

---

# Topographic VAEs learn Equivariant Capsules

---

Anonymous Author(s)

Affiliation  
Address  
email

## Abstract

In this work we seek to bridge the concepts of topographic organization and equivariance in neural networks. To accomplish this, we introduce the Topographic VAE: a novel method for efficiently training deep generative models with topographically organized latent variables. We show that such a model indeed learns to organize its activations according to salient characteristics such as digit class, width, and style on MNIST. Furthermore, through topographic organization over time (i.e. temporal coherence), we demonstrate how predefined latent space transformation operators can be encouraged for observed transformed input sequences – a primitive form of unsupervised learned equivariance. We demonstrate that this model successfully learns sets of equivariant features (i.e. "capsules") directly from sequences and achieves higher likelihood on correspondingly transforming test sequences. Equivariance is verified quantitatively by measuring the commutativity of the inference network and the sequence transformations. Finally, we show the model is capable of learning to be approximately equivariant to complex transformations, expanding upon the capabilities of existing group equivariant neural networks.

## 1 Introduction

Many parts of the brain are organized topographically. Famous examples are the ocular dominance maps and the orientation maps in V1. What is the advantage of such organization and what can we learn from it to develop better inductive biases for deep neural network architectures?

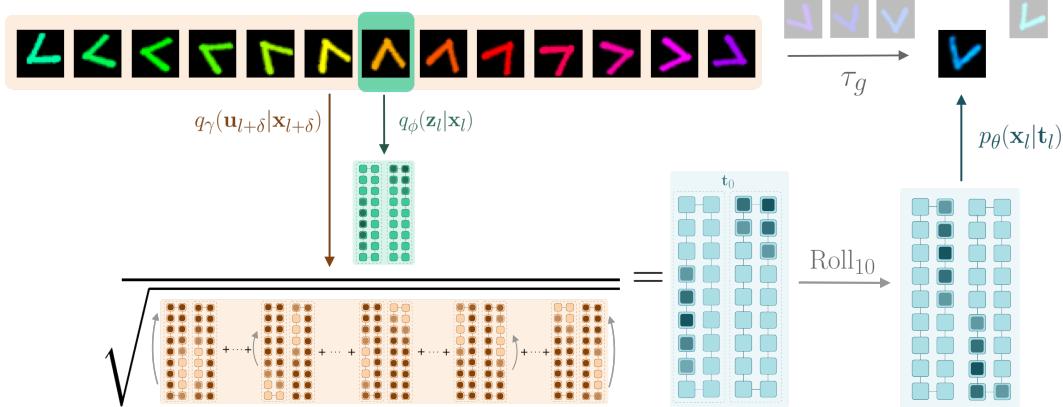


Figure 1: Overview of the Topographic VAE with shifting temporal coherence. The combined color/rotation transformation in input space  $\tau_g$  becomes encoded as a Roll within the equivariant capsule dimension. The model is thus able decode unseen sequence elements by encoding a partial sequence and rolling activations within the capsules. We see this completes a commutative diagram.

20 One potential explanation for the emergence of topographic organization is provided by the principle  
21 of redundancy reduction [1]. In the language of Information Theory, redundancy wastes channel  
22 capacity, and thus to represent information as efficiently as possible, the brain may strive to transform  
23 the input to a neural code where the activations are statistically maximally independent. In the machine  
24 learning literature, this idea resulted in Independent Component Analysis (ICA) which linearly  
25 transforms the input to a new basis where the activities are independent and sparse [2, 12, 28, 41]. It  
26 was soon realized that there are remaining higher order dependencies (such as correlation between  
27 absolute values) that can not be transformed away by a linear transformation. For example, along  
28 edges of an image, linear-ICA components (e.g. gabor filters) still activate in clusters even though the  
29 sign of their activity is unpredictable. This led to new algorithms that explicitly model these remaining  
30 dependencies through a topographic organization of feature activations [26, 42, 43, 51]. Such  
31 topographic models were shown to better match the generative structure of the data, yielding more  
32 efficient representations. Furthermore, the topographically organized features were reminiscent of  
33 pinwheel structures observed in V1, encouraging multiple comparisons with topographic organization  
34 in the biological visual system [27, 29, 39].

35 A second, almost independent body of literature developed the idea of “equivariance” of neural  
36 network feature maps under symmetry transformations. The idea of equivariance is that symmetry  
37 transformations define equivalence classes as the orbits of their transformations, and we wish to  
38 maintain this structure in the deeper layers of a neural network. For instance, for images, asserting a  
39 rotated image contains the same object for all rotations, the transformation of rotation then defines  
40 an orbit where the elements of that orbit can be interpreted as pose or angular orientation. When an  
41 image is processed by a neural network, we want features at different orientations to be able to be  
42 combined to form new features, but we want to ensure the relative pose information between the  
43 features is preserved for all orientations. This has the advantage that the equivalence class of rotations  
44 for the complex composite features is guaranteed to be maintained, allowing for the extraction of  
45 invariant features, a unified pose, and increased data efficiency. Such ideas are reminiscent of the  
46 capsule networks of [20, 21, 46], and indeed formal connections to equivariance have been made  
47 [38]. Interestingly, by explicitly building neural networks to be equivariant, we additionally see  
48 geometric organization of activations into these equivalence classes. The insight of this connection  
49 between topographic organization and equivariance hints at a possibility to encourage approximate  
50 equivariance from an induced topology in feature space.

51 To build a model, we need to ask what mechanisms could induce topographic organization of trans-  
52 formations in biological and artificial neural networks? We have argued that removing dependencies  
53 between latent representations of input data is a possible mechanism, and indeed, as discussed,  
54 there is a rich literature on how to achieve that unsupervised. However, for the more structured  
55 organisation of equivariant capsule representations, the usual approach is to hard-code this structure  
56 into the network, or to encourage it through regularization terms [4, 14]. To achieve this through  
57 *unsupervised learning*, we propose to incorporate another key inductive bias: “temporal coherence”  
58 [17, 23, 48, 52]. The principle of temporal coherence, or “slowness”, asserts that when processing  
59 correlated sequences, we wish for our representations to change smoothly and slowly over space  
60 and time. Thinking of time sequences as symmetry transformations on the input, we desire features  
61 undergoing such transformations to be grouped into equivariant capsules. We therefore suggest that  
62 encouraging slow feature transformations to take place *within a capsule* could induce such grouping  
63 from sequences alone.

64 In the following sections we will explain the details of our Topographic VAE model which lies at  
65 the intersection of topographic organization, equivariance, and temporal coherence, thereby learning  
66 approximate equivariant capsules from sequence data completely unsupervised.

## 67 2 Related Work

68 The history of statistical models upon which this work builds is vast, including sparse coding [41],  
69 Independant Component Analysis (ICA) [2, 12, 28], Slow Feature Analysis (SFA) [49, 52], and  
70 many variants [25, 26]. Most related to this work are topographic generative models including  
71 Generative Topographic Maps [5], Bubbles [24], Topographic ICA [26], and the Topographic Product  
72 of Student’s-t [43, 51]. Prior work on learning equivariant and invariant representations is similarly  
73 vast and also has a deep relationship with these generative models. Specifically, Independant Subspace  
74 Analysis [25], models involving temporal coherence [17, 23, 48, 52], and Adaptive Subspace Self

75 Organizing Maps [33] have all demonstrated the ability to learn invariant feature subspaces and even  
 76 ‘disentangle’ space and time [18]. Our work assumes a similar generative model to these works while  
 77 additionally allowing for efficient estimation of the model through variational inference [31, 45].  
 78 Although our work is not the first to combine Student’s-t distributions and variational inference [6], it  
 79 is the first to provide an efficient method to do so for Topographic Student’s-t distributions.  
 80 Another line of work has focused on constructing neural networks with equivariant representations  
 81 separate from the framework of generative modeling. Analytically equivariant networks such as Group  
 82 Equivariant Neural Networks [11], and other extensions [54, 9, 53, 15, 44, 50, 16] propose to explicitly  
 83 enforce symmetry to group transformations in neural networks through structured weight sharing.  
 84 Other methods propose supervised and self-supervised methods to learn the equivariant or invariant  
 85 structure from pairs of examples or transformations [4, 14, 13]. One related example in this category  
 86 uses a group sparsity regularization term to similarly learn topographic features for the purpose of mod-  
 87 eling invariance [30]. We believe the unsupervised approach presented in this paper provides a valuable  
 88 alternative, principled view of how such structure could be learned in biological neural networks.  
 89 Furthermore, the idea of disentangled representations [3] has also been connected to  
 90 equivariance and representation theory in multiple recent papers [10, 8, 19, 7]. Our work shares a  
 91 fundamental connection to this distributed operator definition of disentanglement, where the slow roll  
 92 of capsule activations can be seen as the latent operator. Recently, the authors of [32] demonstrated  
 93 that incorporating the principle of ‘slowness’ in a variational autoencoders (VAEs) yields the ability  
 94 to learn disentangled representations from natural sequences. While similar in motivation, the  
 95 generative model proposed in [32] is unrelated to topographic organization and equivariance, and  
 96 is more aligned with traditional notions of disentanglement.  
 97 Finally, and importantly, in the neuroscience literature, another popular explanation for topographic  
 98 organization arises as the solution to the ‘wiring length’ minimization problem [34]. Recently,  
 99 models which attempt to incorporate wiring length constraints have been shown to yield topographic  
 100 organization of higher level features, ultimately resembling the ‘face patches’ found in primates  
 101 [37]. Interestingly, the model presented in this paper organizes activity based on the same statistical  
 102 property (local correlation) as the wiring length proxies developed in [37], but from a generative  
 103 modeling perspective, demonstrating a potential generative modeling solution to the same problem.

### 104 3 Background

105 The model in this paper is a first attempt at bridging two yet disjoint classes of models: Topographic  
 106 Generative Models, and Equivariant Neural Networks. In this section, we will provide a brief  
 107 background on these two frameworks.

#### 108 3.1 Topographic Generative models

109 Inspired by Topographic ICA, the class of Topographic Generative models can be understood as gener-  
 110 ative models where the joint distribution over latent variables does not factorize into entirely indepen-  
 111 dent factors, as is commonly done in ICA or VAEs, but instead has a more complex ‘local’ correlation  
 112 structure. The locality is defined by arranging the latent variables into an n-dimensional lattice or grid,  
 113 and organizing variables such that those which are closer together on this grid have greater correlation  
 114 of activities than those which are further apart. In the related literature, activations which are nearby  
 115 in this grid are defined to have higher-order correlation, e.g. correlations of squared activations (aka  
 116 ‘energy’), asserting that all first order correlations are removed by the initial ICA de-mixing matrix.

117 Such generative models can be seen as hierarchical generative models where there exist higher  
 118 level independent ‘variance generating’ variables  $\mathbf{V}$  which are combined locally to generate the  
 119 variances  $\sigma = \phi(\mathbf{W}\mathbf{V})$  of the lower level topographic variables  $\mathbf{T} \sim \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I})$ , for an appropriate  
 120 non-linearity  $\phi$ . The variables  $\mathbf{T}$  are thus independent conditioned on  $\sigma$ . As we will show in the next  
 121 section, the Topographic Product of Student’s-t model can be seen to be a member of this class when  
 122  $\mathbf{V} = \mathbf{U}^2$ ,  $\mathbf{U} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$  and  $\phi(\cdot) = \frac{\sqrt{\nu}}{\sqrt{\cdot}}$ , where  $\nu$  is the number of degrees of freedom.

123 Other related models which can be described under this umbrella include *Independent Subspace*  
 124 *Analysis* (ISA) [25], and certain models which leverage the principle of *temporal coherence* [23].  
 125 ISA can be seen as a special case of a topographic generative model where all variables within a

126 predefined subspace (or ‘capsule’) share a common variance. Temporal coherence can be defined  
 127 as the correlation of energy between time steps for a given variable and has been achieved in  
 128 prior literature [24] by extending the topographic neighborhoods of the TICA model over the time  
 129 dimensions. We will expand upon this idea Section 4.5 to induce latent space operators for observed  
 130 transformations, thereby learning equivariant capsules.

### 131 3.2 Group Equivariant Neural Networks

132 Equivariance is the mathematical notion of symmetry for functions. A function is said to be an  
 133 equivariant map if the result of transforming the input and then computing the function is the same  
 134 as first computing the function and then transforming the output. In other words, the function and  
 135 the transformation commute. Formally,  $f(\tau_\rho[\mathbf{x}]) = \Gamma_\rho[f(\mathbf{x})]$ , where  $\tau$  and  $\Gamma$  denote the (potentially  
 136 different) operators on the domain and co-domain respectively, but are indexed by the same element  $\rho$ .

137 It is well known that convolutional maps in neural networks are translation equivariant, i.e., given  
 138 a translation  $\Gamma_\rho$  (applied to each feature map separately) and a convolutional map  $f(\cdot)$ , we have  
 139  $f(\Gamma_\rho[\mathbf{x}]) = \Gamma_\rho[f(\mathbf{x})]$ . This can be extended to other transformations (e.g. rotation or mirroring)  
 140 using Group convolutions ( $G$ -convolutions) [11]. As a result of the design of  $G$ -convolutions, feature  
 141 maps that are related to each other by a rotation of the filter/input are grouped together. Moreover, a  
 142 rotation of the input results in a transformation (i.e., a permutation and rotation) on the activations of  
 143 each of these groups in the output. Hence, we can think of these equivalence class groups as capsules  
 144 where transformations of the input only cause structured transformations *within* a capsule.

## 145 4 The Generative Model

146 The generative model proposed in this paper is based on the Topographic Product of Student’s-t  
 147 (TPoT) model as developed in [51, 43]. In the following, we will show how a TPoT random variable  
 148 can be constructed from a set of independent univariate standard normal random variables, enabling  
 149 efficient training through variational inference. Subsequently, we will construct a new model where  
 150 topographic neighborhoods are extended over time, introducing temporal coherence and encouraging  
 151 the unsupervised learning of equivariant subspaces we call ‘capsules’.

### 152 4.1 The Product of Student’s-t Model

153 We assume that our observed data is generated by a latent variable model where the joint  
 154 distribution over observed and latent variables  $\mathbf{x}$  and  $\mathbf{t}$  factorizes into the product of the conditional and  
 155 the prior. The prior distribution  $p_{\mathbf{T}}(\mathbf{t})$  is assumed to be a Topographic Product of Student’s-t (TPoT)  
 156 distribution, and we parameterize the conditional distribution with a flexible function approximator:

$$p_{\mathbf{X}, \mathbf{T}}(\mathbf{x}, \mathbf{t}) = p_{\mathbf{X}|\mathbf{T}}(\mathbf{x}|\mathbf{t})p_{\mathbf{T}}(\mathbf{t}) \quad p_{\mathbf{X}|\mathbf{T}}(\mathbf{x}|\mathbf{t}) = p_\theta(\mathbf{x}|g_\theta(\mathbf{t})) \quad p_{\mathbf{T}}(\mathbf{t}) = \text{TPoT}(\mathbf{t}; \nu) \quad (1)$$

157 The goal of training is thus to learn the parameters  $\theta$  such that the marginal distribution of the  
 158 model  $p_\theta(\mathbf{x})$  matches that of the observed data. Unfortunately, the marginal likelihood is generally  
 159 intractable except for all but the simplest choices of  $g_\theta$  and  $p_{\mathbf{T}}$  [42]. Prior work has therefore resorted  
 160 to techniques such as contrastive divergence with Gibbs sampling [51] to train TPoT models as  
 161 energy based models. In the following section, we instead demonstrate how TPoT variables can be  
 162 constructed as a deterministic function of Gaussian random variables, enabling the use of variational  
 163 inference and efficient maximization of the likelihood through the evidence lower bound (ELBO).

