

Seminar Thesis

Learning Disentangled Representations with Semi-Supervised Deep Generative Models

Department of Statistics
Ludwig-Maximilians-Universität München

David B. Hoffmann

Munich, February 6th, 2026



Supervised by M. Sc. Simon Rittel

Abstract

Lorem ipsum dolor sit amet, consectetur adipiscing elit, sed do eiusmod tempor incididunt ut labore et dolore magna aliqua. Ut enim ad minim veniam, quis nostrud exercitation ullamco laboris nisi ut aliquip ex ea commodo consequat. Duis aute irure dolor in reprehenderit in voluptate velit esse cillum dolore eu fugiat nulla pariatur. Excepteur sint occaecat cupidatat non proident, sunt in culpa qui officia deserunt mollit anim id est laborum.

Contents

1	Introduction	1
2	Background and Formulation	2
2.1	Disentanglement in Variational Auto-Encoders	2
2.2	Degrees of Supervision	2
2.3	Semi-Supervised Variational Auto-Encoders	3
2.4	Importance Sampling	4
3	Semi-Supervised Disentanglement with Arbitrary Dependencies	6
3.1	The Generalised Variational Objective	6
3.2	Relation to Kingma’s M2 Model	7
3.3	Graphical Model and Stochastic Implementation	8
3.4	Experiments and Findings	8
3.4.1	Benchmark Classification: MNIST and SVHN	8
3.4.2	Intrinsic Faces: Disentangling Continuous Factors	9
3.4.3	Multi-MNIST: Stochastic Dimensionality and Compositionality	9
4	Extensions and Limitations	11
4.1	Own Criticism	11
4.2	General Critiques	12
4.2.1	Impossibility of Unsupervised Disentanglement	12
4.2.2	The Limitation of Isotropic Priors	13
4.2.3	Unbounded Likelihoods and Mode Collapse	14
4.3	Disentanglement Implications and Limitations	14
4.3.1	Semantic Conflation Problem	14
4.3.2	Breaking the ELBO Bottleneck	15
4.3.3	Differentiating Consistency and Restrictiveness	16
4.4	Extensions in Supervision: From Semi to Weak	17
4.4.1	The Weak Supervision Paradigm	17
4.4.2	Multimodal VAEs and Mutual Supervision	18
4.5	Applications and Integrations in the Literature	18
4.5.1	Domain Adaptation in the Medical Domain	18
4.5.2	Causal Extension and Robustness	19
4.5.3	Catastrophic Forgetting and Generative Replay	19
5	Experiments	21
5.1	Implementation Details	21
5.2	Supervision Weight	21
5.3	Label Corruption	22
6	Discussion	24
6.1	Future Directions	25
7	Conclusion	26

A	Systematic Review Process	V
B	SSVAE Training	V

1 Introduction

2 Background and Formulation

2.1 Disentanglement in Variational Auto-Encoders

Variational Auto-Encoders (VAEs) are generative models that learn a deep latent variable model by maximizing a lower bound on the marginal likelihood of the data (Kingma and Welling, 2014, Rezende et al., 2014). Given a dataset $\mathcal{D} = \{x^{(i)}\}_{i=1}^N$, we assume the data is generated by a random process involving an unobserved continuous random variable z . The generative process consists of a prior distribution $p_\theta(z)$ and a conditional likelihood $p_\theta(x|z)$, typically parametrised by a neural network (the decoder) with parameters θ .

Since the true posterior $p_\theta(z|x)$ is generally intractable, VAEs introduce an approximate posterior $q_\phi(z|x)$, parametrised by a separate neural network (the encoder) with parameters ϕ . The model is trained by maximizing the evidence lower bound (ELBO) on the marginal log-likelihood $\log p_\theta(x)$:

$$\log p_\theta(x) \geq \mathcal{L}_{\text{ELBO}}(x; \theta, \phi) = \mathbb{E}_{q_\phi(z|x)}[\log p_\theta(x|z)] - D_{KL}(q_\phi(z|x) \| p(z)) \quad (1)$$

The first term is the reconstruction error, encouraging the decoder to recover the data from the latent code. The second term is the Kullback-Leibler (KL) divergence, which regularizes the approximate posterior to be close to the prior, typically assumed to be a standard multivariate Gaussian $p(z) = \mathcal{N}(0, I)$.

To allow for backpropagation through the stochastic sampling process, VAEs utilize the reparametrisation trick. Instead of sampling z directly from $q_\phi(z|x) = \mathcal{N}(\mu, \sigma^2)$, we sample an auxiliary noise variable $\epsilon \sim \mathcal{N}(0, I)$ and compute:

$$z = \mu_\phi(x) + \sigma_\phi(x) \odot \epsilon \quad (2)$$

where \odot denotes the element-wise product.

A key goal in representation learning is disentanglement, where individual latent units are sensitive to changes in single generative factors of the data while being invariant to others (Bengio et al., 2013). Formally, a representation is disentangled if it factorizes into independent subspaces corresponding to underlying factors of variation.

Standard VAEs often fail to learn disentangled representations due to the lack of explicit constraints on the latent structure. To address this, the β -VAE framework Higgins et al. (2017) introduces a hyperparameter β to weight the KL divergence term:

$$\mathcal{L}_{\beta\text{-VAE}} = \mathbb{E}_{q_\phi(z|x)}[\log p_\theta(x|z)] - \beta D_{KL}(q_\phi(z|x) \| p(z)) \quad (3)$$

Setting $\beta > 1$ imposes a stronger constraint on the latent bottleneck, encouraging the learned distribution $q_\phi(z|x)$ to align with the isotropic unit Gaussian prior. Since the prior has independent components, this pressure encourages the latent dimensions to become statistically independent, thereby promoting disentanglement, though often at the cost of reconstruction quality.

2.2 Degrees of Supervision

While standard supervised learning relies on a fully labelled dataset $\mathcal{D} = \{(x^{(i)}, y^{(i)})\}_{i=1}^N$, obtaining ground truth labels y is often expensive. Semi-supervised learning addresses

this by leveraging a small set of labelled data \mathcal{D}_S alongside a much larger set of unlabelled data $\mathcal{D}_U = \{x^{(j)}\}_{j=1}^M$. The objective is to utilize the structural information inherent in $p(x)$ from \mathcal{D}_U to improve the estimation of the conditional distribution $p(y|x)$, typically under the assumption that the decision boundary lies in low-density regions of the input space.

In contrast, weak supervision relaxes the requirement for exact ground truth labels, utilizing instead a dataset $\mathcal{D}_W = \{(x^{(i)}, \tilde{y}^{(i)})\}_{i=1}^N$ where \tilde{y} represents a noisy, limited, or heuristic approximation of the true label y . A prominent subset is pair learning, where supervision is provided as pairwise constraints $\mathcal{D}_P = \{(x^{(i)}, x^{(j)}, r^{(ij)})\}$. Here, $r^{(ij)}$ denotes a binary relationship (e.g., must-link or cannot-link) between samples, guiding the model to learn a representation that respects these relative similarities rather than absolute class assignments.

2.3 Semi-Supervised Variational Auto-Encoders

To leverage the generalization benefits of generative modelling for classification tasks, Kingma et al. (2014) extend the VAE framework to the semi-supervised setting. They propose three approaches: the *M1* model, which separates feature learning from downstream classification; the *M2* model, which integrates the class label directly into the latent generative process; and a stacked *M1+M2* architecture. We focus here on the *M2* model, where the data x is generated by both a latent class variable y and a continuous latent variable z .

The model describes the data x as being generated by a latent class variable y and a continuous latent variable z . The generative process factorizes as:

$$p_\theta(x, y, z) = p_\theta(x|y, z)p(y)p(z) \quad (4)$$

where $p(y) = \text{Cat}(y|\pi)$ and $p(z) = \mathcal{N}(z|0, I)$. In this prior specification, y and z are marginally independent. The likelihood $p_\theta(x|y, z)$ is parameterized by a deep neural network that takes both y and z as inputs to generate the observation x .

Crucially, the inference structure introduced to approximate the intractable posterior $p(z, y|x)$ utilizes a specific factorization that introduces a dependency between the latent variables. The recognition model is specified as:

$$q_\phi(z, y|x) = q_\phi(z|x, y)q_\phi(y|x) \quad (5)$$

Here, $q_\phi(y|x)$ is modelled as a categorical distribution, which allows it to be used as a classifier, while $q_\phi(z|x, y)$ is a Gaussian distribution where the parameters (mean and variance) are functions of both the data x and the label y . This dependency structure in the inference network allows the model to learn a class-conditional latent distribution for z , enabling it to effectively separate style z from class y .

In the inference process, however, these variables are coupled. The approximate posterior for the continuous variable is defined as $q_\phi(z|x, y)$. By conditioning on y , the network encourages z to capture residual variations orthogonal to class identity.

