$ except all norms are squared.
 - squaring both sides of the triangle inequality gives an additional cross term on the right that the used inequality omits, without which the inequality does not hold in general.
- 3. first step, p.34: drops the K term, which bounds the decoder Hessian and higher derivatives (in earlier Taylor expansion)
 - this omission is similar to a step in Kumar & Poole (2020) but is not stated, e.g. in Assumption 1.
 - since K is unbounded, any conclusion omitting it without justification is not valid in general.

H THE USE OF LARGE LANGUAGE MODELS (LLMS)

LLMs were used to assist drafting this paper as follows:

- general review for errors, inconsistencies and readability;
- verifying proofs, generating code snippets or identifying code errors;
- creating figure 1.