# ON THE TRANSFERABILITY OF VAE EMBEDDINGS USING RELATIONAL KNOWLEDGE WITH SEMI-SUPERVISION

**Harald Strömfelt**
Department of Computing
Imperial College London
London, SW7 2AZ
h.stromfelt17@imperial.ac.uk

**Luke Dickens**
Department of Information Studies
University College London
London, WC1E 6BT
l.dickens@ucl.ac.uk

**Artur d'Avila Garcez**
Department of Computer Science
City University of London
London, EC1V 0HB
a.garcez@city.ac.uk

**Alessandra Russo**
Department of Computing
Imperial College London
London, SW7 2AZ
a.russo@imperial.ac.uk

November 17, 2020

## ABSTRACT

We propose a new model for relational VAE semi-supervision capable of balancing disentanglement and low complexity modelling of relations with different symbolic properties. We compare the relative benefits of relation-decoder complexity and latent space structure on both inductive and transductive transfer learning. Our results depict a complex picture where enforcing structure on semi-supervised representations can greatly improve zero-shot transductive transfer, but may be less favourable or even impact negatively the capacity for inductive transfer.

## 1 Introduction

When dealing with complex data, the effectiveness of a classifier/predictor is limited by its ability to extract useful information. As such, representations that clearly expose the semantics of the data should then be most amenable to downstream learning [1, 2]. This is often referred to as a challenge of acquiring a *disentangled* representation over the factors of the data [3]. A popular recent trend that has had significant success in this regard uses semi-supervised Variational AutoEncoders (VAE) [4, 5, 6, 7, 8, 9]. Whilst fully unsupervised VAE methods have been shown to require strong inductive bias [10], semi-supervised methods achieve disentanglement by training additional auxiliary tasks that are defined on the factors, alongside the standard VAE objective (see Appendix Eqn. 3).

Recently relation-learning as semi-supervision to VAE representation learning has shown promise in shaping the representations learned [7, 6, 11]. In practice, different relations between data are often interrelated if they are derived from shared underlying factors. We argue that this presents a trade-off between decoder complexity accuracy achievable via highly complex decoders and the value of a latent representation carries over to new data or tasks. As simpler decoders capture fewer independent relationships, they can provide a structural bias towards a beneficial sharing of semantic factors. However, overly simple decoders may only be able to express some global properties of relations and not others, e.g. symmetry, transitivity, etc. We explore this trade-off by investigating the inductive and transductive transfer performance of two relation-decoders: the "Neural Tensor Network" (NTN), a powerful latent factor model (LFM) [12, 13, 14, 15]; and our own novel *Dynamic Comparator* (DC) model with $10\times$ fewer parameters. While our DC decoder has stricter constraints on the expected latent space structure than NTN, it is still sufficiently flexible to express a broad class of global properties on relations. We evaluate these ideas on a variety of tasks using MNIST digit images, our results show that: 1. semi-supervision improves inductive transfer by an appreciable margin (also seen in [10]). 2. Strongly structuring the latent space can degrade the inductive transfer capacity of encodings; and 3. Over 90%

**(a)** Inductive Transfer

**1** Representation Learning Phase
- Train $\beta$-VAE with relation-decoders, using selected context

**2** Inductive Transfer Learning
- Freeze VAE s.t. embeddings are fixed
- Train "Wild Relational Neural Network" on RPM task
- Test WReN on RPM panels from MNIST test set

RPM task instance



(addition)

**(b)** Transductive Transfer

**1** Representation Learning Phase
- Withhold selected subset of digits from **isEqual**
- Include all digits for $\beta$-VAE and any other context relations
- Train $\beta$-VAE with relation-decoders, applying digit exclusion to **isEqual**

**2** (Zero shot) Transductive Transfer Learning
- Freeze all model parameters
- Test **isEqual** performance on held out digits

Baseline Setup - i.e. $\emptyset$ context
- At **1** train $\beta$-VAE unsupervised
- At **2** train **isEqual** on same digit exclusion
- Test **isEqual** on held out digits

Train
**isEqual** trained over:
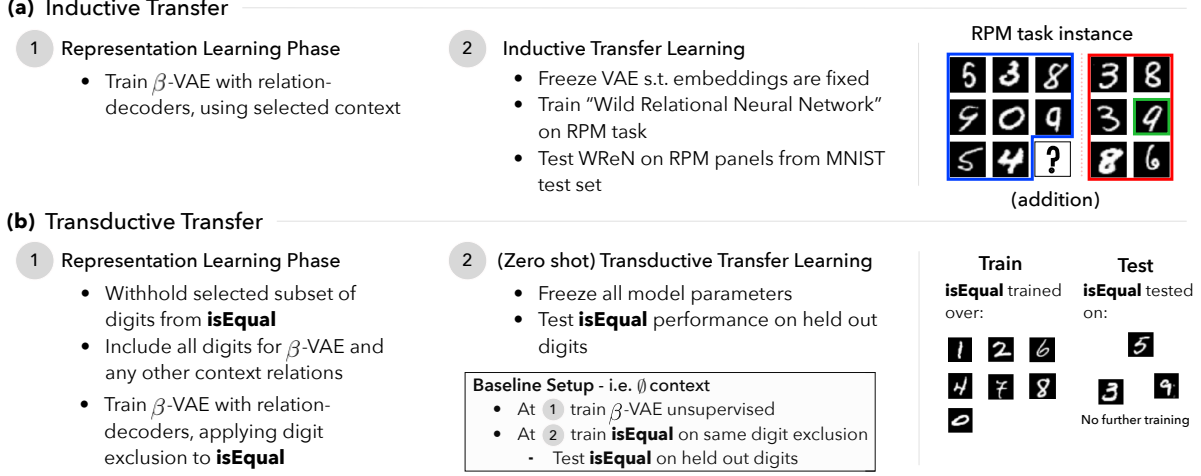
Test
**isEqual** tested on:



No further training

Figure 1: Overview of the transductive and inductive experimental setup used in the evaluation.

zero-shot transductive transfer accuracy on a binary task with semi-supervised VAE representations using our proposed DC model, significantly outperforming no semi-supervision ($< 76\%$) and NTN based semi-supervision ($< 80\%$).

In the following, Section 2 introduces the model, Section 3 presents and then discusses the experimental results. We include additional background on LFM and VAE in Appendix C.

## 2 A simple but flexible relation-decoder

In this paper, we focus our work on *MNIST* [16] and incorporate combinations of the $\{\mathsf{isEqual}, \mathsf{isGreater}, \mathsf{isSuccessor}\}$ binary numeric relations, in order to disentangle the 'number' factor. However, as each can have different symbolic properties (*e.g.* symmetry, transitivity and so on) we are presented with possible trade-off between model flexibility, in terms of the types of relations they can learn, and the degree to which they enforce disentanglement. On the one hand, restrictive relation-decoders that enforce disentanglement may not be able to model each relation type. For instance, Chen and Batmanghelich [6] was able to disentangle digit identity on MNIST using the $\mathsf{isEqual}$ relation - however, the proposed relation-decoder cannot model asymmetrical relations. In contrast, higher capacity decoders may generalise to complex relations but at the expense of instead negatively affecting disentanglement.

To bridge the gap, we propose the following "**D**ynamic **C**omparator" (**DC**) model, a new LFM that encourages disentanglement whilst being able to model (a)symmetric and (non-)transitive relations:

$$f_r^{DC}(\boldsymbol{z}_i, \boldsymbol{z}_j) = a_0 \cdot \underbrace{\sigma\big(\eta_0(\eta_1 - \|\boldsymbol{u} \odot (\boldsymbol{z}_i - \boldsymbol{z}_j + \boldsymbol{b}_\dagger)\|_2^2)\big)}_{f_r^\dagger} + a_1 \cdot \underbrace{\sigma\big((\eta_2 \cdot \boldsymbol{u}^\top (\boldsymbol{z}_i - \boldsymbol{z}_j + \boldsymbol{b}_\ddagger))\big)}_{f_r^\ddagger}. \tag{1}$$

For relation $r$, given $m$-dimensional latent representations $\boldsymbol{z}_i, \boldsymbol{z}_j \in \mathbb{R}^m$ obtained via a VAE encoder: $\boldsymbol{a} = \mathtt{Softmax}(\boldsymbol{A}) \in \mathbb{R}^2$ is an attention weighting between the two functional forms $f_r^\dagger$ and $f_r^\ddagger$; $\boldsymbol{u} = \mathtt{Softmax}(\boldsymbol{U}) \in \mathbb{R}^m$ is an attention mask over the full $m$-dimensional latent space; $\boldsymbol{b}_\dagger, \boldsymbol{b}_\ddagger \in \mathbb{R}^m$ are additional learnable bias terms; and $\eta_0, \eta_1 \in \mathbb{R}^+$ are non-negative and $\eta_2 \in \mathbb{R}$ any-valued scalar terms, respectively. Lastly, $\sigma$ is the sigmoid function used to bound the output to $[0,1]$, $\odot$ denotes element-wise multiplication and $\|\cdot\|_2$ is the $L2$-norm. $f_r^\dagger$ is a generalisation of the relation function from [6] designed for the $\mathsf{isEqual}$ relation and only capable of modelling symmetric 'zero-centred' relations. Firstly, by including $\boldsymbol{b}_\dagger$, the $L2$-norm can depend on the order of $\boldsymbol{z}_i$ and $\boldsymbol{z}_j$, enabling the modelling of asymmetric relations such as $\mathsf{isSuccessor}$. Further, whilst [6] hard-code the relevant subspace as a hyperparameter, we include a learned mask, $\boldsymbol{u}$, which allows the function to 'bind' itself to the relevant latent variable such that latent distance is only calculated on this subspace - this approach was previously done in [7]. In common with [6], $\eta_0$ sets the steepness of the true/false decision boundary and $\eta_1$ is a threshold that sets the width of the relation. $f_r^\ddagger$ generalises [11] whom omit $\boldsymbol{b}_\ddagger$ and set $\eta_2$ to be a non-negative scalar that models confidence in the relation. As such, [11] strictly models one-way ordinal '$>$' relations and critically has fixed predictions at equality (*i.e.* $\sigma(0) = 0.5$). In contrast, the proposed $f_r^\ddagger$ can learn the ordering of the relation and can learn any of $>, \geq, \leq, <$ type relations, such as $\mathsf{isGreater}$.
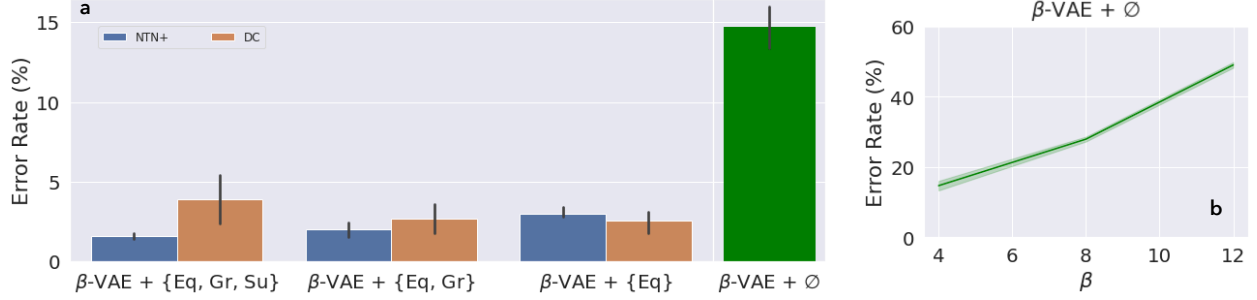
Figure 2: Inductive transfer error rate on our *MNIST* RPM task. (a) demonstrates the performance per context and relation-decoder with $\beta = 4$ and (b) demonstrates $\beta$ effects on an unsupervised $\beta$-VAE. We see a clear shift in performance when semi-supervision is used, with marginally worse performance for **DC**, and negative correlation of performance with $\beta$ increases for $\emptyset$ context.

