# THE SKY IS THE LIMIT WHEN CLUSTERING IS EQUATED WITH DISENTANGLEMENT

#### Anonymous authors

Paper under double-blind review

### ABSTRACT

Disentangled representation learning allows data to be mapped to a latent space where factors of variation can be individually manipulated. These factors define a direct notion of similarity between observations that naturally groups them into clusters with shared factors of variation. While this has been empirically shown to be effective on simple datasets, it is unclear how or when complex real-world data can be disentangled into representations that allow the same degree of manipulation and clustering. To advance the field of disentangled representation learning and clustering, we provide a new theoretical perspective by equating disentanglement with clustering by using factors of variation as a measure of element-wise similarity. This leads to a simple yet important observation: Instead of explicitly clustering the elements of a dataset, we can implicitly cluster them by learning to represent and generate the elements of each cluster. Furthermore, this observation reveals that implicit clusters have a lower bound because (I) explicit clusters are a subset of implicit clusters, and (II) implicit clusters can generate novel elements not present in the finite dataset through combinatorial generalization. Building on these insights, we derive an implicit neural clustering approach based on identifying factors of variation in the latent space. We validate our findings through experiments on synthetic image data and empirical evidence from related stateof-the-art works. This demonstrates the practical relevance of our approach and promising potential for synthesizing complete datasets from limited data, addressing data distribution gaps, improving interpretability in cluster analysis, enhancing SSL and classification tasks, and reducing data storage space.

031 032 033

000

001

002 003 004

005

010 011

012

013

014

015

016

017

018

019

021

025

026

027

028

029

034

# 1 INTRODUCTION

037 Understanding and controlling the underlying factors of variation in data is central to disentangled 038 representation learning (Wang et al., 2022). Disentangled latent spaces not only group similar elements naturally into clusters (Ding et al., 2022) but also allow precise data manipulation when combined with a generator (Higgins et al., 2017). Recent advances in deep generative clustering 040 have shown the potential both to learn disentangled representations and cluster data apoints simulta-041 neously, enabling the generation of high-quality synthetic data (Chen et al., 2016; Mukherjee et al., 042 2019; Lee et al., 2020; Yu & Welch, 2021; Ding et al., 2022). These approaches move away from 043 traditional clustering algorithms, which rely on explicit partitioning based on learned or handcrafted 044 features, towards generative models that leverage disentangled representations. In this context, generative models for controllable image synthesis, such as GANs (Karras et al., 2020; Brock et al., 046 2019) and Diffusion models (Rombach et al., 2022; Croitoru et al., 2023), now produce synthetic 047 images realistic enough to improve downstream tasks like classification (Azizi et al., 2023; Fan 048 et al., 2023) and can help self-supervised learning (SSL) methods learn better general purpose embeddings (Chai et al., 2021; Jahanian et al., 2021; Tian et al., 2024). While prior work on deep generative clustering has focused primarily on improving clustering in the traditional sense, a more 051 implicit approach has the potential to synthesize full datasets and in turn fill gaps in data distributions, improve cluster interpretability, reduce storage needs, or enhance SSL embeddings and 052 classification tasks. This leads to our main research question: What if instead of explicitly clustering the elements of a dataset, we could represent and generate these clusters implicitly?

To answer this question, we put forward a simple sampling method from a disentangled latent space, which we call *Implicit Neural Clustering*. Rather than explicitly assigning data points to clusters, we implicitly represent and generate the clusters with controllable factors of variation.

With Implicit Neural Clustering, we can theoretically 058 equate clustering and generative models with disentanglement, which leads to an implicit neural perspective of clustering. Previous work by Zhao et al. (2020); Ding 060 et al. (2022) has already pointed out that using factors of 061 variation as a measure of similarity will naturally group 062 data into clusters. We move a step further and stress that 063 clusters emerge naturally from controlling factors of vari-064 ation in the latent space, i.e., control over a disentangled 065 semantic latent space inherently dictates cluster member-



Figure 1: Latent traversals of factors of variation changes cluster memberships.

ship (see Figure 1). This shift in perspective provides two key theoretical insights: (I) explicit clusters are a subset of implicit clusters, and (II) implicit clusters have a *lower bound*, which is given by the ability of an encoder to disentangle factors of variation and by the realism as well as combinatorial generalization of a controllable generator.

070 Implicit Neural Clustering samples and modifies encoded disentangled representations of elements 071 in a dataset with atomic group actions to implicitly generate clusters. On top of the group-based 072 definition of disentangled representations by Higgins et al. (2018), we define an atomic group ac-073 tion as a partition of a latent traversal direction. After an atomic group action modifies a latent 074 representation, a generator produces a data point reflecting this change while keeping other factors 075 unchanged. Using disentangling variational autoencoder (VAEs), we show that atomic group actions exist and can be identified using Kernel Density Estimation (KDE) as a partitioning algorithm 076 on each dimension of the disentangled representations of dataset elements separately. Furthermore, 077 probing for atomic group actions leads us to an effective and simple qualitative measure for disentanglement, which is more informative than the commonly-used disentanglement visualization with 079 Hinton Matrices (Eastwood & Williams, 2018; Montero et al., 2022).

We conduct experiments with different unsupervised and semi-supervised VAE-based disentangled
representation learning methods in the well-known dSprites (Higgins et al., 2017), 3DShapes (Kim & Mnih, 2018), and MPI3D real (Gondal et al., 2019) datasets. Our experiments show the validity
of our findings by showing that atomic group actions can be identified in disentangled models and
used to implicitly cluster datasets. This demonstrates that the potential of synthesizing a full dataset,
even from a limited or incomplete dataset, seems promising, which also may help increase the
interpretability of cluster analysis and reduce storage space.

The main contributions of this paper are: (i) We introduce *Implicit Neural Clustering*, a simple sampling procedure that allows to implicitly cluster a dataset. (ii) Based on *Implicit Neural Clustering*, we provide a theoretical analysis of what happens when we equate clustering with disentangled representation learning, which leads to the discovery of a *lower bound* to clustering. (iii) We provide a practical implementation of *Implicit Neural Clustering* for disentangling VAEs and validate our approach through experiments on multiple datasets. (iv) We show an effective qualitative measure for disentanglement, which is more informative than the commonly-used disentanglement visualization with Hinton Matrices Eastwood & Williams (2018); Montero et al. (2022).

096

104 105

# 2 EQUATING CLUSTERING WITH DISENTANGLEMENT

Traditional explicit clustering is defined by a partition function  $C_{sim}$  that partitions an input dataset  $\mathcal{D}$ , under an arbitrary notion of similarity sim, into k clusters.  $C_{sim}$  either maps any  $x \in \mathcal{D}$  to a hard cluster assignment  $C : \mathcal{D} \to \mathbb{N}$  (e.g., k-means) or soft cluster assignment  $C : \mathcal{D} \to \mathbb{R}^k$  (e.g., maximum likelihood). Hard clustering can be defined as applying  $C_{sim}$  on each  $x \in \mathcal{D}$ , which yields k disjoint subsets  $\mathcal{D}_k^{sim} \subseteq \mathcal{D}$ :

$$\mathcal{D} = \bigcup_{k} \mathcal{D}_{k}^{sim} = \bigcup_{k} \left\{ x \mid x \in \mathcal{D} \land C_{sim}(x) = k \right\}$$
(1)

with  $\bigcup_k \mathcal{D}_k^{sim} = \mathcal{D}$  and  $\bigcap_k \mathcal{D}_k^{sim} = \emptyset$ . A different similarity  $sim' \neq sim$  formally defines any arbitrary clustering different from sim over  $\mathcal{D}$ , which is the basis for multi-partition clustering (Galimberti & Soffritti, 2007). In the multi-partition clustering context, explicit clustering w.r.t. sim



Figure 2: We visualize the difference between clustering in the traditional explicit sense and *Implicit Neural Clustering*. Under different measures of similarity, explicit clustering can cluster the dataset correctly in three different partitionings. However, since not all possible combinations between factors of variation are observed in the data, certain combinations are not present in the final clusters because we only *explicitly cluster* the real data. In contrast, *Implicit Neural Clustering* leads to *implicit clusters* that can also include novel crosscombinations not observed in the dataset.

yields only one out of many possible clusterings of the data. For multi-partition clustering, *sim* can
be considered in two ways. One, where *sim* corresponds to clustering over different sub-dimensions
of the feature representation, leading to different clustering partitions (Zhang, 2004; Galimberti &
Soffritti, 2007; Vandewalle, 2020; Rodriguez-Sanchez et al., 2022; Falck et al., 2021; Willetts et al.,
2019; de Chaumaray & Vandewalle, 2023). Two, where based on representation learning, one could
train a different feature extractor for each possible *sim*.

129 130

131

122

## 2.1 DEFINITION OF DISENTANGLED REPRESENTATIONS

Implicit Neural Clustering builds on top of the established symmetry group-based definition for 132 disentangled representations by Higgins et al. (2018). We briefly introduce essential parts needed 133 from this definition. Let G be a symmetry group acting on a set of world states (ground truth factors 134 of variation) W, and let O be a set of observations (e.g., pixel space) and Z the internal agent 135 representation of W. A generative process  $b: W \to O$  leads from world to observation states, and 136 an inference process  $h: O \to Z$  leads from observation to an agent's internal representation of W. 137 In this context, we have a dataset  $\mathcal{D} = \{o_1, ..., o_N\}$  of observations  $o_i \in O$ . We now define the 138 inference process  $h: O \to Z$ , as a parameterized feature extractor  ${}^1 h_{\varphi}: O \to \mathcal{F}$  with parameters 139  $\varphi$ , which yields a disentangled representation  $\mathcal{F}$  of any  $o \in \mathcal{D}$ . Under the assumption that G can be 140 decomposed into a direct product  $G = G_1 \times ... \times G_M$ , the representation  $\mathcal{F}$  is disentangled with 141 respect to G provided that the following conditions are satisfied. (i) There are group actions that act 142 on  $\mathcal{F}, \cdot : G \times \mathcal{F} \to \mathcal{F}$ , (ii) There is a mapping  $d : W \to \mathcal{F}$ , which is equivariant between the actions of G on W and Z:  $g \cdot d(w) = d(g \cdot w), \forall g \in G, \forall w \in W$ , and (iii)  $\mathcal{F}$  decomposes into its factors of 143 variation  $\mathcal{F} = \mathcal{F}_1 \times ... \times \mathcal{F}_M$  so that any  $\mathcal{F}_i$  is only affected by  $G_i$  and invariant to any  $G_j, \forall j \neq i$ . 144 Finally, we assume to have access to a parameterized generator  $G_{\theta} : \mathcal{F} \to O$  with parameters  $\theta$  that 145 transforms samples from the disentangled representation space  $\mathcal{F}$  to the observation space O. 146

147 148

161

## 2.2 FROM EXPLICIT TO IMPLICIT NEURAL CLUSTERING

149 In contrast to *explicitly* clustering a dataset  $\mathcal{D}$  under various sim, Implicit Neural Clustering is de-150 rived from a *disentangled* representation  $\mathcal{F}$  of  $\mathcal{D}$ . As an initial intuition, if we assume that any  $o \in \mathcal{D}$ 151 can be decomposed into its factors of variation, we can impose specific changes to any  $o \in \mathcal{D}$  by 152 modifying the desired parts of the factor in the representation. Figure 2 provides an overview of 153 Implicit Neural Clustering with its main differences to explicit clustering. Following the definition 154 of implicit probabilistic models, *Implicit Neural Clustering* can be defined as a sampling procedure 155 from a disentangled latent space. Different from implicit models where a parameterized generator  $G_{\theta}(\cdot)$  (e.g., GAN) transforms samples from an analytic distribution (e.g., isotropic Gaussian) 156 to synthetic examples (Li & Malik, 2018), Implicit Neural Clustering transforms samples from a 157 disentangled distribution  $\mathcal{F}$  into synthetic clusters. 158

More specifically, for each cluster  $\mathcal{D}_k^{sim}$ , there exists an implicit cluster that can be obtained by sampling from  $\mathcal{F}$  while fixing one respective factor of variation  $\mathcal{F}_i$ . Let  $G = G_1 \times ... \times G_M$ 

<sup>&</sup>lt;sup>1</sup>We change Z to  $\mathcal{F}$  for notation and readability reasons.