For the supervised dataset $\mathcal{D}_S = \{(x^m, y^m)\}$, both x and y are observed. The supervised objective $\mathcal{L}_S(\theta, \phi; x, y)$ consists of two main components: a generative term

which is the ELBO on the joint distribution of x , y , and z as well as a discriminative term weighted by α . By expanding the generative term using the factorization $p_\theta(x, y, z) = p_\theta(x|y, z)p(y)p(z)$, we obtain a reconstruction term which encourages the decoder output to closely match the encoder input, a KL regularization that encourages the style latent z to follow the prior, and the label prior $\log p(y)$:

$$\mathcal{L}_S(\theta, \phi; x, y) = \underbrace{\mathbb{E}_{q_\phi(z|x, y)} \left[\log \frac{p_\theta(x, y, z)}{q_\phi(z|x, y)} \right]}_{\text{Generative (ELBO on joint } x, y)} + \underbrace{\alpha \log q_\phi(y|x)}_{\text{Discriminative}} \quad (6)$$

$$= \mathbb{E}_{q_\phi(z|x, y)} \left[\log p_\theta(x|y, z) + \log p(y) + \log p(z) - \log q_\phi(z|x, y) \right] + \alpha \log q_\phi(y|x) \quad (7)$$

$$= \underbrace{\mathbb{E}_{q_\phi(z|x, y)} [\log p_\theta(x|y, z)]}_{\text{Reconstruction}} - \underbrace{D_{KL}(q_\phi(z|x, y) \| p(z))}_{\text{Style Regularisation}} + \underbrace{\log p(y)}_{\text{Label Prior}} + \underbrace{\alpha \log q_\phi(y|x)}_{\text{Discriminative}} \quad (8)$$

For the unsupervised dataset $\mathcal{D}_U = \{x^n\}$, the label y is treated as a latent variable. We marginalize over all classes to derive the unsupervised bound $\mathcal{L}_U(\theta, \phi; x)$:

$$\mathcal{L}_U(\theta, \phi; x) = \sum_y q_\phi(y|x) \left(\mathbb{E}_{q_\phi(z|x, y)} \left[\log \frac{p_\theta(x, y, z)}{q_\phi(z|x, y)} \right] \right) + \mathcal{H}(q_\phi(y|x)) \quad (9)$$

where $q_\phi(y|x)$ is the recognition network acting as a discriminative classifier.

The final training objective combines these bounds over both datasets:

$$\mathcal{L}(\theta, \phi; \mathcal{D}_U, \mathcal{D}_S) = \sum_{x \in \mathcal{D}_U} \mathcal{L}_U(\theta, \phi; x) + \gamma \sum_{(x, y) \in \mathcal{D}_S} \mathcal{L}_S(\theta, \phi; x, y) \quad (10)$$

Where γ is used to balance the contribution of the supervised data. Note that this additional weighting parameter was not part of the original formulation from Kingma et al. (2014) which solely relied on α .

2.4 Importance Sampling

Importance sampling is a Monte Carlo method used to estimate properties of a target distribution $p(x)$ that is difficult to sample from directly, by utilizing a simpler proposal distribution $q(x)$.

Given a target density $p(x)$ and a function $f(x)$, we wish to estimate the expectation $\mathbb{E}_p[f(x)]$. We introduce a proposal distribution $q(x)$ such that $q(x) > 0$ whenever $p(x)f(x) \neq 0$ (i.e., the support of q covers the support of pf). By multiplying and dividing by $q(x)$, we rewrite the expectation as:

$$\mathbb{E}_{x \sim p}[f(x)] = \int f(x)p(x) dx = \int f(x) \frac{p(x)}{q(x)} q(x) dx = \mathbb{E}_{x \sim q} \left[f(x) \frac{p(x)}{q(x)} \right] \quad (11)$$

The term $w(x) = \frac{p(x)}{q(x)}$ is known as the importance weight. In practice, we draw K samples $\{x^{(k)}\}_{k=1}^K$ from the proposal distribution $q(x)$ and approximate the expectation using the sample mean:

$$\mathbb{E}_p[f(x)] \approx \frac{1}{K} \sum_{k=1}^K f(x^{(k)}) w(x^{(k)}) \quad (12)$$

This standard estimator is unbiased, provided $p(x)$ is a normalized density. However, in many applications, the target distribution is only known up to an intractable normalizing constant Z , i.e., $p(x) = \frac{\tilde{p}(x)}{Z}$. In this scenario, self-normalized importance sampling is used. We calculate unnormalized weights $\tilde{w}(x) = \frac{\tilde{p}(x)}{q(x)}$ and approximate the expectation using the ratio of estimators:

$$\mathbb{E}_p[f(x)] \approx \frac{\sum_{k=1}^K \tilde{w}(x^{(k)}) f(x^{(k)})}{\sum_{k=1}^K \tilde{w}(x^{(k)})} \quad (13)$$

Unlike the standard estimator, the self-normalized estimator is biased for finite K (with bias of order $1/K$), but consistent, converging to the true expectation as $K \rightarrow \infty$.

3 Semi-Supervised Disentanglement with Arbitrary Dependencies

While the $M2$ model by Kingma et al. (2014) provides a powerful framework for semi-supervised learning, it imposes a specific graphical structure: the latent variables y (label) and z (style) are assumed to be independent in the prior, and the recognition model is required to factorise as $q_\phi(y, z|x) = q_\phi(y|x)q_\phi(z|x, y)$. Narayanaswamy et al. (2017) generalise this approach to a broader class of partially-specified graphical models.

In this framework, the modeller specifies a graphical model where a subset of variables y are semantically meaningful (and partially observed), while the remaining variables z are left unstructured to capture nuisance variations. Unlike the $M2$ model, this framework allows for arbitrary dependency structures in both the generative model $p_\theta(x, y, z)$ and the approximate posterior $q_\phi(y, z|x)$. This flexibility enables the modelling of complex relationships, such as hierarchical dependencies or causal structures between labels and style variables, which are often necessary for disentanglement in real world datasets.

3.1 The Generalised Variational Objective

The core challenge in training these generalised models lies in the supervised objective as defined in Equation 6. For the unsupervised term \mathcal{L}_U (Equation 9), standard variational inference applies. However, for the supervised term \mathcal{L}_S , a difficulty arises when the recognition model $q_\phi(y, z|x)$ does not factorise conveniently. Specifically, evaluating the supervised lower bound requires sampling from the conditional posterior $q_\phi(z|x, y)$. For arbitrary graphical structures, this conditional distribution may not be available in closed form or easy to sample from directly.

To address this, Narayanaswamy et al. (2017) propose an importance sampling estimator that alleviates the need to sample from $q_\phi(z|x, y)$ directly. Instead, samples are drawn from an unconditioned proposal distribution, the marginal encoder $q_\phi(z|x)$, and reweighted. They start by rewriting the supervised objective using the fact that $q_\phi(z|x, y)$ factorizes to $\frac{q_\phi(y, z|x)}{q_\phi(y|x)}$ as follows

$$\mathcal{L}_S(\theta, \phi; x, y) = \alpha \log q_\phi(y|x) + \mathbb{E}_{q_\phi(z|x, y)} \left[\log \frac{p_\theta(x, y, z)}{q_\phi(z|x, y)} \right] \quad (14)$$

$$= \alpha \log q_\phi(y|x) + \mathbb{E}_{q_\phi(z|x, y)} \left[\log \frac{p_\theta(x, y, z)}{q_\phi(y, z|x)} + \log q_\phi(y|x) \right] \quad (15)$$

$$= (1 + \alpha) \log q_\phi(y|x) + \mathbb{E}_{q_\phi(z|x, y)} \left[\log \frac{p_\theta(x, y, z)}{q_\phi(y, z|x)} \right] \quad (16)$$

This formulation removes the need to evaluate $q_\phi(z|x, y)$ directly. To approximate the expectation, they employ self-normalised importance sampling as introduced in subsection 2.4. Instead of sampling from the conditional $q_\phi(z|x, y)$, they sample proposals z^s from the unconditioned encoder distribution $q_\phi(z|x)$. The unnormalised importance weights w^s are defined as the ratio of the target joint distribution to the proposal

$$w^s = \frac{q_\phi(y, z^s|x)}{q_\phi(z^s|x)} \quad (17)$$

where $z^s \sim q_\phi(z|x)$. Using these weights, the expectation is approximated as

$$\mathbb{E}_{q_\phi(z|x, y)} \left[\log \frac{p_\theta(x, y, z)}{q_\phi(y, z|x)} \right] \approx \frac{1}{S} \sum_{s=1}^S \frac{w^s}{Z} \log \frac{p_\theta(x, y, z^s)}{q_\phi(y, z^s|x)} \quad (18)$$

where $Z = \frac{1}{S} \sum_{s=1}^S w^s$ is the normalisation constant.

Furthermore, the term $\log q_\phi(y|x)$ is itself approximated using a Monte Carlo estimator of the lower bound normally used in maximum likelihood estimation

$$\log q_\phi(y|x) \geq \mathbb{E}_{q_\phi(z|x)} \left[\log \frac{q_\phi(y, z|x)}{q_\phi(z|x)} \right] \approx \frac{1}{S} \sum_{s=1}^S \log w^s \quad (19)$$

By combining the importance sampling estimator for the expectation (Equation 18) and the lower bound estimator for the log-likelihood (Equation 19), they obtain the final estimator for the supervised objective

$$\hat{\mathcal{L}}_S(\theta, \phi; x, y) := \frac{1}{S} \sum_{s=1}^S \frac{w^s}{Z} \log \frac{p_\theta(x, y, z^s)}{q_\phi(y, z^s|x)} + (1 + \alpha) \log w^s \quad (20)$$

This formulation ensures that the discriminative power of the model is maximised using the same samples z^s and weights w^s for both the generative and discriminative components.

3.2 Relation to Kingma's M2 Model

The framework by Narayanaswamy et al. (2017) generalizes the *M2* model. We can recover the exact *M2* objective by restricting the dependency structure. If we enforce the factorisation:

$$q_\phi(y, z|x) = q_\phi(y|x)q_\phi(z|x, y) \quad (21)$$

then the importance weights simplify to constants with respect to z

$$w^s = \frac{q_\phi(y|x)q_\phi(z^s|x, y)}{q_\phi(z^s|x, y)} = q_\phi(y|x) \quad (22)$$

Substituting this back into the estimator yields the standard *M2* objective. However, the generalised importance sampling formulation allows for alternative factorisations, such as $q_\phi(y, z|x) = q_\phi(y|x, z)q_\phi(z|x)$, where the label is predicted from the latent style z . This inverted structure effectively treats the latent variable z as a sufficient statistic for the label, enforcing a stronger form of disentanglement where z must contain all information necessary to predict y .