Once again, [11] hard-code the sub-space for the ordering, whereas **DC** can discover it via $u$ as we perform a dot product with the mask to calculate the directional $z_i, z_j$ difference on the $u$ mask hyperplane. Importantly, **DC** can learn combinations of semantically similar relations such that they are each calculated using the same latent factors, which can support disentanglement. See Appendix A for further details.

In the next section, we compare the relative transfer performance obtained when using either **NTN+** (a modified version of NTN [17, 18] for $n$-ary relations - see Appendix Eqn. 4) or **DC** relation-decoders.

## 3 Experiments

In this section we compare the quality of representations produced by VAE semi-supervision using either **NTN+** or the proposed **DC**, against those generated by unsupervised $\beta - VAE$, for both inductive and transductive transfer. [1] (see Figure 1). In all experiments, we generate representations for digits from the MNIST dataset and introduce semi-supervision with additional predictions for binary relations isEqual (Eq), isGreater (Gr) and isSuccessor (Su). We consider four configurations of relations for semi-supervision: $\emptyset$ (no supervision), {Eq}, {Eq, Gr} and {Eq, Gr, Su}, and call this the *context*. Our approach could readily be adapted to work with a variety of unsupervised VAE, including those developed for disentanglement [19, 20, 21, 22, 23]. However, we chose to use the $\beta$-VAE due it showing competitive results whilst being straightforward to optimise given that it only has one hyperparameter, $\beta$, which is understood to control disentanglement pressure [3, 24]. This leads to the following joint objective,

$$\ln p_\theta(\boldsymbol{X}, \boldsymbol{Z}) \geq \mathcal{L}_{\beta\text{-VAE}}^{ELBO} - \lambda \underbrace{\mathbb{E}_{r, \boldsymbol{z}_i, \boldsymbol{z}_j, y_{ij} \in \mathcal{T}}[y_{ij} \ln(\hat{y}_{ij}) + (1 - y_{ij}) \ln(1 - \hat{y}_{ij})]}_{\mathcal{L}^{LFM}}, \tag{2}$$

where $\hat{y}_{ij}$ is estimated by $f_r(\boldsymbol{z}_i, \boldsymbol{z}_j)$ and $\mathcal{T}$ is the set of all positive ($y_{ij} = 1$) and negative triples ($y_{ij} = 0$) of the form $(\boldsymbol{z}_i, r, \boldsymbol{z_j})$. $\mathcal{L}_{\beta\text{-VAE}}^{ELBO}$ is the $\beta$-VAE ELBO (see Appendix C) and $\lambda$ is a weighting parameter. All latent representations are sampled according to $\boldsymbol{z} \sim q_\phi(\boldsymbol{Z}|\boldsymbol{X})$ where $q_\phi(\boldsymbol{Z}|\boldsymbol{X})$ is modelled by the VAE encoder. Given a context, we sample two triples per *MNIST* image and randomly select which relation to generate the triple for - this ensures a fixed number of triples between experiments. We direct readers to Appendix C for further details on the $\beta$-VAE and **NTN+**, and Appendix D for additional implementation details.

**Inductive Transfer:** In this setting both the source and target data are the same but the target (downstream) task differs [1]. We follow recent work [25, 26, 27] and create a RPM dataset consisting of $3 \times 3$ *MNIST* image panels, arranged into rows of addition or subtraction. The final row is left incomplete and a downstream reasoner is tasked with using the VAE image-encoding to select the correct tile (see the addition example in Figure 1(a)-right). We explore the downstream performance effects that different forms of semi-supervision have. see Appendix B for a detailed description of the RPM task and downstream reasoner. If digit classification is possible, it would be possible to complete this task by memorizing the addition/subtraction combinations, however including numeric ordering should alleviate the need for memorization. The aim of this experiment is to evaluate the inductive transfer improvement that semi-supervision produces and the relative benefit of regularising for further structure beyond digit identity. **Results:** Figure 2(a) shows the maximum 5000-step moving average test error rate obtained using each context and relation-decoder and (b) demonstrates how $\beta$ settings affect the downstream performance using an unsupervised VAE. **Discussion:** These results are in agreement with recent work that showed semi-supervision improves inductive feature-representation-transfer
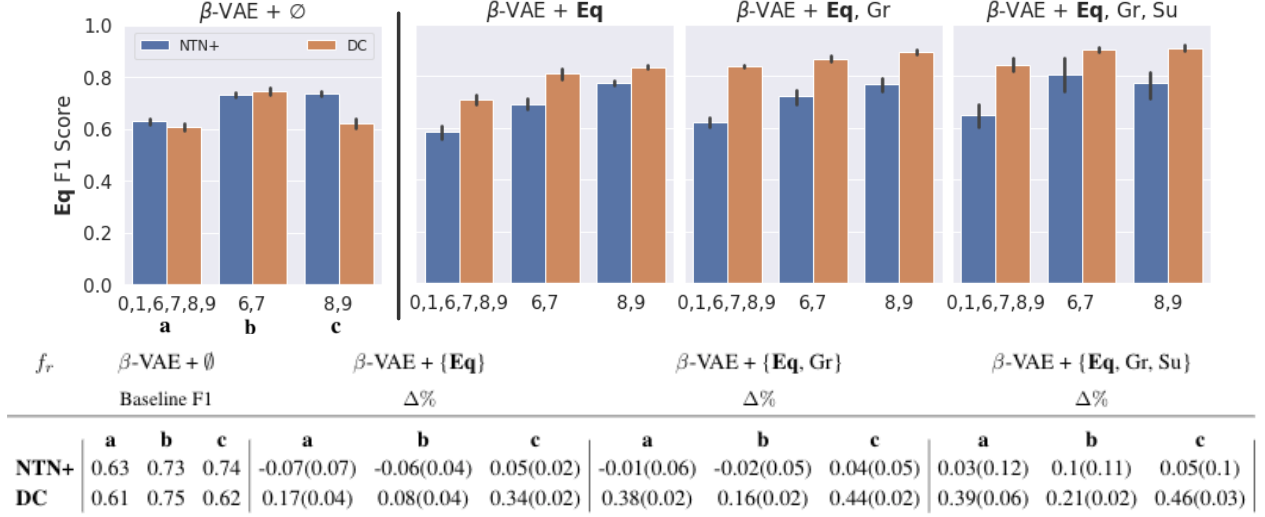
| $f_r$ | $\beta$-VAE + $\emptyset$ | | | $\beta$-VAE + {**Eq**} | | | $\beta$-VAE + {**Eq**, Gr} | | | $\beta$-VAE + {**Eq**, Gr, Su} | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | Baseline F1 | | | $\Delta\%$ | | | $\Delta\%$ | | | $\Delta\%$ | | |
| | **a** | **b** | **c** | **a** | **b** | **c** | **a** | **b** | **c** | **a** | **b** | **c** |
| **NTN+** | 0.63 | 0.73 | 0.74 | -0.07(0.07) | -0.06(0.04) | 0.05(0.02) | -0.01(0.06) | -0.02(0.05) | 0.04(0.05) | 0.03(0.12) | 0.1(0.11) | 0.05(0.1) |
| **DC** | 0.61 | 0.75 | 0.62 | 0.17(0.04) | 0.08(0.04) | 0.34(0.02) | 0.38(0.02) | 0.16(0.02) | 0.44(0.02) | 0.39(0.06) | 0.21(0.02) | 0.46(0.03) |

Figure 3: Zero-shot transductive transfer results showing (**upper**) the isEqual F1 scores on the *held out* subset digit classes (indicated by the horizontal axis groupings); and (**lower**) $\Delta\%$ difference w.r.t. baseline F1 for each relation-decoder and context, with standard deviation in parentheses. Both results obtained with $\beta = 4$ for the VAE. These results indicate that lower complexity decoders perform better at transductive transfer.

[10, 6]. However, in contrast with [25, 26], we observe a negative correlation with $\beta$ increase (Figure 2(b)). In conjunction with an increasing error rate as more relations are included for **DC** but not for **NTN+** (Figure 2(a)), the results indicate that enforcing stronger regularity in the latent embeddings worsens performance. In support of this, Appendix Figure 8 shows that adding more context increases the overall digit factor information captured using **NTN+**. In summary, reduced inductive transfer performance may occur when digit class identity is obscured.

**Transductive Transfer:** In contrast to inductive transfer, here the source and target task data is different but the task itself, in this case isEqual relation prediction, is equivalent in both cases. Concretely, we only train the isEqual relation on a subset of the data by omitting a selection of digits, but show all digits to the VAE and other relations. Hence, in the source domain isEqual only observes a subset of the digits and is then tested on the unseen digits, wherein no further training takes place; as such we test *zero-shot* transductive transfer. With this experiment, we aim to evaluate the amenability of the representations obtained using different contexts and relation-decoders, to the isEqual relation-decoder parameterization learned on the digit subset. **Results:** Figure 3 compares the isEqual F1 test scores on the held out digits, when learned using each relation-decoder. We use $\emptyset$ context as a baseline, wherein we pre-train an unsupervised $\beta$-VAE and post-train isEqual using each relation-decoder on frozen embeddings, with the same digit exclusion strategy. **Discussion:** in the baseline case, each relation-decoder cannot influence the VAE-encodings and so must 'fit' to the frozen latent embeddings that result from the pre-trained $\beta$-VAE. Interestingly, both decoders perform similarly even though they have different complexities. We then observe marked performance increases when including Eq for **DC** but not **NTN+**. This suggests that **DC** is able to impose stronger regularity on the $\beta$-VAE such that the resulting latent embeddings exhibit regularity that is amenable to **DC**, even on untrained digits. As expected, we observe a significant improvement for **DC** when including Gr since it observes the full digit set, for example improving by 39% over the baseline for exclusion setting (**a**). However, **NTN+** does not exhibit the same improvement. This may indicate that, although Gr is symbolically related to Eq, it may not be learned in such a way that this relatedness is captured between the latent embeddings and the Eq versus Gr relation-decoder parameters. Lastly, the increased variance for **NTN+** on context {Eq, Gr, Su} is likely due to it requiring more data to be trained - since we use a fixed triple "quota" that is shared between relation-decoders, adding more relations reduces the total number of triples observed by each relation-decoder. **DC** is by contrast more data efficient, due to it having far fewer parameters to learn. In summary, **DC** outperforms both the baseline and **NTN+** in each exclusion setting, with immediate gains for {Eq} context wherein no relations are trained on the held out digits. This indicates that better transductive transfer performances can be achieved when using relation-decoders that can impose consistent regularity on the $\beta$-VAE with respect to the relation-decoder parameters.