162 be the group actions that act on  $\mathcal{F}$ , and let  $\cdot : G \times \mathcal{F} \to \mathcal{F}$  be the action that changes  $\mathcal{F}$  to the 163 respective factor  $\mathcal{F}_i$ . Given that each factor of variation has of several *atomic* attributes  $\mathcal{F}_{ik}$  (e.g., 164 the class labels of shape or color), we precisely define fixing a factor of variation as follows: Each  $G_i$ 165 consists of *atomic* partitions  $G_i = \{G_{i1}, G_{i2}, ...\}$  that can modify any  $z \in \mathcal{F}$  and are parameterized by a single value  $f \in \mathbb{R}$  or a parameterized distribution  $P(f \mid G_{ik})$ . We further define a function 166 167

 $\stackrel{G_{ik}}{=}$  that yields "true" if an atomic factor of variation  $\mathcal{F}_{ik}$  is present in  $z \in \mathcal{F}$ . 168

When clustering with respect to a factor of variation  $\mathcal{F}_i$ , let  $sim \equiv \mathcal{F}_i$ , equating disentanglement 169 with clustering. Based on Equation 1 and the atomic partitions of a factor of variation, explicit 170 clustering  $C_{\mathcal{F}_i}$  splits  $\mathcal{D}$  into  $|G_i|$  disjoint subsets. In the *implicit* case, together with the generator 171  $G_{\theta}$  and the feature extractor  $h_{\varphi}$ , we can generate each cluster  $\mathcal{D}_{k}^{sim} \approx \mathcal{D}_{if}'$  implicitly to generate a 172 synthetic version  $\mathcal{D}'_i$  of the original dataset D. 173

176

$$\mathcal{D} \approx \mathcal{D}'_i = \bigcup_{f \in G_i} \mathcal{D}'_{if} = \bigcup_{f \in G_i} \{ G_\theta(h_\varphi(o)) \mid o \in \mathcal{D} \land \stackrel{f}{=} (h_\varphi(o)) \}$$
(2)

177 We have that  $\mathcal{D}'_i$  implicitly models  $\mathcal{D}$  with respect to a clustering under a factor of variation  $\mathcal{F}_i$ , under 178 the assumption that the encoder  $h_{\varphi}$  is capable of disentangling and  $G_{\theta}$  is capable of realistically 179 reconstructing the encoded elements. More specifically,  $h_{\varphi}$  must disentangle any  $o \in \mathcal{D}$  w.r.t.  $\mathcal{F}_i$ , so that  $sim \equiv \mathcal{F}_i$ , and  $G_\theta$  must recover an o' from this representation so that  $o' \approx G_\theta(h_\varphi(o))$ . 181 To move beyond *explicit* clusters, we further assume the disentangled representation space  $\mathcal{F}$  to be composable, i.e., we can modify any  $z \in \mathcal{F}$  by acting with the atomic group action  $G_{if}$ , which 182 changes cluster membership under factor  $\mathcal{F}_i$  from any previous  $D'_{il}$  to  $D'_{if}$ ,  $l \neq f$ , or produce a variation of any  $o \in \mathcal{D}$  by acting with the atomic group action  $G_{if}$  of the same cluster on z. 183 Together with the compositionality assumption, Equation 2, the disentanglement assumption on  $h_{\omega}$ , 185 and the reconstruction assumption on  $G_{\theta}$ , we derive *Implicit Neural Clustering*. 186

187 188

189 190

215

 $\mathcal{D} < \approx \mathcal{D}'_i = \bigcup_{f \in G_i} \{ G_{\theta}(z) \mid z \in \{ \cdot (f, z_1), ..., \cdot (f, z_K) \} \sim \mathcal{F} \}$ (3)

where  $(\cdot) \sim \mathcal{F}$  denotes a sampling procedure for each factor of variation  $\mathcal{F}_i$  and K denotes the 191 number of elements to be sampled. For *Implicit Neural Clustering* we fix a factor of variation  $\mathcal{F}_i$ , 192 sample latent representations  $z \in \mathcal{F}$  using the sampling procedure  $(\cdot) \sim \mathcal{F}$ , take/sample a respective 193 group action<sup>2</sup> f of an atomic factor of variation  $\mathcal{F}_i$ , and modify each z accordingly with (f, z). Up 194 to the capabilities of the encoder  $h_{\varphi}$  and generator  $G_{\theta}$ ,  $\mathcal{D} \ll \mathcal{D}'_i$  emphasizes that the set  $\mathcal{D}'_i$ , 195 obtainable from the outlined procedure, can at least *implicitly* represent the original dataset  $\mathcal{D}$  as 196 a lower bound. The lower bound<sup>3</sup> to Equation 3 is given by the reconstruction of the dataset in 197 Equation 2. Therefore, up to the capabilities of the encoder  $h_{\varphi}$  and generator  $G_{\theta}$ , *Implicit Neural Clustering* is able to (a) generate a variety of realistic data by sampling arbitrary data compositions 199 and (b) synthesize novel examples in each cluster not observed in the dataset when  $G_{\theta}$  is capable 200 of combinatorial generalization. Under the respective definition of disentangled representations in Section 2.1 and Equations 1, 2, and 3, the resulting synthetic dataset  $\mathcal{D}'_i$  is partitioned into disjoint 201 subsets w.r.t. a fixed factor of variation  $\mathcal{F}_i \equiv sim$ . In this way, we define an implicit clustering 202 of  $\mathcal{D}$ , where clustering is equated with disentanglement, and clusters are generated implicitly by a 203 generative model controllable by disentangled factors of variation. 204

205 Obtaining atomic group actions in disentangling VAEs. We provide a simple procedure to iden-206 tify atomic group actions in the latent space of disentangling VAEs, which often encode factors of variation in only one dimension l of the representation  $z = (z_1, z_2, ..., z_d) \in \mathbb{R}^d$ , d > M. To 207 obtain the atomic group actions, we first encode the full dataset and then partition each dimension 208 using kernel density estimation (KDE) at local minima of the resulting density estimates<sup>4</sup>. This 209 leads to density-based partitions, which naturally arise in the latent space of disentangling VAEs 210 and each partition is a parameterized probability distribution  $P(f|G_{ik})$  (e.g., uniform or normal), 211

<sup>212</sup> <sup>2</sup>In practice, we would parameterize f with a probability distribution and sample the respective modification 213 for more variety, but a single value, like the mean over all possible values, would also work. 214

<sup>&</sup>lt;sup>3</sup>Extending on this proof sketch, we provide a proof in Appendix A1

<sup>&</sup>lt;sup>4</sup>Any partition algorithm could be used. KDE has the advantage over, e.g., k-means that we do not specify the number of partitions in advance.

which resembles a probability distribution over atomic group actions for  $\mathcal{F}_{ik}$  that we can sample from. However, it is important to point out that these simple partitions are only meaningful when a factor of variation is properly disentangled.

# 219220 Sampling Procedure for Implicit Neural Clustering.

221 Given atomic group actions, the specific sampling procedure Implicit Neural Clustering is defined in lines 5-8 in Algorithm 1. First, we sample a random value from 222 the partition distribution  $f \sim P(f|G_{ik})$ . Next, we sample a random latent  $z \in \mathcal{F}$ . This process can be done in two 224 **Input:** Group actions  $G_i$  for factor of variation  $\mathcal{F}_i$ , ways. (1) Sample only from the 225 Generator  $G_{\theta}$ , number of samples K found partitions (AS) or (2) sample 226 **Output:** Implicit clustering  $\mathcal{D}'$  with respect to  $\mathcal{F}_i$ z from the set of all encoded data-1  $\mathcal{D}'_i \leftarrow \emptyset$ 227 points  $h_{\alpha}(\mathcal{D})$  ( $\neg AS$ ). Afterward, we 2 for  $k in 1..|G_i|$  do 228 act with f on z, i.e., (f, z), which  $D_{ik} \leftarrow \emptyset$ 3 229 modifies z accordingly. Acting can for 1..K do 4 230 also occur with multiple actions at the  $f \sim P(f|G_{ik})$ 5 231 same time (e.g., when a factor is dis $z \sim AS(\mathcal{F})$  or  $\neg AS(\mathcal{F})$ 6 232 entangled across two dimensions). In  $\begin{vmatrix} z' = \cdot (f, z) \text{ or } MDA(f, z) \\ D'_{ik} \leftarrow D'_{ik} \cup \{G_{\theta}(z')\} \end{vmatrix}$ 7 233 the latter case, we denote  $\cdot$  () as Multi-8 Dimensional Action (MDA). Finally, 234 end 9 the process in lines 5-8 is repeated 235  $D'_i \leftarrow D'_i \cup D'_{ik}$ 10 K times for each cluster  $\mathcal{D}_{ik}$ . Note 236 11 end that it is straightforward to achieve an 237 Algorithm 1: Our sampling procedure implicit multi-partition clustering by 238 simply repeating Algorithm 1 with different  $\mathcal{F}_i$ . 239

Identifying meaningful atomic group actions in disentangling VAEs. When ground truth factors of variation are available, we can identify meaningful disentangled atomic group actions by computing if they "uniquely" co-occur with known ground truth atomic factors of variation  $\mathcal{F}_{ij}$ . To this end, we count the frequency of co-occurences between all extracted KDE partitions  $pt_a$ ,  $a \in \mathbb{N}$ and atomic ground truth factors of variation  $\mathcal{F}_{ij}$  of each  $o \in \mathcal{D}$ . This leads to matrix with the factors of variation as rows and the partitions as columns. For each cell [(i, j), a] in row (i, j), we divide the frequency  $freq(\mathcal{F}_{ij}, pt_a)$  by the sum over all frequencies of the row (i, j), which leads

to:  $[(i, j), a] = \frac{freq(\mathcal{F}_{ij}, pt_a)}{\sum_b freq(\mathcal{F}_{ib}, pt_b)}$ . Using this method, we can visualize disentanglement in a qual-

itative manner as shown in Figure 3, where unique co-occurences between a partition and a factor of variation show meaningful disentanglement.

To identify meaningful atomic group actions in the unsupervised case, where no information about ground truth factors of variation is available, one can fall back (i) to a general-purpose zero-shot classifier like CLIP (Radford et al., 2021) to predict if a factor of variation is consistently present in batch of randomly modified images with the same KDE partition, or (ii) use general-purpose embeddings from, e.g., DINOv2 (Oquab et al., 2023) to find a significant image similarity increase when a batch of random images is modified with the same KDE partition.

256 257 258

## **3** EXPERIMENTAL EVALUATION

259 In this evaluation, we systematically increase the number of assumptions made on the encoder  $h_{\varphi}$ 260 and generator  $G_{\varphi}$  that must be fulfilled for *Implicit Neural Clustering* to be valid. In the first set of 261 experiments, we evaluate *Implicit Neural Clustering* with the assumption A.H that an encoder  $h_{\varphi}$ 262 disentangles the factors of variation in atomic partitions. In this setting, we propose an interpretable 263 procedure to quantify and assess qualitatively when  $A \cdot H$  is satisfied. Furthermore, for cases where 264 A.H is satisfied, we quantitatively and qualitatively evaluate the quality of the generated implicit 265 clusters, showing that the assumptions made on the generator  $G_{\varphi}$  regarding realistic sample gener-266 ation (A.G1) and compositionality of atomic partitions for controllable synthesis (A.G2) are valid. Afterward, we evaluate the partition performance using KDE and compare it against traditional 267 clustering algorithms, empirically validating the lower bound to Implicit Neural Clustering. Fur-268 thermore, we perform an ablation study on the sampling procedure regarding sampling type AS and 269 acting type MDA to show the impact of disentanglement on the overall performance. Finally, our



272

275

281

283 284

287

289 290 291

292

293

Figure 4: Co-occurence plot between atomic factors of variation and the dimension-wise partitions

ADA-GVAE, DCI=0.989

3DShapes

β-TCVAE, DCI=0.588

last experiment discusses the validity, strengths, and limitations of *Implicit Neural Clustering* under combinatorial generalization (assumption A.G3) of  $G_{\theta}(\cdot)$  in light of the results of existing works.