### 164 4.2 Constructing the Product of Student’s-t Distribution

165 First, note that a univariate Student’s-t random variable  $T$  with  $\nu$  degrees of freedom can be defined as:

$$T = \frac{Z}{\sqrt{\frac{1}{\nu} \sum_i^{\nu} U_i^2}} \quad \text{with } Z, U_i \sim \mathcal{N}(0, 1) \quad \forall i \quad (2)$$

166 Where  $Z$  and  $\{U_i\}_{i=1}^{\nu}$  are independent standard normal random variables. If  $\mathbf{T}$  is a multidimensional  
 167 Student’s-t random variable, composed of independent  $Z_i$  and  $U_i$ , then  $\mathbf{T} \sim \text{PoT}(\nu)$ , i.e.:

$$\mathbf{T} = \left[ \frac{Z_1}{\sqrt{\frac{1}{\nu} \sum_{i=1}^{\nu} U_i^2}}, \frac{Z_2}{\sqrt{\frac{1}{\nu} \sum_{i=\nu+1}^{2\nu} U_i^2}}, \dots, \frac{Z_n}{\sqrt{\frac{1}{\nu} \sum_{i=(n-1)\nu+1}^{n\nu} U_i^2}} \right] \sim \text{PoT}(\nu) \quad (3)$$

168 Note that the Student's-t variable  $T$  is large when most of the  $\{U_i\}_i$  in its set are small. We can  
 169 therefore think of the  $\{U_i\}_i$  as constraint violations rather than pattern matches: if the input matches  
 170 all constraints  $U_i \approx 0$ , the corresponding  $T$  variables will activate [22].

### 171 4.3 Introducing Topography

172 To make the PoT distribution topographic, we strive to correlate the scales of  $T_j$  which are ‘nearby’ in  
 173 our topographic layout. One way to accomplish this is by *sharing* some  $U_i$ -variables between neighboring  
 174  $T_j$ 's. Formally, we define overlapping neighborhoods  $N(j)$  for each variable  $T_j$  and write:

$$175 \mathbf{T} = \left[ \frac{Z_1}{\sqrt{\frac{1}{\nu} \sum_{i \in N(1)} U_i^2}}, \frac{Z_2}{\sqrt{\frac{1}{\nu} \sum_{i \in N(2)} U_i^2}}, \dots, \frac{Z_n}{\sqrt{\frac{1}{\nu} \sum_{i \in N(n)} U_i^2}} \right] \sim \text{TPoT}(\nu) \quad (4)$$

175 With some abuse of notation, if we define  $\mathbf{W}$  to be the adjacency matrix which defines our  
 176 neighborhood structure,  $\mathbf{U}$  and  $\mathbf{Z}$  to be the vectors of random variables  $U_i$  and  $Z_j$ , we can write  
 177 the above succinctly as:

$$178 \mathbf{T} = \left[ \frac{Z_1}{\sqrt{\frac{1}{\nu} W_1 \mathbf{U}^2}}, \frac{Z_2}{\sqrt{\frac{1}{\nu} W_2 \mathbf{U}^2}}, \dots, \frac{Z_n}{\sqrt{\frac{1}{\nu} W_n \mathbf{U}^2}} \right] = \frac{\mathbf{Z}}{\sqrt{\frac{1}{\nu} \mathbf{W} \mathbf{U}^2}} \sim \text{TPoT}(\nu) \quad (5)$$

178 Due to non-linearities such as ReLUs, it is beneficial to allow the  $Z$  variables to model the mean  
 179 and scale. We found this can be achieved with the following parameterization:  $\mathbf{T} = \frac{\mathbf{Z} - \mu}{\sigma \sqrt{\frac{1}{\nu} \mathbf{W} \mathbf{U}^2}}$ . In  
 180 practice, we found that  $\sigma = \sqrt{\nu}$  often works well, leading to the final simpler form:

$$181 \mathbf{T} = \frac{\mathbf{Z} - \mu}{\sqrt{\mathbf{W} \mathbf{U}^2}} \quad (6)$$

181 Given this construction, we observe that the TPoT generative model can instead be viewed as a latent  
 182 variable model where all random variables are Gaussian and the construction of  $\mathbf{T}$  in Equation 6  
 183 is the first layer of the generative ‘decoder’:  $g_\theta(\mathbf{t}) = g_\theta(\mathbf{u}, \mathbf{z})$ . In Section 5 we then leverage this  
 184 interpretation to show how an approximate posterior for the latent variables  $\mathbf{Z}$  and  $\mathbf{U}$  can be trained  
 185 through variational inference.

### 186 4.4 Capsules as Disjoint Topologies

187 One setting of neighborhood structure  $\mathbf{W}$  which is of par-  
 188 ticular interest is when there exist multiple sets of disjoint  
 189 neighborhoods. Statistically, the variables of two disjoint  
 190 topologies are completely independent. An example of a  
 191 capsule neighborhood structure is shown in Figure 2. The  
 192 idea of independent subspaces has previously been shown  
 193 to learn invariant feature subspaces in the linear setting and  
 194 is present in early work on Independent Subspace Analysis  
 195 [25] and Adaptive Subspace Self Organizing Maps (AS-  
 196 SOM) [33]. It is also very reminiscent of the transformed  
 197 sets of features present in a group equivariant convolu-  
 198 tional neural network. In the next section, we will show  
 199 how temporal coherence can be leveraged to induce the  
 200 encoding of observed transformations into the internal di-  
 201 mensions of such capsules thereby yielding unsupervised  
 202 equivariant capsules.

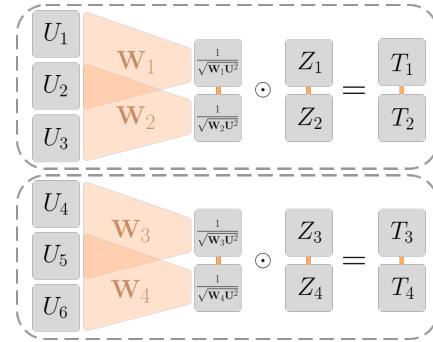


Figure 2: An example of a neighborhood structure which induces disjoint topolo-  
 gies (aka capsules). Lines between variables  $T_i$  indicate that sharing of  $U_i$ , and  
 thus correlation.

### 203 4.5 Temporal Coherence and Equivariance

204 We now describe how the induced topographic organization can be leveraged to learn a basis of  
 205 equivariant capsules for observed transformation sequences. The resulting representation is composed  
 206 of a large set of ‘capsules’ where the dimensions inside the capsule are topographically structured,  
 207 but between the capsules there is independence. To benefit from sequences of input, we encourage  
 208 topographic structure over time between sequentially permuted activations within a capsule, a property  
 209 we refer to as *shifting temporal coherence*. In following we describe how this is formally achieved.

210 **4.5.1 Temporal Coherence**

211 Temporal Coherence can be measured as the correlation of squared activation between time steps.  
 212 One way we can achieve this in our model is by having  $T_j$  share  $U_i$  between time steps. Formally,  
 213 the generative model is identical to Equation 1, factorizing over timesteps denoted by subscript  $l$ , i.e.  
 214  $p_{\mathbf{X}_l, \mathbf{T}_l}(\mathbf{x}_l, \mathbf{t}_l) = p_{\mathbf{X}_l | \mathbf{T}_l}(\mathbf{x}_l | \mathbf{t}_l)p_{\mathbf{T}_l}(\mathbf{t}_l)$ . However,  $\mathbf{T}_l$  is now a function of a sequence  $\{\mathbf{U}_{l+\delta}\}_{\delta=-L}^L$ :

$$\mathbf{T}_l = \frac{\mathbf{Z}_l - \mu}{\sqrt{\mathbf{W} [\mathbf{U}_{l+L}^2; \dots; \mathbf{U}_{l-L}^2]}} \quad (7)$$

215 Where  $[\mathbf{U}_{l+L}^2; \dots; \mathbf{U}_{l-L}^2]$  denotes vertical concatenation of the column vectors  $\mathbf{U}_l$ , and  $L$  can be  
 216 seen as the window size. We see that the choice of  $\mathbf{W}$  now defines correlation structure over time.  
 217 In prior work on temporal coherence [24? ], the grouping over time is such that a given variable  
 218  $T_{l,i}$  has correlated energy with *the same spatial location* ( $i$ ) at a previous time step ( $l-1$ ) (i.e.  
 219  $\text{cov}(T_{l,i}^2, T_{l-1,i}^2) > 0$ ). This can be implemented as:

$$\mathbf{W} [\mathbf{U}_{l+L}^2; \dots; \mathbf{U}_{l-L}^2] = \sum_{\delta=-L}^L \mathbf{W}_\delta \mathbf{U}_{l+\delta}^2 \quad (8)$$

220 Where  $\mathbf{W}_\delta$  defines the topography for a single timestep, and is typically the same for all timesteps.

221 **4.5.2 Equivariance as Shifting Temporal Coherence**

222 In our model, instead of requiring a single location to have correlated energies over a sequence (which  
 223 has been shown to learn invariant features [24]), we would like variables at sequentially permuted  
 224 locations *within a capsule* to have correlated energy between timesteps ( $\text{cov}(T_{l,i}^2, T_{l-1,i-1}^2) > 0$ ).  
 225 Similarly, this can be implemented as:

$$\mathbf{W} [\mathbf{U}_{l+L}^2; \dots; \mathbf{U}_{l-L}^2] = \sum_{\delta=-L}^L \mathbf{W}_\delta \text{Roll}_\delta(\mathbf{U}_{l+\delta}^2) \quad (9)$$

226 Where  $\text{Roll}_\delta(\mathbf{U}_{l+\delta}^2)$  denotes a cyclic permutation of  $\delta$  steps along the capsule dimension. The exact  
 227 implementation, continuous and non-cyclic extensions of Roll can be found in the appendix.

228 **5 Topographic VAE**

229 To train the parameters of the generative model  $\theta$ , we use the above formulation to parameterize an  
 230 approximate posterior for  $\mathbf{t}$  in terms of a deterministic transformation of approximate posteriors over  
 231 simpler Gaussian latent variables  $\mathbf{u}$  and  $\mathbf{z}$ . Explicitly:

$$q_\phi(\mathbf{z}_l | \mathbf{x}_l) = \mathcal{N}(\mathbf{z}_l; \mu_\phi(\mathbf{x}_l), \sigma_\phi(\mathbf{x}_l)\mathbf{I}) \quad p_\theta(\mathbf{x} | g_\theta(\mathbf{t})) = p_\theta(\mathbf{x} | g_\theta(\mathbf{z}, \mathbf{u})) \quad (10)$$

$$q_\gamma(\mathbf{u}_l | \mathbf{x}_l) = \mathcal{N}(\mathbf{u}_l; \mu_\gamma(\mathbf{x}_l), \sigma_\gamma(\mathbf{x}_l)\mathbf{I}) \quad \mathbf{t}_l = \frac{\mathbf{z}_l - \mu}{\sqrt{\mathbf{W} [\mathbf{u}_{l+L}^2; \dots; \mathbf{u}_{l-L}^2]}} \quad (11)$$

232 We denote this model the Topographic VAE (TVAE) and optimize the parameters  $\theta, \phi, \gamma$  (and  $\mu$ )  
 233 through the ELBO, summed over the sequence length  $S$ :

$$\sum_{l=1}^S \mathbb{E}_{Q_{\phi,\gamma}(\mathbf{z}_l, \mathbf{u}_l | \mathbf{x}_l)} ([\log p_\theta(\mathbf{x}_l | g_\theta(\mathbf{t}_l))] - D_{KL}[q_\phi(\mathbf{z}_l | \mathbf{x}_l) || p_{\mathbf{Z}}(\mathbf{z}_l)] - D_{KL}[q_\gamma(\mathbf{u}_l | \mathbf{x}_l) || p_{\mathbf{U}}(\mathbf{u}_l)]) \quad (12)$$

234 where  $Q_{\phi,\gamma}(\mathbf{z}_l, \mathbf{u}_l | \mathbf{x}_l) = q_\phi(\mathbf{z}_l | \mathbf{x}_l) \prod_{\delta=-L}^L q_\gamma(\mathbf{u}_{l+\delta} | \mathbf{x}_{l+\delta})$ . Practically, we observe layer-wise  
 235 weight normalization [47] is required to achieve topographic organization, this is detailed in the  
 236 appendix.

237 **6 Experiments**

238 In the following experiments, we demonstrate the viability of the Topographic VAE as a novel method  
 239 for training deep topographic generative models. Additionally, we quantitatively verify that shifting

240 temporal coherence yields approximately equivariant capsules by computing an ‘equivariance loss’  
 241 and a correlation metric inspired by the disentanglement literature. We show that equivariant capsule  
 242 models yield higher likelihood than baselines on test sequences, and qualitatively support these results  
 243 with visualizations of sequences reconstructed purely from rolled capsule activations.

## 244 6.1 Evaluation Methods

245 As depicted in Figure 1, we make use of *capsule traversals* to qualitatively visualize the transformations  
 246 learned by our network. Simply, these are constructed by encoding a partial sequence into a  $\mathbf{t}_0$   
 247 variable, and decoding sequentially Roll’d copies of this variable. Explicitly, in the top row we show  
 248 the data sequence  $\{\mathbf{x}_l\}_l$ , and in the bottom row we show the decoded sequence:  $\{g_\theta(\text{Roll}_l(\mathbf{t}_0))\}_l$ .