**Concluding remarks:** The results in this paper shed light onto the complex interplay between latent embedding structure and the decoders that are used to perform downstream tasks. In order to obtain transferable latent representations,

we observe that for powerful neural network based downstream learners, stronger regularity is less favourable. On the other hand, our results suggest that we can achieve better transductive transfer results if we enforce regularisation on the representations. This has the potential of encouraging a consistent structure across the latent space which relation-decoders can leverage.

# References

[1] S. J. Pan and Q. Yang, "A survey on transfer learning," *IEEE Trans. on Knowl. and Data Eng.*, vol. 22, no. 10, p. 1345–1359, Oct. 2010.

[2] Y. Bengio, A. Courville, and P. Vincent, "Representation learning: A review and new perspectives," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 35, no. 8, pp. 1798–1828, 2013.

[3] I. Higgins, L. Matthey, A. Pal, C. Burgess, X. Glorot, M. Botvinick, S. Mohamed, and A. Lerchner, "beta-VAE: Learning Basic Visual Concepts with a Constrained Variational Framework," in *5th International Conference on Learning Representations, {ICLR}*, Toulon, France, 2017.

[4] R. Shu, Y. Chen, A. Kumar, S. Ermon, and B. Poole, "Weakly Supervised Disentanglement With Guarantees," in *8th International Conference on Learning Representations, ICLR*, Addis Ababa, Ethiopia, 2020.

[5] F. Locatello, B. Poole, G. Rätsch, B. Schölkopf, O. Bachem, and M. Tschannen, "Weakly-Supervised Disentanglement Without Compromises," *CoRR*, vol. abs/2002.0, 2020.

[6] J. Chen and K. Batmanghelich, "Weakly Supervised Disentanglement by Pairwise Similarities," in *Proceedings of the 32nd AAAI Conference on Artificial Intelligence, AAAI*, New York, NY, USA, 2020.

[7] T. Karaletsos, S. Belongie, and G. Rätsch, "When crowds hold privileges: Bayesian unsupervised representation learning with oracle constraints," in *4th International Conference on Learning Representations, {ICLR}*, San Juan, Puerto Rico, 2016, pp. 1–16.

[8] D. P. Kingma, S. Mohamed, D. J. Rezende, and M. Welling, "Semi-supervised Learning with Deep Generative Models," in *Advances in Neural Information Processing Systems 27: Annual Conference on Neural Information Processing Systems*, Montreal, Quebec, Canada, 2014, pp. 3581—-3589.

[9] Z. Feng, A. Zeng, X. Wang, D. Tao, C. Ke, and M. Song, "Dual swap disentangling," in *Advances in Neural Information Processing Systems 32*, Montreal, Canada, 2018, pp. 5894–5904.

[10] F. Locatello, S. Bauer, M. Lucic, G. Rätsch, S. Gelly, B. Schölkopf, and O. Bachem, "Challenging Common Assumptions in the Unsupervised Learning of Disentangled Representations," in *Proceedings of the 36th International Conference on Machine Learning,{ICML}*, Long Beach, California, USA, 2019, pp. 4114—-4124.

[11] J. Chen and K. Batmanghelich, "Robust ordinal VAE: employing noisy pairwise comparisons for disentanglement," *CoRR*, vol. abs/1910.05898, 2019.

[12] T. Trouillon, É. Gaussier, C. R. Dance, and G. Bouchard, "On inductive abilities of latent factor models for relational learning," *Journal of Artificial Intelligence Research*, vol. 64, pp. 21–53, 2019.

[13] M. Nickel, K. Murphy, V. Tresp, and E. Gabrilovich, "A review of relational machine learning for knowledge graphs," *Proceedings of the IEEE*, vol. 104, no. 1, pp. 11–33, 2016.

[14] Q. Wang, Z. Mao, B. Wang, and L. Guo, "Knowledge graph embedding: A survey of approaches and applications," *IEEE Transactions on Knowledge and Data Engineering*, vol. 29, no. 12, pp. 2724—-2743, 2017.

[15] R. Socher, D. Chen, C. Manning, D. Chen, and A. Ng, "Reasoning With Neural Tensor Networks for Knowledge Base Completion," in *Advances in Neural Information Processing Systems 26: 27th Annual Conference on Neural Information Processing Systems*, 2013, pp. 926–934.

[16] Y. LeCun and C. Cortes, "MNIST handwritten digit database," 2010. [Online]. Available: http://yann.lecun.com/exdb/mnist/

[17] I. Donadello, L. Serafini, and A. d'Avila Garcez, "Logic Tensor Networks for Semantic Image Interpretation," in *Proceedings of the Twenty-Sixth International Joint Conference on Artificial Intelligence*, 2017, pp. 1596—-1602.

[18] L. Serafini and A. D. Garcez, "Logic tensor networks: Deep learning and logical reasoning from data and knowledge," in *Proceedings of the 11th International Workshop on Neural-Symbolic Learning and Reasoning (NeSy'16) co-located with the Joint Multi-Conference on Human-Level Artificial Intelligence {(HLAI} 2016)*, New York, NY, USA, 2016.

[19] R. T. Q. Chen, X. Li, R. B. Grosse, and D. Duvenaud, "Isolating Sources of Disentanglement in Variational Autoencoders," in *Advances in Neural Information Processing Systems 31: Annual Conference on Neural Information Processing Systems*, Montreal, Quebec, Canada, 2018, pp. 2615—-2625.

[20] A. Kumar, P. Sattigeri, and A. Balakrishnan, "Variational inference of disentangled latent concepts from unlabeled observations," in *6th International Conference on Learning Representations, {ICLR}*, Vancouver, BC, Canada, 2018.

[21] H. Kim and A. Mnih, "Disentangling by Factorising," in *Proceedings of the 35th International Conference on Machine Learning, {ICML}*, Stockholm, Sweden, 2018, pp. 2654—-2663.

[22] K. Ridgeway and M. C. Mozer, "Learning Deep Disentangled Embeddings With the F-Statistic Loss," in *Advances in Neural Information Processing Systems 31: Annual Conference on Neural Information Processing Systems*, Montreal, Quebec, Canada, 2018, pp. 185—-194.

[23] C. Eastwood and C. K. I. Williams, "A framework for the quantitative evaluation of disentangled representations," in *6th International Conference on Learning Representations, {ICLR}*, Vancouver, BC, Canada, 2018.

[24] C. P. Burgess, I. Higgins, A. Pal, L. Matthey, N. Watters, G. Desjardins, and A. Lerchner, "Understanding disentangling in $\beta$-VAE," in *Advances in Neural Information Processing Systems 30*, no. Nips, Long Beach, CA, USA, 2017. [Online]. Available: http://arxiv.org/abs/1804.03599

[25] X. Steenbrugge, S. Leroux, T. Verbelen, and B. Dhoedt, "Improving Generalization for Abstract Reasoning Tasks Using Disentangled Feature Representations," in *Neural Information Processing Systems (NeurIPS) Workshop on Relational Representation Learning*, Montreal, Canada, 2018.

[26] S. van Steenkiste, F. Locatello, J. Schmidhuber, and O. Bachem, "Are Disentangled Representations Helpful for Abstract Visual Reasoning?" in *Advances in Neural Information Processing Systems 32: Annual Conference on Neural Information Processing Systems*, Vancouver, BC, Canada, 2019, pp. 14 222—-14 235.

[27] D. G. Barrett, F. Hill, A. Santoro, A. S. Morcos, and T. Lillicrap, "Measuring abstract reasoning in neural networks," *35th International Conference on Machine Learning, ICML 2018*, vol. 10, pp. 7118–7127, 2018.

[28] V. Gutiérrez-Basulto and S. Schockaert, "From Knowledge Graph Embedding to Ontology Embedding? An Analysis of the Compatibility between Vector Space Representations and Rules," in *Principles of Knowledge Representation and Reasoning: Proceedings of the Sixteenth International Conference*, Tempe, Arizona, US, 2018.

[29] J. C. Raven, "Standardization of progressive matrices, 1938," *British Journal of Medical Psychology*, vol. 19, no. 1, pp. 137–150, 1941.

[30] C. Kemp and J. B. Tenenbaum, "The discovery of structural form," *Proceedings of the National Academy of Sciences of the United States of America*, vol. 105, no. 31, pp. 10 687–10 692, 2008.

[31] A. Madsen and A. R. Johansen, "Neural Arithmetic Units," in *8th International Conference on Learning Representations, ICLR*, Addis Ababa, Ethiopia, 2020.

[32] A. Trask, F. Hill, S. E. Reed, J. W. Rae, C. Dyer, and P. Blunsom, "Neural arithmetic logic units," in *Advances in Neural Information Processing Systems 31*, Montreal, Canada, 2018, pp. 8046–8055.

[33] A. Santoro, D. Raposo, D. G. Barrett, M. Malinowski, R. Pascanu, P. Battaglia, and T. Lillicrap, "A simple neural network module for relational reasoning," in *Advances in Neural Information Processing Systems 30*, Long Beach, CA, USA, 2017.

[34] D. P. Kingma and M. Welling, "Auto-Encoding Variational Bayes," in *Proceedings of the 2nd International Conference on Learning Representations*, Banff, Alberta, Canada, 2014.