3.3 Graphical Model and Stochastic Implementation

To perform gradient ascent, the generative and recognition models are mapped onto a stochastic computation graph where each node forms a sub-graph. The generative model assumes a factorization $p_\theta(\mathbf{x}, \mathbf{y}, \mathbf{z}) = p_\theta(\mathbf{x}|\mathbf{y}, \mathbf{z})p(\mathbf{y})p(\mathbf{z})$, while the recognition model utilizes a conditional dependency structure $q_\phi(\mathbf{y}, \mathbf{z}|\mathbf{x}) = q_{\phi_z}(\mathbf{z}|\mathbf{y}, \mathbf{x})q_{\phi_y}(\mathbf{y}|\mathbf{x})$ to disentangle the digit label from the handwriting style.

The implementation handles supervision through three specific node types. First, for the fully supervised variable \mathbf{x} , they compute the likelihood $p_\theta(\mathbf{x}|\mathbf{y}, \mathbf{z}) = \mathcal{N}(\mathbf{x}; \eta_\theta(\mathbf{y}, \mathbf{z}))$, where η_θ is a neural network returning the distribution parameters. Second, the unobserved latent variable \mathbf{z} requires computing both the prior $p(\mathbf{z})$ and the conditional $q_\phi(\mathbf{z}|\mathbf{x}, \mathbf{y})$. This utilizes the Gaussian reparametrisation trick $\mathbf{z} = g(\epsilon, \lambda_{\phi_z}(\mathbf{x}, \mathbf{y}))$ where $\epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ to ensure differentiability. Finally, the partially observed variable \mathbf{y} is treated as observed when labels are available. When unavailable, it is sampled from $q_{\phi_y}(\mathbf{y}|\mathbf{x})$ using a Gumbel-Softmax relaxation to allow for gradient estimation on the discrete distribution.

Narayanaswamy et al. (2017) implement these different node types in a stochastic machine learning framework called ProbTorch which is available at: <https://github.com/probtorch/probtorch>.

3.4 Experiments and Findings

Narayanaswamy et al. (2017) evaluate the proposed framework across three distinct domains to demonstrate its ability to learn disentangled representations under partial supervision and varying structural complexity.

The authors do not provide full details about the implementation of their experiments. However, they do state that they use a linear layers for Modified National Institute of Standards and Technology (MNIST), a convolutional network for Street View House Numbers (SVHN) and the Yale-B datasets and a gated recurrent unit for the Multi-MNIST experiment. For the training process they use Adam with the "default learning rate and momentum correction terms". Narayanaswamy et al. (2017) vaguely state that depending on the dataset a batch size between 100-700 is used. Experiment specific information is mentioned in the following subsections.

3.4.1 Benchmark Classification: MNIST and SVHN

To validate the generalised objective against the standard $M2$ model (Kingma et al., 2014), the authors first applied the framework to the MNIST (Deng, 2012) and SVHN (Netzer et al., 2011) datasets. The graphical model structure mirrored the $M2$ setup, where the latent space consists of a partially observed class label y and an unobserved style variable z . They state that they use $\gamma = 1$ and introduce a supervision rate $\rho = \frac{\gamma M}{N + \gamma M}$. The α parameter is not specified.

On the MNIST dataset classification performance is evaluated for supervised set sizes of 100, 600, 1000 and 3000 out of the 50,000 unsupervised samples, where they marginally outperform the baseline $M2$ model. In the more complex SVHN domain the authors compare classification results for 1000 and 3000 labels out of the 70,000 unsupervised

samples. They find that their Semi-Supervised VAE (SSVAE) performs similar compared to the two-stage $M1+M2$ approach despite using a single-stage training process (Narayanaswamy et al., 2017). Qualitatively, the results demonstrated a clean separation of style and content. By fixing the style latent z and varying the label y , the model generated analogies where the same handwriting style was consistently applied to different digits. Conversely, fixing y and varying z produced variations in stroke width and slant while maintaining the digit identity.

Furthermore, the authors analysed the effect of the supervision weight γ and found that for sparsely labelled datasets ($M \ll N$), over-weighting the supervised term ($\gamma > 1$) significantly improved generalisation on the test set, although excessive weighting eventually led to overfitting. Here they offer conflicting statements about α where the figure description states that $\alpha = 50$ for MNIST and $\alpha = 70$ for SVHN while the discussion of the results of the supervision rate mentions that they used $\alpha = 0.1/\rho$ in line with Kingma et al. (2014). For a more detailed description of experimental results refer to Figure 3 of Narayanaswamy et al. (2017).

3.4.2 Intrinsic Faces: Disentangling Continuous Factors

Moving beyond categorical labels, the authors utilised the Yale-B dataset (Georghiades et al., 2001) to learn disentangled representations of *Identity* (categorical) and *Lighting* (continuous). The generative model defined latent variables for identity, lighting, shading, and reflectance, with partial supervision provided only for identity and lighting direction. The authors use a supervision rate of $\rho = 0.5$ with 6 supervised and 32 unsupervised data points corresponding which implies $\gamma \approx 5.33$. α is set to 10. The authors demonstrate that the model disentangled identity from lighting conditions with qualitative samples where they manipulate lighting direction on a specific identity’s face without altering facial features. Quantitatively, the framework outperformed an existing baseline on the identity classification task, achieving a semi-supervised error rate of 3.5% compared to approximately 30% for the baseline method by Jampani et al. (2015). For the regression of the continuous lighting parameter, the semi-supervised model achieved an angular error of 17.6%, compared to approximately 10% for the fully supervised baseline from Jampani et al. (2015), demonstrating the efficacy of the objective for continuous latents even with limited supervision. For a more detailed description of experimental results refer to Figure 4 of Narayanaswamy et al. (2017).

3.4.3 Multi-MNIST: Stochastic Dimensionality and Compositionality

The final experiment addressed a complex counting task using a custom Multi-MNIST dataset, where each image contained a varying number (1–3) of MNIST digits. This required a model with stochastic dimensionality, where the number of latent variables is itself a random variable K . The authors proposed a recursive generative model that sequentially samples digits and places them on a canvas using Spatial Transformer Networks. For this experiment neither α , γ or ρ are reported by the authors. The model learned to reliably predict the number of digits K and successfully decomposed overlapping digits into their constituent parts and background. Additionally, the authors compared a flat model against a nested model that incorporated a pre-trained MNIST auto-encoder as

a sub-component. This nested configuration demonstrated the framework’s capacity for modular model design and transfer learning. While the flat model achieved higher raw counting accuracy, the nested model still effectively leveraged the pre-learned constituent representations to reconstruct complex multi-digit scenes and decompose inputs into coherent individual digits. For a more detailed description of experimental results refer to Figure 6 of Narayanaswamy et al. (2017).

4 Extensions and Limitations

Several works build on the results from Narayanaswamy et al. (2017) and extend it to specific domains or in a general context. However, there are also some limitations and critiques in the subsequent literature. Both aspects are presented in the following sections. Relevant sources were systematically identified as outlined in Appendix A.

4.1 Own Criticism

In the following we present several issues in presentation Narayanaswamy et al. (2017) as well as unverified claims and minor critiques.

The presentation of the SSVAE framework in Narayanaswamy et al. (2017), while correct, at points leaves the reader guessing at specific details, disrupting the reading flow. As such, some figures do not have captions and are never referenced. Other figures are referenced without explaining their interpretation such as the meaning of shading in the nodes of the computational graph in Figure 1 and 5 of Narayanaswamy et al. (2017). It is never explained that the shading corresponds to the degree to which the variables are supervised. Furthermore, some variables are never defined such as the n in the relationship

$$(n \cdot l) \times r + \epsilon \tag{23}$$

of the supervised variables in the Yale-B dataset (Georghiades et al., 2001) experiment we explain in subsection 3.4.2. The meaning of this variable only becomes clear after reading the cited work of Jampani et al. (2015) which declares it as the normal map.

Further, Narayanaswamy et al. (2017) claim a number of properties and findings that are not shown theoretically or empirically. One example of this is their in the introduction that "a representation that has some factorisable structure, and consistent semantics associated to different parts, is more likely to generalise to a new task". This is not verified by their experiment or any cited source. In a later paper Locatello et al. (2019) investigate this claim and find no evidence to support it as further explained in subsection 4.2.1.

The authors introduce SSVAE as a framework for learning disentangled representations. Narayanaswamy et al. (2017) claim that it provides the user with the ability to precisely specify axes of variation and their dependencies. Specifically, in their Yale-B dataset experiment they claim to learn the relationship between lighting, shading and reflectance from Equation 23. While they qualitatively demonstrate disentangling as well as classification and regression performance, the learning of this relationship is not shown.

While the SSVAE framework generalizes to arbitrary dependency structures, its practical application in real-world settings with multiple partially observed variables reveals significant limitations. The current formulation relies on a binary distinction between unsupervised data \mathcal{D}_U and supervised data \mathcal{D}_S , assuming a consistent set of available labels for the supervised portion. However, in complex real-world datasets involving K distinct supervised variables $\mathbf{y} = \{y_1, \dots, y_K\}$, supervision is often heterogeneous; a datapoint $x^{(n)}$ may possess annotations for a random subset of variables $O^{(n)} \subseteq \mathbf{y}$, while remaining variables $U^{(n)} = \mathbf{y} \setminus O^{(n)}$ are missing.

Under the proposed framework, variables are treated as observed when available and sampled otherwise. Yet, implementing this dynamically in a batched stochastic gradient

descent setting becomes non-trivial. The standard objective presented in Equation 10 breaks down because there is no single \mathcal{L}_S . Instead, the objective function effectively fractures into 2^K potential observation patterns. To handle this, the practitioner must manually specify unique loss components or intricate masking logic for each missingness pattern to correctly toggle between computing the likelihood $p(y_k)$ (for $y_k \in O^{(n)}$) and performing importance sampling using $q_\phi(y_k|x)$ (for $y_k \in U^{(n)}$). This combinatorial explosion necessitates a bespoke and brittle implementation of the stochastic computation graph, undermining the flexibility intended by the general framework.