Experimental Setup. To show the validity and limitations of our concept, we consider 3DShapes (Kim & Mnih, 2018), MPI3D Real (Gondal et al., 2019), and dSprites (Higgins et al., 295 2017), which are widely adopted datasets for disentangled representation learning. Specifically, 296 we evaluate learning disentangled representations in an unsupervised setting with  $\beta$ -TCVAE (Chen 297 et al., 2018), and in a weakly-supervised setting with ADA-GVAE (Locatello et al., 2020) and the 298 architecture of (Montero et al., 2020; 2022), which we refer to by Comp-WAE. We use these models 299 as a means to show the validity of our concept due their strong disentanglement performance w.r.t. 300 the DCI disentanglement metric (Eastwood & Williams, 2018) and their shown compositionality ca-301 pabilities. The DCI disentanglement metric measures the degree of capturing at most one generative 302 factor for each latent variable. In all experiments<sup>5</sup>, all models are trained from scratch on a single 303 46GB RTX A6000.

304 **E0 w/** A.H: 305

We propose a simple but effective qualitative 306 evaluation procedure for our very strict disen-307 tanglement requirement into atomic partitions. 308 As presented in Figure 3, we propose a visu-309 alization scheme that is more informative than 310 Hinton Matrices eastwood2018framework, montero2022lost). Compared to a Hinton 311 Matrix (left), our method provides much more 312 details on disentanglement without the need for 313 a classifier, which makes it a complementary 314 visualization tool to Hinton Matrices and 315 DCI (proposed in (Eastwood & Williams, 316 2018)) for assessing disentanglement. Figure 4 317 compares several unsupervised models with 318



MPI3D-real

 $\beta$ -TCVAE, DCI=0.311

Comp-WAE, DCI=0.99

0.10

0.06

0.04

0.02



Figure 3: Comparing a Hinton Matrix against atomic group actions. Atomic group actions in VAEs are well disentangled and can be used to control the generation. weakly-supervised approaches under DCI and shows the corresponding qualitative co-occurrence plot between found atomic partitions and ground truth factors of variation. We find overall that we are always able to consistently generate the implicit clusters based on factors of variation that can be seen in Figure 4 as clean and unique partitions. Furthermore, a lower DCI score indicates worse disentanglement that correlates with our qualitative measure.

319

320

<sup>322</sup> 323

<sup>&</sup>lt;sup>5</sup>Code, models, and all details to reproduce our experiments will be publicly available upon acceptance.

			$F1 \text{ Macro}(\uparrow)$												
Dataset	Approach	DCI	$\mathcal{F}_1$	$\mathcal{F}_2$	$ \mathcal{F}_3 $	$ \mathcal{F}_4 $	$\mathcal{F}_5$	$\mathcal{F}_6$	$\mathcal{F}_7$						
dSprites	CompWAE	0.999	0.422	0.945	0.118	0.619	0.698	-	-						
uspines	Oracle	-	0.99	0.99	0.92	0.65	0.70	-	-						
3DShapes	ADA-GVAE	0.99	0.73	0.64	0.46	0.95	0.96	0.98	-						
JDShapes	Oracle	-	1.0	1.0	1.0	1.0	1.0	1.0	-	_					
MPI3D	CompWAE	0.99	1.0	0.49	0.98	1.0	1.0	0.28	0.85	-					
WII ISD	Oracle	-	1.0	0.98	0.99	1.0	1.0	0.90	0.99						
			Color												
	Shape														
											2				
- 👞 👞															
				10.00											
	-										6				
			<b>T</b>				Size								
											Ŀ				
				100											

Table 1: Quality of generated samples in different datasets with different models.

Figure 5: Exemplary *Implicit Neural Clustering* of 3DShapes. Each row represents random samples for some atomic factor of variation. Each row is the result of applying atomic group actions we extracted from the disentangled representation space to random samples.

From experiment E0, we find that unsupervised methods do not satisfy the necessary disentanglement requirement for our theoretical assumptions. For this reason, the remaining experiments of this section are performed only with weakly-supervised approaches, showing the validity of our theoretical findings.

352 E1 w/ A.H, A.G1, A.G2: Quantifying Realism of the Generated Samples. To evaluate the 353 quality of the generated samples, we train a CNN as an oracle on the images of the real dataset to 354 predict which factor of variation is present in each image. For each  $\mathcal{F}_i$ , we (1) split the dataset into 355 training and test sets (random 0.8/0.2), (2) train on the training split of the real data for 15 epochs 356 to predict  $\mathcal{F}_i$ , (3) generate an *Implicit Neural Clustering* of the dataset w.r.t. a factor of variation 357  $\mathcal{F}_i$  with 10,000 samples for each atomic factor, and (4) evaluate the classifier on the real test split 358 as well as on the synthetically generated dataset. In step (3), we run the sampling procedure from Algorithm 1 with MDA and  $\neg AS$  as components. We found this combination to perform the 359 best (see ablation study in Experiment E4). The results are shown in Table 1, where we report the 360 macro F1 score over all atomic factors of each  $\mathcal{F}_i$ . For very low scores, some implicit clusters can be 361 generated well while others fail (full confusion matrices are provided in the supplementary material). 362 We notice that MPI3D and Shapes3D can be implicitly clustered much better than dSprites. To 363 summarize, most of the generated samples can be predicted accurately compared to the real data, 364 showing that the generated samples are (1) realistic because there is only small drop in performance on the oracle performance, and (2) can be synthesized by acting with the atomic partitions, showing 366 the assumption of compositionality is fulfilled. The quantitative results relate to the qualitative 367 evaluation procedure in Figure 4, in which the factors of variation that are "atomically" disentangled 368 provide the highest F1 scores. Analogously, the factors that do not have unique co-occurences with 369 atomic partitions or that span multiple dimensions exhibit lower quality in the generated samples.

370 E2 w/ A.H, A.G1, A.G2: Qualitative Evaluation of Implicit Neural Clustering. Figure 5 shows 371 three coherent synthetic multi-partition clusters of 3DShapes with respect to shape, color, and size 372 that we have *implicitly* clustered with our concept. We obtained these results by applying to random 373 samples the atomic group actions that we were able to extract from the disentangled representations 374 of the underlying dataset. In Figure 6 we present the modifications of arbitrary samples based on 375 the atomic group actions, which demonstrates that randomly generated samples can be modified to the desired atomic factor the variation. It is relevant to note that although we are able to find 376 atomic group actions for all ground truth labels, not all of them are invariant to the other factors 377 of variation, which demonstrates a limitation in the disentanglement for some factors of variation.

349

350

351

324



Figure 6: For random generated samples, atomic group actions specifically modify a certain factor of variation, e.g., change the object shape, color, or size. Exemplary demonstration in 3DShapes.

Table 2: Comparing partitioning performance of KDE against k-Means and GMM.

ACC,NMI (†)	J	<b>F</b> <sub>4</sub>	J	r <sub>5</sub>	${\cal F}$	6
<i>k</i> -means (requires number of partitions)	0.87	0.92	0.75	0.80	1.0	1.0
GMM (requires number of partitions)	0.87	0.92	0.75	0.82	0.99	0.99
KDE (non-parametric)	0.99	0.97	1.0	1.0	0.99	0.99

393

394

382

384

In the supplementary material, we provide the full implicit clusters for all three datasets and latent traversals with atomic group actions for MPI3D and 3DShapes.

E3: On Clustering Performance and the Lower Bound. Using the 3DShapes and the bestperforming factors from Table 1, we compare KDE against k-means and a Gaussian Mixture Model 396 (GMM) to show its effectiveness in partitioning atomic group actions in each dimension. At the 397 same time, this experiment allows to evaluate the clustering performance and empirically validate 398 the *lower bound* to clustering given by Equation 2. When a factor of variation is used as the measure 399 of similarity, then under a high degree of disentanglement, implicit clusters are equivalent to explicit 400 clustering. We set the number of clusters in k-means and GMM to the ground truth number of atomic 401 factors corresponding to each dimension, which we identify with our qualitative co-occurrence dis-402 entanglement measure. KDE does not require any number of clusters as a parameter. We use the 403 commonly used Purity (ACC) and Normalized Mutual Information (NMI) as metrics. The results 404 are shown in Table 2, where we can see that using KDE to partition a dimension outperforms both 405 k-Means and GMM in all factors. In addition, another advantage of using our KDE approach is 406 that it is non-paramatric, i.e., does not require specifying the number of clusters in advance. Together with the realism of generated samples in Table 1, the implicit clusters in Figure 5, and the 407 partitioning performance in Table 2, it becomes evident that we can implicitly cluster a dataset by 408 controllable atomic factors of variation. That is, because (i) implicit clusters in 2 are equivalent to 409 explicit clustering due to reconstruction (Equation 2), and (ii) randomly generated samples for each 410 cluster always include the factor of variation (Equation 3). These results empirically show the *lower* 411 bound for clustering. 412

E4: Ablation Study, Sampling Procedure. We test two different kinds of sampling strategies using 413 the 3DShapes dataset. First, sampling only from the atomic partitions (AS) or using random samples 414 from the encoded dataset ( $\neg AS$ ). Second, acting with multi-dimensional group actions (MDA) or 415 only acting with a single group action ( $\neg$  MDA). Table 3 shows the results for each combination. 416 We can see that  $\neg AS$  together with MDA lead to the best results. It is expected for AS to exhibit 417 low performance, since not all factors factors of variation are perfectly disentangled. Given the 418 qualitative disentanglement results, it is also expected that MDA should perform better, since some 419 factors are disentangled across multiple dimensions while still being unique combinations. This 420 result highlights that compositionality can even apply to the tested models when atomic actions 421 span multiple dimensions at the same time.

422 Note on our Results and Combinatorial Generalization (w/ A.H, A.G1, A.G2, A.G3). Among 423 other works, Montero et al. (Montero et al., 2022) have shown that combinatorial generalization 424 can be achieved in some special cases with learned disentangled representations. Despite these 425 special cases, there is no theoretical guarantee that proper disentanglement leads to combinatorial 426 generalization or that these special cases will transfer across different kinds of models. Dividing 427 the problem in two, i.e., disentangling first, and training a separate generator on the disentangled 428 representations afterwards potentially leads to better results. Building on the theory of our concept and the empirical results of our experiments, we can straightforwardly apply our concept in a setting 429 where  $G_{\theta}$  has learned combinatorial generalization (A.G3). In this setting, our experiments close 430 the loop in the conceptual illustration of our concept in Figure 2, where  $G_{\theta}$  will "fill-in-the-blanks" 431 by synthesizing cross-combinations between factors of variation not seen in the dataset. In the same

432 Table 3: Ablation study regarding the impact of sampling procedure on the performance of the 433 generated implicit clusters. We evaluate predicting the factor of variation using classifiers trained on 434 the real data in all generated implicit clusters.

				F1 Ma	$cro(\uparrow)$			
w/ AS	w/ MDA	$\mathcal{F}_1$	$\mathcal{F}_2$	$\mathcal{F}_3$	$\mathcal{F}_4$	$\mathcal{F}_5$	$\mathcal{F}_6$	
		0.43	0.39	0.52	0.88	0.98	0.98	
		0.73	0.64	0.46	0.95	0.96	0.98	
$\checkmark$		0.36	0.44	0.20	0.74	0.91	0.81	
	$\checkmark$	0.68	0.56	0.38	0.74	0.95	0.79	
Grou	nd Truth ${\cal D}$	1.0	1.0	1.0	1.0	1.0	1.0	

context, our concept can be applied to any kind of models trained on ground truth factors of variation that have learned combinatorial generalization (e.g., with conditional or composable diffusion 444 models (Okawa et al., 2024; Liu et al., 2022)) to implicitly cluster the data. Further experiments to 445 validate the aforementioned claims are provided in Appendix C and Appendix D. However, even in this relaxed setting, where the ground truth factors of variation are known in advance (which resembles an optimal encoder  $h^*$ ), there is still no theoretical guarantee for combinatorial generalization 448 in  $G_{\theta}$ , when trained on the factors of variation obtained by an optimal encoder  $h^*$ . Therefore, in line 449 with many previous works (e.g., (Montero et al., 2020; 2022; Okawa et al., 2024; Wiedemer et al., 450 2024b), it is important to emphasize that disentanglement and compositionality does not imply combinatorial generalization. We stress that our concept is constrained by the current limitations faced 452 by the field in relation to combinatorial generalization.