[35] A. Paszke, S. Gross, F. Massa, A. Lerer, J. Bradbury, G. Chanan, T. Killeen, Z. Lin, N. Gimelshein, L. Antiga, A. Desmaison, A. Kopf, E. Yang, Z. DeVito, M. Raison, A. Tejani, S. Chilamkurthy, B. Steiner, L. Fang, J. Bai, and S. Chintala, "Pytorch: An imperative style, high-performance deep learning library," in *Advances in Neural Information Processing Systems 32*, H. Wallach, H. Larochelle, A. Beygelzimer, F. d'Alché-Buc, E. Fox, and R. Garnett, Eds. Curran Associates, Inc., 2019, pp. 8024–8035. [Online]. Available: http://papers.neurips.cc/paper/9015-pytorch-an-imperative-style-high-performance-deep-learning-library.pdf

Table 1: Parameters learned for each *MNIST* relation included for representation learning. Parameter names refer to Eqn. 1. In particular, $\boldsymbol{a} = [a_0 \quad a_1]$ weightings indicate which functional form is used to model the relation. These values are complimented by the relation function parameter visualisation given by Figure 4.

| Relation | Param. | Value |
|---|---|---|
| isEqual | $\eta_2$ | [-0.] |
| isEqual | $\mathbf{b}_\ddagger$ | [-2.31] |
| isEqual | $\mathbf{b}_\dagger$ | [-0.] |
| isEqual | $\eta_0$ | [190] |
| isEqual | $\eta_1$ | [-0.04] |
| isEqual | $\mathbf{a}$ | [0.01 0.99] |
| isGreaterThan | $\eta_2$ | [70] |
| isGreaterThan | $\mathbf{b}_\ddagger$ | [-0.18] |
| isGreaterThan | $\mathbf{b}_\dagger$ | [0.38] |
| isGreaterThan | $\eta_0$ | [0.] |
| isGreaterThan | $\eta_1$ | [2.24] |
| isGreaterThan | $\mathbf{a}$ | [0.99 0.01] |
| isSuccessor | $\eta_2$ | [-100] |
| isSuccessor | $\mathbf{b}_\ddagger$ | [-2.65] |
| isSuccessor | $\mathbf{b}_\dagger$ | [-0.35] |
| isSuccessor | $\eta_0$ | [170] |
| isSuccessor | $\eta_1$ | [1.17] |
| isSuccessor | $\mathbf{a}$ | [0.01 0.99] |

# A    Relation-decoder case study

In this section, we expose what is learned, in terms of the parameterization and resulting latent embeddings, when training each relation-decoder alongside a $\beta$-VAE. For DC, we include mask visualisations ($\boldsymbol{u}$ in Eqn. 1) to examine which latent subspace is used to calculate each relation, as well as the full parameterization of DC for *MNIST*. We then include an exploration of each latent dimension against ground truth factors for both *MNIST* and the benchmark *dSprites* dataset, which consists of grey-scale images of hearts, square and ovals; each varying across scale, orientation and position [3] For both datasets, Mutual Information Gap (MIG) scores are calculated - the MIG score calculates the normalised difference between the two latent dimensions that share the greatest mutual information w.r.t. each ground truth factor [19]. When quoted as a single value it is averaged across ground truth factors:

$$\frac{1}{K} \sum_{k=1}^{K} \frac{1}{H(v_k)} \big( I_i(z_{j^{(k)}} ; v_k) - \max_{j \neq j^{(k)}} I_i(z_j ; v_k) \big)$$

where $K$ is the total number of ground truth factors, $I_i(\cdot ; \cdot)$ is the mutual information between two random variables for input $\boldsymbol{x}^i$, $H(\cdot)$ is the entropy of a given random variable which acts as a normalisation term and $z_j^{(k)}$ is the latent factor that maximises the mutual information with factor $v_k$. Note that the $k$th factor gives a large contribution, if one dimension of $\boldsymbol{z}$, say $z_{j(k)}$, can explain a large amount of the variation in $v_k$, while other dimensions explain little. For *MNIST*, we use the 'digit' factor (K=1) and for *dSprites*, we follow [19] and present the MIG score as an average over scale, orientation, y-position, x-position factors ($K = 5$).

## A.1    Dynamic Comparator

To begin with, Table 1 presents an example of the learned DC parameters after training $\beta$-VAE + {Eq, Gr, Su} on *MNIST*. Figure 4 then includes true/false region visualisation for each relation-decoder to enable visualisation of the function itself - this is possible since the masks (shown by Figure 5) select 1-dimensional subspaces, namely $\boldsymbol{z}^l$ with $l = 6$, so we can plot the relation-decoder output over $\boldsymbol{z}_i^6$ versus $\boldsymbol{z}_j^6$.

As shown by Figure 4, each of isEqual, isGreater and isSuccessor are modelled differently, in accordance with their dissimilar symbolic properties (see Table 2). Firstly, the transitive and asymmetrical isGreater relation uses $f_r^\ddagger$ and is thus modelled by a step function around $(\boldsymbol{z}_i^6 - \boldsymbol{z}_j^6) > 0$. Here, a non-zero $\boldsymbol{b}_\ddagger$ ensures that $f_{\text{isGreater}}^{DC}(\boldsymbol{z}_i^6, \boldsymbol{z}_j^6) \approx 0$ if $\boldsymbol{z}_i^6 = \boldsymbol{z}_j^6$. isSuccessor, which is non-transitive and asymmetrical, is modelled as a relative distance based function, through the use of $f^\dagger$. However, unlike isEqual, which is by contrast symmetric and thus invariant to input
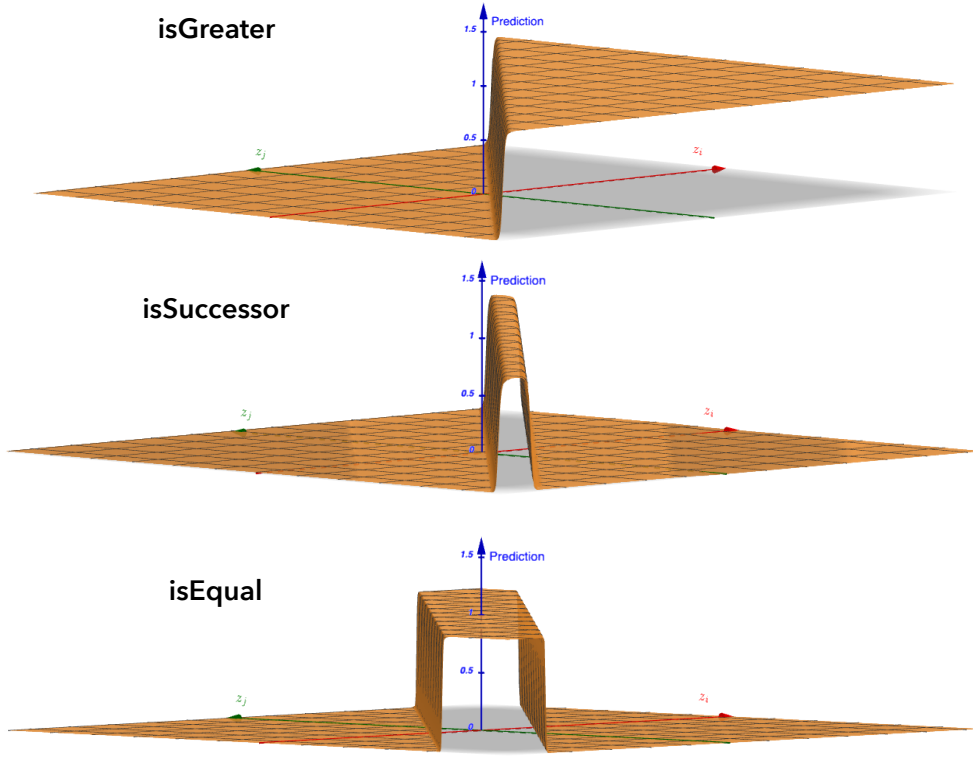
Figure 4: Visualisation of **DC** parameters given by Table 1. $x$- and $y$-axis (red and green) correspond with $z_i^6$ and $z_j^6$ respectively, whilst the $z$-axis (blue) gives the relation-decoder output. We can see that isEqual is learned as a symmetric function, whilst isGreater and isSuccessor are asymmetric. For isGreater, we can see that non-zero $b_\ddagger$ ensures that $f_{\text{isGreater}}^{DC}(z_i^6, z_j^6) \approx 0$ if $z_i^6 = z_j^6$. In terms of transitivity, we can see that isGreater is globally transitive whilst isEqual and isSuccessor can only model localised forms of transitivity.



Figure 5: Example of the $\boldsymbol{u} = \texttt{Softmax}(\boldsymbol{U})$ masks learned by **DC** on *MNIST* for context $\beta$-VAE + {Eq, Gr, Su}. For each relation, **DC** relation-decoders end up using the same latent subspace, namely $z^6$.

9

Table 2: Relations relevant to the 'numeric' semantic factor associated with *MNIST* and their symbolic properties. Other than recursion, all are covered by the proposed relation-decoder.

| Relation | Arity | Symbolic Attributes |
|---|---|---|
| $\mathsf{isZero}(x_i), \ldots, \mathsf{isNine}(x_i)$ | unary | classifier |
| $\mathsf{isEqual}(x_i, x_j)$ | binary | symmetrical, transitive |
| $\mathsf{isGreater}(x_i, x_j), \mathsf{isLess}(x_i, x_j)$ | binary | asymmetrical, transitive |
| $\mathsf{isSuccessor}(x_i, x_j), \mathsf{isPredecessor}(x_i, x_j)$ | binary | asymmetrical, non-transitive |
| $\mathsf{isEven}(x_i), \mathsf{isOdd}(x_i)$ | unary | recursively defined |

ordering, $\mathsf{isSuccessor}$ includes an offset offset via $\boldsymbol{b}_\dagger \neq 0$ and sets a narrow channel-width via $\eta_1$. This leads to $\mathsf{isSuccessor}(z_i^6, z_j^6) \implies \mathsf{isEqual}(z_i^6, z_j^6 - \boldsymbol{b}_\dagger)$. For each relation, decision thresholds are set to be steep, using $\eta_0$ for $\mathsf{isEqual}$ and $\mathsf{isSuccessor}$, and $\eta_2$ for $\mathsf{isGreater}$.

An important nuance between each relation is their strictness over transitivity. Whilst $f_r^\ddagger$ will produce 'global' transitivity since it outputs 1 for all $z_i^6 > z_k^6$. On the other hand, by incorporating a distance measure, it is important that all triples that are true under $f_r^\dagger$ will maintain a maximum distance between head and tail within the $\eta_1$ channel width. It therefore makes sense that $\mathsf{isSuccessor}$ learns a small $\eta_1$ - if it was larger, it is possible that $\mathsf{isSuccessor}$ would demonstrate transitivity across two digits. We refer to this as 'localised-transitivity' since transitivity will depend on the distances between inputs included in any transitive clause.

Finally, to evaluate how well DC can learn a common mask between all semantically related relations, when we include $K > 1$ ground truth factors, Figure 6 shows the learned masks over a set of factors for dSprites. We can clearly see that the latent space is divided between each factor, such that each semantically related relation-decoder calculates any truth-values using only the corresponding subspace. **DC** relation-decoders are trained for binary relations: {isSameX, isRight, isLeft, isSameY, isAbove, isBelow, isSameScale, isBigger, isSmaller, isSameShape} and unary relations {isHeart, isOval, isSquare}. Unary relations are learned by setting $\boldsymbol{z}_j = \boldsymbol{0}$ and $\boldsymbol{a} = [1 \quad 0]$. We indeed see that each set of semantically similar set of relations learn similar masks.