Additionally, while it is considered a good practice to reproduce experiments for methods presented in other papers instead of directly copying their results. This is not done by Narayanaswamy et al. (2017). For their MNIST/SVHN and Yale-B experiments (outlined in subsubsection 3.4.1 and subsubsection 3.4.2 respectively) they merely copy experimental result values from Kingma et al. (2014) and Jampani et al. (2015) instead of confirming the performance in their own experiments. Moreover, the definition of hyperparameters such as α and γ is not always clear as alluded to in subsection 3.4. Another minor error is the wrong column description in Figure 6 of Narayanaswamy et al. (2017) which states that the table on the bottom right reports count error, while the table actually contains accuracy values as correctly pointed out in the corresponding caption.

Lastly, in the same experiments, Narayanaswamy et al. (2017) do not mention which exact values they report. They state that they report error rates but not which metric they use. This complicates comparability to other methods and isolated interpretation of their results.

4.2 General Critiques

In the following the semi-supervised disentanglement approach from Narayanaswamy et al. (2017) is confirmed by results from Locatello et al. (2019), which demonstrate the impossibility of stable unsupervised disentanglement. Casale et al. (2018) offer a general critique of the independent and identically distributed (i.i.d.) assumption in VAEs directly mentioning Narayanaswamy et al. (2017).

4.2.1 Impossibility of Unsupervised Disentanglement

In their paper, "Challenging Common Assumptions in the Unsupervised Learning of Disentangled Representations", Locatello et al. (2019) presented a proof that unsupervised disentanglement is fundamentally impossible without inductive biases on both the model and the data.

The proof relies on the geometry of the latent space and the properties of the Gaussian prior used in VAEs. Assuming a ground-truth generative process $x = f(z)$ where $z \sim P(z)$, then if $P(z)$ is a standard multivariate Gaussian, it is rotationally symmetric. That is, for any orthogonal rotation matrix R , the distribution of Rz is identical to the distribution of z .

Given a new generative function $f'(z) = f(R^T z)$. The marginal distribution of observations $p(x)$ produced by this new model is identical to the original:

$$p(x) = \int p(x|z)p(z)dz = \int p(x|R^T z)p(z)dz$$

Since the unsupervised ELBO only depends on the marginal likelihood of the data $p(x)$, it cannot distinguish between the true disentangled model $f(z)$ and the entangled model $f'(z)$.

In their experiments Locatello et al. (2019) "do not find any evidence that the considered models can be used to reliably learn disentangled representations in an unsupervised manner as random seeds and hyperparameters seem to matter more than the model choice. Furthermore, good trained models seemingly cannot be identified without access to ground-truth labels even if we are allowed to transfer good hyperparameter values across data sets". Furthermore, they state that for the considered models and data sets, "increased disentanglement does not seem to lead to a decreased sample complexity of learning for downstream tasks".

This result confirms the necessity of some form of inductive bias for stable latent representations as proposed by Narayanaswamy et al. (2017). Narayanaswamy et al. (2017) anchor the variable y with labels, however, the variable z is essentially learned in an unsupervised manner relative to the style factors. Locatello et al. (2019)'s proof implies that there is no guarantee that z will factorize into meaningful sub-components purely based on the VAE objective. The model is free to learn any rotation of these style factors that explains the data.

Locatello et al. (2019) further highlight that they couldn't find any evidence that disentangled representation lead to better downstream performance, shifting the focus more towards interpretability and fairness. While Narayanaswamy et al. (2017) claim the contrary, namely that consistent and factorisable representation are more likely to generalise to a new task, they never formally prove or demonstrate this empirically in their paper.

4.2.2 The Limitation of Isotropic Priors

Casale et al. (2018) critique the use of overly simplistic i.i.d. Gaussian priors. While Narayanaswamy et al. (2017) move away from such priors by composing latent graphical models with deep likelihoods, Casale et al. (2018) argue that the framework still often relies on conditional independence assumptions that may induce excessive regularization. Casale et al. (2018) point out that for many important datasets—such as time-series of images, medical scans of the same patient, or rotated views of an object—the samples exhibit structured correlations that are better captured through the prior's covariance.

Using a prior that ignores these sample-to-sample correlations is a model misspecification that forces the encoder to discard the correlation structure. Casale et al. (2018) introduced the Gaussian Process Prior VAE (GPPVAE), which replaces the independent prior with a Gaussian Process:

$$\mathbf{z} \sim \mathcal{GP}(0, K(X, X'))$$

Here, the covariance kernel K explicitly models the correlation between samples (e.g., temporal proximity or identity). This allows the model to disentangle "object identity" from "view" by leveraging the kernel structure, a form of disentanglement that Narayanaswamy et al. (2017)'s graphical approach typically handles through semi-supervised labels and specific message-passing variational families.

Esmaili et al. (2019) also address the limitations of the isotropic Gaussian prior, arguing that it fails to exert sufficient regularizing pressure to force effective disentanglement. To

resolve the common problem of disentangling discrete factors of variation from continuous variables—a task that remains problematic for many contemporary approaches—they propose a structured decomposition of the ELBO objective. This modification allows for explicit control over the relative levels of disentanglement within different groups of latent variables. However, their analysis also highlights a significant optimization hurdle: discrete variables tend to exhibit higher likelihood values than continuous variables, which can introduce a bias that skews the optimization process in semi-supervised generative models.

4.2.3 Unbounded Likelihoods and Mode Collapse

Mattei and Frellsen (2018) investigated the maximum likelihood estimation for Deep Latent Variable Models (DLVMs), the class of models to which Narayanaswamy et al. (2017) semi-supervised VAE belongs. They proved that for continuous data, if the variance of the decoder’s output distribution is learned without constraints, the likelihood function is unbounded.

Consider a decoder that outputs a Gaussian distribution $\mathcal{N}(\mu_\theta(z), \sigma_\theta^2(z))$. If the model can map a specific latent point z_i to exactly match a data point x_i (i.e., $\mu_\theta(z_i) \approx x_i$) and simultaneously drive the variance $\sigma_\theta^2(z_i)$ toward zero, the likelihood density $p(x_i | z_i)$ approaches infinity:

$$p(x_i | z_i) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(x_i - \mu)^2}{2\sigma^2}\right) \rightarrow \infty \text{ as } \sigma \rightarrow 0$$

This creates a "hole" in the optimization landscape where the model collapses to storing data points (overfitting) rather than learning a generalized generative process. This is a severe limitation for Narayanaswamy et al. (2017)s framework when applied to continuous data, as it can lead to numerical instability and meaningless latent representations where the "style" z degenerates into a lookup table index rather than a semantic factor. Mattei and Frellsen (2018) suggest constraints on σ or specific architectural choices to ensure the existence of generalized maximum likelihood estimates.

4.3 Disentanglement Implications and Limitations

In the following, disentanglement and its limitations are presented.

4.3.1 Semantic Conflation Problem

Narayanaswamy et al. (2017)’s model feeds the label y and the latent z into the decoder: $p(x | y, z)$. The assumption is that y handles the class-specific information and z handles the rest. Joy et al. (2021) argue that this rigid separation is flawed because labels often imply continuous characteristics. Specifically they state that "Originally motivated by the desiderata of semi-supervised classification, each label is given a corresponding latent variable of the same type (e.g. categorical), whose value is fixed to that of the label when the label is observed and imputed by the encoder when it is not. Though natural, we argue that this assumption is not just unnecessary but actively harmful from a representation-learning perspective, particularly in the context of performing manipulations. To allow

manipulations, we want to learn latent factors that capture the characteristic information associated with a label, which is typically much richer than just the label value itself.”

They provide an example pertaining to a dataset of face. Consider a face with the label ”Young”. Joy et al. (2021) argue that age is not inherently discrete and that youth correlates with a number of continuous features such as smooth skin, specific facial structures, and hair density. In a standard VAE, the decoder sees $y = 1$ (”Young”) and learns to generate these features. Consequently, the latent z is absolved of the responsibility to encode the degree of youth or the specific way youth manifests in that image. This is semantic conflation: the discrete label ”steals” the semantic content from the continuous latent space.

They state that this leads to a failure in disentanglement: if we want to manipulate the ”age” of a face continuously, we cannot do so easily because the age information is locked inside the discrete variable y . Conversely, if we change the label y from ”Young” to ”Old” while keeping z fixed, the style z (which might encode ”smiling”) might be ignored or misinterpreted by the decoder because, so Joy et al. (2021) argue, the correlations between style and age are broken.

To resolve this, they proposed the Characteristic Capturing VAE (CCVAE) which builds on the SSVAE proposed by Narayanaswamy et al. (2017). They radically altered the graphical model. Instead of conditioning the generation on the label, they condition the inference of specific latents on the label. Specifically, they split the latent space into: (1) characteristic latents (z_c) which capture information correlated with the label and (2) salient latents (z_s) which capture the residual information.

The generative model is $p(x | z_c, z_s)$. Note that y is not an input to the decoder. Instead, an auxiliary classifier ensures that z_c is predictive of y : $\mathbb{E}[\log p(y | z_c)]$. By forcing z_c to predict y but feeding z_c (not y) to the decoder, the model must encode the semantic content of the label into the continuous space z_c . This avoids conflation. The label y acts as a supervisor for the structure of the latent space rather than a crutch for the decoder. Joy et al. (2021) demonstrated that this leads to superior performance in ”attribute traversal.” For example, dealing with the attribute ”Young,” the CCVAE could generate smooth transitions from young to old faces by interpolating in z_c , whereas the baseline models could only flip the binary switch y , resulting in discontinuous and often lower-quality transitions.