249 To measure equivariance quantitatively, we measure an *equivariance error* similar to [14]. The  
 250 equivariance error can be seen as the difference between traversing the two distinct paths of the  
 251 commutative diagram, and provides some measure of how precisely the function and the transform  
 252 commute. Formally, for a sequence of length  $S$ , and  $\hat{\mathbf{t}} = \mathbf{t} / \|\mathbf{t}\|_2$ , the error is defined as:

$$\mathcal{E}_{eq}(\{\mathbf{t}_l\}_{l=1}^S) = \sum_{l=1}^{S-1} \sum_{\delta=1}^{S-l} \|\text{Roll}_\delta(\hat{\mathbf{t}}_l) - \hat{\mathbf{t}}_{l+\delta}\|_1 \quad (13)$$

253 Additionally, inspired by existing disentanglement metrics, we measure the degree to which observed  
 254 transformations in capsule space are correlated with input transformations by introducing a new  
 255 metric we call  $\text{CapCorr}_y$ . Simply, this metric computes correlation between the amount of observed  
 256 roll of a capsule’s activation at two timesteps  $l$  and  $l + \delta$ , and the shift of the ground truth generative  
 257 factors  $y_l$  in that same time. Formally, for correlation coefficient  $\text{Corr}$ :

$$\text{CapCorr}(\mathbf{t}_l, \mathbf{t}_{l+\delta}, y_l, y_{l+\delta}) = \text{Corr}(\text{argmax}[\mathbf{t}_l \star \mathbf{t}_{l+\delta}], |y_l - y_{l+\delta}|) \quad (14)$$

258 Where  $\star$  is discrete cross-correlation across the capsule dimension, and the correlation coefficient  
 259 is computed across the entire dataset. We see the argmax of the cross-correlation is an estimate  
 260 of the degree to which a capsule activation has shifted from time  $l$  to  $l + \delta$ . To extend this to  
 261 multiple capsules, we can replace the argmax function with the mode of the argmax computed  
 262 for all capsules. We provide additional details and extensions of this metric in the appendix. For  
 263 measuring capsule-metrics on baseline model which do not naturally have capsules, we simply divide  
 264 the latent space into a fixed set of corresponding capsules and capsule dimensions.

## 265 6.2 Topographic VAE without Temporal Coherence

266 To validate the TVAE is capable of learning topographically organized representations with deep neural networks,  
 267 we first perform experiments on a Topographic VAE without Temporal Coherence. The model is constructed as in  
 268 Equations 10 and 11 with  $L = 0$ , and is trained to maximize Equation 12. We fix  $\mathbf{W}$  such that globally the latent  
 269 variables are arranged in a grid on a 2-dimensional torus (a single capsule), and locally  $\mathbf{W}$  sums over 5x5 2D groups  
 270 of variables. In this setting,  $\mathbf{W}$  can be easily implemented as 2D convolution with a 5x5 kernel of 1’s, stride 1, and  
 271 cyclic padding. We see that training the model with 3-layer MLP’s for the encoders and decoder indeed yields  
 272 a 2D topographic organization of higher level features. In Figure 3, we show the maximum activating image for each  
 273 final layer neuron of the capsule, plotted as a flattened torus. We see that the neurons become arranged according  
 274 to class, orientation, width, and other learned features.



Figure 3: Maximum activating images for a Topographic VAE trained with a 2D torus topography on MNIST.

## 283 6.3 Learning Equivariant Capsules

284 In the remaining experiments, we provide evidence that the Topographic VAE can be leveraged to  
 285 learn equivariant capsules by incorporating shifting temporal coherence into a 1D baseline topographic

286 model. We compare against two baselines: standard normal VAEs (without weight normalization),  
 287 and models that have non-shifting ‘stationary’ temporal coherence (denoted ‘BubbleVAE’ [24]).

288 In all experiments we use a 3-layer MLP with ReLU activations for both encoders and the decoder.  
 289 We arrange the latent space into 15 circular capsules each of 15-dimensions for dSprites [40], and 18  
 290 circular capsules each of 18-dimensions for MNIST [36]. Example sequences  $\{\mathbf{x}_l\}_{l=1}^S$  are formed by  
 291 taking a random initial example, and sequentially transforming it according to one of the available  
 292 transformations: (X-Pos, Y-Pos, Orientation, Scale) for dSprites, and (Color, Scale, Orientation)  
 293 for MNIST. All transformation sequences are cyclic such that when the maximum transformation  
 294 parameter is reached, the subsequent value returns to the minimum. We denote the length of a full  
 295 transformation sequence by  $S$ , and the time-extent of the induced temporal coherence (i.e. the length  
 296 of the input sequence) by  $2L$ . For simplicity, both datasets are constructed such that the sequence  
 297 length  $S$  equals the capsule dimension (for dSprites this involves taking a subset of the full dataset  
 298 and looping the scale 3-times for a scale-sequence). Exact details and subsets are in the appendix.

299 In Figure 4, we show the capsule traversals for TVAE models with  $L \approx \frac{1}{3}S$ . We see that despite the  
 300  $t_0$  variable encoding only  $\frac{2}{3}$  of the sequence, the remainder of the transformation sequence can be  
 301 decoded nearly perfectly by permuting the activation through the full capsule – implying the model has  
 302 learned to be equivariant to full sequences while only observing partial sequences per training point.  
 303 Furthermore, we see that the model is able to successfully learn all transformations simultaneously for  
 304 the respective datasets. Capsule traversals for the non-equivariant baselines, as well as TVAEs with  
 305 smaller values of  $L$  (which learn only equivariance to partial sequences) are shown in the appendix.  
 306 We note that the capsule traversal plotted in Figure 1 uses a separate dataset where color and rotation  
 307 transformations are combined simultaneously. This is intended solely to demonstrate the ability of  
 308 the TVAE to learn complex transformations, but such combined transformations are not present in the  
 309 remainder of the experiments presented in this section. Additional results involving complex learned  
 310 transformations (such as perspective transforms) are shown in the appendix.

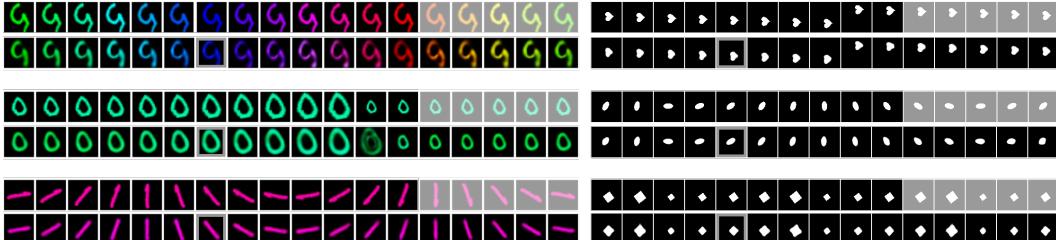


Figure 4: Capsule Traversals for TVAE models on dSprites and MNIST. The top rows show the encoded sequences (with greyed-out images held-out), and the bottom rows show the images generated from by decoding sequentially rolled copies of the initial activation  $t_0$  (indicated by a grey border).

311 In Table 1 we measure the equivariance error and log-likelihood (reported in nats) of the test data  
 312 under our trained MNIST models as estimated by importance sampling with 10 samples. We observe  
 313 that models which incorporate temporal coherence (BubbleVAE and TVAE with  $L > 0$ ) achieve  
 314 higher log-likelihood than the baseline VAE, while the TVAE models with shifting temporal coherence  
 315 achieve the highest likelihood and the lowest equivariance error. Surprisingly, we also see that the  
 316 TVAE with no temporal coherence ( $L = 0$ ) still achieves significantly lower equivariance error than  
 317 the VAE baseline. By examining the capsule traversals of this model (see appendix), we posit that  
 318 this is due to the induced smoothness of the representation within capsules.

Table 1: Log Likelihood and Equivariance Error on MNIST for different settings of temporal coherence length  $L$  relative to sequence length  $S$ . Mean  $\pm$  std. over 3 random initializations.

Model ( $L, S$ )	TVAE $L = \frac{1}{2}S$	TVAE $L = \frac{5}{36}S$	TVAE $L = 0$	BubbleVAE $L = \frac{5}{36}S$	VAE $L = 0$
$\log p(\mathbf{x}) \uparrow$	-188.5 $\pm$ 1.0	<b>-184.7</b> $\pm$ 0.7	-217.3 $\pm$ 1.3	190.0 $\pm$ 0.3	-189.0 $\pm$ 0.8
$\mathcal{E}_{eq} \downarrow$	<b>386.7</b> $\pm$ 1.0	3149.6 $\pm$ 20.7	3047.8 $\pm$ 59.2	2522.5 $\pm$ 21.3	13273.9 $\pm$ 0.5

319 To further understand how capsules transform for observed input transformations, in Table 2 we  
 320 measure  $\mathcal{E}_{eq}$  and the CapCorr metric on the dSprites dataset for the four proposed transformations.

321 We see that the TVAE with  $L \geq \frac{1}{3}S$  achieves perfect correlation – implying the learned representation  
 322 indeed permutes cyclically within capsules for observed transformation sequences. Further, we see  
 323 that this correlation gradually decreases as  $L$  decreases, eventually reaching the same level as the  
 324 baselines. We see that contrary to the equivariance loss, the correlation of transformations for  $L = 0$   
 325 is not higher than the standard VAE. We believe this to be due to the fundamental difference between  
 326 the metrics.  $\mathcal{E}_{eq}$  measures continuous L1 similarity which is still low when a representation is locally  
 327 smooth (even if the change of the representation does not exactly follow the observed transformation),  
 328 whereas CapCorr more strictly measures the correspondence between the transformation of the input  
 329 and the transformation of the representation. Additionally,  $\mathcal{E}_{eq}$  is misleadingly low for invariant  
 330 representations, whereas CapCorr strictly measures equivariance. Interestingly, we also see the  
 331 BubbleVAE achieves even lower correlation than the VAE baseline, agreeing with the assumption  
 that such a model should learn invariant features.

Table 2: Equivariance error ( $\mathcal{E}_{eq} \downarrow$ ) and correlation of observed capsule roll with ground truth factor shift (CapCorr  $\uparrow$ ) for the dSprites dataset. Mean  $\pm$  standard deviation over 3 random initializations.

Model ( $L, S$ )	TVAE $L = \frac{1}{2}S$	TVAE $L = \frac{1}{3}S$	TVAE $L = \frac{1}{6}S$	TVAE $L = 0$	BubbleVAE $L = \frac{1}{6}S$	VAE $L = 0$
CapCorr <sub>X</sub> $\uparrow$	<b>1.0</b> $\pm$ 0	<b>1.0</b> $\pm$ 0	0.68 $\pm$ 0.02	0.17 $\pm$ 0.03	0.07 $\pm$ 0.02	0.24 $\pm$ 0.01
CapCorr <sub>Y</sub> $\uparrow$	<b>1.0</b> $\pm$ 0	<b>1.0</b> $\pm$ 0	0.69 $\pm$ 0.02	0.17 $\pm$ 0.04	0.08 $\pm$ 0.01	0.25 $\pm$ 0.01
CapCorr <sub>O</sub> $\uparrow$	<b>1.0</b> $\pm$ 0	<b>1.0</b> $\pm$ 0	0.57 $\pm$ 0.01	0.11 $\pm$ 0.02	0.05 $\pm$ 0.01	0.09 $\pm$ 0.01
CapCorr <sub>S</sub> $\uparrow$	<b>1.0</b> $\pm$ 0	<b>1.0</b> $\pm$ 0	0.42 $\pm$ 0.01	0.52 $\pm$ 0.01	0.26 $\pm$ 0.00	0.47 $\pm$ 0.00
$\mathcal{E}_{eq} \downarrow$	<b>131</b> $\pm$ 0	964 $\pm$ 9	2313 $\pm$ 33	1075 $\pm$ 20	1397 $\pm$ 29	6929 $\pm$ 1

332

## 333 7 Future Work & Limitations

334 The model presented in this work has a number of limitations in its existing form which we believe to  
 335 be interesting directions for future research. Foremost, the model is challenging to compare directly  
 336 with existing disentanglement and equivariance literature since it requires an input sequence which  
 337 determines the transformations reachable through the capsule roll. We believe this limitation could  
 338 be overcome in multiple manners, either by increasing the capsule dimension beyond 1D circles, or  
 339 by increasing the sparsity between capsules such that individual capsules begin to specialize on trans-  
 340 formations. Related to this, we note the temporal coherence proposed in our model is not ‘causal’ ( $t_0$   
 341 depends on future  $x_l$ ). However this may be easily remedied when desired, and future work may only  
 342 be a function of past inputs. We provide more details on these proposals in the supplementary material.  
 343 We additionally note that a priori definition of topographic structure may be a burden to model  
 344 developers. While true, we know that the construction of appropriate priors is always a challenging  
 345 task in latent variable models, and we observe that our proposed TVAE achieves strong performance  
 346 even with improper specification. Furthermore, in future work, we believe adding learned flexibility  
 347 to the parameters  $\mathbf{W}$  may alleviate some of this burden.  
 348 Finally, we note that while this work does demonstrate improved log-likelihood and equivariance  
 349 error, the study is inherently preliminary and does not examine all important benefits of equivariant  
 350 representations. Specifically, further study of this model in terms of the sample complexity, semi-  
 351 supervised classification accuracy, and invariance through capsule pooling would be enlightening.

## 352 8 Conclusion

353 In the above work we introduce the Topographic Variational Autoencoder as a method to train deep  
 354 topographic generative models, and show how topography can be leveraged to learn equivariant sets  
 355 of features, a.k.a. capsules, directly from sequences of data with no other supervision. Ultimately, we  
 356 believe these results may shine some light on how biological brains could hard-wire themselves to  
 357 more effectively learn representations with equivariant capsule structure. In terms of broader impact,  
 358 the learned transformations of our model could be used to generate more realistic transformations of  
 359 ‘deepfakes’, enhancing disinformation. Given that the model learns *approximate* equivariance, we  
 360 caution against the over-reliance on equivariant properties as these have no formal guarantees yet.

