

1 We thank all reviewers for their valuable and constructive feedback. The consensus appears to be that the paper is well  
2 written, well motivated and easy to follow (R1, R2). We also appreciate the assessment of our work as “*innovative and*  
3 *unusual in the field*” (R1) and “*a very valuable contribution*” (R2).

#### 4 Questions and Comments of Reviewer 1

- 5 • **Resolution of the "3Dprinted-realistic" dataset.** Re-rendering and releasing higher resolution images is an  
6 important future work we plan to address. It is however beyond the scope of this paper due to the significant  
7 engineering challenges it entails. To illustrate the point, the rendering of 1m high resolution (512x512) images for  
8 the ‘3Dprinted-realistic’ image dataset, would roughly take 140 days of continuous computation using our existing  
9 infrastructure. Moreover, SOTA disentanglement learning algorithms [37] (Locatello et al., ICML 2019) still use  
10 64x64 resolution, such that already the provided 256x256 resolution is a significant improvement and for a resolution  
11 of 512x512 we directly provide the real-world dataset.
- 12 • **Downsampling of real-world images for transfer learning.** Thanks for pointing this out. We will make clear in  
13 the final version that in all experiments we used images with resolution 64x64. This resolution is used in recent  
14 large-scale evaluations and by SOTA disentanglement learning algorithms [37](Locatello et al., ICML 2019).
- 15 • **Link to code or the dataset.** Unfortunately, the NeurIPS’ policy does not allow to provide any external links  
16 in rebuttal documents. However, we guarantee to make all datasets and trained models publicly available before  
17 publication.

#### 18 Questions and Comments of Reviwer 2

- 19 • **"Is disentanglement helpful in general? and can authors comment on tasks associated with real-world disentan-**  
20 **glement?"** Thanks for pointing this out. We agree that currently there are multiple research directions for disentangled  
21 representation learning. One of them is the suggested investigation of the usefulness of disentangled representations.  
22 Our work is encouraging that and allowing to investigate the effectiveness of disentangled representations with access  
23 to ground truth labels on real-world data.  
24 Prior work e.g. [20] (Higgins et al. (ICML 2017) already indicated that disentanglement is useful for reinforcement  
25 learning. We believe that applying these ideas on the real-world platform is an interesting direction and we will  
26 update the future work section.  
27 As pointed out by the reviewer, the effort to find a thorough answer to these questions is significantly beyond the  
28 scope of one paper. However, we are already looking into replicating the analyses of the provided reference (van  
29 Steenkiste et al., arxiv 2019)<sup>1</sup> and a new fairness approach (Locatello et al. arxiv 2019)<sup>2</sup>, both illustrating the benefits  
30 of disentanglement, on our real-world data set.
- 31 • **"The authors state that the training on low resolution images results in the instability therefore random seed and**  
32 **hyper-parameters being more important than the model. Could the authors clarify this?"** Our intention was not  
33 to suggest that this instability is necessarily due to the low resolution of those datasets. We agree that this is slightly  
34 misleading and we will rephrase this sentence in the final version.
- 35 • **Used architecture.**  
36 All the methods use the same convolutional encoder and decoder where the latent dimension is fixed to be 10.  
37 *Encoder : input(64 × 64 × 3) → 4 × (4 × 4conv, 32ReLU, stride2) → FC1(256) → FC2(10)*  
38 *Decoder : 10 → FC1(256), ReLU → FC2(4 × 4 × 64), ReLU → 3 × (4 × 4upconv, 64ReLU, stride2) →*  
39 *4 × 4upconv, 3, stride2*  
40 This can be considered the standard architecture and is used in e.g. [4,11,19,38]. We will make sure to add this  
41 information in the final paper.
- 42 • **"Could the authors maybe comment on the difficulty of disentanglement (or reconstruction) for each individual**  
43 **feature? (which features are harder to disentangle?) I would be also interested to see which features cause the**  
44 **most difference in image space."** Empirically, we observed that some of the best trained models were able to  
45 disentangle, though imperfectly, the factors of camera height, background colors, object sizes and the motions along  
46 the first and second degrees of freedom. In contrast, they performed poorly in disentangling the factors of some  
47 object shapes (for e.g. pyramid and cone) and some colors (for e.g. olive and brown). This may well be because of  
48 less pixel variation in the respective factors. The features of camera height and background color cause the most  
49 difference (maximum L2 distance) in image space. Similarly the object positions at different joint configurations (not  
50 the consecutive frames) also have big L2 distance which may explain why the models focus more on learning these  
51 factors. We will add a detailed analysis to the final version.

<sup>1</sup>van Steenkiste, S., Locatello, F., Schmidhuber, J., and Bachem, O. Are disentangled representations helpful for abstract visual reasoning? arXiv preprint arXiv:1905.12506, 2019

<sup>2</sup>Locatello, F., Abbati, G., Rainforth, T., Bauer, S., Schölkopf, B. and Bachem, O. On the Fairness of Disentangled Representations. arXiv preprint arXiv:1905.13662. 2019.