## A.2 Limitations of DC with respect to symbolic attributes of relations

As shown by Table 2, different relations can have different symbolic properties. This is important to consider when defining an LFM relation-decoder, since each relation-decoder will define a region of input-space which outputs true/false [28] - we require that these various geometric spaces, which are themselves the result of the relation-decoder parameterization, can accommodate the different symbolic attributes. The proposed Dynamic Comparator relation-decoder (Eqn. 1) can accommodate the majority of Table 2, where for simple unary attribute relations, we set $\boldsymbol{z}_j = \boldsymbol{0}$. However, it cannot model recursively defined relations such as,

$$\mathsf{isOdd}(a) \implies \mathsf{isPredecessor}(b, a) \wedge \mathsf{isEven}(b).$$

Whilst this can be learned as a true false step function on a separate latent dimension (using $f^\ddagger$), we cannot learn both $\mathsf{isEven}/\mathsf{isOdd}$ such that they share a common latent subspace with $\mathsf{isEqual}$. This is a core motivation of **DC**. We have attempted to use periodic function components to model recursively defined relations, as this could accommodate same-dimension encoding, but in practise training became unstable.

## A.3 Latent space structure comparison with NTN+

This section exemplifies the latent embedding structure changes induced when using each relation-decoder. Firstly, we present MIG scores on *MNIST* for each relation-decoder and, for further perspective, when no semi-supervision is used. Table 3 shows the MIG scores for each context and relation-decoder pairing, when using different values of $\beta$ and Figure 7 presents the normalised mutual information of each latent dimension w.r.t. the digit factor. We include the latter as this provides more insight into how the ground truth factor is being encoded in the latent space. Looking at both of these results, some key observations are: 1. Overall, DC achieves the highest MIG scores and shows the greatest disparity of digit factor information being encoded by each latent dimension. 2. as a general rule, increasing $\beta$ seems to positively influence MIG scores, this is especially the case for DC, where it seems to regularise out digit factor information from all but one dimension at the extreme case; and 3. although there is no great difference in the average MIG scores between NTN+ semi-supervision and no supervision, we do observe a per-dimension increase in digit factor information being encoded across the latent dimensions. In summary, we observe a marked increase in digit factor information being encoded by relation-decoders when included versus when performing fully unsupervised
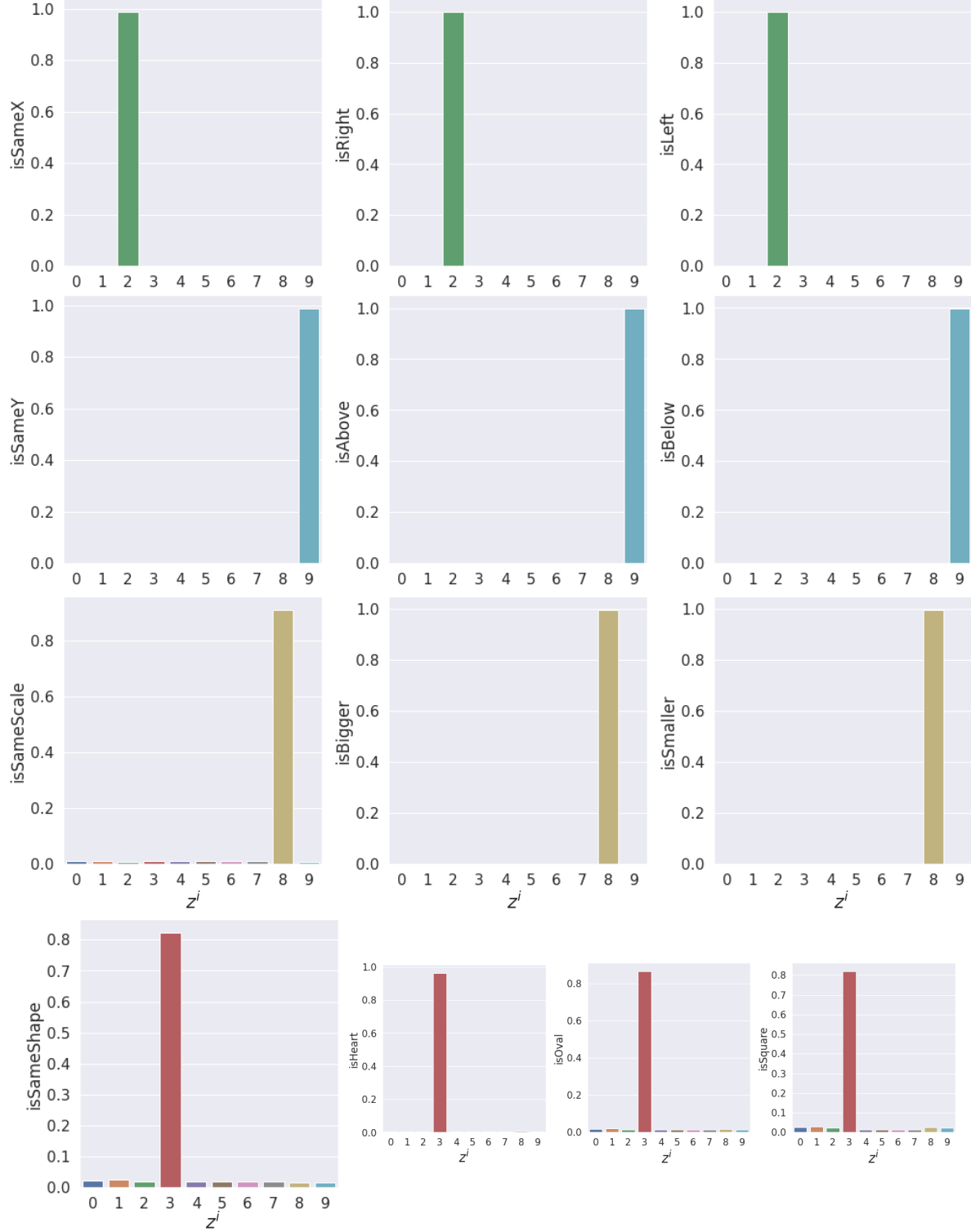
Figure 6: Example of the $\boldsymbol{u} = \mathrm{Softmax}(\boldsymbol{U})$ masks learned by our LFM on *dSprites*. y-axis labels give the relations trained. All 'same' and categorical relations used $a_0 = 1, a_1 = 0$ with $b_\dagger = 0$, meaning they learn a symmetric transitive relations as expected. Comparative relations on the other hand used $a_0 = 0, a_1 = 1$, thereby learning a transitive asymmetrical relation.

Table 3: Mean and standard deviation MIG scores reported for each relation-decoder and context pairing. Results are included for three $\beta$ settings: 4, 8 and 12.

| $f_r$ | $\beta$ | $\beta$-VAE + $\emptyset$ | $\beta$-VAE + Eq | $\beta$-VAE + Eq, Gr | $\beta$-VAE + Eq, Gr, Su |
|---|---|---|---|---|---|
| **-** | 4 | 0.02(0.01) | - | - | - |
|  | 8 | 0.02(0.01) | - | - | - |
|  | 12 | 0.04(0.02) | - | - | - |
| **NTN+** | 4 | - | 0.04(0.03) | 0.03(0.03) | 0.02(0.01) |
|  | 8 | - | 0.02(0.02) | 0.02(0.01) | 0.04(0.02) |
|  | 12 | - | 0.08(0.01) | 0.11(0.02) | 0.03(0.01) |
| **DC** | 4 | - | 0.05(0.02) | 0.15(0.18) | 0.31(0.26) |
|  | 8 | - | 0.07(0.05) | 0.21(0.07) | 0.37(0.06) |
|  | 12 | - | 0.08(0.02) | 0.2(0.11) | 0.53(0.26) |



Figure 7: Violin plots showing the normalised mutual information scores for each latent dimension $z^i$ ($i \in 0, \ldots, 9$) w.r.t. the *MNIST* 'digit' factor. Results have been ordered to improve objective clarity over the common spectra that each context and relation-decoder setting induces. We report results for three $\beta$ settings indicated by their colour codings as follows: 4-blue, 8-orange and 12-green.

12

representation learning. This explains, at least partly, the significant inductive transfer performance increase observed when adding in semi-supervision. However, it is clear that DC is able to extract and disentangle the digit factor information into the fewest latent dimensions, particularly as more context is included for semi-supervision. This is in contrast to the inductive transfer performance, where we observed an evident decrease in performance, but, interestingly, transductive transfer was observed to improve. It therefore seems that this increased regularity w.r.t. the digit factor is beneficial for transductive performance, but not for inductive transfer.

### A.4  Further latent space investigations

In this section we provide further visualisations to substantiate any claims regarding the levels of semantic 'structure' that each relation-decoder induces when employed for $\beta$-VAE semi-supervision. These results aid in understanding the interplay between semantic structure and representation transferability. Figures 8 and 9 present latent dimension versus ground truth digit class scatter plots, for the best/worst performing DC and NTN+ relation-decoder experiments, respectively.

In each case, the mutual information 'spectra' across the latent dimensions are also included to demonstrate how digit class disentanglement affects transfer performance. In summary, we again see that the improved disentanglement of the digit factor obtained by DC positively improves transductive transfer, but has negatively affects on inductive transfer.