361 **References**

- 362 [1] Horace B Barlow et al. Possible principles underlying the transformation of sensory messages.  
363 *Sensory communication*, 1(01), 1961.
- 364 [2] Anthony J. Bell and Terrence J. Sejnowski. An Information-Maximization Approach to Blind  
365 Separation and Blind Deconvolution. *Neural Computation*, 7(6):1129–1159, 11 1995.
- 366 [3] Yoshua Bengio, Aaron Courville, and Pascal Vincent. Representation learning: A review and  
367 new perspectives. *IEEE transactions on pattern analysis and machine intelligence*, 35(8):1798–  
368 1828, 2013.
- 369 [4] Gregory W. Benton, Marc Finzi, Pavel Izmailov, and Andrew Gordon Wilson. Learning  
370 invariances in neural networks. *CoRR*, abs/2010.11882, 2020.
- 371 [5] Christopher Bishop, Markus Svensen, and Christopher Williams. Gtm: The generative topo-  
372 graphic mapping. *Neural Computation*, 10:215–234, 05 1997.
- 373 [6] Benedikt Boenninghoff, Steffen Zeiler, Robert M. Nickel, and Dorothea Kolossa. Vari-  
374 ational autoencoder with embedded student- $t$  mixture model for authorship attribution. *ArXiv*,  
375 abs/2005.13930, 2020.
- 376 [7] Diane Bouchacourt, Mark Ibrahim, and Stéphane Deny. Addressing the topological defects of  
377 disentanglement via distributed operators. *ArXiv*, abs/2102.05623, 2021.
- 378 [8] Taco Cohen and M. Welling. Transformation properties of learned visual representations. *CoRR*,  
379 abs/1412.7659, 2015.
- 380 [9] Taco Cohen and M. Welling. Steerable cnns. *ArXiv*, abs/1612.08498, 2017.
- 381 [10] Taco Cohen and Max Welling. Learning the irreducible representations of commutative lie  
382 groups. In *Proceedings of the 31st International Conference on Machine Learning*, volume 32  
383 of *Proceedings of Machine Learning Research*, pages 1755–1763, Beijing, China, 22–24 Jun  
384 2014. PMLR.
- 385 [11] Taco Cohen and Max Welling. Group equivariant convolutional networks. In *International  
386 conference on machine learning*, pages 2990–2999, 2016.
- 387 [12] Pierre Comon. Independent component analysis, a new concept? *Signal processing*, 36(3):287–  
388 314, 1994.
- 389 [13] Marissa Connor, Gregory Canal, and Christopher Rozell. Variational autoencoder with learned  
390 latent structure. In *Proceedings of The 24th International Conference on Artificial Intelligence  
and Statistics*, volume 130 of *Proceedings of Machine Learning Research*, pages 2359–2367.  
391 PMLR, 13–15 Apr 2021.
- 393 [14] Nichita Diaconu and Daniel E. Worrall. Learning to convolve: A generalized weight-tying  
394 approach. *CoRR*, abs/1905.04663, 2019.
- 395 [15] Marc Finzi, Samuel Stanton, Pavel Izmailov, and Andrew Gordon Wilson. Generalizing  
396 convolutional neural networks for equivariance to lie groups on arbitrary continuous data.  
397 In *Proceedings of the 37th International Conference on Machine Learning*, volume 119 of  
398 *Proceedings of Machine Learning Research*, pages 3165–3176. PMLR, 13–18 Jul 2020.
- 399 [16] Marc Finzi, M. Welling, and Andrew Gordon Wilson. A practical method for constructing  
400 equivariant multilayer perceptrons for arbitrary matrix groups. *ArXiv*, abs/2104.09459, 2021.
- 401 [17] Peter Földiák. Learning invariance from transformation sequences. *Neural Computation*,  
402 3:194–200, 06 1991.
- 403 [18] Will Grathwohl and Aaron Wilson. Disentangling space and time in video with hierarchical  
404 variational auto-encoders. *CoRR*, abs/1612.04440, 2016.
- 405 [19] Irina Higgins, David Amos, David Pfau, Sebastien Racaniere, Loic Matthey, Danilo Rezende,  
406 and Alexander Lerchner. Towards a definition of disentangled representations. *ArXiv*,  
407 abs/1812.02230, 2018.

- 408 [20] Geoffrey E. Hinton, Alex Krizhevsky, and Sida D. Wang. Transforming auto-encoders. In  
 409 Timo Honkela, Włodzisław Duch, Mark Girolami, and Samuel Kaski, editors, *Artificial Neural*  
 410 *Networks and Machine Learning – ICANN 2011*, pages 44–51, Berlin, Heidelberg, 2011.  
 411 Springer Berlin Heidelberg.
- 412 [21] Geoffrey E Hinton, Sara Sabour, and Nicholas Frosst. Matrix capsules with EM routing. In  
 413 *International Conference on Learning Representations*, 2018.
- 414 [22] Geoffrey E. Hinton and Yee-Whye Teh. Discovering multiple constraints that are frequently  
 415 approximately satisfied. In *Proceedings of the Seventeenth Conference on Uncertainty in*  
 416 *Artificial Intelligence*, UAI'01, page 227–234, 2001.
- 417 [23] Jarmo Hurri and Aapo Hyvärinen. Simple-Cell-Like Receptive Fields Maximize Temporal  
 418 Coherence in Natural Video. *Neural Computation*, 15(3):663–691, 03 2003.
- 419 [24] A. Hyvärinen, J. Hurri, and Jaakko J. Väyrynen. A unifying framework for natural image  
 420 statistics: spatiotemporal activity bubbles. *Neurocomputing*, 58-60:801–806, 2004.
- 421 [25] Aapo Hyvärinen and Patrik Hoyer. Emergence of phase-and shift-invariant features by decom-  
 422 position of natural images into independent feature subspaces. *Neural computation*, 12(7):1705–  
 423 1720, 2000.
- 424 [26] Aapo Hyvärinen, Patrik O Hoyer, and Mika Inki. Topographic independent component analysis.  
 425 *Neural computation*, 13(7):1527–1558, 2001.
- 426 [27] Aapo Hyvärinen, Jarmo Hurri, and Patrick O Hoyer. *Natural image statistics: A probabilistic*  
 427 *approach to early computational vision.*, volume 39. Springer Science & Business Media, 2009.
- 428 [28] Aapo Hyvärinen and Erkki Oja. Independent component analysis: algorithms and applications.  
 429 *Neural networks*, 13(4-5):411–430, 2000.
- 430 [29] Aapo Hyvärinen and Patrik O. Hoyer. A two-layer sparse coding model learns simple and  
 431 complex cell receptive fields and topography from natural images. *Vision Research*, 41(18):2413–  
 432 2423, 2001.
- 433 [30] Koray Kavukcuoglu, Marc’Aurelio Ranzato, Rob Fergus, and Yann LeCun. Learning invariant  
 434 features through topographic filter maps. In *2009 IEEE Conference on Computer Vision and*  
 435 *Pattern Recognition*, pages 1605–1612. IEEE, 2009.
- 436 [31] Diederik P Kingma and Max Welling. Auto-encoding variational bayes. *arXiv preprint*  
 437 *arXiv:1312.6114*, 2013.
- 438 [32] David Klindt, Lukas Schott, Yash Sharma, Ivan Ustyuzhaninov, Wieland Brendel, Matthias  
 439 Bethge, and Dylan Paiton. Towards nonlinear disentanglement in natural data with temporal  
 440 sparse coding. *ArXiv*, abs/2007.10930, 2021.
- 441 [33] Teuvo Kohonen. Emergence of invariant-feature detectors in the adaptive-subspace self-  
 442 organizing map. *Biological cybernetics*, 75(4):281–291, 1996.
- 443 [34] Alexei A Koulakov and Dmitri B Chklovskii. Orientation preference patterns in mammalian  
 444 visual cortex: a wire length minimization approach. *Neuron*, 29(2):519–527, 2001.
- 445 [35] Quoc Le, Alexandre Karpenko, Jiquan Ngiam, and Andrew Ng. Ica with reconstruction cost for  
 446 efficient overcomplete feature learning. In J. Shawe-Taylor, R. Zemel, P. Bartlett, F. Pereira,  
 447 and K. Q. Weinberger, editors, *Advances in Neural Information Processing Systems*, volume 24.  
 448 Curran Associates, Inc., 2011.
- 449 [36] Yann LeCun, Corinna Cortes, and CJ Burges. Mnist handwritten digit database. *ATT Labs*  
 450 *[Online]. Available: <http://yann.lecun.com/exdb/mnist>*, 2, 2010.
- 451 [37] Hyodong Lee, Eshed Margalit, Kamila M. Jozwik, Michael A. Cohen, Nancy Kanwisher, Daniel  
 452 L. K. Yamins, and James J. DiCarlo. Topographic deep artificial neural networks reproduce the  
 453 hallmarks of the primate inferior temporal cortex face processing network. *bioRxiv*, 07/2020  
 454 2020.

- 455 [38] Jan Eric Lenssen, Matthias Fey, and Pascal Libuschewski. Group equivariant capsule networks.  
 456 *arXiv preprint arXiv:1806.05086*, 2018.
- 457 [39] Libo Ma and Liqing Zhang. Overcomplete topographic independent component analysis.  
 458 *Neurocomputing*, 71(10-12):2217–2223, 2008.
- 459 [40] Loic Matthey, Irina Higgins, Demis Hassabis, and Alexander Lerchner. dsprites: Disentangle-  
 460 ment testing sprites dataset. <https://github.com/deepmind/dsprites-dataset/>, 2017.
- 461 [41] Bruno A Olshausen and David J Field. Sparse coding with an overcomplete basis set: A strategy  
 462 employed by V1? *Vision research*, 37(23):3311–3325, 1997.
- 463 [42] Simon Osindero, Max Welling, and Geoffrey E. Hinton. Topographic Product Models Applied  
 464 to Natural Scene Statistics. *Neural Computation*, 18(2):381–414, 02 2006.
- 465 [43] Simon Kayode Osindero. *Contrastive Topographic Models*. PhD thesis, University of London,  
 466 2004.
- 467 [44] Siamak Ravanbakhsh, Jeff Schneider, and Barnabas Poczos. Equivariance through parameter-  
 468 sharing. *ArXiv*, abs/1702.08389, 2017.
- 469 [45] Danilo Jimenez Rezende, Shakir Mohamed, and Daan Wierstra. Stochastic backpropagation  
 470 and approximate inference in deep generative models. *ArXiv*, abs/1401.4082, 2014.
- 471 [46] Sara Sabour, Nicholas Frosst, and Geoffrey E Hinton. Dynamic routing between capsules. *arXiv  
 472 preprint arXiv:1710.09829*, 2017.
- 473 [47] Tim Salimans and Diederik P. Kingma. Weight normalization: A simple reparameterization to  
 474 accelerate training of deep neural networks. *ArXiv*, abs/1602.07868, 2016.
- 475 [48] James V. Stone. Learning Perceptually Salient Visual Parameters Using Spatiotemporal Smooth-  
 476 ness Constraints. *Neural Computation*, 8(7):1463–1492, 10 1996.
- 477 [49] Richard Turner and Maneesh Sahani. A maximum-likelihood interpretation for slow feature  
 478 analysis. *Neural computation*, 19:1022–38, 05 2007.
- 479 [50] Maurice Weiler, Mario Geiger, Max Welling, Wouter Boomsma, and Taco Cohen. 3d steerable  
 480 cnns: Learning rotationally equivariant features in volumetric data. *ArXiv*, abs/1807.02547,  
 481 2018.
- 482 [51] Max Welling, Simon Osindero, and Geoffrey E Hinton. Learning sparse topographic represen-  
 483 tations with products of student-t distributions. In *Advances in neural information processing  
 484 systems*, pages 1383–1390, 2003.
- 485 [52] Laurenz Wiskott and Terrence J Sejnowski. Slow feature analysis: Unsupervised learning of  
 486 invariances. *Neural computation*, 14(4):715–770, 2002.
- 487 [53] Daniel Worrall and Max Welling. Deep scale-spaces: Equivariance over scale. In H. Wallach,  
 488 H. Larochelle, A. Beygelzimer, F. d'Alché-Buc, E. Fox, and R. Garnett, editors, *Advances in  
 489 Neural Information Processing Systems*, volume 32. Curran Associates, Inc., 2019.
- 490 [54] Daniel E. Worrall, Stephan J. Garbin, Daniyar Turmukhambetov, and Gabriel J. Brostow.  
 491 Harmonic networks: Deep translation and rotation equivariance. *ArXiv*, abs/1612.04642, 2017.

492    **Checklist**

- 493    1. For all authors...
- 494    (a) Do the main claims made in the abstract and introduction accurately reflect the pa-  
495    per's contributions and scope? Yes, see Figure 1 for an example of a learned highly  
496    complex transformation, and Sections 5 and 6 for the proposed model and empirical  
497    demonstration of claimed contributions.
- 498    (b) Did you describe the limitations of your work? Yes, see Section 7
- 499    (c) Did you discuss any potential negative societal impacts of your work? Yes, see Section  
500    ??
- 501    (d) Have you read the ethics review guidelines and ensured that your paper conforms to  
502    them? Yes.
- 503    2. If you are including theoretical results...
- 504    (a) Did you state the full set of assumptions of all theoretical results? N/A
- 505    (b) Did you include complete proofs of all theoretical results? N/A
- 506    3. If you ran experiments...
- 507    (a) Did you include the code, data, and instructions needed to reproduce the main experi-  
508    mental results (either in the supplemental material or as a URL)? Yes, see supplementary  
509    material.
- 510    (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they  
511    were chosen)? Yes, see supplementary material.
- 512    (c) Did you report error bars (e.g., with respect to the random seed after running experi-  
513    ments multiple times)? Yes, see Tables 1 and 2.
- 514    (d) Did you include the total amount of compute and the type of resources used (e.g., type  
515    of GPUs, internal cluster, or cloud provider)? Yes, see supplementary material.
- 516    4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
- 517    (a) If your work uses existing assets, did you cite the creators? Yes, see Section 6
- 518    (b) Did you mention the license of the assets? Yes, see supplementary material.
- 519    (c) Did you include any new assets either in the supplemental material or as a URL? Yes,  
520    see supplementary material.
- 521    (d) Did you discuss whether and how consent was obtained from people whose data you're  
522    using/curating? N/A
- 523    (e) Did you discuss whether the data you are using/curating contains personally identifiable  
524    information or offensive content? N/A
- 525    5. If you used crowdsourcing or conducted research with human subjects...
- 526    (a) Did you include the full text of instructions given to participants and screenshots, if  
527    applicable? N/A
- 528    (b) Did you describe any potential participant risks, with links to Institutional Review  
529    Board (IRB) approvals, if applicable? N/A
- 530    (c) Did you include the estimated hourly wage paid to participants and the total amount  
531    spent on participant compensation? N/A

532 **A Experiment Details**

533 The code for reproducing all experiments in this paper can be found in the following GitHub  
534 repository: <https://github.com/ReallyAnonNeurips2021/TopographicVAE>

535 **A.1 Optimizer Parameters**

536 Given the differences between the training procedures of the model presented in Section 6.2, and those  
537 in Section 6.3, the optimizer parameters for the two settings differed slightly. The 2D Topographic  
538 VAE without Temporal Coherence presented in Figure 3 was trained with stochastic gradient descent  
539 on batches of size 128, using a learning rate of  $1 \times 10^{-3}$ , and standard momentum of 0.9. All  
540 models in Section 6.3 were trained with stochastic gradient descent on batches of size 8 (due to each  
541 batch-example being a length 15 or 18 sequence), using a learning rate of  $1 \times 10^{-4}$ , and standard  
542 momentum of 0.9. All models were trained for 100 epochs. As mentioned in Section 5, we applied  
543 layer-wise weight normalization [47] to all topographic models (TVAE & BubbleVAE), but not to the  
544 baseline VAE as we found its likelihood was significantly higher without weight normalization. The  
545 details of the weight normalization method are in Section A.8 below.

546 **A.2 Initialization**

547 All weights of the models were initialized with uniformly random samples from  $U(-\frac{1}{\sqrt{m}}, \frac{1}{\sqrt{m}})$ ,  
548 where  $m$  is the number of input units. For all topographic models including BubbleVAE,  $\mu$  was  
549 initialized to a large value (30.0) as this was observed to increase the speed of convergence and was  
550 sometimes necessary for observed topographic organization in deeper models. For the 2D topographic  
551 model in Figure 2,  $\mu$  was initialized to 5.