Nie et al. (2020) underline this critique by and state that ”it still remains unclear how the use of supervision impacts the disentanglement learning”. Which in part is also due to a lack of relevant metrics (Adel et al., 2018), making it hard to quantify how well disentanglement works and how the semi-supervised framework of Narayanaswamy et al. (2017) impacts it.

4.3.2 Breaking the ELBO Bottleneck

Feng et al. (2021) identified a specific optimization pathology in loss function of SSVAEs (Narayanaswamy et al., 2017), termed the *ELBO Bottleneck* and proposed Smooth-ELBO Optimal InTerpolation VAE (SHOT-VAE) as an extension of SSVAE to mitigate this limitation.

In SSVAEs, the objective combines the ELBO and a classification loss. Feng et al. (2021) observed that as training progresses, the ELBO term often plateaus before the inference

accuracy is maximized. The standard ELBO does not strictly penalize misclassifications in the latent space as long as the reconstruction is good. This leads to a disconnect: the model might reconstruct the image of a "9" perfectly well even if the latent code z and label y are slightly misaligned or ambiguous.

In the SSVAEs framework proposed by Narayanaswamy et al. (2017), the latent space is explicitly partitioned into a discrete label y and a continuous style variable z . While this separation is designed to "anchor" semantics, Feng et al. (2021) identify a fundamental optimization pathology in this joint objective: the *ELBO Bottleneck*. Because the generative term $p(x|y, z)$ can achieve high likelihood by absorbing label-relevant information into the "unspecified" style variable z , the model often reaches an ELBO plateau where reconstruction is near-perfect despite poor classification accuracy. This indicates that the inductive bias intended by Narayanaswamy et al. (2017) is often bypassed during training, as the standard ELBO lacks the gradient pressure to prevent z from "leaking" the information that should be exclusive to y .

Feng et al. (2021) introduce SHOT-VAE with two key innovations (1) Smooth-ELBO which is an approximation that integrates the label predictive loss directly into the ELBO derivation, rather than treating it as an auxiliary loss. This aligns the generative and discriminative objectives more tightly. (2) Optimal Interpolation for which they utilized data augmentation in the latent space (mixup) to fill the gaps between class clusters. By enforcing linearity in the latent space (interpolating between a "1" and a "7" should yield a semantic blend), they break the ELBO bottleneck and force the encoder to learn a more robust, disentangled structure.

Empirically, SHOT-VAE achieved significant error rate reductions on CIFAR-10 (6.11% vs baseline) compared to standard SSVAEs from Kingma et al. (2014), demonstrating that optimization dynamics are as critical as model architecture for disentanglement.

4.3.3 Differentiating Consistency and Restrictiveness

Shu et al. (2020) introduce a refined definition of disentanglement that is made up of consistency and restrictiveness. Here consistency is the degree to which a representation is deterministic with respect to the ground-truth factors. If we fix the ground-truth factor (e.g., color), the latent code should ideally be constant. On the other hand, restrictiveness is the degree to which a single dimension of the representation encodes only one ground-truth factor. This prevents a single latent dimension from encoding both color and shape. Shu et al. (2020) state that Narayanaswamy et al. (2017) falsely claim that their method leads to disentangled results through semi-supervision. Instead it merely creates consistent representation which are not necessarily restrictive and hence not guaranteed to be disentangled. However, they relativise this by conceding that on real world data consistency and restrictiveness are often strongly correlated.

These definitions provide a finer granularity than the mutual information metrics used previously. Critically, observational metrics (like the Mutual Information Gap (MIG) (Chen et al., 2018) used in Locatello et al. (2019)) often fail to capture these causal properties accurately. The differentiation of disentanglement introduced by Shu et al. (2020) paves the way for researchers to formally prove which types of weak supervision (e.g., restricted labeling, match-pairing) guarantee consistency or restrictiveness. This represents a theoretical maturation from Narayanaswamy et al. (2017)'s reliance on empirical

validation via reconstruction.

4.4 Extensions in Supervision: From Semi to Weak

In response to the impossibility results from Locatello et al. (2019), the research community pivoted. If pure unsupervised disentanglement is impossible, and full supervision is expensive, what is the minimal signal required? Locatello et al. (2020) proposed Weakly-Supervised Disentanglement as a robust extension to the semi-supervised paradigm. Joy et al. (2022) further extend the concept of weakly-supervised VAEs to multimodal learning and inference.

4.4.1 The Weak Supervision Paradigm

While Locatello et al. (2019) highlight the need for an inductive bias to create stable representations. Several works (Lin et al., 2020, Bouchacourt et al., 2018, Ke et al., 2024, Kim and Mnih, 2018) also critique the use of labels in the framework of Narayanaswamy et al. (2017). These criticisms are summarized well by Kim and Mnih (2018) who note that "semi-supervised approaches that require implicit or explicit knowledge about the true underlying factors of the data have excelled at disentangling" referring to Narayanaswamy et al. (2017). However, they also point out that "ideally we would like to learn these in an unsupervised manner, due to the following reasons: 1. Humans are able to learn factors of variation unsupervised [Perry et al. (2010)]. 2. Labels are costly as obtaining them requires a human in the loop. 3. Labels assigned by humans might be inconsistent or leave out the factors that are difficult for humans to identify".

Following these critiques and moving away from semi-supervision a number of weak-supervision frameworks emerged.

In "Weakly-Supervised Disentanglement Without Compromises," Locatello et al. (2020) model observations not as independent samples, but as pairs (x_1, x_2) that share at least one underlying factor of variation, even if the label of that factor is unknown.

For example, in a video of a moving arm, two adjacent frames might share the same "background color" and "object identity" but differ in "arm position." The weak label here is simply the knowledge that some factors are shared, without specifying what they are (unlike Narayanaswamy et al. (2017) who require the explicit label $y = \text{"Digit 7"}$).

They proved theoretically that knowing how many factors have changed (or stayed the same) between pairs is sufficient to guarantee disentanglement. This relaxes the requirement of Narayanaswamy et al. (2017) for expensive annotations for parts of the data while maintaining the capability of supervised disentanglement, only requiring a data collection process that provides paired samples (e.g., temporal coherence in video, or multi-camera setups).

Chen and Batmanghelich (2020) also critique the necessity for labels Narayanaswamy et al. (2017) and take a similar approach to Locatello et al. (2020), but focus more on the geometric structure of the latent space, utilizing a ranking-based loss to ensure that pairwise similarities in the data are preserved as proximity in specific latent dimensions. While Locatello et al. (2020) provide formal identifiability guarantees using *rank-1* pairs, where exactly one factor is varied, Chen and Batmanghelich (2020) rely on more flexible binary same/different judgments. This allows the model to learn from supervision that

indicates whether any factor is shared, without requiring the explicit knowledge of which specific factor differs between the pair.

Yang and Yao (2019) extend these principles to hand pose estimation and image synthesis, claiming the ability to learn interpretable disentangled representations without the necessity for additional weak labels. However, their framework still fundamentally relies on explicit pose annotations to supervise the partitioning of the latent space into pose and appearance components. This highlights a recurring theme in the literature where, despite claims of minimizing supervision, the underlying disentanglement mechanism remains dependent on the primary task’s labels to maintain structural integrity bringing us back to the impossibility theorem from Locatello et al. (2019).

4.4.2 Multimodal VAEs and Mutual Supervision

Joy et al. (2022), in "Learning Multimodal VAEs through Mutual Supervision," extended the semi-supervised framework to scenarios with multiple high-dimensional modalities (e.g., Image + Text).

In Narayanaswamy et al. (2017), supervision comes from a low-dimensional label y . In Multimodal VAEs, supervision comes from another rich modality. Joy et al. (2022) utilized a Product-of-Experts (PoE) aggregation. The joint posterior is approximated as the product of individual posteriors:

$$q(z|x_{\text{img}}, x_{\text{text}}) \propto q(z|x_{\text{img}}) \cdot q(z|x_{\text{text}})$$

This structure enforces a form of inter-modality disentanglement. Similar to Locatello et al. (2020), the shared information (semantics) must be encoded in the intersection of the posteriors, while modality-specific noise is filtered out. This mutual supervision allows the model to learn robust representations even when one modality is missing, extending the partial-specification idea of Narayanaswamy et al. (2017) to complex, unstructured labels.

4.5 Applications and Integrations in the Literature

While multiple papers simply use SSVAE as a baseline in their experiments (Gordon et al., 2019, Kulinski and Inouye, 2023), other apply various aspects of Narayanaswamy et al. (2017)s framework to their method. Vaze et al. (2023) critiques the prevalent use of synthetic data in the literature, several papers apply SSVAE to real world problems. This section discusses several works that apply the concepts of SSVAEs or use it as a baseline in their benchmark, which provides perspective into the utility of the framework proposed by Narayanaswamy et al. (2017) in real world conditions.

4.5.1 Domain Adaptation in the Medical Domain

Yang et al. (2019) apply the framework of SSVAEs to the critical problem of unsupervised domain adaptation in medical imaging (e.g., segmenting livers in CT vs. MRI scans). Building on the latent decomposition proposed by Narayanaswamy et al. (2017), Yang et al. (2019)s framework utilizes separate encoders (E_c content code and E_s style

code) to partition representations into a shared content manifold for anatomical preservation and a domain-specific style subspace for modality features. This architecture facilitates many-to-many cross-domain translation via style-code swapping and adversarial alignment, forcing the content space to converge on a modality-agnostic representation suitable for unsupervised liver segmentation.

By enforcing that the content code is shared while the style code is domain-specific, they could train a segmenter on labeled CT scans and apply it to unlabeled MRI scans by swapping the style codes. This effectively treats the "domain" (CT/MRI) as the specified factor y in Narayanaswamy et al. (2017)'s framework, but extends the mechanism to allow for image-to-image translation via the disentangled latent space.