#### 453 454

451

443

446

447

#### 455 456

4

DISCUSSION

457 Our experiments show that *Implicit Neural Clustering* has a *lower bound* and is particularly lim-458 ited by the assumptions on  $h_{\alpha}$  regarding disentanglement (A.H), and  $G_{\theta}(\cdot)$  regarding realism of 459 synthetic samples (A.G1), compositionality (A.G2), and combinatorial generalization (A.G3). Es-460 pecially for disentanglement, we notice a huge gap between unsupervised and weakly-supervised 461 approaches. While combinatorial generalization is feasible in relaxed synthetic settings and to some extent in real-world data (Wiedemer et al., 2024a; Montero et al., 2022; Okawa et al., 2024), effec-462 tive methods for learning disentangled representations and achieving combinatorial generalization 463 from complex real-world data remain elusive. In real-world tasks, general-purpose embeddings 464 learned through SSL methods like SimCLR (Chen et al., 2020) or DINOv2 (Caron et al., 2021; 465 Oquab et al., 2023) can learn effective representations that disentangle real-world data to limited 466 extent. Empirical evidence of disentanglement in SSL representations, as shown by Bordes et al. 467 (2022), demonstrates that training a generator on SSL features allows for concept swapping in rep-468 resentations, producing samples that reflect these changes. While our approach naturally extends 469 to SSL representations and would make it applicable to datasets like CIFAR-10 (Krizhevsky et al., 470 2009) or ImageNet (Deng et al., 2009), atomic group actions do not exist in the simple form as in 471 disentangling VAEs. New methods to identify potential subspaces in these representations may lead to new insights. A different way to apply our method on real-world data would be through finding 472 a way that effectively partitions the interpretable directions learned by latent navigators from GANs 473 (e.g., (Voynov & Babenko, 2020; Georgopoulos et al., 2022)), or Diffusion Models (e.g., (Yang 474 et al., 2023)). Finally, it is easy to see that our sampling procedure can be easily applied to gen-475 erative models trained on ground truth factors of variation to synthesize datasets. We point to an 476 important closely related work by (Okawa et al., 2024), where a conditional diffusion model is 477 trained on ground truth factors of variation and combinatorial generalization is achieved. In this 478 context, our work provides valuable insights on a research question by Jahanian et al. (2021): "If we 479 have good enough generative models, do we still need datasets?" With Implicit Neural Clustering, 480 we can potentially generate any realistic synthetic variations of a dataset with a corresponding class 481 label, fill gaps in its distribution, and it can be a basis to replace datasets in order to save valuable cost 482 for storage and acquisition of data. Applications of *Implicit Neural Clustering* are not only limited 483 to datasets of images but can also be applied to completely different kinds of data, such as natural language, videos, or time series. Negative societal impact can occur when a model can achieve com-484 binatorial generalization under "full" disentanglement for, e.g, DeepFakes in the imaging or video 485 domain.

# 486 5 RELATED WORK

488 Existing *related works* on disentangled representation learning, deep generative clustering, and con-489 trollable image generation have shown the following points. (i) Factors of variation are embedded in 490 single (e.g., (Locatello et al., 2019; 2020; Wang et al., 2022)) or multiple dimensions (e.g., (Bordes 491 et al., 2022; Falck et al., 2021)) of a disentangled latent space, which can be learned with disen-492 tangled representation learning approaches that are unsupervised (e.g., VAE-based (Higgins et al., 2017; Kim & Mnih, 2018; Locatello et al., 2019; Falck et al., 2021), from pre-trained generative 493 models (Ren et al., 2022; Yang et al., 2023), deep-clustering (Mukherjee et al., 2019; Lee et al., 494 2020; Yu & Welch, 2021; Ding et al., 2022; Zhao et al., 2020)), or (weakly) supervised (e.g., (Hris-495 tov et al., 2018; Locatello et al., 2020; Montero et al., 2020; 2022; Wang et al., 2022)), (ii) Deep 496 generative clustering approaches (e.g., (Mukherjee et al., 2019; Lee et al., 2020; Yu & Welch, 2021; 497 Ding et al., 2022)) can simultaneously learn a disentangled latent space, cluster assignments, and 498 allows controllable generation of elements for each cluster with disentanglement. (iii) Learned dis-499 entangled representations are often composable and can be used to control the factors of variation 500 in images using latent traversal or the recombination/swapping of different latent dimensions be-501 tween images (e.g., (Bordes et al., 2022; Wang et al., 2022; Falck et al., 2021; Montero et al., 2020; 502 2022)). (iv) Both generative and disentangled representation learning models can learn combinato-503 rial generalization in rare synthetic and real-word settings (e.g., (Okawa et al., 2024; Montero et al., 504 2020; 2022; Wiedemer et al., 2024b)), which allows to synthesize novel cross-combinations between factors of variation not observed in the data. (v) Depending on the degree of disentanglement, 505 meaningful directions to traverse a disentangled latent space to control images can be straightfor-506 ward one-dimensional and linear (e.g., (Higgins et al., 2017; 2018)), or non-linear multi-dimensional 507 traversals can be learned from disentangled latent spaces in VAEs (e.g., (Ren et al., 2022; Yang et al., 508 2023)), GANs (e.g., (Voynov & Babenko, 2020; Georgopoulos et al., 2022)), or Diffusion Models 509 (e.g., (Yang et al., 2023)) in the form of a navigator. 510

Based on the points above, existing methods can already effectively control and sample data for
implicit cluster generation. However, they primarily focus on improving clustering performance,
disentanglement performance, controllable generation, realism, or combinatorial generalization in
isolation. Instead of treating these problems individually, with *Implicit Neural Clustering*, our work
takes a new unified perspective by equating disentanglement with clustering and generative models.
This reveals that existing methods are inherently limited by a clustering *lower bound*, given by
disentanglement capabilities, realism of generated samples, and combinatorial generalization.

517 518 519

# 6 CONCLUSION

520 In this paper, we present Implicit Neural Clustering, a sampling method for generating clusters im-521 plicitly through disentangled representations. Through theoretical analysis and empirical validation, 522 we show that equating disentangled representation learning with clustering and generative mod-523 els reveals that this method has *lower bound*, governed by the degree of disentanglement, realism 524 of generated samples, and combinatorial generalization in generative models. This lower bound 525 of Implicit Neural Clustering highlights strong potential for relevant future applications, such as 526 implicitly generating clusters driven by factors of variation in real datasets, synthesizing complete 527 datasets from limited data, improving interpretability in cluster analysis, enhancing SSL and clas-528 sification tasks, and reducing data storage needs. At last, in line with many prior works, we also 529 underscore the importance of focusing on combinatorial generalization in future research.

- 530 531
- 532
- 533
- 534
- 535
- 536
- 537
- 538

540 541	Reproducibility Statement
542	We commit to ensuring the reproducibility of our work as follows:
543	• All of our implementations, trained model checkpoints, hyperparameters, and all pacessary
544	• All of our implementations, trained model checkpoints, hyperparameters, and an necessary details to reproduce our empirical and qualitative results will be publicly available upon
545 546	acceptance.
547	• The source code of models by related work and the datasets used in this work are all pub-
548	licly available.
549	• The hardware setup used for our experiments (46GB NVIDIA RTX A6000 GPU) is de-
550	scribed.
551	
552	References
553	Shakoofah Azizi Simon Komhlith Chitwan Saharia Mahammad Narouzi and David I Eleat Sun
554 555 556	thetic data from diffusion models improves imagenet classification. <i>Transactions on Machine Learning Research</i> , 2023. ISSN 2835-8856. URL https://openreview.net/forum?
557	Id=DIRSOX JYPM.
558	Florian Bordes, Randall Balestriero, and Pascal Vincent. High fidelity visualization of what your
559 560	self-supervised representation knows about. <i>Transactions on Machine Learning Research</i> , 2022. ISSN 2835-8856. URL https://openreview.net/forum?id=urfWb7VjmL.
561	Andrew Brock Leff Donahue and Karen Simonyan Large scale GAN training for high fidelity
562	natural image synthesis. In International Conference on Learning Representations, 2019. URL
563	https://openreview.net/forum?id=B1xsqj09Fm.
564	Mathilda Caron, Hugo Touwron, Ishan Misro, Harvá Jágou, Julian Mairol, Piotr Rojanowski, and
565 566 567	Armand Joulin. Emerging properties in self-supervised vision transformers. In <i>Proceedings of the IEEE/CVF international conference on computer vision</i> , pp. 9650–9660, 2021.
568	Lucy Chai Jun Van Zhu Eli Shaahtman Dhillin Isala and Diahard Zhang. Encambling with door
569 570	generative views. In <i>Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern</i> <i>Recognition</i> pp. 14997–15007, 2021
571	Recognition, pp. 17997-18007, 2021.
572 573 574	Ricky TQ Chen, Xuechen Li, Roger B Grosse, and David K Duvenaud. Isolating sources of disen- tanglement in variational autoencoders. <i>Advances in neural information processing systems</i> , 31, 2018.
575 576 577	Ting Chen, Simon Kornblith, Mohammad Norouzi, and Geoffrey Hinton. A simple framework for contrastive learning of visual representations. In <i>International conference on machine learning</i> , pp. 1597–1607. PMLR, 2020.
578 579 580 581	Xi Chen, Yan Duan, Rein Houthooft, John Schulman, Ilya Sutskever, and Pieter Abbeel. Info- gan: Interpretable representation learning by information maximizing generative adversarial nets. <i>Advances in neural information processing systems</i> , 29, 2016.
582 583	Florinel-Alin Croitoru, Vlad Hondru, Radu Tudor Ionescu, and Mubarak Shah. Diffusion models in vision: A survey. <i>IEEE Transactions on Pattern Analysis and Machine Intelligence</i> , 2023.
584 585	Marie du Roy de Chaumaray and Vincent Vandewalle. Non-parametric multi-partitions clustering. <i>arXiv preprint arXiv:2301.02422</i> , 2023.
587 588 589	Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hi- erarchical image database. In 2009 IEEE conference on computer vision and pattern recognition, pp. 248–255. Ieee, 2009.
590 591	Fei Ding, Yin Yang, and Feng Luo. Clustering by directly disentangling latent space. In 2022 IEEE International Conference on Image Processing (ICIP), pp. 341–345. IEEE, 2022.
592 593	Cian Eastwood and Christopher KI Williams. A framework for the quantitative evaluation of disen- tangled representations. In <i>International conference on learning representations</i> , 2018.