552 **A.3 Model Architectures**

553 All models presented in this paper make use of the same 3-Layer MLP for parameterizing the encoders  
554 and decoders. Specifically, the model is constructed as 3 fully connected layers with ReLU activations  
555 in-between the layers. For MNIST, the layers of both the  $\mathbf{u}$  and  $\mathbf{z}$  encoders have (972, 648, 648)  
556 output units each for the first, second, and third layers respectively. The 648 units in the third layer  
557 are divided into two sets to compute the mean and log standard deviation of the respective  $u$ 's and  
558  $z$ 's, yielding 324  $t$  variables. This is then divided into 18 capsules, each of 18 dimensions. The layers  
559 of the decoder have (648, 972, 2352) output units respectively. For dSprites, both encoder layers  
560 have output sizes (674, 450, 450), where the resulting 225  $t$  variables are divided into 15 capsules,  
561 each of 15 dimensions. The decoder layers then have output sizes (450, 675, 4096). We note the  
562 non-topographic VAE baselines make use of only a single encoder for the Gaussian variable  $\mathbf{z}$  (as  $\mathbf{u}$   
563 is not needed), and do not incorporate a  $\mu$  parameter.

564 **A.4 Choice of  $\mathbf{W}$**

565 For all topographic models (TVAE and BubbleVAE) in Section 6.3, the global topographic organiza-  
566 tion afforded by  $\mathbf{W}$  was fixed to a set of 1-D tori ('circular capsules') as depicted in Figure 1. The  
567 model presented in Section 6.2 organizes its variables as a single 2-D torus. Practically, multiplication  
568 by  $\mathbf{W}$  was performed by convolution over the appropriate dimensions (time & capsule dimension)  
569 with a kernel of all 1's, taking advantage of circular padding to achieve toroidal structure.

570 **A.5 Choice of  $\mathbf{W}_\delta$**

571 The choice of  $\mathbf{W}_\delta$  determines the local topographic structure within a single timestep. For all TVAE  
572 models with  $L > 0$  we experimented with local neighborhood sizes (denoted  $K$ ) of 3 units (effective  
573 kernel size 3 in the capsule dimension), and 1 unit (no neighborhood). For MNIST it was observed  
574 that  $K = 3$  performed best, while  $K = 1$  worked best for dSprites. This is likely due to the slower,  
575 smoother, and more overlapping transformations constructed on MNIST, whereas our subset of  
576 dSprites contained non-smooth transformations where the overlap between successive images was  
577 smaller (e.g. due to sub-setting, see Section A.10), which made larger neighborhood sizes  $K > 1$   
578 less fitting. For TVAE models with  $L = 0$ ,  $\mathbf{W}_\delta = \mathbf{W}$  was fixed to sum over neighborhoods of size  
579  $K = 9$  for MNIST and  $K = 10$  for dSprites. These values were chosen arbitrarily to be close to half

580 of the capsule size. For BubbleVAE models, the extent of topographic organization in the capsule  
 581 dimension was set to be equal to the organization in time dimension  $K = 2L$ . We additionally  
 582 experimented with  $K = 3$  to match the TVAE, but observed this performed worse in on all metrics.

583 **A.6 Choice of  $L$**

584 The choice of  $L$  determines the extent of temporal coherence where  $2L$  equals the input sequence  
 585 length. For Table 1, we experimented with values of  $L$  in the set  $\{0, \frac{5}{36}S, \frac{1}{4}S, \frac{1}{2}S\}$  for both the TVAE  
 586 and BubbleVAE. Both the BubbleVAE and TVAE achieved highest likelihoods at  $L = \frac{5}{36}S$ , and  
 587 TVAE achieved lowest equivariance error at  $L = \frac{1}{2}S$ . We additionally included TVAE experiments  
 588 with  $L = \frac{13}{36}S$  for purposes of visualization in Figures 1 and 4. For Table 2, we experimented with  
 589 values of  $L$  in the set  $\{0, \frac{1}{6}S, \frac{4}{15}S, \frac{1}{3}S, \frac{2}{5}S, \frac{1}{2}S\}$  for both TVAE and BubbleVAE, and presented a  
 590 broad selection in the table. The results of all models are shown in Section B below.

591 **A.7 Hyperparameter Selection**

592 Hyperparameters such as learning rate, batch size, number of capsules, capsule size, and ultimately  
 593 model architecture were chosen to allow for quick training on limited resources and were not tuned  
 594 significantly. Since it was conceptually simpler to have an equal number of capsule dimensions and  
 595 sequence elements, this limited the number of capsules we could then train efficiently. In Section C.1  
 596 we explain how a model with fewer capsule dimensions than sequence elements could be constructed.

597 **A.8 Weight Normalization**

598 As noted in Section 5, it was found that some form of weight normalization at each layer was necessary  
 599 in order to achieve topographic organization. This is consistent with prior work on overcomplete  
 600 Topographic Product of Student’s-t models [51] and other overcomplete ICA methods [28] where  
 601 the same technique has been used to prevent degenerate solutions. In the ICA literature, a *Projected*  
 602 *Gradient Descent* method was used to enforce this constraint [35]. In our work, we observe *Weight*  
 603 *Normalization* [47], achieved by reparameterization, works equally well while frequently obtaining  
 604 higher log-likelihoods. Simply, *Weight Normalization* is achieved by reparameterizing the weight  
 605 vectors (rows  $\mathbf{w}_i$  of  $\mathbf{W}$ ) as the product of a normalized vector  $\frac{\mathbf{v}_i}{\|\mathbf{v}_i\|_2}$  and a ‘rescaling’ scalar  $g_i$ , i.e:

$$\mathbf{w}_i = g_i \cdot \frac{\mathbf{v}_i}{\|\mathbf{v}_i\|_2} \quad (15)$$

606 In the table below we show a comparison of Projected Gradient Descent (PGD) and weight normaliza-  
 607 zation (WN) for the TVAE ( $L = \frac{5}{36}S$ ) and VAE ( $L = 0$ ) on the MNIST Dataset.

Table 3: Comparison of weight normalization methods for TVAE ( $L = \frac{5}{36}S$ ) and baseline VAE. PGD  
 stands for Projected Gradient Descent [35], while WN stands for Weight Normalization [47]. Mean  
 $\pm$  std. over 3 random initializations.

Model Norm. Method	TVAE PGD	TVAE WN	VAE PGD	VAE WN	VAE None
$\log p(\mathbf{x}) \uparrow$	$-189.0 \pm 0.1$	$-184.9 \pm 0.8$	$-209.1 \pm 1.1$	$-192.5 \pm 0.5$	$-189.0 \pm 0.8$
$\mathcal{E}_{eq} \downarrow$	$3196.7 \pm 5.0$	$3137.9 \pm 6.4$	$13262.0 \pm 0.5$	$13273.2 \pm 1.4$	$13273.9 \pm 0.5$

607

608 **A.9 MNIST Transformations**

609 The first set of experiments presented in this paper are based on the MNIST dataset [36] (MIT  
 610 Licence). For Section 6.2 (Figure 3) an MNIST training set of 48,000 images was used, while the  
 611 standard test set of 10,000 images was used to compute the maximum activating image. For Section  
 612 6.3 (Figure 4 and Table 1), sequences of MNIST images were created by picking a random training  
 613 image (with a random transformation ‘pose’) and successively transforming it according to one of  
 614 the 3 available transformations (e.g. only one attribute is changed per sequence). The available  
 615 transformations consisted of rotation, color (hue rotation), and scale with increments of 20-degrees for

616 rotation and color, and 3.66% increments for scale. Since scale is inherently non-cyclic, the bounds of  
 617 the transformation were set at 60% and 126%, and the transformations were constructed to be periodic  
 618 such that once scale reached 126%, the next element was at 60% scale. The final sequences were thus  
 619 constructed to be 18 images long, where each element in the batch had an independently randomly  
 620 chosen transformation. Again, the likelihood  $\log p(\mathbf{x})$  and equivariance error  $\mathcal{E}_{eq}$  were computed on  
 621 the held-out 10,000 example test set, where the same random transformation sequences were applied.

## 622 A.10 dSprites Transformations

623 The second set of experiments presented in this paper are based on the dSprites dataset [40] (Apache-  
 624 2.0 License). To reduce computational complexity of this dataset, we took a subset of the dataset which  
 625 consisted of all 3 shapes, the largest 5 scales, and every other example from the first 30 orientations,  
 626 x-positions, and y-positions. The resulting dataset thus had 50,625 total images (3 shapes, 5 scales, 15  
 627 orientations, 15 x-positions, 15 y-positions), compared to the original 737,280 images. To construct  
 628 sequences, we followed the same procedure as for MNIST, whereby first a random example and  
 629 transformation were chosen, and a sequence of 15 images was constructed where only the chosen  
 630 transformation was applied successively. We define the transformations available for sequences as  
 631 scale, orientation, x-position, and y-position, omitting shape since smooth shape transforms are not  
 632 present in the dSprites dataset. Again, we define all transformations to be cyclic such that once the  
 633 15th element is reached, the 1st element follows. For scale transformations, we simply loop over all 5  
 634 scales 3 times per sequence. We observe that although these sequences do not match the latent priors  
 635 exactly, the models still train relatively well, implying some degree of robustness.

## 636 A.11 Capsule Correlation Metric (CapCorr)

637 Here we define CapCorr more precisely as it is implemented in our work. First, we denote the  
 638 ground truth transformation parameter of the sequence at timestep  $l$  as  $y_l$  (e.g. the rotation angle at  
 639 timestep  $l$  for a rotation sequence), and the corresponding activation at time  $l$  as  $t_l$ . Next, to get an  
 640 arbitrary starting point, we let  $l = \Omega$  denote the timestep when  $y_l$  is at its canonical position (e.g.  
 641 rotation angle 0, x-position 0, or scale 1). We see  $\Omega$  is not necessarily 0 since the first timestep of  
 642 each sequence ( $l = 0$ ) is a randomly transformed example. Then, we observe that we can measure  
 643 the approximate observed roll in the capsule dimension between time 0 and  $\Omega$  as a ‘phase shift’ by  
 644 computing the index of the maximum value of a discrete (periodic) cross-correlation of  $t_\Omega$  and  $t_0$ :

$$\text{ObservedRoll}(\mathbf{t}_\Omega, \mathbf{t}_0) = \text{argmax} [\mathbf{t}_\Omega \star \mathbf{t}_0] \quad (16)$$

645 Where  $\star$  is discrete (periodic) cross-correlation across the (cyclic) capsule dimension and  $\text{argmax}$   
 646 is also subsequently performed over the capsule dimension. Then, the CapCorr metric for a single  
 647 capsule is given as:

$$\text{CapCorr}(\mathbf{t}_\Omega, \mathbf{t}_0, y_\Omega, y_0) = \text{Corr} (\text{ObservedRoll}(\mathbf{t}_\Omega, \mathbf{t}_0), |y_\Omega - y_0|) \quad (17)$$

648 Where the correlation coefficient  $\text{Corr}$  is then computed across all examples for the entire dataset. In  
 649 our experiments we use the Pearson correlation coefficient for  $\text{Corr}$ . We thus see this metric is the  
 650 correlation of the estimated observed capsule roll with the shift in ground truth generative factors,  
 651 which is equal to 1 when the model is perfectly equivariant. To extend this definition to multiple  
 652 capsules, we estimate  $\text{ObservedRoll}$  for each capsule separately, and then correlate the mode of all  
 653  $\text{ObservedRoll}$  values with the true shift in ground truth generative factors. We see empirically that  
 654 the  $\text{ObservedRolls}$  for all capsules are almost always identical (i.e. all capsules roll simultaneously  
 655 for each transformation), therefore computing the mode does not destroy significant information.  
 656 Finally, for transformation sequences which have multiple timesteps where  $y_l$  is at the canonical  
 657 position (e.g. scale transformations on dSprites where scale is looped 3 times), we select  $l = \Omega$  to  
 658 be the one from this possible set which yields the minimal absolute distance between  $|y_\Omega - y_0|$  and  
 659  $\text{ObservedRoll}(\mathbf{t}_\Omega, \mathbf{t}_0)$ .

## 660 A.12 Definition of Roll for Capsules

661 As stated in Section 4.5.2,  $\text{Roll}_\delta(\mathbf{u})$ , is defined as a cyclic permutation of  $\delta$  steps along the capsule  
 662 dimension of  $\mathbf{u}$ . Explicitly, if  $\mathbf{u}$  is divided into  $C$  capsules each with  $D$  dimensions, the  $\text{Roll}_\delta$   
 663 operation can be written as:

$$\begin{aligned} \text{Roll}_\delta(\mathbf{u}) &= \text{Roll}_\delta ([u_1, u_2, \dots, u_{C \cdot D}]) \\ &= [u_D, u_1, \dots, u_{D-1}, u_{2 \cdot D}, u_{D+1}, \dots, u_{2 \cdot D-1}, u_{3 \cdot D}, \dots, \dots, u_{C \cdot D-1}] \end{aligned} \quad (18)$$

664 **B Extended Results**

665 In this section we include an extended version of Tables 1 & 2, showing all tested settings of the  
 666 TVAE & BubbleVAE. Additionally, in Figure 5 below, we plot the average (over all transformations)  
 667 CapCorr metric for all settings of  $L$  as shown in Table 5. We observe the TVAE achieves perfect  
 correlation ( $\text{CapCorr} = 1$ ) for  $L \geq \frac{1}{3}$ , and steadily decreasing correlation for lower values of  $L$ .

Table 4: Log Likelihood and Equivariance Error on MNIST for all models tested. Mean  $\pm$  std. over 3 random initializations.

Model	TVAE	TVAE	TVAE	TVAE	TVAE
$(L, S)$	$L = \frac{1}{2}S$	$L = \frac{13}{36}S$	$L = \frac{1}{4}S$	$L = \frac{5}{36}S$	$L = 0$
$K$	$K = 3$	$K = 3$	$K = 3$	$K = 3$	$K = 9$
$\log p(\mathbf{x}) \uparrow$	$-188.5 \pm 1.1$	$-187.1 \pm 1.2$	$-186.9 \pm 0.6$	$-184.7 \pm 0.7$	$-217.3 \pm 1.3$
$\mathcal{E}_{eq} \downarrow$	<b><math>386.7 \pm 1.1</math></b>	$990.6 \pm 10.7$	$2045.7 \pm 3.3$	$3149.6 \pm 20.8$	$3047.8 \pm 59.2$

Model	BubbleVAE	BubbleVAE	BubbleVAE	BubbleVAE	VAE
$(L, S)$	$L = \frac{1}{2}S$	$L = \frac{1}{4}S$	$L = \frac{5}{36}S$	$L = \frac{5}{36}S$	$L = 0$
$K$	$K = 2L$	$K = 2L$	$K = 2L$	$K = 3$	$K = 1$
$\log p(\mathbf{x}) \uparrow$	$-196.3 \pm 1.1$	$-203.2 \pm 2.1$	$-190.0 \pm 0.3$	$-190.4 \pm 0.9$	$-189.0 \pm 0.8$
$\mathcal{E}_{eq} \downarrow$	$2066.5 \pm 596.7$	$1005.3 \pm 12.7$	$2522.5 \pm 21.3$	$3277.0 \pm 15.6$	$13273.9 \pm 0.5$

Table 5: Equivariance error and CapCorr for all models tested on the dSprites dataset. Mean  $\pm$  standard deviation over 3 random initializations.