Biswal et al. (2021) extend the semi-supervised generative framework to the medical domain through the development of EVA, a model designed for the generation of longitudinal electronic health records. Their approach builds upon the SSVAE architecture by introducing hierarchically factorized latent variables to better capture the complex dependencies inherent in temporal patient data. In this hierarchical structure, the upper-level latent variables are shared across the entire population to represent common clinical characteristics and global trends, while the lower-level variables remain patient-specific to capture individual disease trajectories and variations. By incorporating this multi-level factorization, the model ensures semantic consistency across the population while simultaneously allowing for the fine-grained representation of individual patient histories. Li et al. (2018) work on "multi-objective de novo drug design with conditional graph generative models" and mention SSVAEs as a possible extension to improve their work

4.5.2 Causal Extension and Robustness

Zhang et al. (2020) provide a causal extension to the graphical latent variable framework, applying it to the study of neural network robustness. By adopting a setup similar to the semi-supervised generative models in Narayanaswamy et al. (2017), they frame the relationship between latent factors and observed data through a structural causal model (SCM). This causal lens allows for a formal investigation into how perturbations and distribution shifts propagate through the model's representations. Their work demonstrates that leveraging a graphical latent model is not only useful for disentanglement but is also critical for analyzing and improving the robustness of neural networks against adversarial or environmental interventions.

4.5.3 Catastrophic Forgetting and Generative Replay

To address catastrophic forgetting in deep learning models trained sequentially, Ye and Bors (2020) introduced the Lifelong VAEGAN, which combines the inference capabilities of VAEs with the high-quality generation of Generative Adversarial Networks (GANs). Instead of a standard teacher-student setup, the model employs a generative replay mechanism and a two-step "wake-dreaming" optimization. In the "wake" phase, the generator is updated to approximate the distribution of both current data and replayed samples from previous tasks. In the "dreaming" phase, the model maximizes the log-likelihood of these samples to learn shared, disentangled latent representations across multiple domains.

Building on Narayanaswamy et al. (2017), Ye and Bors (2020) extend their Lifelong VAEGAN framework to the supervised setting and directly compare to SSVAE in their experiments where it has the second lowest classification error after the stacked generative semi-supervised model (M1+M2) model from Kingma et al. (2014).

In a later paper, Ye and Bors (2022) extend the SSVAE framework into the domain of continual learning by integrating it within a lifelong teacher-student network. Their method builds upon the generative principles established by Narayanaswamy et al. (2017) to facilitate knowledge retention and transfer across sequential learning tasks, utilizing the original model as a primary baseline for experimental comparison. Notably, while the authors formally cite Narayanaswamy et al. (2017), they attribute the architectural foundations to the Kingma et al. (2014) model in their technical descriptions, reflecting the shared heritage of these semi-supervised generative approaches. This extension demonstrates how the structured latent spaces of SSVAEs can be leveraged to mitigate catastrophic forgetting in dynamic data environments. Further, they find that the SSVAE architecture performs well in an experiment comparing the semi-supervised classification results on MNIST data, when considering MNIST to MNIST-Fashion lifelong learning. Specifically, they compare ten architectures of which SSVAE as the second lowest error.

5 Experiments

5.1 Implementation Details

Following previous work, our experiments were conducted with the MNIST dataset Deng (2012) using the digit label for partial supervision. We use the same fully connected feed forward architecture as Narayanaswamy et al. (2017), with Rectified Linear Units (ReLU) Agarap (2018) activations. The SSVAE architecture used in the following experiments is detailed in Figure 1.

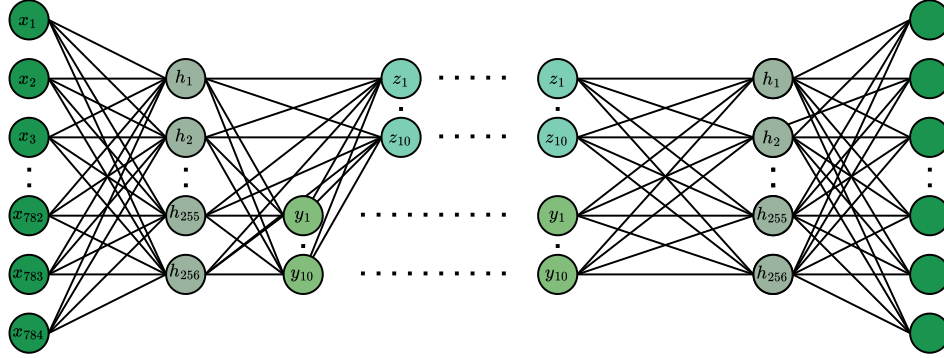


Figure 1: SSVAE architecture used throughout all experiments. The network has $28 \times 28 = 784$ inputs matching the MNIST input, 256 hidden neurons, 10 style variables and 10 digit neurons.

Following Narayanaswamy et al. (2017) we use the Adam optimizer (Kingma and Ba, 2015) with default parameters, a learning rate of 10^{-3} and a batch size of 128. Unlike the original paper we train for 40 instead of 200 epochs due to computational constraints. As shown in Appendix B this is sufficient for all trained models to converge. Each experiment is repeated with 10 different random seeds. Experiments were run on a Nvidia Tesla P40 GPU and implemented with the ProbTorch package introduced by Narayanaswamy et al. (2017). To reproduce the experiments presented in this paper see the code at: <https://github.com/davidbhoffmann/ssvae>.

5.2 Supervision Weight

While Narayanaswamy et al. (2017) explore the supervision rate ρ and therefore indirectly γ , the supervision weight α is not explored by the authors. Furthermore, the authors make conflicting statements about the value of α in the MNIST experiments as critiqued in subsection 4.1. In the Yale-B experiment α is set to a different value without justification and in the Multi-MNIST experiments the value of α is never mentioned. While Narayanaswamy et al. (2017) change its value, indicating its importance, they leave the reader guessing at the effect of α .

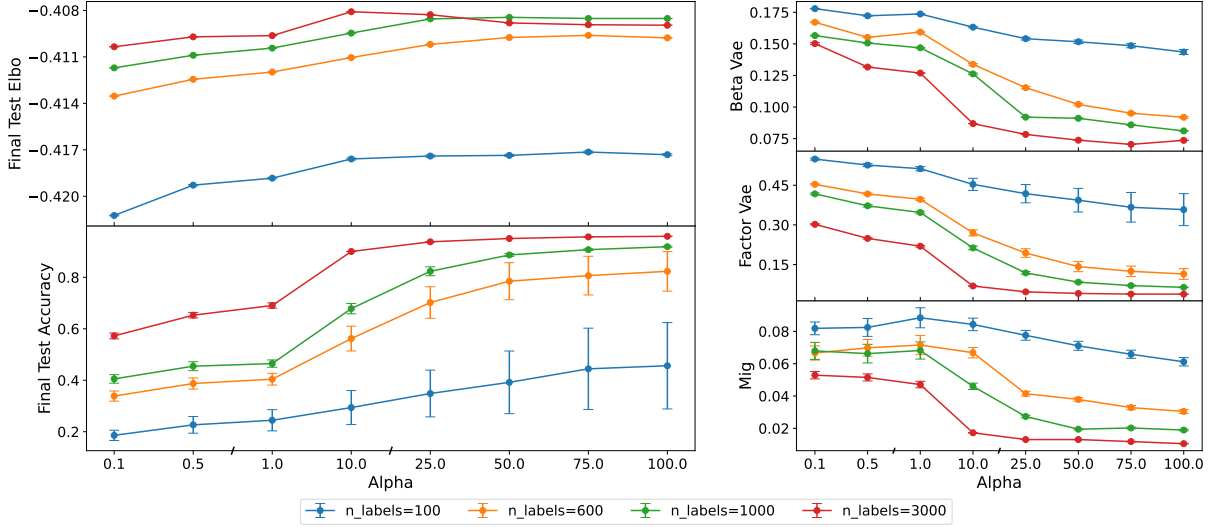


Figure 2: **Supervision Weight** experiment results with varying α for supervised set sizes of 100, 600, 1000 and 3000. Error bars indicate variance across 10 random seeds. **Top Left:** final test ELBO for varying α remains relatively stable. **Bottom Left:** final test accuracy of y improves for higher values of α . **Right:** disentanglement metrics (β -VAE Higgins et al. (2017), Factor-VAE Kim and Mnih (2018), Mutual Information Gap Chen et al. (2018)) decreases with increasing α across all three metrics.

To address this, we explore how robust the SSVAE framework is with respect to the supervision weight α which we evaluate for eight different values: 0.1, 0.5, 1, 10, 25, 50, 75 and 100. The remaining configurations follow the setup explained in subsection 5.1.

Figure 2 shows that classification accuracy of the semi-supervised y variable is sensitive to the supervision weight α . As expected increasing α leads to better classification performance, although the effect of increasing the value beyond 10 only leads to marginal improvements. While we did not report final training accuracy in the experiments, we can still observe there is no overfitting for high α values to the point of decreasing test accuracy.

The final test ELBO remains relatively constant with respect to α , which can also be seen in Appendix B. A minor increase of the ELBO can be observed with increasing α . The results further show that for the unsupervised z , stronger supervision decreases disentanglement.

5.3 Label Corruption

The most frequent critique identified in our review process was the use of labels (Lin et al., 2020, Bouchacourt et al., 2018, Ke et al., 2024, Kim and Mnih, 2018) as explained in subsubsection 4.4.1. One aspect of these critiques is the potential for incorrect human annotations.

