

1 We thank reviewers for their insightful and positive feedback! We are encouraged that they find our motivation and  
2 idea novel (R1), important (R1, R4), of value to the community (R1, R4), and has high potential significance (R1). We  
3 are glad they found our paper does not have major flaw (R2), well written (R2, R4), theoretically sound (R1, R4), our  
4 empirical evaluation well designed and thoroughly performed (R1, R2, R4), and the experimental results promising (R1,  
5 R2). We address the reviewer’s comments below and will incorporate all feedback.

6 **[R1] GNM distinguishes between purely symbolic and purely distributed representation while AIR and SPACE  
7 use it conservatively:** Not true. GNM does not require purely symbolic representations. In the structured representation  
8 layer, we are actually using hybrid (symbolic +distributed) representation which is exactly the same as SPACE and  
9 AIR, i.e., distributed for  $z^{\text{what}}$  and symbolic for the others. Yes, the distributed part is still groundable. We will make  
10 this point clearer. **Comparison to SPACE is missing:** SPACE uses independent priors for symbolic priors and thus *by*  
11 *design* cannot model the distribution of global scene structure. See also a relevant answer for R4. **Reporting model  
12 sizes and hyperparameters:** Yes, we use comparable model sizes. We will include the model sizes, hyperparameters,  
13 and training curve in the supplementary material. We will also release our code. **Increasing model size of GENESIS  
14 would make it better?** In our experiments, GENESIS used 4 times more parameters (14M) than our model (3M).  
15 This is because GENESIS uses additional conv-layers followed by MLP layers for the encoder and decoder of the  
16 segmentation networks which our model does not need. **Why not compare to existing datasets, e.g. CLEVR:** We  
17 agree that to compare on existing datasets can be better if possible. However, we found that objects in existing datasets  
18 are rather independently generated without complex spatial and inter-object relations and thus do not reveal required  
19 global scene structures. This makes it difficult to evaluate the learned global structure, which is of our main interest.  
20 For follow-up works, we will release our datasets. **Why GENESIS cannot obtain digit-wise decomposition for  
21 multi-MNIST scenes:** Although we have worked hard to make GENESIS decompose the multi-MNIST scene, we  
22 could not observe the decomposition for this dataset. We believe that our result is valid. First, similar to our result, the  
23 results in Figure 3 and Figure 4 in the GENESIS paper, exhibit strong color-based segmentation bias (e.g., the circles  
24 in the wall). We actually also observed that in our Arrow-Room dataset, GENESIS often grouped different objects  
25 into the same segmentation when they have the same color. A careful hyperparameter tuning was required to prevent  
26 this. We believe that this is a reasonable thing to observe in GENESIS because it does not use the *locality* principle  
27 of objectness. Second, it is indeed interesting to know that R1 has observed the decomposition on a similar dataset.  
28 It would have been great if we were able to access the result to compare with ours. However, given that there is no  
29 published result about the multi-MNIST task using GENESIS including R1’s experience, our result is the first and  
30 the only accessible result on which we can make any reasonable decision. **Strong claim that all prior methods have  
31 failed:** We fully agree and thanks for pointing this. We will tone down this and the sentence “this ability is not properly  
32 studied”. **Other modeling possibilities:** We agree that it is indeed an interesting suggestion. The SPACE paper shows  
33 that the auto-regressive approach does not scale gracefully and iterative inference (e.g., of IODINE) is known to be  
34 computationally expensive. However, as pointed, we note that our proposed framework is more generally applicable  
35 beyond the SPACE representation style. **Interaction MLP** is a standard MLP that takes the CNN feature map as input  
36 and mixes (interacts) globally. Lacking this, standard ConvDRAW do not mix information between features spatially  
37 distant. **Dual representation:** both  $z^s$  and  $z^g$  are inferred from the same input  $x$ . So they are different representations  
38 encoding the same information. **Beta parameter effect** is studied because some VAE-family models (ConvDRAW in  
39 our case) work better with a different beta. We thank all the **additional feedback** which will make our paper clearer.

40 **[R2] Why the particular form of the symbolic representation:** In this work, we have taken, as an example imple-  
41 mentation of the symbolic inference model  $q(z^s|x)$ , the SPACE model, and thus uses the same symbolic representation  
42 of theirs. However, the proposed framework is more general and not limited to the way SPACE represents an object. It  
43 can be applied to make any symbolic or hybrid (symbol + distributed) representation conform a proper distribution.  
44 **More visualization would be helpful although Figure 5 is “already really nice”:** We thank you for suggesting this.  
45 We will add more visualization if we come up with an additional idea for that. **Why not categorical variable for  $z^{\text{what}}$ :**  
46 because the categorical value cannot model within-category variation.

47 **[R4] Need the evaluation of structured representation learning ability.** Thanks for suggesting this. We initially  
48 thought that we do not need this experiment because our structured representation learning is basically the same as  
49 SPACE (with a single segment for the background modeling). However, following the suggestion, we performed the  
50 suggested experiments by comparing ours and SPACE. We indeed obtained a similar performance between ours and  
51 SPACE. Specifically, for average precision IOU on the MNIST-10 dataset, we obtained 0.459 and 0.453 for our model  
52 and SPACE, respectively. And for classification accuracy on the recognized  $z^{\text{what}}$  representation, we obtained 0.984 and  
53 0.980 for ours and SPACE, respectively. The SPACE paper demonstrates its advantages against MoNET and IODINE  
54 in some settings and without loss of generality, we can extend that result to our model. We will add this result to  
55 the revision. **Manual checking data:** As suggested, we will publish the corresponding data we used for the manual  
56 investigation. **FID Score:** We will consider adding FID score. However, as agreed, it may not affect the main result of  
57 the paper. **Missing reference:** Thanks! We will add the suggested paper.