594 595 596	Fabian Falck, Haoting Zhang, Matthew Willetts, George Nicholson, Christopher Yau, and Chris C Holmes. Multi-facet clustering variational autoencoders. <i>Advances in Neural Information Processing Systems</i> , 34:8676–8690, 2021.
597 598 599	Lijie Fan, Kaifeng Chen, Dilip Krishnan, Dina Katabi, Phillip Isola, and Yonglong Tian. Scaling laws of synthetic images for model training for now. <i>arXiv preprint arXiv:2312.04567</i> , 2023.
600 601	Giuliano Galimberti and Gabriele Soffritti. Model-based methods to identify multiple cluster struc- tures in a data set. <i>Computational statistics &amp; data analysis</i> , 52(1):520–536, 2007.
602 603 604 605	Markos Georgopoulos, James Oldfield, Grigorios G Chrysos, and Yannis Panagakis. Cluster-guided image synthesis with unconditional models. In <i>Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition</i> , pp. 11543–11552, 2022.
606 607 608 609	Muhammad Waleed Gondal, Manuel Wuthrich, Djordje Miladinovic, Francesco Locatello, Martin Breidt, Valentin Volchkov, Joel Akpo, Olivier Bachem, Bernhard Schölkopf, and Stefan Bauer. On the transfer of inductive bias from simulation to the real world: a new disentanglement dataset. <i>Advances in Neural Information Processing Systems</i> , 32, 2019.
610 611 612	Irina Higgins, Loic Matthey, Arka Pal, Christopher P Burgess, Xavier Glorot, Matthew M Botvinick, Shakir Mohamed, and Alexander Lerchner. beta-vae: Learning basic visual concepts with a constrained variational framework. <i>ICLR (Poster)</i> , 3, 2017.
613 614 615	Irina Higgins, David Amos, David Pfau, Sebastien Racaniere, Loic Matthey, Danilo Rezende, and Alexander Lerchner. Towards a definition of disentangled representations. <i>arXiv preprint arXiv:1812.02230</i> , 2018.
617 618	Jonathan Ho and Stefano Ermon. Generative adversarial imitation learning. Advances in neural information processing systems, 29, 2016.
619 620 621	Yordan Hristov, Alex Lascarides, and Subramanian Ramamoorthy. Interpretable latent spaces for learning from demonstration. In <i>Conference on Robot Learning</i> , pp. 957–968. PMLR, 2018.
622 623 624	Ali Jahanian, Xavier Puig, Yonglong Tian, and Phillip Isola. Generative models as a data source for multiview representation learning. In <i>International Conference on Learning Representations</i> , 2021.
625 626 627 628	Justin Johnson, Bharath Hariharan, Laurens Van Der Maaten, Li Fei-Fei, C Lawrence Zitnick, and Ross Girshick. Clevr: A diagnostic dataset for compositional language and elementary visual reasoning. In <i>Proceedings of the IEEE conference on computer vision and pattern recognition</i> , pp. 2901–2910, 2017.
629 630	Tero Karras, Samuli Laine, Miika Aittala, Janne Hellsten, Jaakko Lehtinen, and Timo Aila. Analyz- ing and improving the image quality of StyleGAN. In <i>Proc. CVPR</i> , 2020.
632 633	Hyunjik Kim and Andriy Mnih. Disentangling by factorising. In International conference on ma- chine learning, pp. 2649–2658. PMLR, 2018.
634 635 636	Thomas Kreutz, Max Mühlhäuser, and Alejandro Sanchez Guinea. Common sense initialization of mixture density networks for motion planning with overestimated number of components. In <i>The Second Tiny Papers Track at ICLR 2024</i> , 2024.
638 639	Alex Krizhevsky, Geoffrey Hinton, et al. Learning multiple layers of features from tiny images. 2009.
640 641 642	Wonkwang Lee, Donggyun Kim, Seunghoon Hong, and Honglak Lee. High-fidelity synthesis with disentangled representation. In <i>Computer Vision–ECCV 2020: 16th European Conference, Glas-</i> gow, UK, August 23–28, 2020, Proceedings, Part XXVI 16, pp. 157–174. Springer, 2020.
643 644 645	Ke Li and Jitendra Malik. Implicit maximum likelihood estimation. <i>arXiv preprint arXiv:1809.09087</i> , 2018.
646 647	Nan Liu, Shuang Li, Yilun Du, Antonio Torralba, and Joshua B Tenenbaum. Compositional visual generation with composable diffusion models. In <i>European Conference on Computer Vision</i> , pp. 423–439. Springer, 2022.

- Francesco Locatello, Stefan Bauer, Mario Lucic, Gunnar Raetsch, Sylvain Gelly, Bernhard Schölkopf, and Olivier Bachem. Challenging common assumptions in the unsupervised learning of disentangled representations. In *international conference on machine learning*, pp. 4114–4124.
   PMLR, 2019.
- Francesco Locatello, Ben Poole, Gunnar Rätsch, Bernhard Schölkopf, Olivier Bachem, and Michael Tschannen. Weakly-supervised disentanglement without compromises. In *International conference on machine learning*, pp. 6348–6359. PMLR, 2020.
- Milton L. Montero, Jeffrey Bowers, Rui Ponte Costa, Casimir JH Ludwig, and Gaurav Malhotra. Lost in latent space: Examining failures of disentangled models at combinatorial generalisation. In Alice H. Oh, Alekh Agarwal, Danielle Belgrave, and Kyunghyun Cho (eds.), Advances in Neural Information Processing Systems, 2022. URL https://openreview.net/forum? id=7yUxTNWyQGf.
- Milton Llera Montero, Casimir JH Ludwig, Rui Ponte Costa, Gaurav Malhotra, and Jeffrey Bowers. The role of disentanglement in generalisation. In *International Conference on Learning Representations*, 2020.
- Sudipto Mukherjee, Himanshu Asnani, Eugene Lin, and Sreeram Kannan. Clustergan: Latent space
   clustering in generative adversarial networks. In *Proceedings of the AAAI conference on artificial intelligence*, volume 33, pp. 4610–4617, 2019.
- Maya Okawa, Ekdeep S Lubana, Robert Dick, and Hidenori Tanaka. Compositional abilities emerge
   multiplicatively: Exploring diffusion models on a synthetic task. *Advances in Neural Information Processing Systems*, 36, 2024.
- Maxime Oquab, Timothée Darcet, Theo Moutakanni, Huy V. Vo, Marc Szafraniec, Vasil Khalidov, Pierre Fernandez, Daniel Haziza, Francisco Massa, Alaaeldin El-Nouby, Russell Howes, Po-Yao Huang, Hu Xu, Vasu Sharma, Shang-Wen Li, Wojciech Galuba, Mike Rabbat, Mido Assran, Nicolas Ballas, Gabriel Synnaeve, Ishan Misra, Herve Jegou, Julien Mairal, Patrick Labatut, Armand Joulin, and Piotr Bojanowski. Dinov2: Learning robust visual features without supervision, 2023.
- Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In *International conference on machine learning*, pp. 8748–8763. PMLR, 2021.
- Kuanchi Ren, Tao Yang, Yuwang Wang, and Wenjun Zeng. Learning disentangled representation by exploiting pretrained generative models: A contrastive learning view. In *International Conference on Learning Representations*, 2022. URL https://openreview.net/forum?id= j-63FSNc05a.
  - Fernando Rodriguez-Sanchez, Concha Bielza, and Pedro Larrañaga. Multipartition clustering of mixed data with bayesian networks. *International Journal of Intelligent Systems*, 37(3):2188– 2218, 2022.
- Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High resolution image synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF confer- ence on computer vision and pattern recognition*, pp. 10684–10695, 2022.
- Yonglong Tian, Lijie Fan, Phillip Isola, Huiwen Chang, and Dilip Krishnan. Stablerep: Synthetic
   images from text-to-image models make strong visual representation learners. *Advances in Neural Information Processing Systems*, 36, 2024.
- <sup>697</sup> Vincent Vandewalle. Multi-partitions subspace clustering. *Mathematics*, 8(4):597, 2020.

687

688

689

- Andrey Voynov and Artem Babenko. Unsupervised discovery of interpretable directions in the gan latent space. In *International conference on machine learning*, pp. 9786–9796. PMLR, 2020.
- 701 Xin Wang, Hong Chen, Si'ao Tang, Zihao Wu, and Wenwu Zhu. Disentangled representation learning. *arXiv preprint arXiv:2211.11695*, 2022.

702 703 704 705 706	Thaddäus Wiedemer, Jack Brady, Alexander Panfilov, Attila Juhos, Matthias Bethge, and Wieland Brendel. Provable compositional generalization for object-centric learning. In <i>The Twelfth International Conference on Learning Representations</i> , 2024a. URL https://openreview.net/forum?id=7VPTUWkiDQ.
707 708 709	Thaddäus Wiedemer, Prasanna Mayilvahanan, Matthias Bethge, and Wieland Brendel. Compositional generalization from first principles. Advances in Neural Information Processing Systems, 36, 2024b.
710 711	Matthew Willetts, Stephen Roberts, and Chris Holmes. Disentangling to cluster: Gaussian mixture variational ladder autoencoders. <i>arXiv preprint arXiv:1909.11501</i> , 2019.
712 713 714 715	Tao Yang, Yuwang Wang, Yan Lu, and Nanning Zheng. Disdiff: Unsupervised disentanglement of diffusion probabilistic models. In <i>Thirty-seventh Conference on Neural Information Processing Systems</i> , 2023. URL https://openreview.net/forum?id=3ofe0lpwQP.
716 717	Hengshi Yu and Joshua D Welch. Michigan: sampling from disentangled representations of single- cell data using generative adversarial networks. <i>Genome biology</i> , 22(1):158, 2021.
718 719 720	Nevin L Zhang. Hierarchical latent class models for cluster analysis. <i>The Journal of Machine Learning Research</i> , 5:697–723, 2004.
721 722 723	Junjie Zhao, Donghuan Lu, Kai Ma, Yu Zhang, and Yefeng Zheng. Deep image clustering with category-style representation. In <i>Computer Vision–ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part XIV 16</i> , pp. 54–70. Springer, 2020.
724 725 726 727	Bolei Zhou, Xiaogang Wang, and Xiaoou Tang. Understanding collective crowd behaviors: Learn- ing a mixture model of dynamic pedestrian-agents. In 2012 IEEE Conference on Computer Vision and Pattern Recognition, pp. 2871–2878. IEEE, 2012.
728 729 730	
731	
733	
734	
735	
736	
737	
738	
739	
740	
741	
742	
744	
745	
746	
747	
748	
749	
750	
751	
752	
753	
754	
755	

#### 756 APPENDICES

#### **PROOF FOR LOWER BOUND** А

760 We now provide a proof that Equation 2 is a lower bound for implicit clustering, which is given by explicit clustering. 762

763 *Proof.* Suppose we cluster  $\mathcal{D}$  by sim, then we obtain k disjoint subsets of  $\mathcal{D}$  (Equation 1). When 764  $sim \equiv \mathcal{F}_i$  holds, Equation 2 formalizes the generation of any partitioning of the data. Under an 765 encoder that perfectly disentangles any datapoint into its factors of variation in single dimensions, 766 we can transform any  $o \in \mathcal{D}$  to a latent representation that include its factors of variation  $z = h_{\varphi}(o)$ . 767 Given that the factors are perfectly disentangled, we can then group all  $o \in \mathcal{D}$  by the respective 768 dimension representing the respective factor of variation  $\mathcal{F}_i$  and obtain the same disjoint subsets as with explicit clustering under  $C_{sim}$ . If instead of explicitly grouping the elements o, we reconstruct 769 them using a generator  $G_{\theta}(\cdot)$ , i.e.,  $o^* = G_{\theta}(h_{\varphi}(o))$ , we specifically end up with Equation 2: 770

758

759

761

$$\mathcal{D} = \bigcup_{k} \mathcal{D}_{k}^{sim} = \bigcup_{k} \left\{ x \mid x \in \mathcal{D} \land C_{sim}(x) = k \right\} = \bigcup_{f \in G_{i}} \left\{ o \mid o \in \mathcal{D} \land \stackrel{f}{=} h_{\varphi}(o) \right\}$$
(4)

777

$$\approx \bigcup_{f \in G_i} \{ G_\theta(h_\varphi(o)) \mid o \in \mathcal{D} \land \stackrel{f}{=} h_\varphi(o) \} = \bigcup_{f \in G_i} \mathcal{D}'_{if} = \mathcal{D}' \approx \mathcal{D}$$
(5)

778 Up to the reconstruction capabilities of the generator  $G_{\theta}$ , the reconstructed dataset  $\mathcal{D}'$  is a synthetic 779 version of  $\mathcal{D}$ , which under a perfect generator would be equivalent, i.e.,  $\mathcal{D} \equiv \mathcal{D}'$ . However, because 780 we can not assume a perfect reconstruction from the latent representation z, we write the synthetic 781 version of  $\mathcal{D}'$  is an appropriate realistic synthetic version of  $\mathcal{D}$ , i.e.,  $\mathcal{D} \approx \mathcal{D}'$ . 782

Suppose that we can modify the factors of variation of an object o in its disentangled representation 783  $z = h_{\omega}(o)$  with a function mod, and the generator  $G_{\theta}$  creates realistic synthetic elements  $o' = b_{\omega}(o)$ 784  $G_{\theta}(mod(h_{\varphi}(o)))$ . In this case, any minor modification yields a new object o' different from o, 785 effectively extending the cardinality and variety of samples in  $\mathcal{D}'$ , which makes Equation 2 a lower 786 bound to Implicit Neural Clustering, given by explicit clustering. More specifically: 787

792

 $\mathcal{D} = \bigcup_k \mathcal{D}_k^{sim} < \approx \bigcup_{f \in G_i} \mathcal{D}'_{if} = \mathcal{D}'$ f )

(6)

$$= \bigcup_{f \in G_i} \{G_\theta(h_\varphi(o)) \mid o \in \mathcal{D} \land \stackrel{j}{=} h_\varphi(o)\} \cup \{G_\theta(mod(h_\varphi(o))) \mid o \in \mathcal{D} \land \stackrel{j}{=} mod(h_\varphi(o)\}$$
(7)

Assuming the factors of variation F for any  $x \in \mathcal{D}$  can be obtained with an encoder h, derived ana-793 lytically, or are given by annotations, we can train  $G_{\theta}(\cdot)$  to synthesize samples from the underlying 794 distribution that we can group naturally by a factor of variation. A more general formulation that 795 encompasses compositional and combinatorial generalization is given in Equation 3. When we can 796 partition or factorize the underlying generative factors into respective atomic partitions and compo-797 sitionality emerges in  $G_{\theta}(\cdot)$ , synthesis of known cross combinations between factors of variation 798 is possible, which results in (a) dataset reconstruction (Equation 2) and (b) controllable synthesis, 799 i.e., we can now move beyond only reconstructing the dataset, but can also specifically control the 800 synthesis of new examples. Finally, when  $G_{\theta}(\cdot)$  also learns to generalize to combinations that are 801 not in the data distribution, we move beyond the lower bound given by Equation 2 with Equation 3, 802 where arbitrary novel cross-combinations between factors of variation can be synthesized. 