To investigate the sensitivity of SSVAEs to potentially noisy labels, we use the setup as explained above and introduce a random corrupt labels into the training process. Specifically, in each epoch a specified amount of labels is flipped to a random class. We

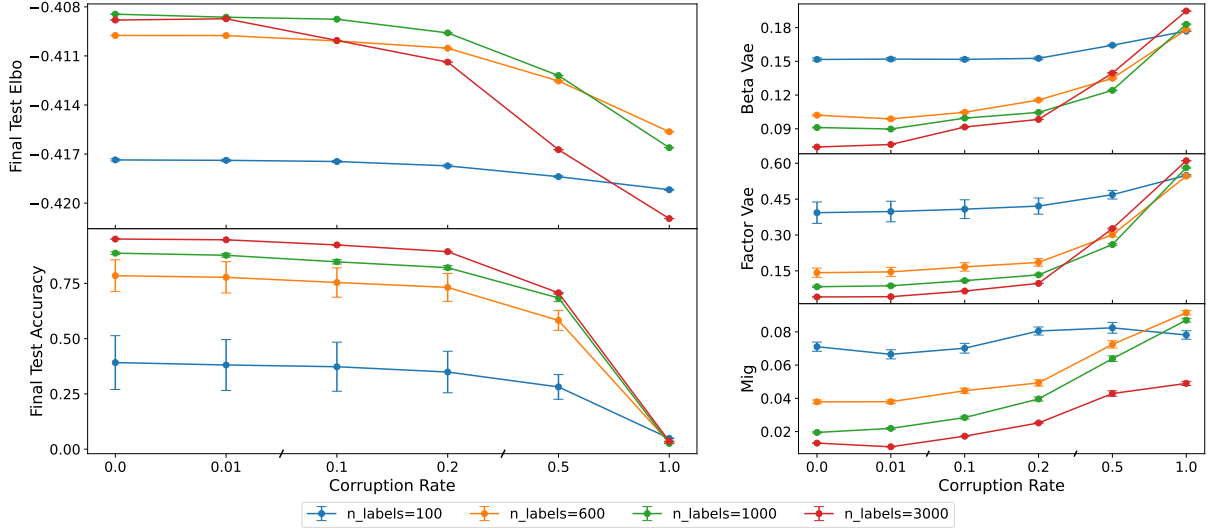


Figure 3: **Label Corruption** experiment results with varying label corruption rate for supervised set sizes of 100, 600, 1000 and 3000. Error bars indicate variance across 10 random seeds. **Top Left:** final test ELBO for varying noise levels remains relatively stable. **Bottom Left:** while the final test accuracy of y decreases for higher corruption rates it remains relatively high for up to 20% noise. **Right:** disentanglement (as measured through β -VAE Higgins et al. (2017), Factor-VAE Kim and Mnih (2018), Mutual Information Gap Chen et al. (2018)) increase with a higher corruption rate across all three metrics.

evaluate the SSVAEs for corruption rates of 0.0, 0.01, 0.1, 0.2, 0.5 and 1. The remaining configurations follow the setup explained in subsection 5.1.

The result of this experiment is shown in Figure 3. As expected We find that classification accuracy of the latent y variable decreases as the corruption rate of the labels increases. However, the accuracy does remain relatively stable up to a corruption rate of 20%.

Similar to the the supervision weight experiments, the final test ELBO remains relatively constant. However, a marginal ELBO drop with increased corruption rates is noticeable across configurations. The results further show that for the unsupervised z , more label corruption in y decreases disentanglement.

6 Discussion

In this work, we review the SSVAEs by Narayanaswamy et al. (2017), a unified framework for learning disentangled representations by integrating deep generative models with partially-specified probabilistic graphical models. Unlike standard variational autoencoders that often assume a flat latent space, their approach allows for the injection of domain knowledge through structured dependencies while retaining the flexibility of neural networks for function approximation. By allowing for dependency structures they extend previous work on semi-supervised variational auto-encoders from Kingma et al. (2014).

While most literature builds on the original SSVAE paper from Kingma et al. (2014) (with XX citation son google scholar add reference), there is a also a considerably body of derivative work for Narayanaswamy et al. (2017) (with XX citations on google scholar). However, we only found few works that deeply discuss and critique the method specifically. While Locatello et al. (2019) validate the use of an inductive bias for disentanglement they also critique the use of labels. This is generally the most frequent critique in the literature (Lin et al., 2020, Bouchacourt et al., 2018, Ke et al., 2024, Kim and Mnih, 2018). Other papers critique the effect of supervision on disentanglement and mroe general aspects which are not specific to SSVAEs. We do however find that none of the critiques are substantiate with empirical evidence specific to SSVAEs.

In our experiments, we investigate the critique of labels being potentially noisy, harming the disentanglement process. While we find that noisy labels decrease classification accuracy of the semi-supervised y latent, we also find that a small share of corrupt labels (up to 20%)only marginally affects performance. We further investigate the effect varying the supervision weight α on the SSVAE and find that it improves classification of the y latent whitout leading to overfitting even for values up to 100 and do not affect the final test ELBO. For both experiments we find that altering the strenght of supervision, through α or weakening the supervision signal through label corruption, has an effect on the accuracy and the disentanglement of z . ohh semanitic conflation problem and my experiments

While the framework generalizes to arbitrary dependency structures, its practical application in real-world settings with multiple partially observed variables reveals significant limitations. The current formulation relies on a binary distinction between unsupervised data \mathcal{D} and supervised data \mathcal{D}^{sup} , assuming a consistent set of available labels for the supervised portion[cite: 116]. However, in complex real-world datasets involving K distinct supervised variables $\mathbf{y} = \{y_1, \dots, y_K\}$, supervision is often heterogeneous; a datapoint may possess annotations for a random subset of variables while others are missing. Under the proposed framework, variables are treated as observed when available and sampled otherwise. Implementing this dynamically in a batched setting becomes non-trivial, as the objective function effectively fractures into 2^K potential observation patterns. This combinatorial explosion necessitates a bespoke and brittle implementation of the stochastic computation graph, undermining the flexibility intended by the general framework.

6.1 Future Directions

The current implementation relies on static graph definitions, but the natural evolution of this work lies in the integration with probabilistic programming languages. Probabilistic programs would allow for more expressive models involving recursion and control flow, which are currently limited by the static nature of the underlying computation graph. Future work will explore amortizing inference over these dynamic structures, enabling the learning of disentangled representations for increasingly complex, structured data domains such as video or language.

7 Conclusion

A concise summary of contents and results

The trajectory of research stemming from "Learning Disentangled Representations with Semi-Supervised Deep Generative Models" illustrates the rigorous self-correction mechanism of the scientific community. Siddharth et al. (2017) correctly identified that supervision was the key to unlocking interpretable representations, effectively predicting the impossibility results that would later prove purely unsupervised disentanglement futile.

However, the specific mechanism they proposed—a simple partitioning of the latent space into label y and style z —proved to be structurally insufficient for complex real-world data. It suffered from semantic conflation (Joy et al.), lacked theoretical guarantees for the unspecified variables (Locatello et al.), and was prone to optimization failures (Feng et al., Mattei et al.).

The extensions detailed in this report represent a maturation of the field. We have moved from simple semi-supervision to weakly supervised causal mechanisms, from static models to lifelong learning agents, and from deterministic encoders to uncertainty-aware Bayesian frameworks. The legacy of Siddharth et al. lies not just in their specific architecture, but in shifting the paradigm towards principled, supervision-guided structure learning.

A Systematic Review Process

To identify relevant extension and limitations of Narayanaswamy et al. (2017) in the literature, the results of several retrieval methods were synthesized. Note that the papers identified through the different methods had some overlap. For any paper to be considered it had to fulfil one of the following selection criteria: (1) directly critique Narayanaswamy et al. (2017), (2) cite the paper in the context of a more general critique, (3) extend the SSVAE framework or (4) use it in an experiment.

Firstly, Connected Papers (<https://www.connectedpapers.com>) lead to a list of 10 derivative works, of which only one matched the selection criteria.

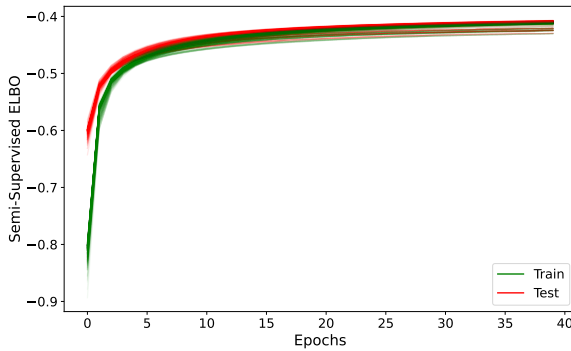
Secondly, Google Scholar (<https://scholar.google.com>) was used to find results relevant literature in the list of 445 papers which cite Narayanaswamy et al. (2017). The first 150 citing papers sorted by relevance were checked for the selection criteria. For the first 80 results, all paragraphs that mentioned Narayanaswamy et al. (2017) were reviewed in detail and for the remaining 70 papers only papers that mentioned one of the following keywords in the title were considered: *semi-supervised*, *weak-supervision* and *label*. With this method 21 relevant papers were identified.

Additionally, the search functionality of Google Scholar was used to find results for the search terms *semi-supervised disentanglement* and *supervised disentanglement*. For each term the first 10 results were reviewed in detail leading to a total of 4 relevant papers.

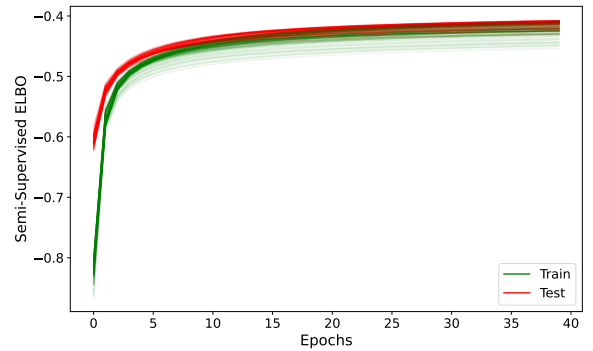
Lastly, Gemini-3-Pro was used to generate a report with relevant extensions and critiques which lead to 5 relevant papers.

