

1 We thank the reviewers for their encouraging feedback. We will revise our writing as suggested (R1, R2 and R4), and
2 discuss the prior work mentioned by the reviewers in the final manuscript (R1, R2 and R4).

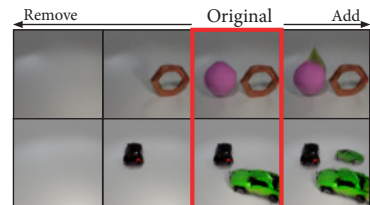
3 BlockGAN is a generative model that learns 3D **object-aware** scene representations using only *unlabelled* images.
4 We show that BlockGAN works with both synthetic datasets with simple backgrounds and real images with complex
5 background and lighting. We evaluated BlockGAN on 64×64 images, which are common for this line of work [1,2].
6 Due to resource constraints, we did not train BlockGAN on images with higher resolution, but instead chose to focus on
7 pushing the complexity, both in terms of number of objects and, especially, in terms of texture and cluttered background,
8 of the datasets. We believe that this is the first demonstration that deep 3D **object** representations can be learnt directly
9 from natural images without any template geometry, pretrained object detector, or multi-view input.

10 **R1: Claims:** Although pose is an input to our model, no GT pose labels were used for training. Hence, we maintain
11 our claim that we do *not* need any pose *supervision*. In addition to synthetic images with a simple background, we
12 also train BlockGAN on the real CAR dataset with complex, natural backgrounds (see Figures 5 and 6). Changing the
13 background object, in this case, changes not only lighting but also colour and texture.

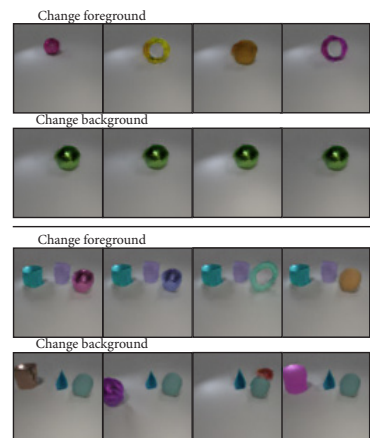
14 **Learnt renderer:** Yes, it is category-specific. We discuss a shared renderer as future work in line 131 of the supplement.
15 **Objects' appearance interaction:** As the foreground object moves, its appearance and shadow move accordingly,
16 depending on where the object is in relation to the camera view (specularity) and lighting positions (shadows) – this can
17 be observed most clearly in the animated results. Indeed, we leave more complex effects, such as inter-object reflection
18 between foreground objects, as future works as discussed in line 272 in the paper.

19 **3D vs 2D convolutions:** Our goal is to perform 3D transformations on deep 3D
20 features (including the background). Performing scene combination in 3D allows
21 representing geometry and appearance independent of camera specification.

22 **Removing objects:** We show adding and removing objects on the right →
23 **ShapeNet dataset:** ShapeNet contains only a limited number of textured models,
24 most of which are low quality, e.g., no specularity. Note that our SYNTH-CHAIR
25 dataset contains ShapeNet chairs with high-quality textures from PhotoShape.



26 **R2:** In the final version of the paper, we will also add a discussion on the projec-
27 tion/depth composition method, in addition to the rendering function. **Number of**
28 **objects:** On the right, we show BlockGAN with 2 foreground (FG) object gener-
29 ators trained with images containing 1 or 3 FG objects. **1 object (top):** Changing
30 either FG object changes the object's appearance and pose; changing the back-
31 ground works as expected. **3 objects (bottom):** Changing one FG object changes
32 one object as expected; changing the background changes one FG object and the
33 background. **Compositing function:** We perform max pooling across objects. This
34 does not require any learning, and, more importantly, is agnostic to the number of
35 inputs, allowing any number of objects to be added at test time. We have set up an
36 experiment with a learned linear weight per voxel, as suggested by R2; however,
37 the training collapsed, and we could not get it to work during the rebuttal period.



38 **R4: Performance of voxel grids:** Voxel grids can be more memory efficient when
39 adopting warping fields [Neural Volumes, SIGGRAPH 2019], a multi-resolution
40 strategy [Lighthouse, CVPR 2020] or sparsity [Neural Sparse Voxel Fields, arXiv
41 2020]. Moreover, HoloGAN [ICCV 2019] showed that voxel grids with low spatial resolution but high feature dimension
42 can be very expressive. Note that the choice of voxel grid does not affect whether shape and appearance can be separated
43 – both voxel grids [HoloGAN] and implicit functions [Texture Fields, CVPR 2020] can separate shape and appearance.
44 **Image encoding:** Many GAN models (apart from ALI and BiGAN) lack an inference (image encoding) mechanism.
45 However, recent work on training image encoders, such as Image2StyleGAN [ICCV 2019], exploit the representations
46 learnt by GANs to great effect. We look forward to extending BlockGAN towards this direction in future work.

47 **Alternative representations:** We thank the reviewer for pointing out [3], which we will cite and discuss. However, this
48 work requires predefined 3D mesh templates for each category (which are not always available) and a pretrained Mask
49 R-CNN to detect objects in images, while BlockGAN learns to disentangle and represent objects using only unlabelled
50 2D images. Scene Representation Networks [NeurIPS 2019] need many images with labelled poses to learn a good
51 representation for each scene. More importantly, only Visual Object Networks [NeurIPS 2018] and HoloGAN [ICCV
52 2019], which both use voxel grids, are trained successfully in an unsupervised manner, and can work across multiple
53 scenes. This makes voxel grids a reasonable and effective choice to achieve the goals of our paper.

54 **References:** [1] Investigating object compositionality in Generative Adversarial Networks. Neural Networks 2020.
55 [2] Towards Unsupervised Learning of Generative Models for 3D Controllable Image Synthesis, CVPR 2020. [3]
56 3D-Aware Scene Manipulation via Inverse Graphics. NeurIPS 2018.