803 804

805

#### В ADDITIONAL DETAILS FOR EXPERIMENTAL SETUP AND DESIGN

806 We will provide additional training details on hyperparameters, code, and further setups upon ac-807 ceptance of this publication. 808

The abbreviations for each factor of variation  $\mathcal{F}_i$  in each dataset as used in our experiments is based 809 on the specification of the datasets.



(a) Ground Truth

(b) Relaxed Implicit Neural Clustering

Figure 7: Relaxed implicit neural clustering of crowd trajectory data (right) effectively approximates the underlying trajectory distribution (left)

# C APPLICATION OF IMPLICIT NEURAL CLUSTERING TO OTHER DOMAINS

826 Multi-partition clusters exist in many different domains. Different from the hard challenge of learn-827 ing the factors of variation with an encoder  $h_{\varphi}$ , they are often provided in ground truth or can be 828 obtained analytically. In such contexts, where factors of variation might be already available, can 829 be obtained analytically, or with a zero-shot classifier like CLIP, Implicit Neural Clustering is also applicable when we relax the assumption on having an encoder  $h_{\varphi}$ . To show this, in the following 830 experiment, we relax the assumption that there exists an encoder  $h_{\varphi}$  that extracts factors of variation 831 from a given dataset  $\mathcal{D}$  to show that the compositionality and combinatorial generalization require-832 ments for our approach can be fulfilled in different applications. More specifically, we assume that 833 factors of variation have been obtained by, e.g., applying an "optimal" encoder  $h^*$  (e.g., a human, 834 an analytical relaxation, or a zero-shot classifier Radford et al. (2021)) on D. In such cases, Implicit 835 *Neural Clustering* is applicable to generate implicit clusters resembling the conceptual overview of 836 our approach in Figure 2. While this relaxation might seem trivial, we can effectively show a novel 837 application of disentangled representations, while also providing a new perspective on clustering. 838

839 **Experimental Setup and Models.** To show different applications of *Implicit Neural Clustering* 840 with analytical or provided ground truth factors of variation and different domains, we evaluate our 841 approach with the following datasets. (i) the Grand Central Station (GC) Zhou et al. (2012) dataset 842 that consists of time series that resemble trajectories of pedestrians traversing a public train station. In GC, we can vary the start and goal position of the agent to traverse the underlying environment, 843 which passed to a generator as continuous inputs. In this dataset, we train a goal-conditioned policy 844 with behavior cloning as the generator  $G_{\theta}$  Kreutz et al. (2024). (ii) the CLEVR Relations Johnson 845 et al. (2017) dataset, which allows us to vary three ground truth factors of variation, the number 846 of objects, as well as their respective X and Y coordinates. We use a pre-trained compositional 847 diffusion model Liu et al. (2022) provided by the authors<sup>6</sup> that is conditioned on natural language 848 prompts. In this way, we show an application to text-to-image generative models. (iii) a synthetic 849 dataset of simple shapes provided by Okawa et al. (2024) (in the remainder referred to as Simple-850 Shapes). SimpleShapes has shape, color, and size as factors of variation, which are used as continu-851 ous input to the generator. We train a conditional diffusion model <sup>7</sup> based on Okawa et al. (2024) in 852 SimpleShapes on the provided ground truth factors of variation. All models are trained from scratch with a single 46GB RTX A6000. We will provide additional training details on hyperparameters, 853 code, and further setups upon acceptance of this publication. 854

855 856

857

820

821

822 823 824

825

C.1 ANALYTICAL FACTORS OF VARIATION — SEQUENTIAL DECISION MAKING.

Sequential decision making tasks, such as motion planning or navigation, can be relaxed analytically into several factors of variation. For instance, start and goal positions influence the outcome of a trajectory. When varying these two positions as factors of variation, a policy trained on a dataset of expert demonstrations will then "fill-in-the-blanks" and generate a trajectory that follows the

 <sup>&</sup>lt;sup>6</sup>https://energy-based-model.github.io/Compositional-Visual-Generation-with-Composable-Diffusion-Models/

<sup>&</sup>lt;sup>7</sup>public github repository https://github.com/phys-ai/concept\_graphs/\$



....

(b) Implicit clusters for each atomic X coordinate partitions

Figure 8: Relaxed implicit neural clustering of object compositions, where one can implicitly cluster the data based on the relaxed representation of the factors of variation, i.e., discretized x and y coordinates

ground truth distribution. In this setting, we can apply *Implicit Neural Clustering* where start andgoal positions can be considered factors of variation.

We evaluate *Implicit Neural Clustering* on the grand central station dataset (GC) Zhou et al. (2012), where we train a goal-conditioned policy with behavior cloning as the generator  $G_{\theta}$ . Analogous to images, we partition the factors of variation in each dimension separately, i.e., start and goal positions using DBSCAN as a partition algorithm, and only keep the top k number of (start, goal) pairs from the real dataset as factors of variation. Given the respective positions at each start and goal, we compute the parameters of a normal distribution for each of these sets, which serve as the parameters to sample from atomic start and goal partitions for *Implicit Neural Clustering*.

890 Figure 7 visualizes an implicit clustering into clusters of (start,goal) combinations and shows the 891 results of the clusters of ground truth trajectories (a) and their corresponding implicit clusters (b). 892 In this experiment, start and goal positions are considered factors of variation that are varied and 893 the motion planner learns to generate samples that approximate the original dataset. We control the 894 generation by varying the respective factors of variation and generate paths according to Equation 3. The atomic partitions in this context correspond to pairs of (start, goal) partitions in an euclidean 895 space. In comparison to the ground truth clusters, variations of the paths are generated that mimic the 896 ground truth distribution of expert demonstrations. In the same context, completely new scenarios 897 can be synthesized by algorithms that would allow the model to learn react to the environment, such 898 as GAIL Ho & Ermon (2016). 899

Analytical Factors of Variation — Compositional Image Synthesis. Similar to relaxing start and goal positions for motion planning as factors of variation, placement of objects in an image is also a straightforward relaxation in the euclidean space. In this experiment, we show an application of our approach to compositional generation of images by relaxing several factors of variation required for object composition. More specifically, we relax the placement (x,y) placement coordinates and the number of objects, which gives three analytical factors of variation.

Figure 8a and Figure 8b show exemplary implicit clusters on the CLEVR Relations dataset. We can control the generation by varying the respective factors of variation and generating images according to Equation 3. We emphasize that the underlying model fails to generate coherent clusters near the distribution boundaries, highlight the need for better OOD generalization even in "simple" placement tasks. However, we want to highlight that this kind of task can as well be expressed under Equation 3 as an *Implicit Neural Clustering*.

912

872 873 874

875 876

877

878

879 880 881

913 914

# D ADDITIONAL COMMENTS AND EXPERIMENTS ON COMBINATORIAL GENERALIZATION

915 916

In the relaxed setting with  $h^*$ , an application to diffusion models with known factors of variation that shows combinatorial generalization under Equation 3 can be given based on the work by Okawa



Figure 9: Partitions on x and y with all cross combinations.

et al. (2024). In this experiment, the factors of variation can be naturally partitioned into atomic partitions due to their discrete nature. We reproduce their experiments on combinatorial generalization by training a conditional diffusion model to generate images conditioned on restricted factors of variation. We show the explicit clusters compared to the implicit clusters that can be generated in Figure 10. We can see that the model learns combinatorial generalization to generate small spheres, big blue rectangles, and small red and blue rectangles. Note how this experiment mimics our conceptual illustration of *Implicit Neural Clustering* in Figure 2, showing strong empirical evidence for the correctness of our definition. Finally, we want to highlight that Okawa et al. have rigorously tested when diffusion models achieved combinatorial generalization in this synthetic dataset Okawa et al. (2024). Their overall experiments provide empirical evidence for the validity and practicality of our approach while satisfying all of our conditions (*A.H*, *A.G*1, *A.G*2, *A.G*3) in the special case of relaxation to  $h^*$ .

# E ADDITIONAL RESULTS FOR DSPRITES

We provide the full implicit multi-partition clustering of dSprites.

970 971

952

953 954

955

956

957

958

959

960

961

962

963

964

965 966 967



Figure 10: Relaxed implicit neural clustering of shapes (right) effectively approximates and generalizes to novel cross combinations of the underlying ground truth distribution (left). This experiment provides empirical evidence that diffusion models trained on synthetic data can satisfy the most difficult part of moving beyond the data distribution with *Implicit Neural Clustering*. In this example, the generative model generalizes to being able to synthesize elements not seen in the training dataset, hence filling in the gaps in the data distribution as previously illustrated in Figure 2.