Model	TVAE	TVAE	TVAE	TVAE	TVAE	TVAE
$(L, S)$	$L = \frac{1}{2}S$	$L = \frac{2}{5}S$	$L = \frac{1}{3}S$	$L = \frac{4}{15}S$	$L = \frac{1}{6}S$	$L = 0$
$\text{CapCorr}_X \uparrow$	<b><math>1.0 \pm 0</math></b>	<b><math>1.0 \pm 0</math></b>	<b><math>1.0 \pm 0</math></b>	$0.96 \pm 0.01$	$0.68 \pm 0.02$	$0.17 \pm 0.03$
$\text{CapCorr}_Y \uparrow$	<b><math>1.0 \pm 0</math></b>	<b><math>1.0 \pm 0</math></b>	<b><math>1.0 \pm 0</math></b>	$0.97 \pm 0.01$	$0.69 \pm 0.02$	$0.17 \pm 0.04$
$\text{CapCorr}_O \uparrow$	<b><math>1.0 \pm 0</math></b>	<b><math>1.0 \pm 0</math></b>	<b><math>1.0 \pm 0</math></b>	$0.90 \pm 0.01$	$0.57 \pm 0.01$	$0.11 \pm 0.02$
$\text{CapCorr}_S \uparrow$	<b><math>1.0 \pm 0</math></b>	<b><math>1.0 \pm 0</math></b>	<b><math>1.0 \pm 0</math></b>	$0.97 \pm 0.00$	$0.42 \pm 0.01$	$0.52 \pm 0.01$
$\mathcal{E}_{eq} \downarrow$	<b><math>131 \pm 0</math></b>	$679 \pm 4$	$964 \pm 9$	$1298 \pm 5$	$2313 \pm 33$	$1075 \pm 20$

Model	BubbleVAE	BubbleVAE	BubbleVAE	BubbleVAE	BubbleVAE	VAE
$(L, S)$	$L = \frac{1}{2}S$	$L = \frac{2}{5}S$	$L = \frac{1}{3}S$	$L = \frac{4}{15}S$	$L = \frac{1}{6}S$	$L = 0$
$\text{CapCorr}_X \uparrow$	$0.17 \pm 0.01$	$0.14 \pm 0.01$	$0.13 \pm 0.02$	$0.12 \pm 0.03$	$0.07 \pm 0.02$	$0.24 \pm 0.01$
$\text{CapCorr}_Y \uparrow$	$0.16 \pm 0.01$	$0.14 \pm 0.00$	$0.12 \pm 0.01$	$0.12 \pm 0.02$	$0.08 \pm 0.01$	$0.25 \pm 0.00$
$\text{CapCorr}_O \uparrow$	$0.12 \pm 0.03$	$0.13 \pm 0.01$	$0.10 \pm 0.01$	$0.09 \pm 0.01$	$0.05 \pm 0.01$	$0.09 \pm 0.01$
$\text{CapCorr}_S \uparrow$	$0.52 \pm 0.00$	$0.53 \pm 0.00$	$0.52 \pm 0.02$	$0.49 \pm 0.03$	$0.26 \pm 0.00$	$0.47 \pm 0.00$
$\mathcal{E}_{eq} \downarrow$	$3543 \pm 441$	$6922 \pm 6$	$2037 \pm 924$	$2420 \pm 1192$	$1397 \pm 29$	$6929 \pm 1$

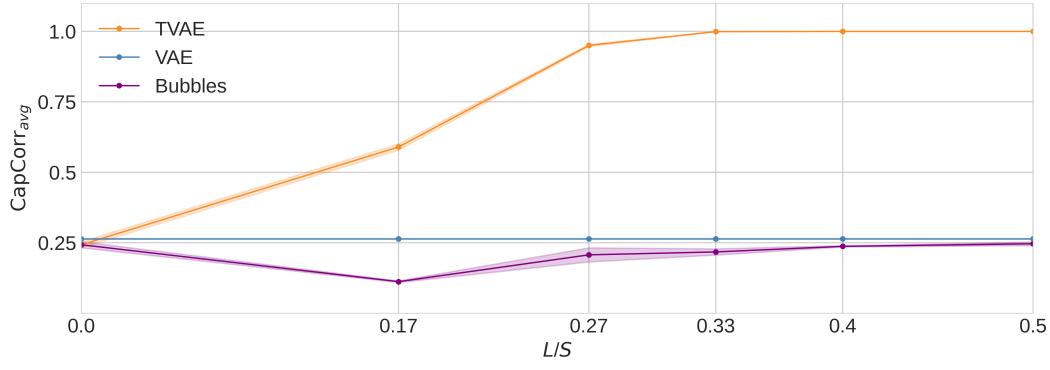


Figure 5: Average CapCorr over all transformations vs. Temporal Coherence Length  $L$ .

669 **C Proposed Model Extensions**

670 **C.1 Extensions to Roll & CapCorr**

671 The Roll operation can be seen as defining the speed at which  $\mathbf{t}$  transforms corresponding to an  
 672 observed transformation. For example, with Roll defined as in Section A.12 above, we implicitly  
 673 assume that for each observed timestep, we would like the representation  $\mathbf{t}$  to cyclically permute  
 674 1-unit within the capsule. For this to match the observed data, it requires the model to have an  
 675 equal number of capsule dimensions and sequence elements. If we wish to reduce the size of our  
 676 representation, we need a ‘partial permutation’ for each observed transformation. For a single capsule  
 677 with  $D$  elements, an example of a simple linear version of such a partial permutation (for  $0 < \alpha \leq 1$ )  
 678 can be implemented as:

$$\text{Roll}_\alpha(\mathbf{u}) = [\alpha u_D + (1 - \alpha)u_1, \alpha u_1 + (1 - \alpha)u_2, \dots, \alpha u_{D-1} + (1 - \alpha)u_D] \quad (19)$$

679 A slightly more principled partial roll for periodic signals could also be achieved by performing a  
 680 phase shift of the signal in Fourier space, and performing the inverse Fourier transform to obtain the  
 681 resulting rolled signal. To extend the CapCorr metric to similarly allow for partial Rolls, we see  
 682 that we can simply redefine the ObservedRoll (originally given by discrete cross-correlation) to be  
 683 given by the argmax of the inner product of a sequentially partially rolled activation with the initial  
 684 activation  $\mathbf{t}_\Omega$ . Formally:

$$\text{ObservedRoll}(\mathbf{t}_\Omega, \mathbf{t}_0) = \text{argmax} [\mathbf{t}_\Omega \cdot \text{Roll}_0(\mathbf{t}_0), \mathbf{t}_\Omega \cdot \text{Roll}_\alpha(\mathbf{t}_0), \dots, \mathbf{t}_\Omega \cdot \text{Roll}_{D-\alpha}(\mathbf{t}_0)] \quad (20)$$

685 **C.2 Non-Cyclic Capsules**

686 We can also see that there is nothing beyond convenience which inherently requires the capsules to  
 687 be circular (i.e. have periodic boundary conditions). To implement linear capsules, we propose one  
 688 solution is to add  $L$  additional  $U_i$  variables to both the left and right boundaries of each capsule. In  
 689 this way, the vector  $\mathbf{U}$  is larger than the vector  $\mathbf{Z}$  and can be seen as a ‘padded’ version, where the  
 690 padding is composed of independant random variables. Additionally, the transformation sequences  
 691 can then be padded on both sides by replicating the first and final elements  $L$  times. The construction  
 692 of  $\mathbf{T}$  variables is then performed identically as in Equations 7 and 9. The choice of a cyclic or  
 693 non-cyclic Roll operation which can be defined as filling the boundaries with 0 since these values  
 694 will not be used as part of the computation.

695 **C.3 Multi-dimensional Temporally Coherent Capsules**

696 In consideration of transformations which may naturally live in multiple dimensions, we wish to  
 697 extend the original model to support multi-dimensional capsules. Such multi-dimensional capsules  
 698 could additionally support more well-defined ‘disentanglement’ of transformations by encouraging  
 699 each transformation to be axis-aligned with one dimension of each capsule. We see that in the  
 700 non-temporally coherent case ( $L = 0$ ), the model can easily be extended to capsules of multiple  
 701 dimensions through multi-dimensional neighborhoods. An example of a model with 2-dimensional  
 702 neighborhoods is presented in Figure 3. However, when considering shifting temporal coherence  
 703 as we defined in Section 6.3, it is not clear how the shift operator or the neighborhoods should be  
 704 defined for higher dimensional capsules. In this section we propose to modify the definitions of  $\mathbf{T}$  in  
 705 Equations 7 and 9 with an extension resembling ‘group sparsity’ in the denominator.

706 First, we again assume that each input sequence is an observation of a single transformation at a time.  
 707 Formally, the multi-dimensional capsules are then constructed by arranging  $\mathbf{U}$  into a  $D$  dimensional  
 708 lattice. In such a model, we desire to roll and sum only along a single axis of the lattice for a given  
 709 sequence. Incorporating this into the construction of  $\mathbf{T}$  yields the following:

$$\mathbf{T}_l = \frac{\mathbf{Z}_l - \mu}{\sum_{d=1}^D \sqrt{\mathbf{W}^d [\mathbf{U}_{l+L}^2; \dots; \mathbf{U}_{l-L}^2]}} = \frac{\mathbf{Z}_l - \mu}{\sum_{d=1}^D \sqrt{\sum_{\delta=-L}^L \mathbf{W}_\delta^d \text{Roll}_\delta^d(\mathbf{U}_{l+\delta}^2)}} \quad (21)$$

710 Where  $\mathbf{W}_\delta^d$  refers to a matrix which sums locally along the  $d^{th}$  dimension of each capsule, and not at  
 711 all along the others, and similarly  $\text{Roll}_\delta^d$  rolls only along the  $d^{th}$  dimension. In practice we observe  
 712 such models can indeed disentangle up to 2 distinct transformations, but become more challenging  
 713 to optimize for higher dimensions. We believe this is potentially due to the exponential growth in  
 714 capsule size with increasing dimension, but leave further exploration to future work.

715 **C.4 Causal Temporal Coherence**

716 Finally, as noted in the Conclusion, the sequence model in this paper is not ‘causal’, meaning that  
 717 each variable  $\mathbf{T}_l$  requires variables from future timesteps in the sequence ( $\mathbf{U}_{l+\delta}$  for  $\delta > 0$ ). Although  
 718 for the purpose of learning equivariance in practice this may not be an issue, it may be relevant for  
 719 some online learning applications. We can modify Equations 7 and 9 by changing the matrix  $\mathbf{W}$   
 720 (implemented as convolution) to a causal convolution (i.e. masking out  $\mathbf{W}_\delta$  for  $\delta > 0$ ). Formally:

$$\mathbf{T}_l = \frac{\mathbf{Z}_l - \mu}{\sqrt{\mathbf{W} [\mathbf{U}_l^2; \dots; \mathbf{U}_{l-L}^2]}} = \frac{\mathbf{Z}_l - \mu}{\sqrt{\sum_{\delta=-L}^0 \mathbf{W}_\delta \text{Roll}_\delta(\mathbf{U}_{l+\delta}^2)}} \quad (22)$$

721 In a causal setting, it is also likely the transformations are no longer assumed to be circular. We thus  
 722 refer the reader to Section C.2 above on non-circular capsules, which can be combined with Equation  
 723 22, to achieve such a model.

724 **D Capsule Traversals**

725 In this section we provide a set of 12 capsule traversals for each of the models presented in main  
 726 text. The traversals are randomly selected such that all transformations (and dSprites shapes) are  
 727 shown evenly. Unlike the main section, we additionally include a middle row which shows the direct  
 728 reconstruction of the input without any rolling (i.e.  $\{g_\theta(\mathbf{t}_l)\}_l$ ). We find the direct reconstructions  
 729 valuable to determine if poor traversals are due to bad reconstructions (low  $\log p_\theta(\mathbf{x}|\mathbf{t})$ ) or a lack of  
 730 equivariance ( $\mathcal{E}_{eq}$ ). For example, with the baseline VAE models, we see that the reconstructions in  
 731 the middle row are accurate for the full sequence, while the capsule traversals obtained by sequentially  
 732 rolling the initial activation (shown in the bottom row) are nothing like the input transformation (top  
 733 row). In all traversals, the left-most image corresponds to  $\mathbf{t}_0$ , and thus input sequences of length  
 734  $2L$  cover both the left and right edges when  $L > 0$ . Additionally in the figure captions we include  
 735 the value of  $K$ , the topographic neighborhood size within a single timestep, to make the difference  
 736 between the TVAE and VAE at  $L = 0$  clear.

737 Finally, in Figures 18 & 19 at the end of the section, we include capsule traversals for models trained  
 738 on MNIST with more complex transformations such as combined color & rotation (as shown in  
 739 Figure 1), and combined color & perspective transforms. These models were trained in an identical  
 740 manner to the other MNIST models, with the same architecture, only changing the transformation  
 741 sequences of the dataset.