For completeness, Figures 12 and 13 present visualisations of each latent dimension w.r.t. each ground truth factor, on *dSprites* when training on same relations as in Figure 6. We can see that DC produces clearer correlations with each ground truth factor than NTN+.
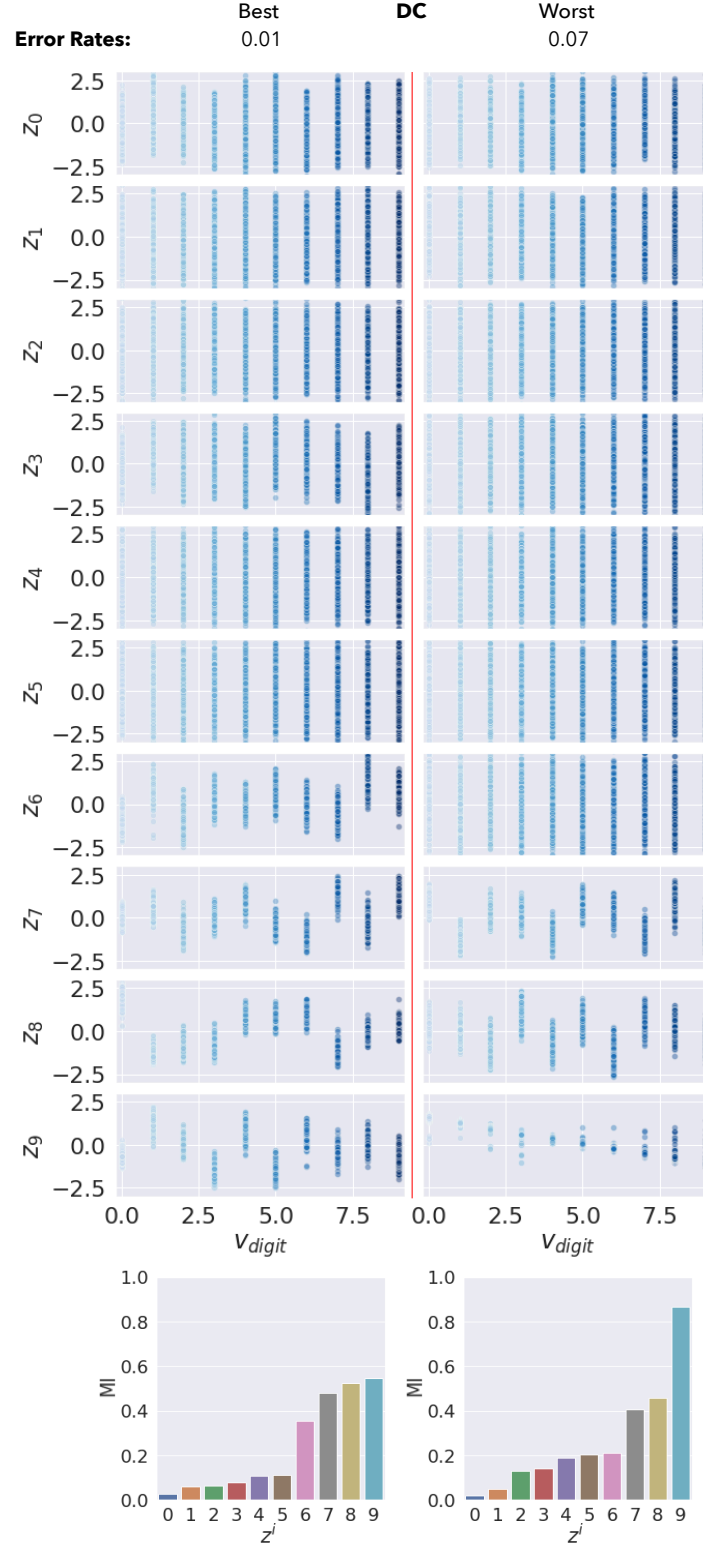
Figure 8: Latent dimension versus digit factor visualisations (top) and normalised mutual information (bottom) when applying semi-supervision with a DC relation-decoder, presented for the best (left) and worse (right) case inductive transfer experiments. Numeric performance values are given in the top row.
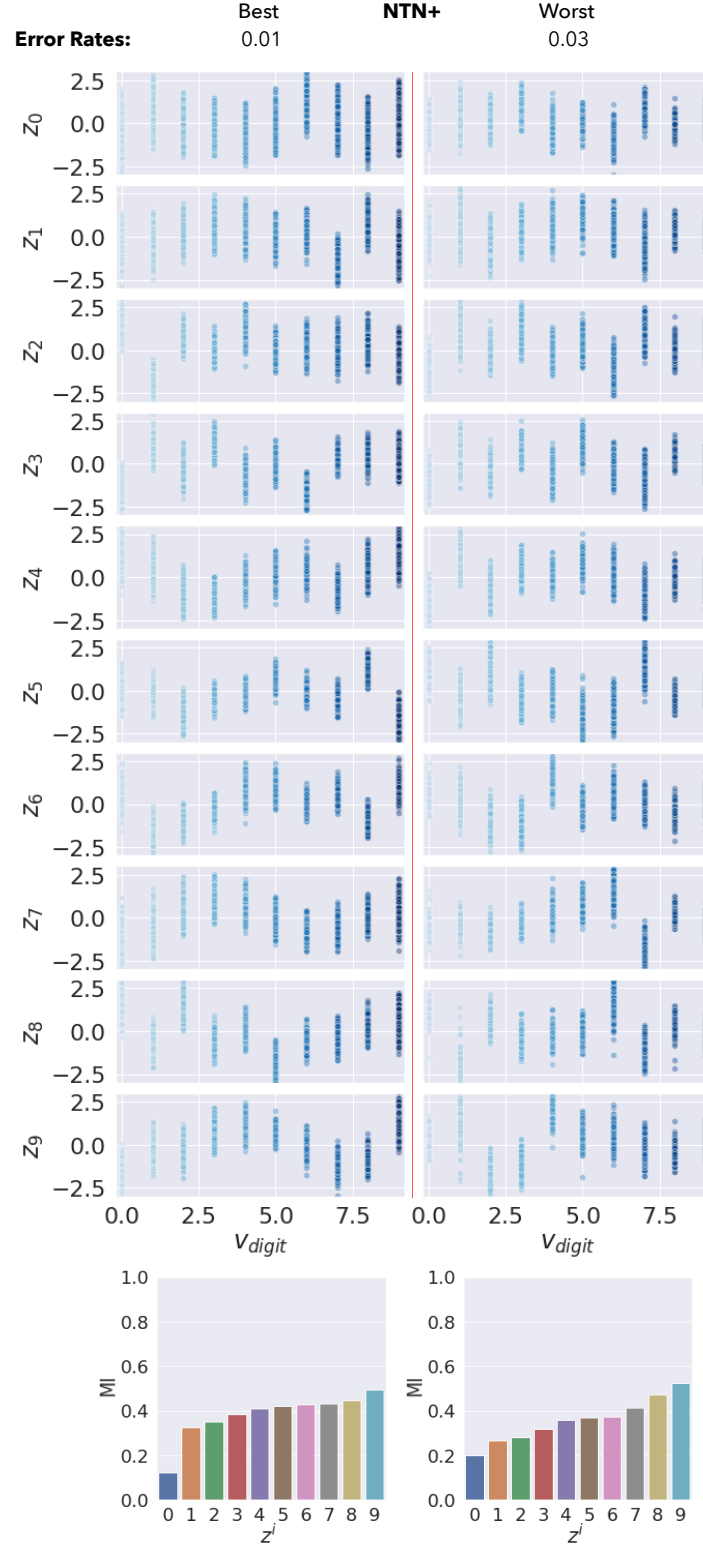
Figure 9: Latent dimension versus digit factor visualisations (top) and normalised mutual information (bottom) when applying semi-supervision with a NTN+ relation-decoder,, presented for the best (left) and worse (right) case inductive transfer experiments. Numeric performance values are given in the top row.
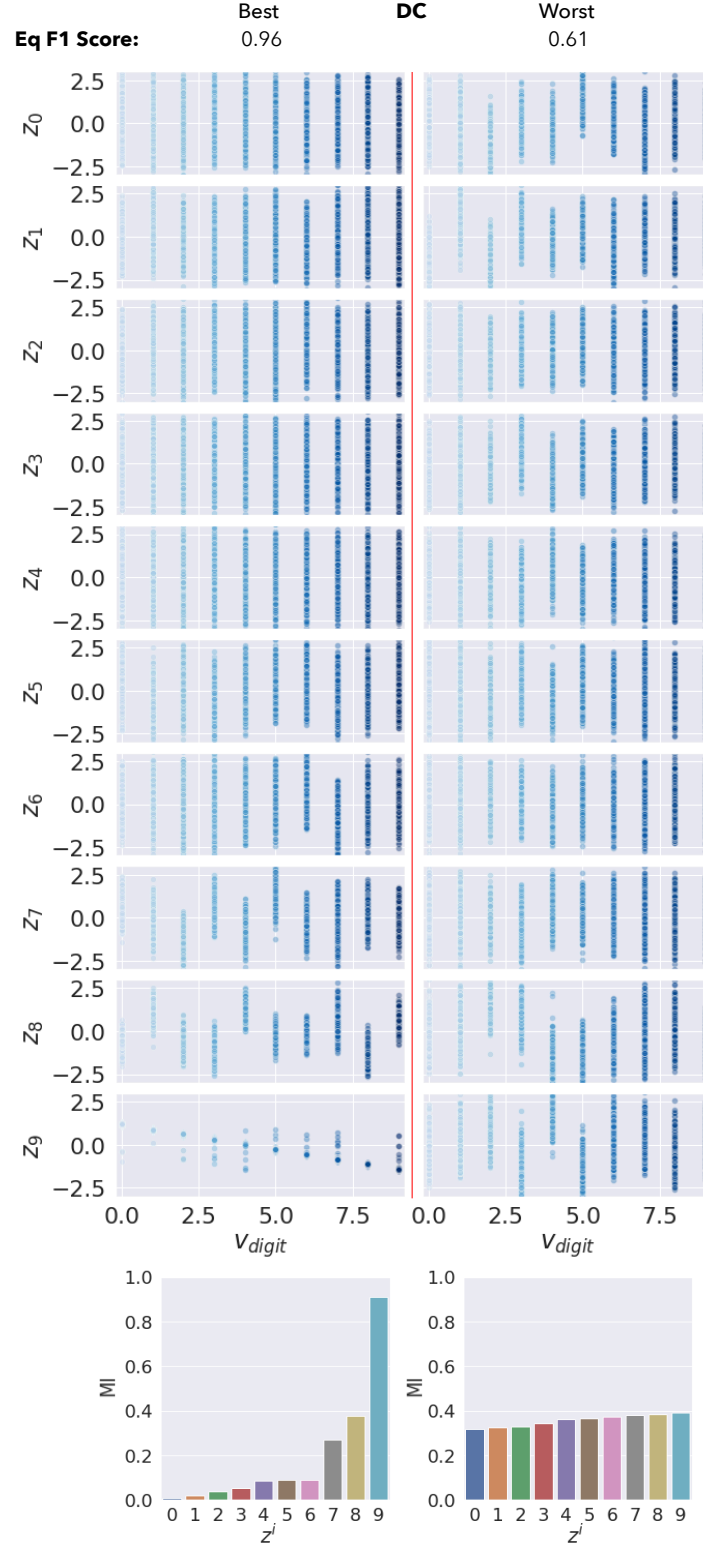
Figure 10: Latent dimension versus digit factor visualisations (top) and normalised mutual information (bottom) when applying semi-supervision with a DC relation-decoder, presented for the best (left) and worse (right) case transductive transfer experiments. Numeric performance values are given in the top row.