1134	۰.	•			•						•				•	-		1	•		l
1135		,	•	•	•	•	•	•	•		•	•	•	•			•	•	•		
1136	٠	•	٠	•	•	•	٠	•	·	•	•	•	•	•	•	۰,	•	•	•	•	
1137	•	*	٠	١	•	•	•	٠	•	•	٠	١	4		•	٠	•	•	-	•	1
1138	•	•	•	•	8	•	•	•	•	•	٩	•	٠	•	٠	•	•	•	٠	٠	
1139	÷	÷	•	<u>`</u>	•	•	•	•	•	•	•	•	<u>·</u>	•	-		•	-	•	÷	
1140	•	È	Ť	•	•	_		,	,	•	•	÷	-	-	-	•	•		<u>.</u>	•	1
1141	·	ī	•	·	٠	•	•	٠	•	۶	•	٠	•	·	•	٠	*	·	•	•	
11/2	١	٠	٠	•	•		2	•		18	•	•	٠	•	•	٠	•	•	٠		Ì
1143	•	•	•	•	٠	•	•	•	•	۶	۰	٠	•	*	•	•	٠	•	•	•	
1144	-		•	•	*	•	<u> </u>	-		•		•	÷	•	•	*	•		•	•	
11/5		۰	÷.	•	,	•	•		•	•	_	•	*	•	÷	•	•	•		•	
11/6	·		•	٠	•	٠	٠	·	·	٠	۲	•	•	٠		•	•	•	•	٠	
1147	•	•	•	٠	•	•	٠	٠	•	*		•	٠	٠	¥	•	•	•	•	•	1
1147	·	•	•		•		•	•	٠	•	٠	•	•	٠	٩	·	•	•	•	•	
1140	•	-		<u>•</u>	•	-		*	•	•	•		<u> </u>	•	•	•	•		•	·	
1150	•	, •	•	•		4	•		4	•		•	•	•	·	•	•	-	•	•	
1150	•	٠	•	•	,	•	•	٠	•	•	•	,	•	٠	•	•	•	•	•	·	
1151		•	4	•	٠	-			•	•	•	¥	٠		•	•	•	٠	•	ŀ	Ì
1152	٠	٠	•	٠	٠	•	•	•	•	•	۰	•	•	٠	•	•		•	•	٠	
1153	•			•	•	•			•	•	•	-	•	•		•	•		•		1
1154	•	•	-	•	•	•	•	•	•	•	•	•	•		•		•	•		•	1
1155		•	•	٠	,	•		٠	•	٠	٠	·		•	•	•	•	•	·	•	
1156	•	٠	•	٠	١			•	•		•	۲	•	٩	٠		•	•	٠	•	ĺ
1157	٠	۰	•	•	•	٠	*	•	٠	•	•	•	•	•	•	•	•	·	•	٠	
1158		•		•	•	•	•	•	•	<u> </u>	•	•	•	•	•	•	*	Ľ	٠	•	
1159		•	Ì	*		•	•	•	•		•	5	•	÷	÷		•	•	•		
1160							_		_		_			_						-	1
							1.	п		1		7	1								
1161	_	_		_	_		(a)	) R	ano	dor	n S	San	npl	es	_						
1161 1162	١	•	*	•	<mark>ار</mark>	•	(a) •	) R	ano •	doı •	n S •	San	npl •	es •	<b>`</b>	-	•	۲ د	•	4	<b>I</b>
1161 1162 1163	۰ ۰	• • •	* , ,	•	<u>م</u>	• •	(a) • •	) R	an • •	doı • •	n S • •	San • •	npl • •	es • •	、 、 -	•	11 ^ /	1 4 11	• • •	•	
1161 1162 1163 1164	۹ ۱۰ ۱۰	• • •	* > *	• • •	• • •	• • •	(a) • •	) R • •	ano • •	don • •	n S • •	San • •	npl • •	es • •	、 、 ・	• • •	# * ·	/ < %	۰ ۲	۲ ۱ ۱	
1161 1162 1163 1164 1165	• • •	• • • •	* > * *	• • •	* * *		(a) • • •	) R • •	an • •	don • •	n S • • •	San • •	npl • •	es • •	۲ ۲	• • •	* * ·	/ 4 8 /	• • •	۲ ۱ ۱ ۴	
1161 1162 1163 1164 1165 1166	• • • •	• • • • •	* > * *	• • • • • • • •	• • • •	<ul> <li>.</li> <li>.&lt;</li></ul>	(a) • • • •	) R - - - - -	and • • •	don • • •	n S • • •	San • • •	npl	es • • • •	۲ ۱ ۱	• • • •	* * ·	/ 4 / / /	• • •	۲ ۱ ۱ ۲ ۲	
1161 1162 1163 1164 1165 1166 1167	• • • • • • • • •	• • • • • • •	* > * * *	• • • • • • •	• • • • •	<ul> <li>.</li> <li>.&lt;</li></ul>	(a) • • • •	) R	and • • •	doi • • • •	n ( •	San	npl	es • • • • •	<ul> <li>.</li> <li>.&lt;</li></ul>	<ul> <li>.</li> <li>.&lt;</li></ul>	* / / *	/ ( 1 - - - - - - - - - - - - -	• • • •	<ul> <li></li> <li><td></td></li></ul>	
1161 1162 1163 1164 1165 1166 1167 1168	• • • • •	• / / / / / / /	* * * * * *	• • • • •	<ul> <li>.</li> <li>.&lt;</li></ul>	<ul> <li>.</li> <li>.&lt;</li></ul>	(a) • • • •	) R • • • • •		doi	n \$ • • • • • •	San M M M M M M	npl • • • • • • •	es • • • • • •	<ul> <li>.</li> <li>.</li></ul>	*	1 / / 1 1 1 1 1 1 1 1 1 1 1 1 1	/ 8 / / / / / / / / / / / / / / / / / /	> + - - - - - -	<ul> <li>C</li> <li>J</li> <li>H</li> <li>H</li> <li>V</li> <li>H</li> <li>V</li> <li>N</li> </ul>	
1161 1162 1163 1164 1165 1166 1167 1168 1169	• • • • •	•	* > * * * * * * * * * * * * *	• • • • •	<ul> <li></li> &lt;</ul>	<pre>/ / / / / / / / / / / / / / / / / / /</pre>	(a) • • • •	) R 2 2 3 3 4 3 3 4 3 3 3 4 3 3 4 3 4 3 3 4 3 4 3 4 3 4 3 4 3 4 3 4 3 4 3 4 3 4 3 4 3 4 3 4 3 4 3 4 3 4 3 4 3 4 5 4 5	and		n \$ • • • • • • •	San - - - - - - - - - - - - -	npl	es · · · · · · · · · ·	<ul> <li></li> &lt;</ul>	<ul> <li>.</li> <li>.&lt;</li></ul>	1 - - - - - - - - - - - - -	/ 1 2 4 1 2 4 1 1 1 1 1 1 1 1 1 1 1 1 1 1	<ul> <li></li> <li><!--</td--><td><ul> <li>.</li> <li>.</li></ul></td><td></td></li></ul>	<ul> <li>.</li> <li>.</li></ul>	
1161 1162 1163 1164 1165 1166 1167 1168 1169 1170			* > - - - - - - - - - - - - -	* * * * * * *	* • • • • • • • • •	<ul> <li></li> &lt;</ul>	(a) • • • • • • • • • • • • •	) R - - - - - - - - - - - - - - - - - - -	and		n (	San 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2	npl	es	<ul> <li></li> &lt;</ul>	<ul> <li>.</li> <li>.&lt;</li></ul>			* *	<ul> <li>C</li> <li>A</li> <li>B</li> <li>C</li> <li>A</li> <li>C</li> <li>C</li></ul>	
1161 1162 1163 1164 1165 1166 1167 1168 1169 1170					<ul> <li></li> &lt;</ul>	<ul> <li></li> &lt;</ul>	(a) - - - - - - - - - - - - -	) R 2 2 3 3 4 3 3 4 3 3 4 3 3 4 3 4 3 4 3 4	and - - - - - - - - - - - - -		n (	San San San San San San San San	npl 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2	es	<ul> <li></li> &lt;</ul>					<ul> <li></li> &lt;</ul>	
1161 1162 1163 1164 1165 1166 1167 1168 1169 1170 1171					<ul> <li></li> &lt;</ul>	<ul> <li>.</li> <li>.&lt;</li></ul>	(a) - - - - - - - - - - - - - - - - - - -	) R	and - - - - - - - - - - - - -		n (	San San San San San San San San	npl • • • • • • • • • • • • • • • • • • •	es • • • • • • • • • • • • • • • • • • •	<ul> <li></li> &lt;</ul>	<ul> <li>.</li> <li>.&lt;</li></ul>			<ul> <li></li> &lt;</ul>	<ul> <li>.</li> <li>.</li></ul>	
1161 1162 1163 1164 1165 1166 1167 1168 1169 1170 1171 1172							(a) - - - - - - - - - - - - -	) R - - - - - - - - - - - - - - - - - - -	and		n (* * * * * * * * * *	San 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2		es					<ul> <li></li> &lt;</ul>	<ul> <li>.</li> <li>.</li></ul>	
1161 1162 1163 1164 1165 1166 1167 1168 1169 1170 1171 1172 1173								R         I <td< td=""><td>and - - - - - - - - - - - - -</td><td></td><td></td><td>San San San San San San San San</td><td></td><td>es</td><td><ul> <li></li> &lt;</ul></td><td></td><td></td><td></td><td></td><td></td><td></td></td<>	and - - - - - - - - - - - - -			San San San San San San San San		es	<ul> <li></li> &lt;</ul>						
1161 1162 1163 1164 1165 1166 1167 1168 1169 1170 1171 1172 1173 1173												San San San San San San San San	npl								
1161 1162 1163 1164 1165 1166 1167 1168 1169 1170 1171 1172 1173 1174								R         I <td< td=""><td>and - - - - - - - - - - - - -</td><td></td><td>n S - - - - - - - - - - - - -</td><td>San San San San San San San San</td><td>npl</td><td>es</td><td></td><td></td><td></td><td></td><td></td><td></td><td></td></td<>	and - - - - - - - - - - - - -		n S - - - - - - - - - - - - -	San San San San San San San San	npl	es							
1161 1162 1163 1164 1165 1166 1167 1168 1169 1170 1171 1172 1173 1174 1175 1176 1177							(a) - - - - - - - - - - - - -	R 2 2 3 3 3 3 3 4 3 3 3 4 3 3 3 4 3 3 3 4 3 3 3 4 3 3 4 3 4 3 4 3 4 3 4 3 4 3 4 3 4 3 4 3 4 3 4 3 4 5 4 5	and - - - - - - - - - - - - -			San San San San San San San San	npl	es							
1161 1162 1163 1164 1165 1166 1167 1168 1169 1170 1171 1172 1173 1174 1175 1176							(a) - - - - - - - - - - - - -	R 2 2 3 3 3 3 3 3 4 3 3 3 3 4 3 3 3 4 3 3 3 4 3 3 3 4 3 3 4 3 3 4 3 4 3 4 3 4 3 4 3 4 3 4 3 4 3 4 3 4 3 4 5 4 5	and - - - - - - - - - - - - -		n 5 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1	San	npl	es							
1161 1162 1163 1164 1165 1166 1167 1168 1169 1170 1171 1172 1173 1174 1175 1176 1177 1178 1179								R 2 2 3 3 3 3 3 3 4 3 3 3 3 4 3 3 3 4 3 3 3 4 3 3 4 5 4 5	and - - - - - - - - - - - - -			San									
1161 1162 1163 1164 1165 1166 1167 1168 1169 1170 1171 1172 1173 1174 1175 1176 1177 1178 1179 1180									and 					es							
1161         1162         1163         1164         1165         1166         1167         1168         1169         1170         1171         1172         1173         1174         1175         1176         1177         1178         1179         1181								R - - - - - - - - - - - - -	and 				npl								
1161         1162         1163         1164         1165         1166         1167         1168         1169         1170         1171         1172         1173         1174         1175         1176         1177         1178         1179         1180         1181								R	and - - - - - - - - - - - - -				npl								
1161         1162         1163         1164         1165         1166         1167         1168         1169         1170         1171         1172         1173         1174         1175         1176         1177         1178         1180         1181         1182         1183									and - - - - - - - - - - - - -												
1161         1162         1163         1164         1165         1166         1167         1168         1169         1170         1171         1172         1173         1174         1175         1176         1177         1178         1179         1180         1181         1183         1184																					
1161         1162         1163         1164         1165         1166         1167         1168         1169         1170         1171         1172         1173         1174         1175         1176         1177         1178         1179         1180         1181         1182         1183         1184         1185																					
1161         1162         1163         1164         1165         1166         1167         1168         1169         1170         1171         1172         1173         1174         1175         1176         1177         1178         1179         1180         1181         1182         1183         1184         1185         1186																					