B SSVAE Training

Figure 4a and Figure 4b show the training curve plots of all runs from both the supervision factor (alpha) and corruption rate experiments. The figures show that all runs converge without overfitting.



(a) All training runs of the **supervision factor α experiment** in subsection 5.2



(b) All training runs of the **corruption rate experiment** in subsection 5.3

References

- Adel, T., Ghahramani, Z. and Weller, A. (2018). Discovering interpretable representations for both deep generative and discriminative models, *in* J. Dy and A. Krause (eds), *Proceedings of the 35th International Conference on Machine Learning*, Vol. 80 of *Proceedings of Machine Learning Research*, PMLR, pp. 50–59.
- Agarap, A. F. (2018). Deep learning using rectified linear units (relu), *ArXiv Preprint*.
- Bengio, Y., Courville, A. and Vincent, P. (2013). Representation learning: A review and new perspectives, *IEEE transactions on pattern analysis and machine intelligence* **35**(8): 1798–1828.
- Biswal, S., Ghosh, S., Duke, J., Malin, B., Stewart, W., Xiao, C. and Sun, J. (2021). Eva: Generating longitudinal electronic health records using conditional variational autoencoders, *Machine Learning for Healthcare Conference*, PMLR, pp. 260–282.
- Bouchacourt, D., Tomioka, R. and Nowozin, S. (2018). Multi-level variational autoencoder: Learning disentangled representations from grouped observations, *Proceedings of the AAAI Conference on Artificial Intelligence*, Vol. 32.
- Casale, F. P., Dalca, A. V., Saglietti, L., Listgarten, J. and Fusi, N. (2018). Gaussian process prior variational autoencoders, *Advances in Neural Information Processing Systems (NeurIPS 2018)*.
- Chen, J. and Batmanghelich, K. (2020). Weakly supervised disentanglement by pairwise similarities, *Proceedings of the AAAI Conference on Artificial Intelligence*, Vol. 34, pp. 3495–3502.
- Chen, R. T., Li, X., Grosse, R. B. and Duvenaud, D. K. (2018). Isolating sources of disentanglement in variational autoencoders, *Advances in neural information processing systems* **31**.
- Deng, L. (2012). The mnist database of handwritten digit images for machine learning research, *IEEE Signal Processing Magazine* **29**(6): 141–142.
- Esmaili, B., Wu, H., Jain, S., Bozkurt, A., Siddharth, N., Paige, B., Brooks, D. H., Dy, J. and Meent, J.-W. (2019). Structured disentangled representations, *The 22nd International Conference on Artificial Intelligence and Statistics*, PMLR, pp. 2525–2534.
- Feng, H.-Z., Kong, K., Chen, M., Zhang, T., Zhu, M. and Chen, W. (2021). Shot-vae: Semi-supervised deep generative models with label-aware elbo approximations, *Proceedings of the AAAI Conference on Artificial Intelligence* **35**(8): 7413–7421.
- Georghiades, A., Belhumeur, P. and Kriegman, D. (2001). From few to many: Illumination cone models for face recognition under variable lighting and pose, *IEEE Trans. Pattern Anal. Mach. Intelligence* **23**(6): 643–660.

- Gordon, J., Bronskill, J., Bauer, M., Nowozin, S. and Turner, R. E. (2019). Meta-learning probabilistic inference for prediction, *International Conference on Learning Representations (ICLR)*.
- Higgins, I., Matthey, L., Pal, A., Burgess, C. P., Glorot, X., Botvinick, M. M., Mohamed, S. and Lerchner, A. (2017). Beta-vae: Learning basic visual concepts with a constrained variational framework, *International Conference on Learning Representations*.
- Jampani, V., Eslami, S. M. A., Tarlow, D., Kohli, P. and Winn, J. (2015). Consensus Message Passing for Layered Graphical Models, in G. Lebanon and S. V. N. Vishwanathan (eds), *Proceedings of the Eighteenth International Conference on Artificial Intelligence and Statistics*, Vol. 38 of *Proceedings of Machine Learning Research*, PMLR, San Diego, California, USA, pp. 425–433.
URL: <https://proceedings.mlr.press/v38/jampani15.html>
- Joy, T., Schmon, S. M., Torr, P. H. S., Siddharth, N. and Rainforth, T. (2021). Capturing label characteristics in vaes, *International Conference on Learning Representations (ICLR)*.
- Joy, T., Shi, Y., Torr, P. H., Rainforth, T., Schmon, S. M. and Siddharth, N. (2022). Learning multimodal vaes through mutual supervision, *ICLR*.
- Ke, Q., Jing, X., Woźniak, M., Xu, S., Liang, Y. and Zheng, J. (2024). Apgvae: Adaptive disentangled representation learning with the graph-based structure information, *Inf. Sci.* **657**(C).
- Kim, H. and Mnih, A. (2018). Disentangling by factorising, *International Conference on Machine Learning*.
- Kingma, D. P. and Ba, J. (2015). Adam: A method for stochastic optimization, *International Conference on Learning Representations (ICLR)*.
- Kingma, D. P., Rezende, D. J., Mohamed, S. and Welling, M. (2014). Semi-supervised Learning with Deep Generative Models, *Advances in Neural Information Processing Systems*, Vol. 27, Curran Associates, Inc.
- Kingma, D. P. and Welling, M. (2014). Auto-Encoding Variational Bayes, *2nd International Conference on Learning Representations*.
- Kulinski, S. and Inouye, D. I. (2023). Towards explaining distribution shifts, *Proceedings of the 40th International Conference on Machine Learning (ICML)*.
- Li, Y., Zhang, L. and Liu, Z. (2018). Multi-objective de novo drug design with conditional graph generative model, *Journal of Cheminformatics* **10**(1): 33.
- Lin, Z., Thekumparampil, K., Fanti, G. and Oh, S. (2020). InfoGAN-CR and Model-Centrality: Self-supervised model training and selection for disentangling GANs, in H. D. III and A. Singh (eds), *Proceedings of the 37th International Conference on*

- Machine Learning*, Vol. 119 of *Proceedings of Machine Learning Research*, PMLR, pp. 6127–6139.
- Locatello, F., Bauer, S., Lucic, M., Gelly, S., Schoelkopf, B. and Bachem, O. (2019). Challenging common assumptions in the unsupervised learning of disentangled representations, *Proceedings of the 36th International Conference on Machine Learning (ICML)*.
- Locatello, F., Poole, B., Raetsch, G., Schoelkopf, B., Bachem, O. and Tschannen, M. (2020). Weakly-supervised disentanglement without compromises, *Proceedings of the 37th International Conference on Machine Learning (ICML)*.
- Mattei, P.-A. and Frellsen, J. (2018). Leveraging the exact likelihood of deep latent variable models, *Advances in Neural Information Processing Systems* **31**.
- Narayanaswamy, S., Paige, B., van de Meent, J.-W., Desmaison, A., Goodman, N., Kohli, P., Wood, F. and Torr, P. (2017). Learning Disentangled Representations with Semi-Supervised Deep Generative Models, *Advances in Neural Information Processing Systems*, Vol. 30, Curran Associates, Inc.
- Netzer, Y., Wang, T., Coates, A., Bissacco, A., Wu, B., Ng, A. Y. et al. (2011). Reading digits in natural images with unsupervised feature learning, *NIPS workshop on deep learning and unsupervised feature learning*, Granada, p. 7.
- Nie, W., Karras, T., Garg, A., Debnath, S., Patney, A., Patel, A. and Anandkumar, A. (2020). Semi-supervised StyleGAN for disentanglement learning, in H. D. III and A. Singh (eds), *Proceedings of the 37th International Conference on Machine Learning*, Vol. 119 of *Proceedings of Machine Learning Research*, PMLR, pp. 7360–7369.
- Perry, G., Rolls, E. and Stringer, S. (2010). Continuous transformation learning of translation invariant representations, *Experimental brain research* **204**(2): 255–270.
- Rezende, D. J., Mohamed, S. and Wierstra, D. (2014). Stochastic backpropagation and approximate inference in deep generative models, *International conference on machine learning* pp. 1278–1286.
- Shu, R., Chen, Y., Kumar, A., Ermon, S. and Poole, B. (2020). Weakly supervised disentanglement with guarantees, *International Conference on Learning Representations (ICLR)*.
- Vaze, S., Vedaldi, A. and Zisserman, A. (2023). No representation rules them all in category discovery, *Advances in Neural Information Processing Systems (NeurIPS)*.
- Yang, J., Dvornek, N. C., Zhang, F., Chapiro, J., Lin, M. and Duncan, J. S. (2019). Unsupervised domain adaptation via disentangled representations: Application to cross-modality liver segmentation, *International Conference on Medical Image Computing and Computer-Assisted Intervention*, Springer, pp. 255–263.

- Yang, L. and Yao, A. (2019). Disentangling latent hands for image synthesis and pose estimation, *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*.
- Ye, F. and Bors, A. G. (2020). Learning latent representations across multiple data domains using lifelong vaegan, *European Conference on Computer Vision (ECCV)*.
- Ye, F. and Bors, A. G. (2022). Lifelong teacher-student network learning, *IEEE Transactions on Pattern Analysis and Machine Intelligence* **44**(10): 6280–6296.
- Zhang, C., Zhang, K. and Li, Y. (2020). A causal view on robustness of neural networks, *Advances in Neural Information Processing Systems (NeurIPS)*.

Declaration of authorship

I hereby declare that the report submitted is my own unaided work. All direct or indirect sources used are acknowledged as references. I am aware that the Thesis in digital form can be examined for the use of unauthorized aid and in order to determine whether the report as a whole or parts incorporated in it may be deemed as plagiarism. For the comparison of my work with existing sources I agree that it shall be entered in a database where it shall also remain after examination, to enable comparison with future Theses submitted. Further rights of reproduction and usage, however, are not granted here. This paper was not previously presented to another examination board and has not been published.

Location, date

Name