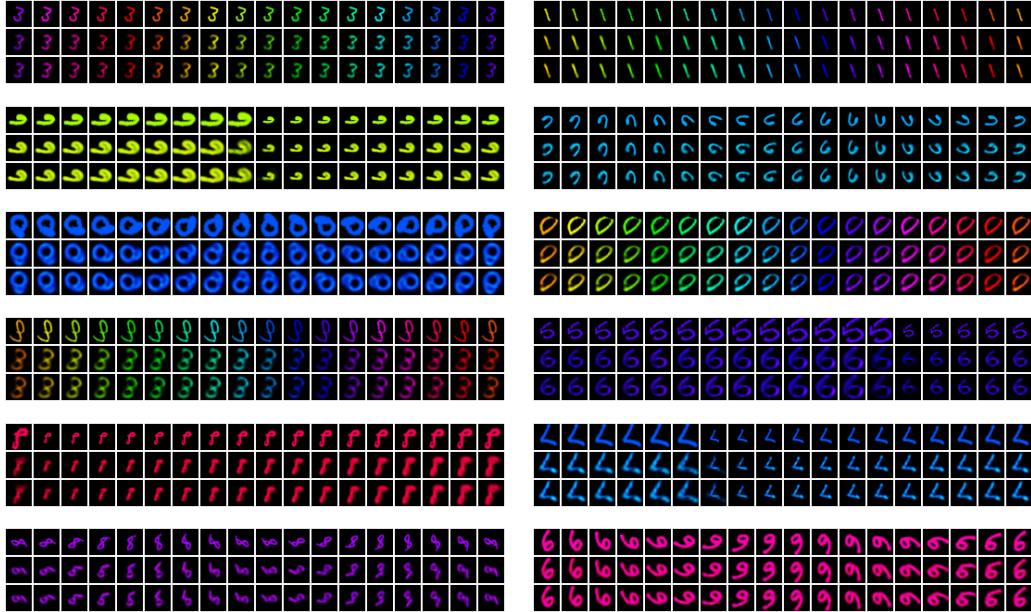


Figure 6: MNIST TVAE  $L = \frac{1}{2}S$ ,  $K = 3$

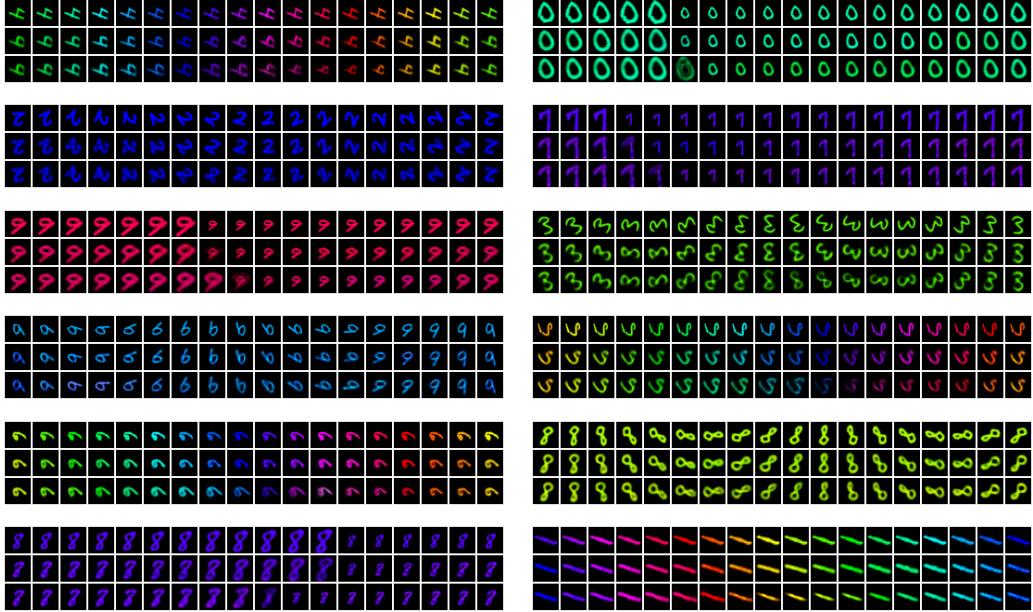


Figure 7: MNIST TVAE  $L = \frac{13}{36}S$ ,  $K = 3$

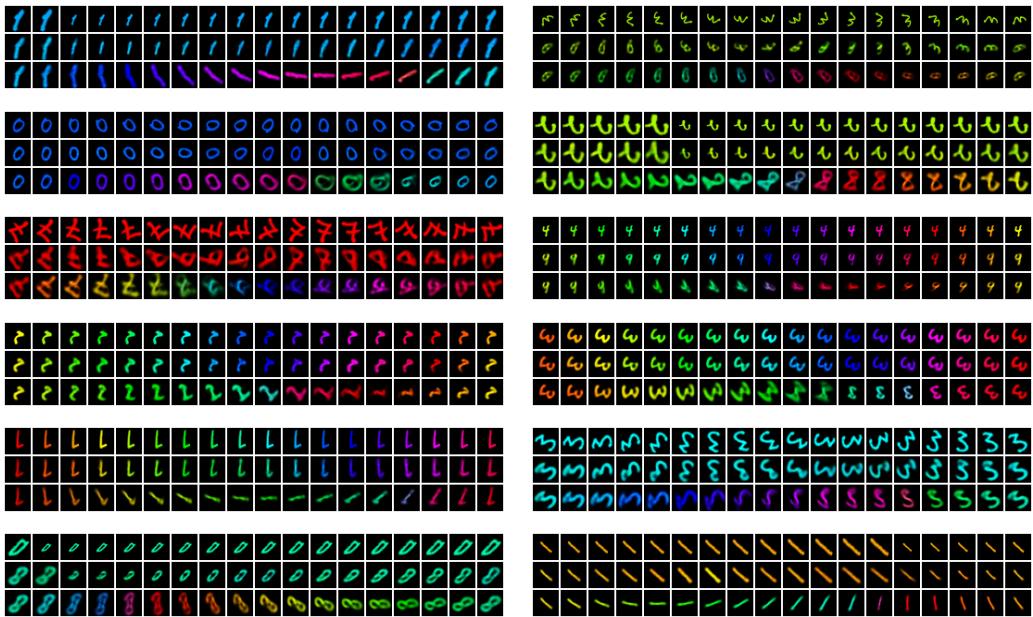


Figure 8: MNIST TVAE  $L = \frac{5}{36}S$ ,  $K = 3$

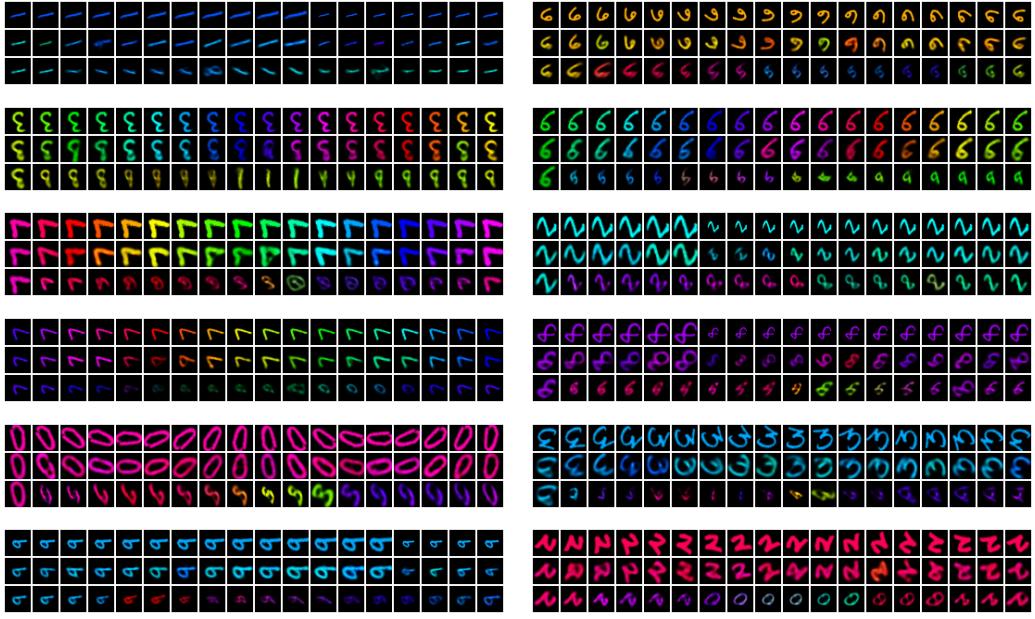


Figure 9: MNIST TVAE  $L = 0, K = 9$

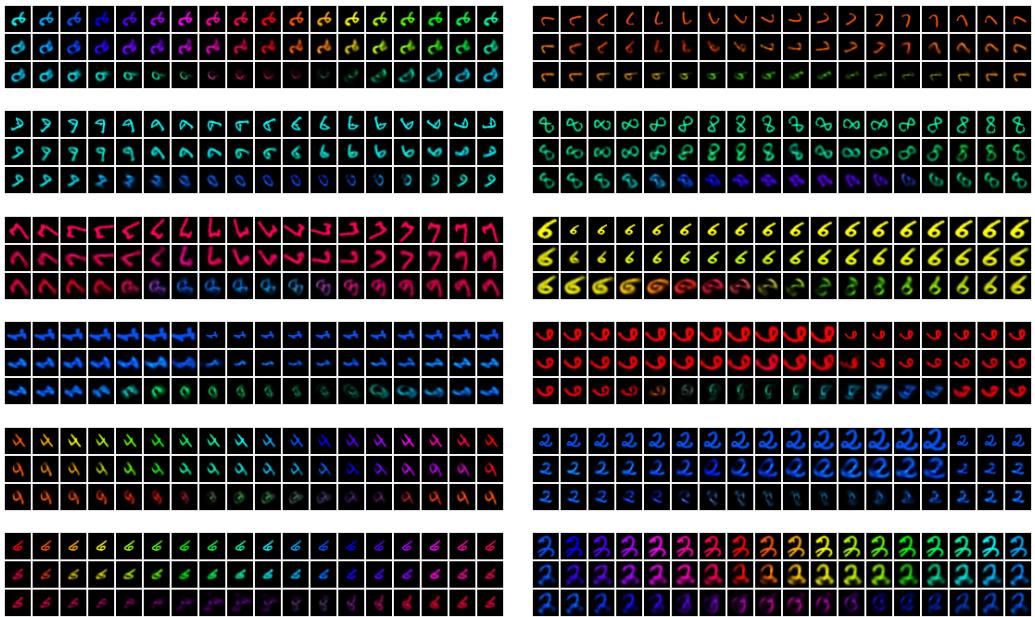


Figure 10: MNIST BubbleVAE  $L = \frac{5}{36}S, K = 2L$

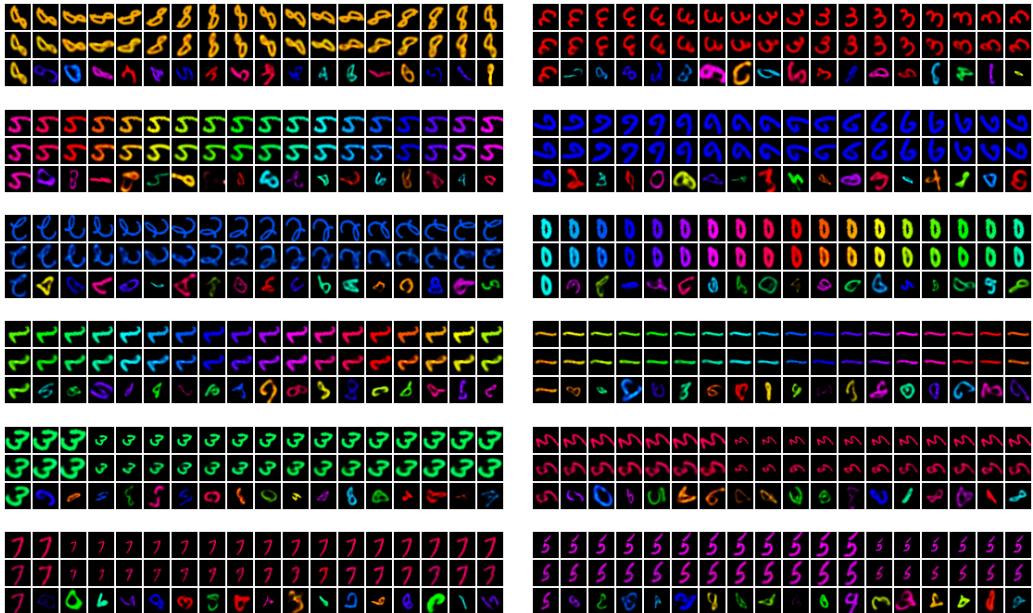


Figure 11: MNIST VAE  $L = 0, K = 1$

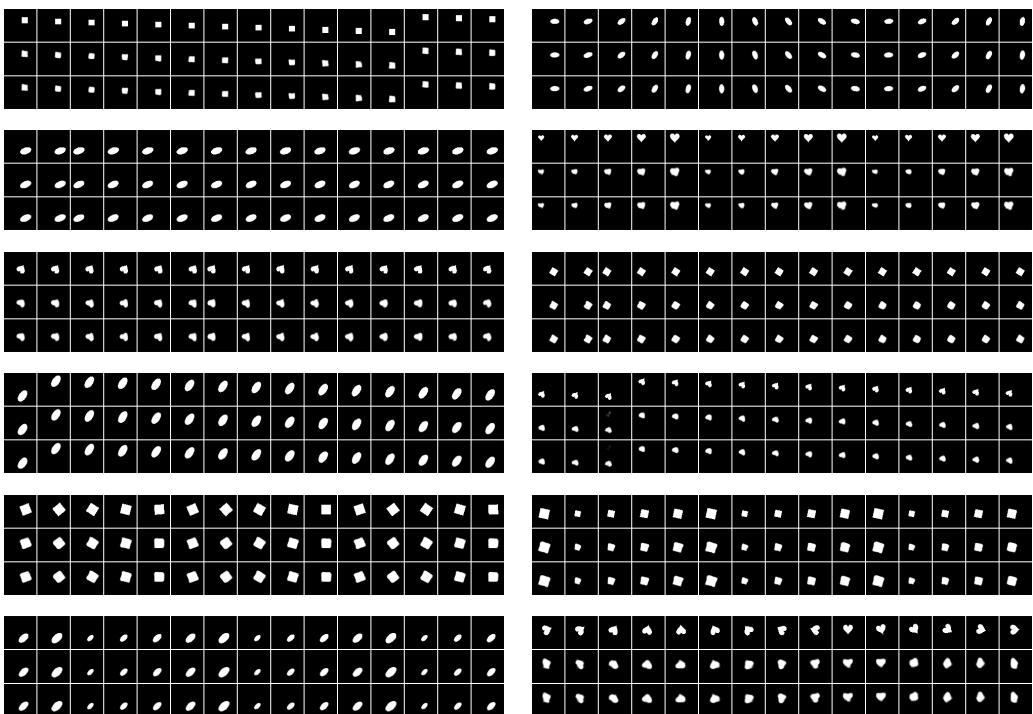


Figure 12: dSprites TVAE  $L = \frac{1}{2}S, K = 1$

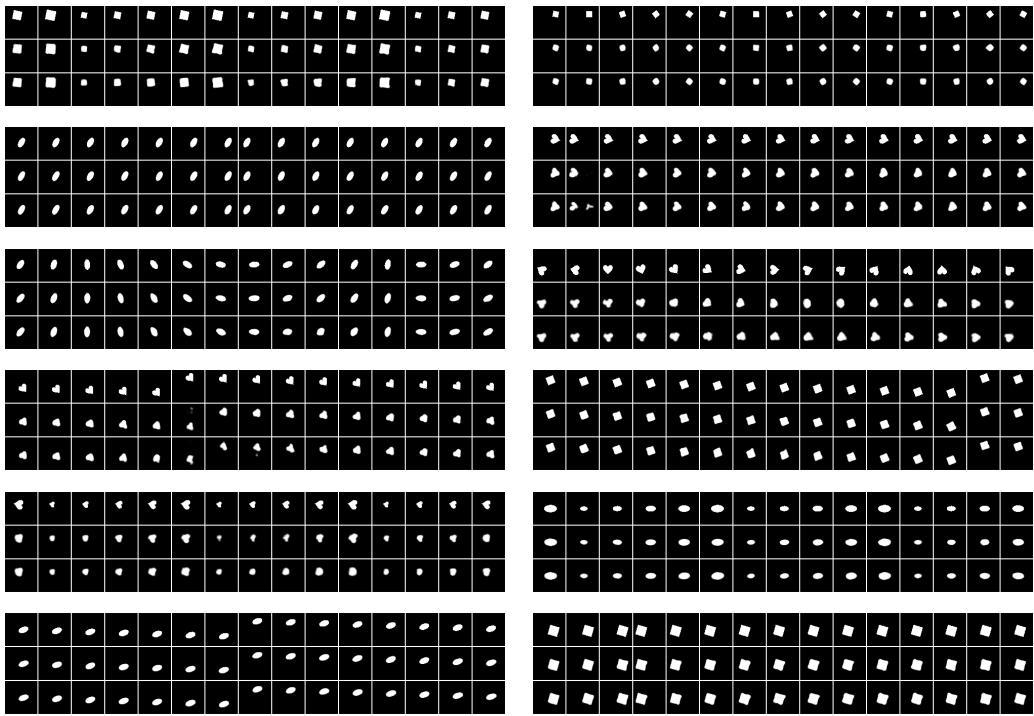


