1   We thank the reviewers for their valuable feedback and encouraging reception of our work. We're glad they found
2   our 3D object-centric model of videos to be highly significant, of interest to the NeurIPS community, and potentially
3   impactful. We now address some points raised in the reviews; we will of course incorporate the other suggestions.

4   **R1: videos are short; illustrate 3D with different cameras.** We have now trained a model on 12-frame (ROOMS)
5   videos, showing our method can scale to significantly longer sequences. Fig. A shows generated images, objects,
6   and depths, which still show coherent scenes. Fig. B shows reconstructions; the rows are original, reconstructed,
7   reconstructed with higher camera angle (to better show the 3D structure), and objects with 3D bounding boxes (from
8   that angle). The additional images to the right use even more extreme viewing angles, revealing the 3D layout clearly.

9   **R1: low resolution & fidelity.** Training video models is expensive, particularly in 3D. We consider the resolutions
10  used (96x72 & 80x80) a good trade-off between computation and quality; they are also comparable with similar works
11  [10,11,21]. Our results show O3V-voxel produces higher-quality videos than the state-of-the-art (SCALOR) on the
12  complex (TRAFFIC) dataset (Tab. 2b, Fig. 5 & S6). Moreover, our method is the first that can address our tasks and
13  setting; it is natural that scope remains for improving visual fidelity and resolution in future work.

14  **R1,2,4: complexity of datasets.** Our claim of handling more visually-complex videos than prior work was meant with
15  reference to the state-of-the-art in generative object-centric video modelling (SCALOR) [21]; we'll clarify this. [21] is
16  demonstrated only on objects of near-uniform color, without shading nor perspective/3D effects. O3V successfully
17  models videos containing all these effects, while SCALOR fails to do so (Fig. S6 & Tab. 1). While we do only
18  demonstrate O3V on synthetic videos (**R2**), our (TRAFFIC) dataset contains significantly more complex videos than
19  state-of-the-art (**R4**)—and so we already go a significant way towards bridging the gap to natural videos.

20  **R4: limited number of object slots $G$.** Many similar models have the same limitation [21,10,29], but we agree it
21  would be interesting to lift it. Note $G$ is rather large in our case (tens of objects).

22  **R1: correctness/clarity of L86-87.** Our intended meaning is simply that the grid is 3D, with its cells placed in 3D
23  world-space, rather than being a 2D grid in the space of the image. Such a 2D grid is used in the generative models
24  [6,21], inspired by YOLO [34], which produces detections based on a grid of cells. We'll clarify in the camera-ready.

25  **R1: overly strong claims; inductive biases are 3D supervision.** We note that numerous related models (e.g. SQAIR,
26  SCALOR, IODINE), regarded by the community as unsupervised, have similar inductive biases (priors on object size,
27  speed, uniformity of color), albeit in 2D. We therefore respectfully disagree with this characterisation. Of course, we do
28  agree that inductive biases make learning easier (indeed, possible!). We will qualify the claim of compositionality at
29  L36 to note this refers to the generative model itself, not the encoder used to train it.

30  **R4: explain differing generation quality.** We'll expand the discussion in Sec. 5.2. O3V-mesh has poor FID on
31  (TRAFFIC) as it is prone to local optima where cars are not tracked correctly—see the images in the supplementary.
32  GENESIS has good FID on (ROOMS) as color segmentation (which it readily exploits) is a strong cue here.

33  **R1,2: use of FID.** We agree KID is a more-principled generation metric than FID; we used FID for consistency with
34  the works we compare to. Following **R1**'s suggestion, we re-ran our generation evaluation using KID (see Tab. A). We
35  see that all statements made in the paper regarding relative quality of methods remain true—in fact, FID & KID are
36  highly correlated on our datasets. Note (**R2**) that in calculating FID, the ground-truth feature distribution is defined by
37  synthetic images, so the use of a 'natural image network' is reasonable (and, indeed, common [10,28]).

38  **R2,3: several components of O3V (e.g. mesh renderer) are known techniques.** This is true for very many works,
39  and there is clearly value in showing that a novel model using some existing techniques can achieve state-of-the-art
40  results. Importantly, we also examine different variants (i.e. mesh vs. voxel representations), and discuss their strengths.

41  **R1: authors do not show scene generation results from prior works.** These are in supplementary figures S1–S6.

42  **R3: motivation for discrete grid of objects.** We'll expand the discussion at L90. A grid of objects ensures gradients
43  with respect to object locations are non-zero, even when the current predictions are poor, as at least one candidate object
44  should be near enough each true object. Note that the actual object locations are continuous, as they are offset by a
45  vector $\Delta_g$ from the cell center.

46  **R1: compare to GQN and SRN.** These perform novel-view synthesis, but cannot sample new scenes *a priori*. They
47  also cannot perform segmentation/tracking, as they lack a representation of separate objects and assume a static scene.

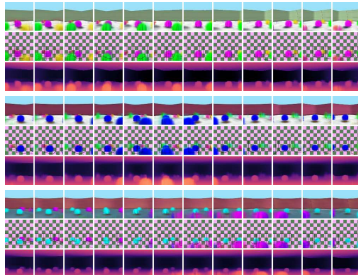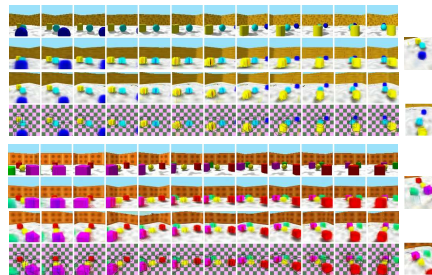|           | (ROOMS) | (TRAFFIC) |
|-----------|---------|-----------|
| MONet     | 0.151   | 0.305     |
| GENESIS   | **0.083** | 0.257   |
| SCALOR    | 0.148   | 0.272     |
| O3V-voxel | 0.108   | **0.157** |
| O3V-mesh  | 0.106   | 0.345     |

Table A: KID scores



Fig. A: 12-frame generations



Fig. B: 12-frame reconstructions