Figure 11: Latent dimension versus digit factor visualisations (top) and normalised mutual information (bottom) when applying semi-supervision with a NTN+ relation-decoder, presented for the best (left) and worse (right) case transductive transfer experiments. Numeric performance values are given in the top row.
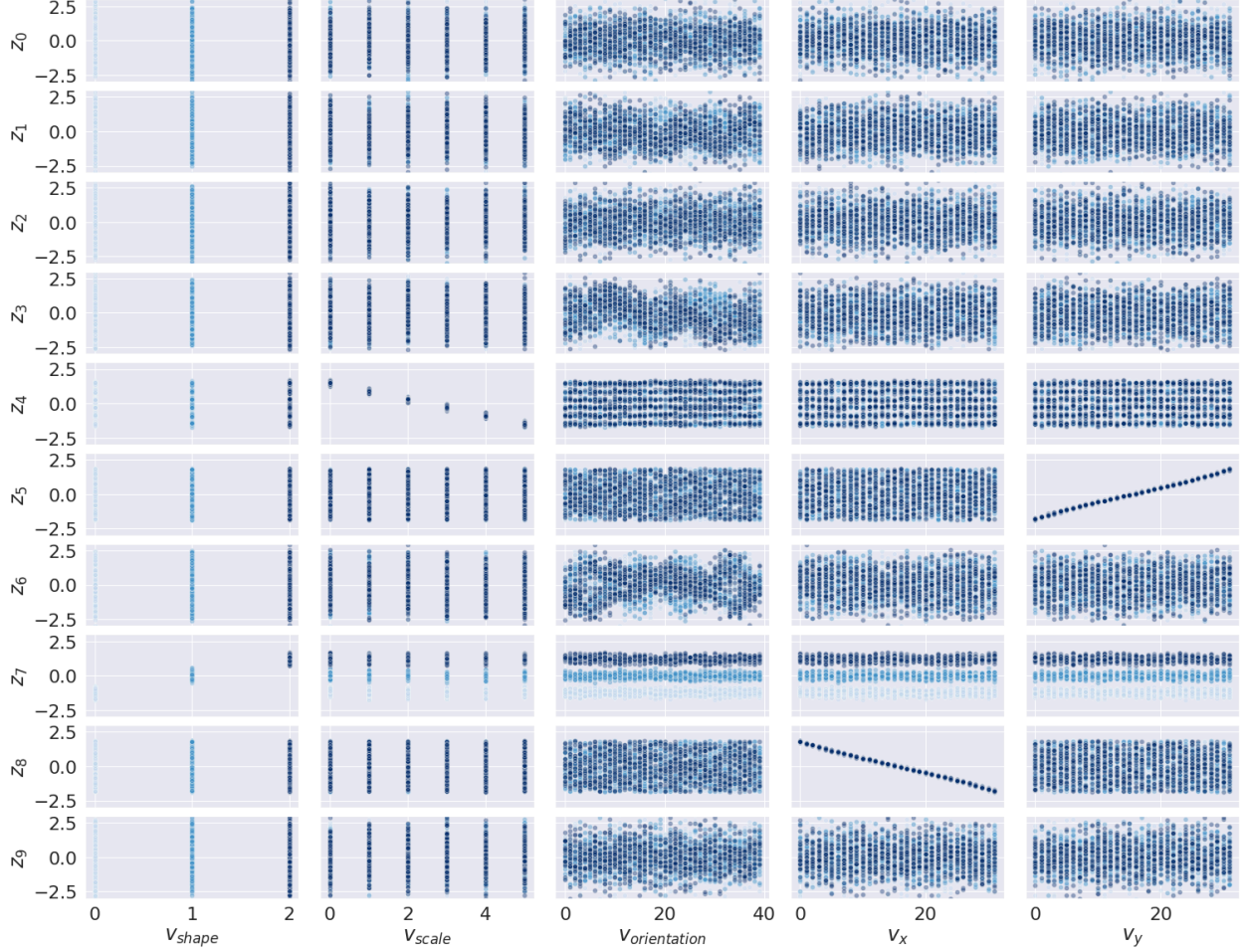
Figure 12: Latent dimension versus ground truth factor visualisations for *dSprites*, when using a DC relation-decoder and training on each relation shown by Figure 6.
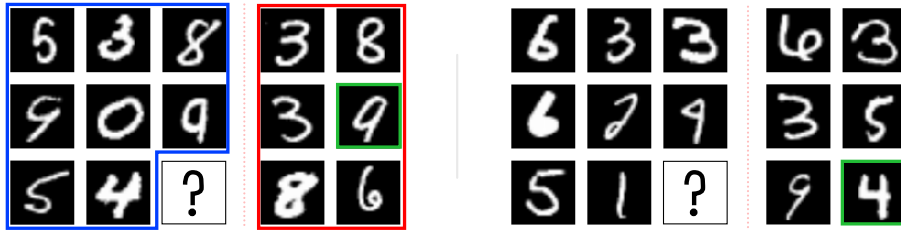


Figure 14: Two examples of our RPM-like *MNIST*-arithmetic tasks, where a reasoner must identify and perform addition (left) or subtraction (right) over *MNIST* digits. Each RPM instance is constructed by a set of tiles corresponding to the context (blue), question (**?**), answer set (red) and answer (green)

## B   Abstract Reasoning over Non-Visual Semantics

This section provides further details regarding the inductive task setup used in the main text.

Recently, abstract reasoning tasks inspired by Raven's Progressive Matrices (RPM) - a well-established measure of non-verbal intelligence [29] - have been used to illustrate the generalisation capability of disentangled representations [25, 26]. Neural networks designed specifically for the RPM task were used in Barrett et al. [27]. In an RPM task, the
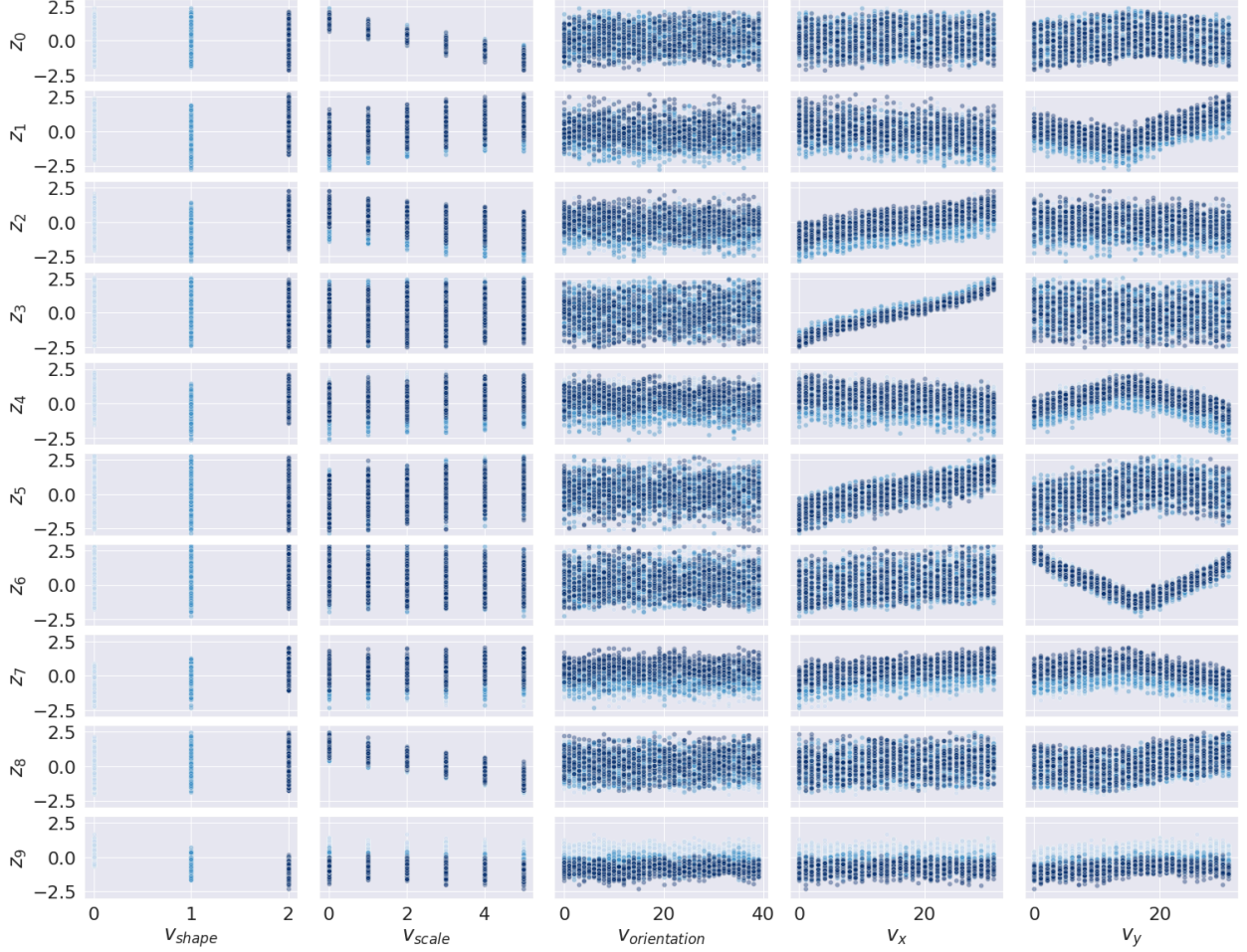
Figure 13: Latent dimension versus ground truth factor visualisations for *dSprites*, when using a NTN+ relation-decoder and training on each relation shown by Figure 6.

learner is presented with a panel made of sequences of context tiles, following one or more relational consistencies, and a final question tile. The learner is tasked with identifying the underlying consistency patterns by selecting the missing question tile from a set of possible answers. This testing mechanism is useful here as it tests the learner's ability to identify semantic relationships in the data [30] and it is thought that a disentangled representation should express clearly the relevant relationships upon which the reasoning behind the RPM task is to be constructed. We construct an RPM-inspired task based on the recent work using *MNIST* data for arithmetic [31, 32]. For each RPM panel, the downstream learner is required first to identify the mathematical operation being performed, and then to apply it to the final panel in order to select the correct answer. See Figure 14 for an example.

**Reasoning Model:** In order to compute the downstream reasoning task, we use a "Wild Relational Network" (WReN), a purpose-built architecture designed for RPM tasks [27, 33]. A WReN leverages a previous relational architecture [33] in order to compute pairwise interaction embeddings for each context-to-context tile pairing and context-to-answer. The model uses two shared neural networks, $g_\theta$ and $f_\phi$, one for interaction learning and another for overall scoring, which ensures that the same reasoning method is performed for each possible answer. In the standard WReN, tile image embeddings are acquired by way of a CNN feature extraction module. However, it was found that a VAE-obtained disentangled representation can lead to improvements, for example on sample complexity [25, 26]. In this paper, we follow the same procedure: we compare representations obtained from a VAE trained with and without a semi-supervision on our *MNIST* RPM task evaluated using a WReN downstream reasoner but with fewer parameters than was used in [26]. WReN model details are provided in Section D.