(b) w/ Atomic Group Actions

Figure 14: dSprites: Implicit 22usters for object X coordinate

							٠	•				•	٠	•		•	٠	•	•	•
1189		•	•	÷	•	•	•		•	•			•	•		•	٠	٠		•
1190		•	•	•	٠	٠	٠	٠	•	٠		•	*		٠		•	•	٠	
1191	•	•		`		٠	•	•	•	•	*	`	•	٠	•	٠	a	•	•	•
1192	•	•	È	•			•	*	•		·	÷	•	•	•		•	•	-	
1193	÷	•	•	÷	•	•	•	÷		•	, ,	<i>•</i>			•	•	•	•		
1194	•	4	•	÷	•		٠	•	·	•	٠		•	•	•	•	*	•	·	
1195	•	٠	٠	•	٠		•	٠	٠	٠		•	٠		•	÷	•	•		•
1196	٠	٠	•	٠	٠	•	•	•	•	•	,	•	•	٠	•	•	•	•	•	•
1197	-	•	•	•	•	•	•	Ľ	•		Ŀ	÷	•	•	Ŀ		÷	•	•	•
1198	•	•	•	•		•	•	*			È	•	-		•	•	•		÷	4
1199		•	•	•	•	•	•	•		•	·	•	•	•	٠		•	·	•	•
1200			٠	٠	•	٠	•	•	•	٠	•	٠	•	•	•		•	•	•	•
1201	•	٠	•	•	٠	۰	•	•	•	•	•	•		•	•	•	•	'	•	•
1202	•	•	•	•	•	-	•	•	•		-	-	•	•	•	•	•	-	•	•
1203	•	•	、	•	•		-	•		۴	•	•	_	•	÷	•	÷	•	·	
1204		٠	٠		•	•	•	•	,	٠	Ţ	·	,	•	•	•	•	•	•	•
1204		4	٠	٠	'		•	•	٠	٠	,	•	٠	٠	•	•	۶.	•	۰	•
1206		•	`	•	·	•	•	•		٠	•	•	•	•	Ľ	٠	•	<u>'</u>	•	•
1200	•	-	•	•	•			•	•	•	•	*	•		÷	•	È	È	•	
1207		•	•	•	•	÷	÷		•		•	•	•	•			•	È	•	•
1200	٠	•	•	٠	•	٠	•	•	۰	•	·			•	٠	۲	•	٠	•	•
1210	•	•	•	•		*	•	٠	•	•	•	•	•	٠	•	•	*	•	•	*
1011	٠	•	•	<u> </u>	•	•	•	•	•	•	•	•	•	•	'	•	•	-	•	*
1010	•	•		•	•	-			+	-	÷	•	•	•	•			•	-	•
1010	-			•	÷	•		•	•	٠	•	、		·	•	•	•		•	•
1213	*	•	•	٠	•	٠	٠	•	٠	٠	•	•	•		•	٠	*	•		•
1015							(a	) R	an	doi	n S	San	nnl	es						
1215	•	•	۰	•	•	•		•			•	<b>_</b> `	•		•	٩	٠	•	•	•
1210	·	•	•	•	•	•	•	•	•	•	·	•	•	•	·	•	·	•	-	•
1010	•	•	•	•	•	٠	1	*	*	۲	۴	•	*	*	4	۴	•	•	۰	*
1210	÷	•	•	<u>`</u>	•	•	<u> </u>	<u>'</u>	<u>`</u>	•	1	Ľ	•	•	-	•	4	•	-	•
1213		· ·	•				•	•		•	•	-	•		Ì	•	•			•
1990	•	•			٠	-					_		•	•	_				٠	
1220	•	`	•	•	•	•	`	٠	۲	•	•				`	٠	٠	•	•	•
1220 1221 1222	*	` '	• •	•	• • •	•	` *	*	•	•	•	•	•	•	` •	•	•	•	* • •	•
1220 1221 1222 1223	•	、 ・ ・	• •	•	• • •	•	* •	*	•	•	*	•	۶ ۲	•	` •	• • •	* * •	• • • •	* • •	
1220 1221 1222 1223 1224	•	、 ・ ・ ・	• • • • •	•	• • • •	•	• • •	* * *	• • •	• • •	• • •	• • •	> > 4	# • •	、 ・ ・	* * *	• • •	• • • •	* • • •	
1220 1221 1222 1223 1224 1225	•	、 ・ ・ ・ ・	• • • •	• • • •	• • • • •	) ) ) ) ) ) )	* * * *	• • • •	• • • •	• • • •	* * * *	• • • •	> • •	# • • 4	、 ・ ・ ・	* * * *	• • • •	*	* • • • •	
1220 1221 1222 1223 1224 1225 1226	•	<ul> <li>.</li> <li>.&lt;</li></ul>	• • • • •	* * * *	• • • • •	•	* * * *	* * * *	• • • •	• • • • •	• • • • •	• • • •	> > • •	# • • • •	` * * *	* * * * *	• • • • • •	<ul> <li>*</li> <li>*&lt;</li></ul>	* • • • •	
1220 1221 1222 1223 1224 1225 1226 1227	•••••	<ul> <li>.</li> <li>.&lt;</li></ul>	• • • •			•	* * * *	· • • •	• • • • •	• • • • •		· • • •	> - - - -	* * * * *	` • • •	* * * * *	· • • •		* • • • •	
1220 1221 1222 1223 1224 1225 1226 1227 1228		<ul> <li>.</li> <li>.&lt;</li></ul>										· • • • • •	· · ·		` • • • •				• • • • • •	
1220 1221 1222 1223 1224 1225 1226 1227 1228 1229		<ul> <li>.</li> <li>.&lt;</li></ul>															· • • • • • • • • • • • • • • • • • • •			
1220 1221 1222 1223 1224 1225 1226 1227 1228 1229 1230		<ul> <li></li> &lt;</ul>																		
1220 1221 1222 1223 1224 1225 1226 1227 1228 1229 1230		>							· · · · · · · · ·											
1220 1221 1222 1223 1224 1225 1226 1227 1228 1229 1230 1231																				
1220 1221 1222 1223 1224 1225 1226 1227 1228 1229 1230 1231 1232																				
1220 1221 1222 1223 1224 1225 1226 1227 1228 1229 1230 1231 1232 1232																				
1220 1221 1222 1223 1224 1225 1226 1227 1228 1229 1230 1231 1232 1233 1234																				
1220 1221 1222 1223 1224 1225 1226 1227 1228 1229 1230 1231 1232 1233 1234 1235																				
1220 1221 1222 1223 1224 1225 1226 1227 1228 1229 1230 1231 1232 1233 1234 1235 1236																				
1220 1221 1222 1223 1224 1225 1226 1227 1228 1229 1230 1231 1232 1233 1234 1235 1236 1237 1238																				
1220 1221 1222 1223 1224 1225 1226 1227 1228 1229 1230 1231 1232 1233 1234 1235 1236 1237 1238 1239																				
1220 1221 1222 1223 1224 1225 1226 1227 1228 1229 1230 1231 1232 1233 1234 1235 1236 1237 1238 1239																				
1220 1221 1222 1223 1224 1225 1226 1227 1228 1229 1230 1231 1232 1233 1234 1235 1236 1237 1238 1239 1240																				

(b) w/ Atomic Group Actions

Figure 15: dSprites: Implicit 2Busters for object Y coordinate



















Figure 23: 3DShapes: Latent traversal with atomic group actions, wall color





Figure 25: 3DShapes: Latent traversal with atomic group actions, object size





	$+\hat{\cap}$
	$\uparrow$ + + + $\uparrow$ $\uparrow$ $\uparrow$ + $\uparrow$ $\uparrow$ $\uparrow$ + + + $\uparrow$ $\uparrow$ $\uparrow$ + $\uparrow$ $\uparrow$ $\uparrow$
	ナナナテナマ アナデナ ナナナナキ マアナデナ
	(a) Random Samples (b) W/ Atomic Group Actions
	Figure 28: MPI3D: Implicit clusters for object color
G	Additional Results on MPI3D
U	
We p	provide the full implicit multi-partition clustering of MPI3D.
<b>C</b> 1	
G.I	IMPLICIT NEURAL CLUSTERINGS









2160			~	<u> </u>	(		$\left( \right)$								(					
2161	4		4		-		. *	+		1	4	-	••		•	<u>_</u>				1.0
2162		<b></b> _			1		1	÷	1	3	* 8	1	41	1.	•1		-			- in
2163		-	.1			•	4				3.		• 1				~		<b>.</b>	
2164			•		•		*	*	•					-			÷			
2165	1	4				4	1	1		1	•	•	1			1	1	1	•	•1
2166	-	1		1	1		7	1	-								•		1	
2167					$\sim$					•					$\sim$					
2168	4		× 4		1994 - C.	1 A.	- 4	- X.	•		4	8	<b>*</b> •	4	- 19		- 0			ţ.
2109	d.	-				1		-	1	-	8		7		-	1				J.
2170			*	$\frown$							-					$\sim$	*			
2171		•	~	$\sim$				•			-		*					0	•	•
2172	1	1		4	1	1		1	1		1			1	1				-3	•
2173	1.	-		$\frown$		1		-	-			-			-					- F
2175	-		-		4	P		$\frown$			0		<u> </u>							-
2176	4	1	•				÷			4		4	•			+	+		- <b>*</b> -	4
2177	T		1		4			1	i	1	1		1					•	1	
2178		-			5		-							-					<b>n</b>	
2179		• —				•				•	Þ								2	•
2180	Ţ	ð.*	1	4	- <b>k</b>	8	3	7				- d		- A.,		6	4		4	+
2181	1	-1	) O			1	8 -	1	1								Ĩ		1	
2182		- 4	+			•	-	4				4	-			•				
2183						*		•	•		$\sim$	*				•		•	•	
2184	)	Å.	1			4	1	1			1	- 3	1			3			1	
2185		- I	1																	
2186	1				-			$\sim$			<b>^</b>			-	-				$\sim$	
2187	4	1		1		1. A.		4	- 0			1 A.	- <b>P</b>	- 2		19 - 19 - 19 - 19 - 19 - 19 - 19 - 19 -		4	4	
2188	8.	Ì					4	1		*							i		1	1
2189				-			•1			$\frown$	4		1							
2190				*									-			-		-		-
2102	1	4	1				-1	•	4	1	4	4	4		0				4	1
2193	1	7	Ţ	1		-3		1		*						1			1	1
2194	-						J.	-		-						<u>م</u>		-	-	-
2195	4		$\sim$		•		<i>•</i>	¢			4. •	•	$\sim$	~	•			•	_	
2196			)	- <b>h</b> e	3		T	1.	*	1			1	3	1		1		1	
2197	-	-		-E			1	1												
2198	÷ ,				-					-		-			-	$\sim$		-	$\sim$	-
2199	+	-0	4			1. A		•	- 45	1		-0				1	•	•	4	*
2200	J	1		4		2	Ø1	1			Ì			4	1	1			1	
2201	+	$\overline{1}$	Ĩ	-		-1-	1	F								-		$\widehat{}$	-	
2202								_											-	
2203	4		- 8		<b>-</b> 1			6	4			$\sim$	- A	1-				4		
2204	7	1	$\widehat{}$			1	dr.									-				
2205		4						•		-		_			C					
2206		•		9	2			÷						•	2					
2207	1	3	1		-		1	$\uparrow$		3	- A	1	1				*	1		1
2208	3*						1.										$\widehat{}$		-	
2203								-					•			_		_	-	
2211	1	+	- 1				1	1	1				- 6	1	0				1	- 2
2212			(	a) Ra	andoi	n Sai	nples	5					(b) w	/ Ato	omic	Grou	p Ac	tions		
2213																				

Figure 33: MPI3D: Implicit clusters for horizontal axis



2268							
2269							
2270							
2271	ک	ک	ک	ک	٩	ک	60
2272							
2273							
2274		•					
2275					N N		
2276	چ	چ	چ	چ	۹		چ
2277							
2278							
2279							
2280	da		do l	da		de la	
2281		· · ·			~		*
2282							
2283							
2284						1	1
2285							
2286							
2287	-		-	-	-	-	-
2288			10			10 M	
2289							
2290	۲	۹	۲	۲	۲	۲	۲
2201							

Figure 35: MPI3D: Latent traversal with atomic group actions, object color

G.2 LATENT TRAVERSALS WITH ATOMIC GROUP ACTIONS

2323		
2324		
2325		
2326		
2327		
2328		
2329		
2330	<b>A</b>	<b>_</b>
2331	T	T
2332		
2333		
2334		
2335		
2336		
2330		
2338		
2330		
2340		
2340	4.6	1.
2340	(SHD)	6.5
2342		
2343		
2344		
2345		
2340		
2347		
2340		
2345	6.00	6.00
2350		
2001		
2002		
2353		
2334		
2300		
2330		
2337		<u>a</u>
2300		
2359		
2300		
2301	-	-
2362		
2303		
2304	100	100
2300		2.0
2300		
2001		
2000	Figure 36: MPI3D: Latant traversal with	h atomic group actions, object shape
2009	rigute 50. Wir 15D. Latent traversal with	n atomic group actions, object shape
2370		
2011		
2312		
2313		
2375		
2010		



Figure 37: MPI3D: Latent traversal with atomic group actions, object size



Figure 38: MPI3D: Latent traversal with atomic group actions, camera height



Figure 39: MPI3D: Latent traversal with atomic group actions, background color