Figure 13: dSprites TVAE  $L = \frac{1}{3}S$ ,  $K = 1$

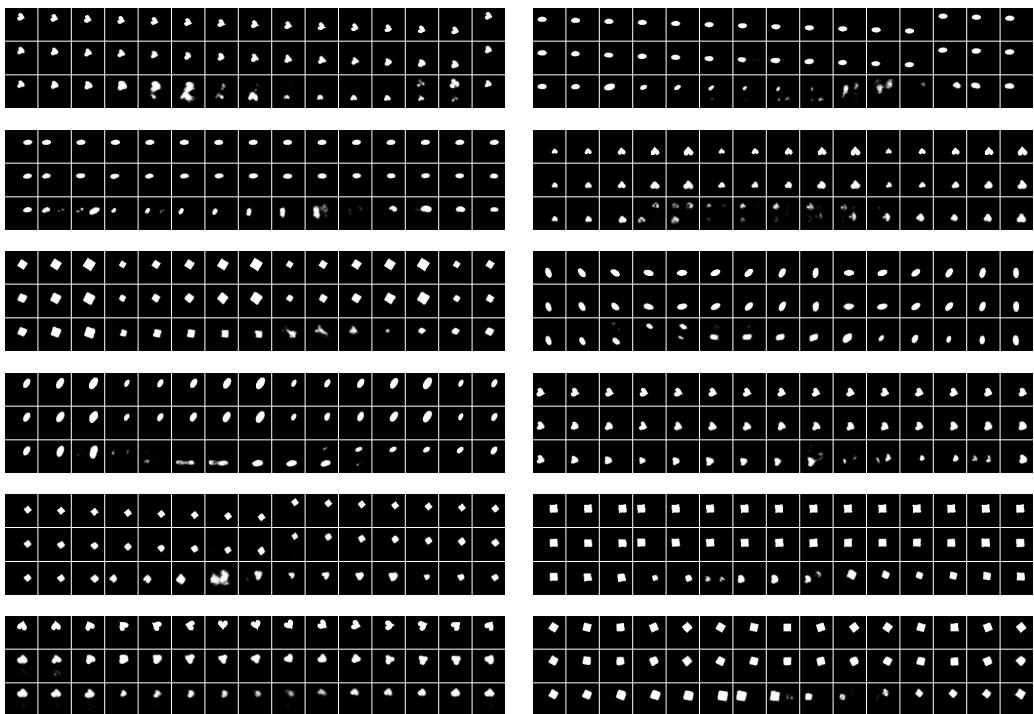


Figure 14: dSprites TVAE  $L = \frac{1}{6}S$ ,  $K = 1$

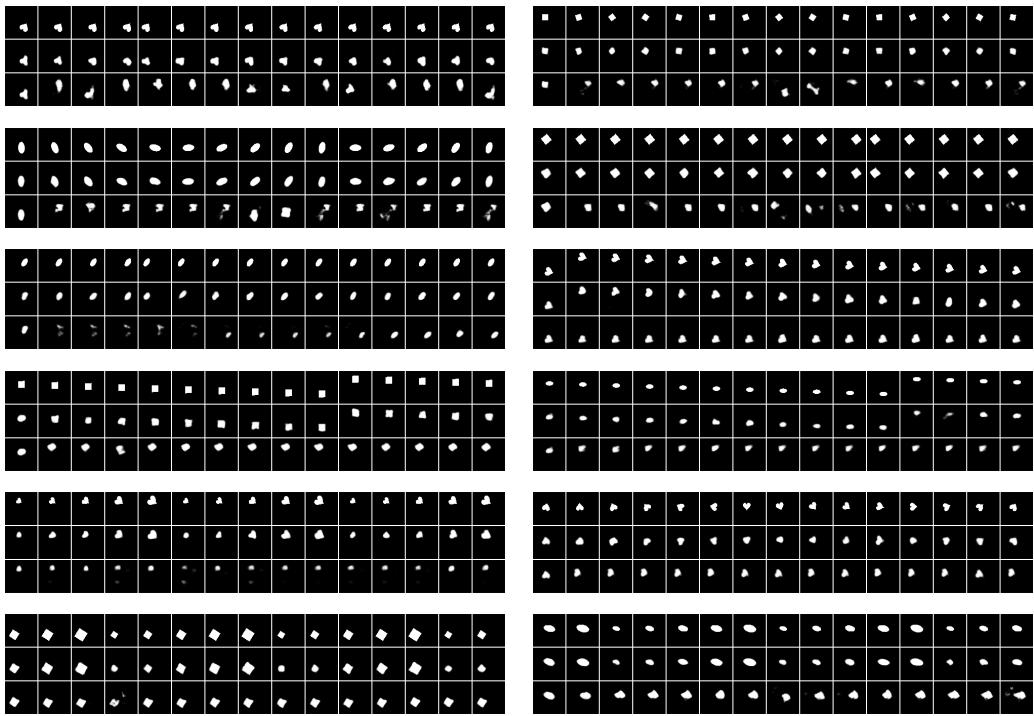


Figure 15: dSprites TVAE  $L = 0, K = 10$

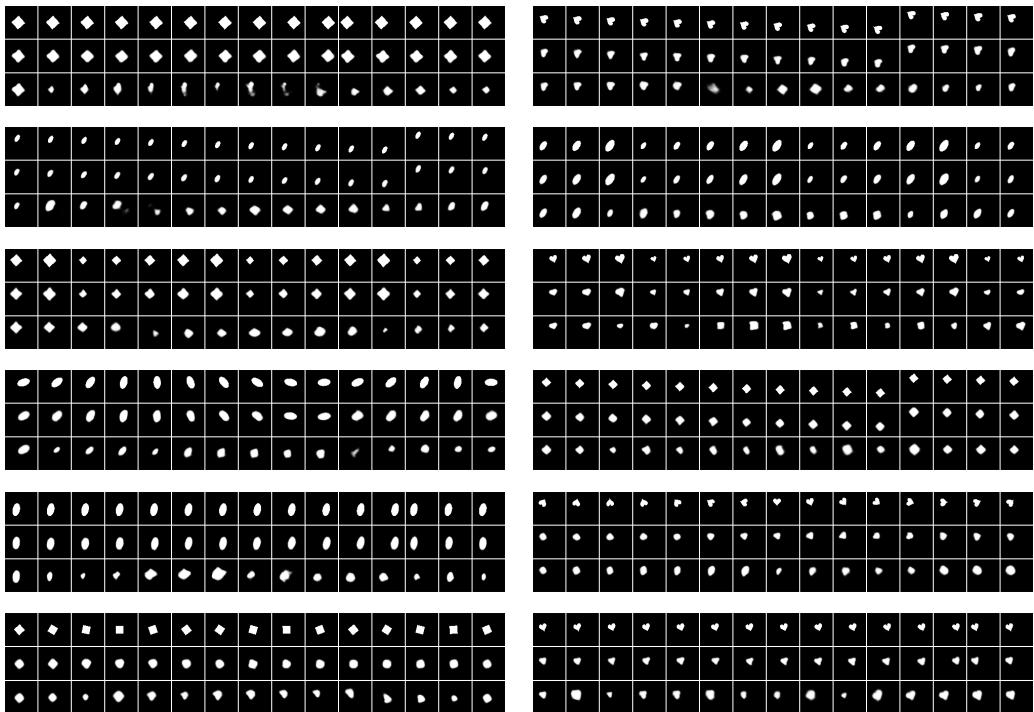


Figure 16: dSprites BubbleVAE  $L = \frac{1}{6}S, K = 2L$

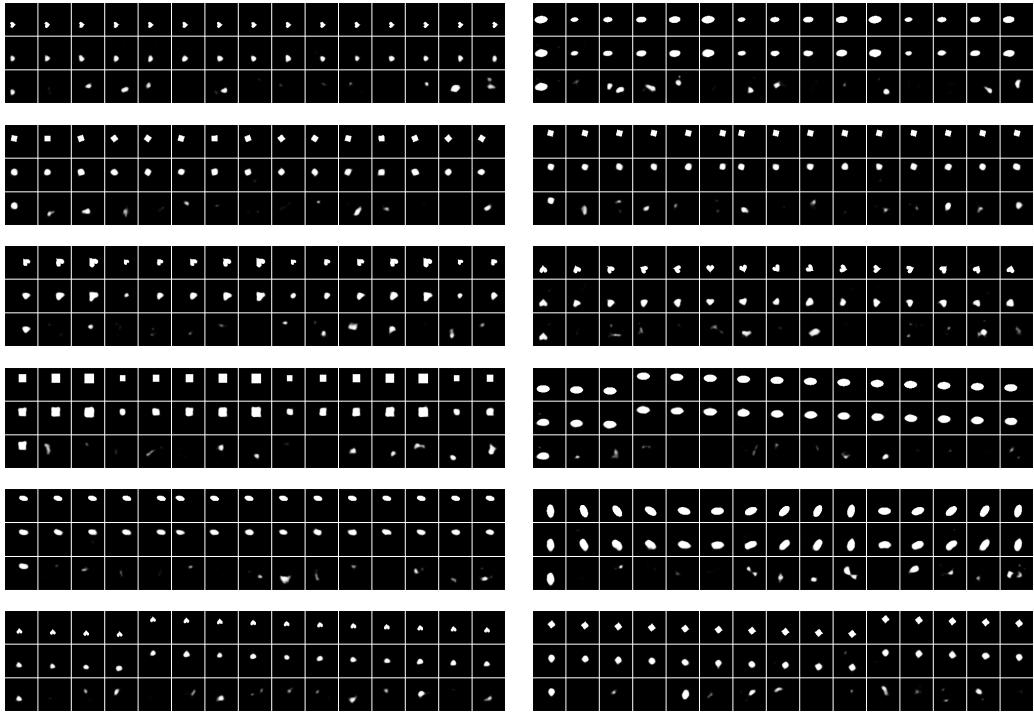


Figure 17: dSprites VAE  $L = 0, K = 1$

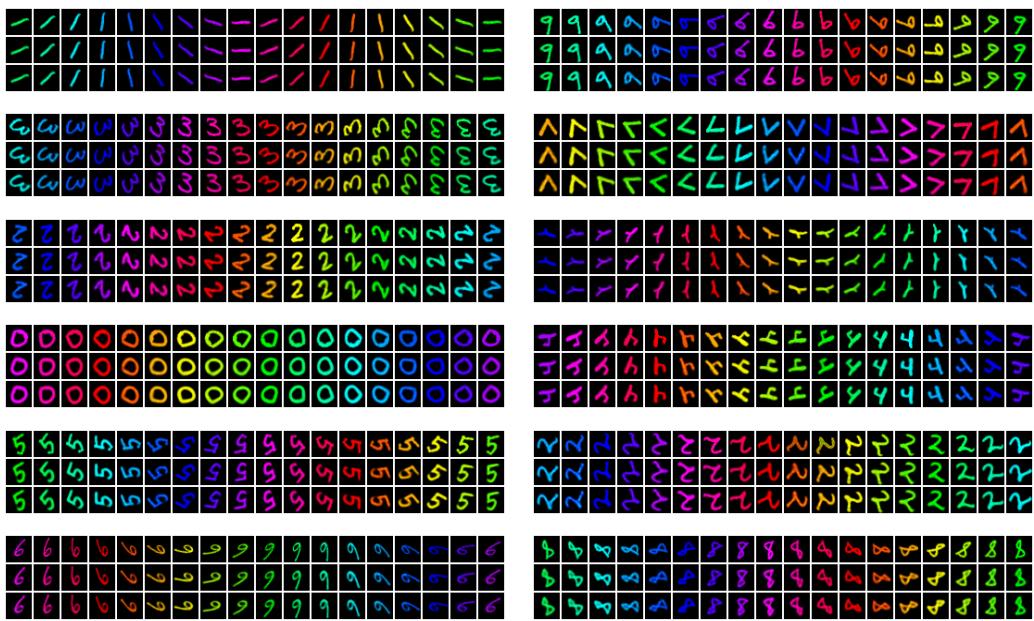


Figure 18: Combined Color & Rotation MNIST TVAE  $L = \frac{13}{36}S, K = 3$

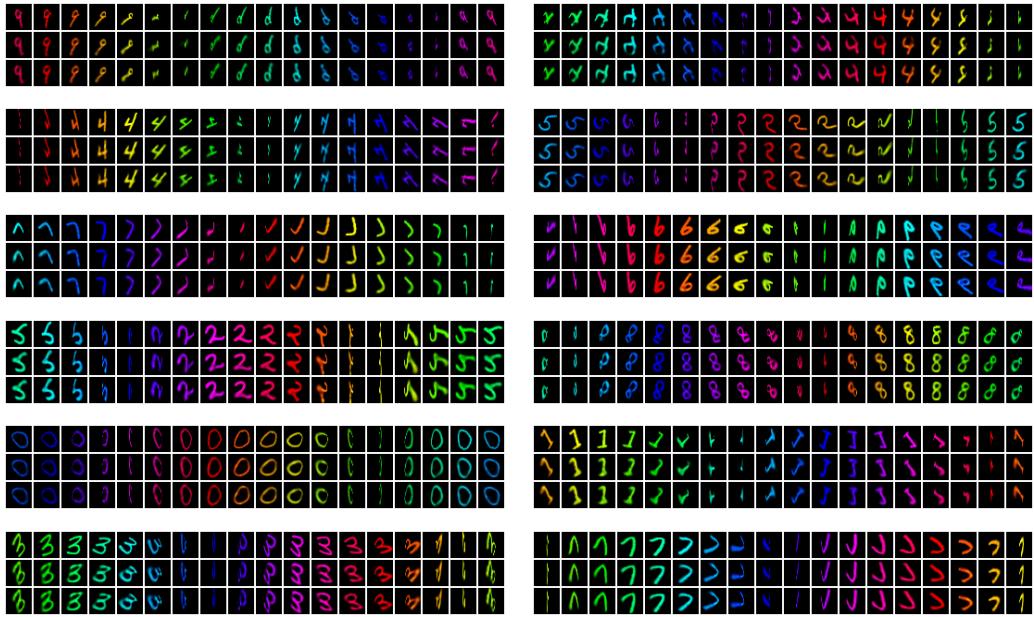


Figure 19: Combined Color & Perspective MNIST TVAE  $L = \frac{13}{36}S$ ,  $K = 3$