## C   Background Theory

**Variational AutoEncoders (VAE) -** The VAE is derived by introducing an approximate posterior $q_\phi(\boldsymbol{Z}|\boldsymbol{X})$, from which a lower bound (commonly referred to as the Evidence LOwer Bound (ELBO)) on the true marginal $\ln p_\theta(\boldsymbol{X})$ can be obtained by using Jensen's inequality [34]. The VAE maximises the log-probability by maximising this lower bound, given by:

$$\mathcal{L}^{ELBO}_{\beta\text{-VAE}} = \mathbb{E}_{q_\phi(\boldsymbol{Z}|\boldsymbol{X})}[\log p_\theta(\boldsymbol{X}|\boldsymbol{Z})] - \beta D_{KL}(q_\phi(\boldsymbol{Z}|\boldsymbol{X})\|p_\theta(\boldsymbol{Z})), \tag{3}$$

where $q_\phi(\boldsymbol{Z}|\boldsymbol{X})$ is the approximate posterior, typically modelled as a neural network encoder with parameters $\phi$. Similarly $p_\theta(\boldsymbol{X}|\boldsymbol{Z})$ is modelled as a decoder with parameters $\boldsymbol{\theta}$ and is calculated as a Monte Carlo estimation. A reparameterization trick is used to enable differentiation through this term (see [34]). In the $\beta$-VAE [3, 24], an additional $\beta$ scalar hyperparameter was added as it was found to influence disentanglement through stronger distribution matching pressure with the isotropic zero-mean Gaussian prior $p_\theta(\boldsymbol{Z})$. When $\beta = 1$ we obtain the standard VAE objective [34].

**Latent Factor Models (LFM) -** LFMs are a technique for knowledge graph embedding, where latent representations for data are learned by jointly optimising for them alongside parameterized relation-decoders. These methods are often applied to link prediction, where data that hold similar relations will have similar latent representations, and thus computing held out relations on the entities should produce the correct true/false scoring (see [12, 13, 14] for further details). The importance of LFMs is that the relation-decoders parameters and latent representations together provide the semantics. However, unlike in disentanglement, attention is rarely given to the semantic value of each dimension of the latent space.

**NTN+ -** In this paper, we use a modified Neural Tensor Network (NTN) back-end [15] given by:

$$f_r(\boldsymbol{z}_0, \ldots, \boldsymbol{z}_n) = \sigma(f_r^{\text{NTN+}}(\boldsymbol{z}')), \tag{4}$$

$$f_r^{\text{NTN+}}(\boldsymbol{z}') = \boldsymbol{u}_r^\top[\tanh(\boldsymbol{z}'^\top \boldsymbol{M}_r \boldsymbol{z}' + \boldsymbol{V}_r \boldsymbol{z}_c + \boldsymbol{b}_r)]$$

$$\boldsymbol{z}_c \in \mathbb{R}^{nm}, \boldsymbol{M}_r \in \mathbb{R}^{k \times nm \times nm}; \boldsymbol{V}_r \in \mathbb{R}^{k \times nm}; \boldsymbol{b}_r \in \mathbb{R}^k, \boldsymbol{u}_r \in \mathbb{R}^k,$$

where, $\sigma$ is a sigmoid function used to bound the output of $f_r^{NTN+}$ to $[0,1]$ (interpreted as the truth-value of a predicate in a many-valued logic) and $\boldsymbol{z}' = (\boldsymbol{z}_0; \cdots; \boldsymbol{z}_n)$ is a concatenation of the relation-decoder's arguments with $m$-dimensional latent embeddings $\boldsymbol{z}_0, \cdots, \boldsymbol{z}_n \in \mathbb{R}^m$. The original NTN does not apply the sigmoid nor the concatenation operation, since it was strictly defined for binary relations, whereas NTN+ can accommodate $n$-ary relations and can model boolean relations without requiring any additional arbitrary true/false thresholding. The only hyperparameter to consider is $k$ which controls the model's capacity [15] - in all experiments, we set this to 1.

## D Model Details

Our experiments were implemented using PyTorch [35]. For all models, we use an Adam optimiser with the same parameters for all relation-decoders and VAE models. These are: learning rate of 0.0001, betas $= (0.9, 0.999)$, $\epsilon = 1 \times 10^{-8}$. No weight decay is used. In all experiments, we repeat hyperparameter configurations with 5 restarts. Furthermore, for *MNIST*, all datasets are pre-sampled and shared across experiments. These include a $60,000 : 10,000$ train and test *MNIST* data split, with corresponding knowledge graphs and RPM datasets. For any *dSprites* experiments, we randomly produce a 8:2 train-test split and produce triples from combinations of the inputs in each sampled image batch. In all experiments (both *MNIST* and *dSprites*), we use an image batch size of 64.

### D.1 VAE configuration

In all representation learning experiments, we use a $\beta$-VAE trained for 300,000 steps, following accepted practise from [10, 25]. The encoder-decoder model parameters are given in Table 4 - we include the model configurations used for both *MNIST* and *dSprites* datasets.

Table 4: Specification of our $\beta$-VAE encoder and decoder model parameters, for both $28\times28$ (top) and $64\times64$ (bottom) size input data. I: Input channels, O: Output channels, K: Kernel size, S: Stride, P: Padding, A: Activation

| Encoder | Decoder |
|---|---|
| Input: $28 \times 28 \times N_C$ | Input: $\mathbb{R}^{10}$ |
| **Layer_ID ; I ; O ; K ; S ; P ; A** | **Layer_ID ; Num Nodes : In - Out ; A** |
| Conv2d_1 ; $N_C$ ; 32 ; $4 \times 4$ ; 2 ; 1 ; ReLU | FC_z ; 10 - 144 ; ReLU |
| Conv2d_2 ; 32 ; 32 ; $4 \times 4$ ; 2 ; 1 ; ReLU | FC_z_mu ; 144 - 576 ; ReLU |
| Conv2d_3 ; 32 ; 64 ; $3 \times 3$ ; 2 ; 1 ; ReLU | |
| Conv2d_4 ; 64 ; 64 ; $2 \times 2$ ; 2 ; 1 ; ReLU | **Layer_ID ; I ; O ; K ; S ; P ; A** |
| | UpConv2d_1 ; 64 ; 64 ; $2 \times 2$ ; 2 ; 1 ; ReLU |
| **Layer_ID ; Num Nodes : In - Out ; A** | UpConv2d_1 ; 64 ; 32 ; $3 \times 3$ ; 2 ; 1 ; ReLU |
| FC_z ; 576 - 144 ; ReLU | UpConv2d_1 ; 32 ; 32 ; $4 \times 4$ ; 2 ; 1 ; ReLU |
| FC_z_mu ; 144 - 10 ; None | UpConv2d_1 ; 32 ; $N_C$ ; $4 \times 4$ ; 2 ; 1 ; Sigmoid |
| FC_z_logvar ; 144 - 10 ; None | |

| Encoder | Decoder |
|---|---|
| Input: $64 \times 64 \times N_C$ | Input: $\mathbb{R}^{10}$ |
| **Layer_ID ; I ; O ; K ; S ; P ; A** | **Layer_ID ; Num Nodes : In - Out ; A** |
| Conv2d_1 ; $N_C$ ; 32 ; $4 \times 4$ ; 2 ; 1 ; ReLU | FC_z ; 10 - 256 ; ReLU |
| Conv2d_2 ; 32 ; 32 ; $4 \times 4$ ; 2 ; 1 ; ReLU | FC_z_mu ; 256 - 1024 ; ReLU |
| Conv2d_3 ; 32 ; 64 ; $4 \times 4$ ; 2 ; 1 ; ReLU | |
| Conv2d_4 ; 64 ; 64 ; $4 \times 4$ ; 2 ; 1 ; ReLU | **Layer_ID ; I ; O ; K ; S ; P ; A** |
| | UpConv2d_1 ; 64 ; 64 ; $4 \times 4$ ; 2 ; 1 ; ReLU |
| **Layer_ID ; Num Nodes : In - Out ; A** | UpConv2d_1 ; 64 ; 32 ; $4 \times 4$ ; 2 ; 1 ; ReLU |
| FC_z ; 1024 - 256 ; ReLU | UpConv2d_1 ; 32 ; 32 ; $4 \times 4$ ; 2 ; 1 ; ReLU |
| FC_z_mu ; 256 - 10 ; None | UpConv2d_1 ; 32 ; $N_C$ ; $4 \times 4$ ; 2 ; 1 ; Sigmoid |
| FC_z_logvar ; 256 - 10 ; None | |

### D.2 WReN module

For abstract reasoning tasks, we follow the setup of [25], training for 100,000 steps and a batch size of 32, and testing 100 RPM new samples after each 1000 steps. As in [25], we ensure that each RPM instance is new, meaning that the training set consists of $3.2\times10^6$ samples. In this paper, each RPM panel consists of eight context tiles $C = \{c_1, \ldots, c_8\}$ with six possible answer tiles $A = \{a_1, \ldots, a_6\}$. For each possible answer, the WReN reasoning module computes:

$$WReN(a_k, C) = f_\phi\big( \sum_{\boldsymbol{z}_i, \boldsymbol{z}_j \in Z} g_\theta(h_\gamma(\boldsymbol{z}_i), h_\gamma(\boldsymbol{z}_j))\big), \tag{5}$$

$$Z = \{\psi(c_1), \ldots \psi(c_8) \cup \psi(a_k)\}, \quad \psi(\cdot) = CNN(\cdot)\,(\vee\,\psi_{enc}(\cdot)), \tag{6}$$

where $f_\theta$ is a scoring multilayer perceptron (MLP), which takes as input an aggregation over inter-tile interactions, as computed by the relation network function $g_\theta$. As in [25, 26], instead of using a CNN feature extractor as in the original

model, we replace the initial CNN feature extractor [27] with pre-trained representations (*i.e.* $z_i, z_j$) taken from the VAE bottleneck as extracted by the VAE encoder: $\psi_{enc}$. Finally, $h_\gamma$ serves the purpose of incorporating positional features into each tile, where each tile has its feature vector concatenated with a one-hot position vector. This is then passed through a single-layer fully connected MLP to obtain each tiles' feature vector. Note that all answer tiles are considered to be at the final position (position 9) of the panel. See [27, 33] for more information on the overall WReN architecture. The WReN model parameters are provided in Table 5, where the reduced parameters sizes, used in the paper's main text experiments, are given in parenthesis.

Table 5: Specification of our the WReN model parameters. A: Activation

**Scoring function** $f_\phi$
Input: $28 \times 28 \times N_C$

| Layer_ID ; Num Nodes : In - Out ; A |
| --- |
| FC_1 ; 64 ; 64 ; ReLU |
| FC_2 ; 64 ; 64 ; ReLU |
| – Drop-out layer $p = 0.5$ – |
| FC_4 ; 64 ; 1 ; None |

**Relation net** $g_\theta$
Input: concatenation $[\mathbb{R}^{64}; \mathbb{R}^{64}]$

| Layer_ID ; Num Nodes : In ; Out ; A |
| --- |
| FC_1 ; 2*64 ; 128 ; ReLU |
| FC_2 ; 128 ; 128 ; ReLU |
| FC_3 ; 128 ; 128 ; ReLU |
| FC_4 ; 128 ; 64 ; ReLU |

**Tile-position feature projection** $h_\gamma$
Input: $\mathbb{R}^{10} + \mathbb{R}^9$ position one-hot vector

| Layer_ID ; Num Nodes : In ; Out ; A |
| --- |
| FC_1 ; 10+9 ; 64 ; ReLU |