

1 **Overview.** We would like to thank all reviewers for their thorough reviews and helpful feedback/suggestions.

2 **Generalization Experiments.** Three of four reviewers asked about generalization to novel objects and scenes. To
 3 address this, we trained MulMON, IODINE and GQN on CLE-Aug. Then, we compared their performance on CLE-MV
 4 and 2 new datasets—Black-Aug and UnseenShape. Black-Aug contains the CLE-Aug objects but only in single, unseen
 5 colour (black). This tests the models’ ability perform segmentation without colour differences/cues. UnseenShape
 6 contains only novel objects that are not in the CLE-Aug dataset—cups, cars, spheres, and diamonds. This directly
 7 tests generalization capabilities. Both datasets contain 30 scenes (more for final paper), each with 10 views. Table 1
 8 shows that 1) all models generalize well to novel scenes, 2) MulMON still performs best for all tasks but observation
 9 prediction—where GQN does slightly better due to its more direct prediction procedure (features → layout vs. features
 10 → objects → layout), 3) MulMON can indeed understand the composition of novel objects in novel scenes—impressive
 11 novel-view predictions (observations and segmentations) and disentanglement.

Tasks	Models	CLE-Aug (train)	CLE-MV	Black-Aug	UnseenShape
Seg. (mIoU)	IODINE	0.51 ± 0.001*	0.61 ± 0.002	0.50 ± 0.006	0.51 ± 0.004
	MulMON	0.71 ± 0.000	0.71 ± 0.004	0.67 ± 0.002	0.64 ± 0.004
Pred.Obs (RMSE)	GQN	0.15 ± 0.000	0.15 ± 0.001	0.24 ± 0.003	0.17 ± 0.002
	MulMON	0.07 ± 0.000	0.16 ± 0.002	0.26 ± 0.002	0.21 ± 0.006
Disent. (D,C,I)	IODINE	0.54, 0.48, 0.21	0.14, 0.12, 0.26	0.2, 0.26, 0.27	0.13, 0.12, 0.26
	MulMON	0.63, 0.54, 0.68	0.52, 0.48, 0.63	0.55, 0.55, 0.66	0.5, 0.47, 0.67
Pred.Seg (mIoU)	MulMON	0.69 ± 0.001	0.71 ± 0.004	0.68 ± 0.005	0.60 ± 0.005

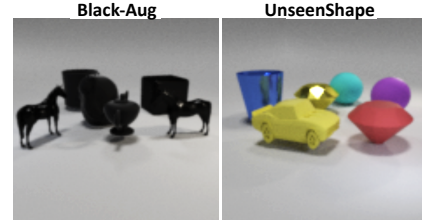


Table 1: Models’ generalization performance. *Paper correction. **Figure 1:** Samples from novel-scene datasets.

12 **Writing.** We would like to thank the reviewers for feedback on writing style, including both structuring and typos,
 13 as well as related work. In particular, we would like to thank R1 for detailed suggestions to improve the abstract,
 14 introduction and related work sections. These have all been incorporated, and have greatly improved the paper’s clarity.

15 **R4: ‘the contribution in the learning method seems not enough. Most of the method are based on the previous work
 16 (IODINE)’.** IODINE’s approximation of the MOSV posterior cannot maintain object correspondence or “matching”
 17 across multiple views, and thus, it is not a feasible solution to the MOMV problem. To address this, we make two
 18 contributions in the learning method itself: 1) we side-step the object matching problem by iterating over multiple views,
 19 each time using the previous iteration’s approximate posterior as the new prior—see eq. 4; 2) we introduce a training
 20 procedure that forces the model to make use of these iterative updates, aggregating spatial information across multiple
 21 views. Specifically, by randomly partitioning the views of a scene into two subsets, we can ask the model to predict
 22 the views in one (novel viewpoint-queried generation) having observed the views in the other (scene learning)—see
 23 eq. 5. This forces the model to use the iterative updates to aggregate spatial information across multiple views, as it
 24 needs to form a complete 3D scene understanding in order to perform well. If instead we had asked the model to simply
 25 reconstruct the observed view, as IODINE does, it would completely overwrite the prior on each iteration—as it would
 26 not need to aggregate spatial information across multiple views in order to perform well.

27 **R1: ‘results for scene factorization, novel view predictions and disentanglement... only 1, 2 and 3 observations’.**
 28 Figure 3 in Appendix E plots novel-viewpoint observation and segmentation results as a function of the number of
 29 views T (1-10). However, we appreciate R1’s suggestion of adding baselines (IODINE, GQN) to these analyses, and
 30 also analyzing both segmentation and disentanglement results as a function of the number of views T . We will definitely
 31 include these in the final paper. **‘Can the trained model be used to generate random scenes?’** Yes, we can generate
 32 random scenes by composing independently-sampled objects. However, to focus on forming accurate, disentangled
 33 representations of multi-object scenes, we must assume objects are i.i.d. and thus ignore inter-object correlations—e.g.
 34 two objects can appear at the same location. Nonetheless, we will include random samples in the final paper. **‘not clear
 35 ... how many views the model sees during training’.** Using a fixed number of observations could harm the model’s
 36 robustness at test time. Thus, the number of observations given to the model is sampled at every training step—from the
 37 range [1, 6] for the GQN-Jaco dataset, and from [1, 5] for the CLEVR-based datasets. Then, the model is asked to predict
 38 the remaining, unseen views. Note that GQN-Jaco dataset has a total of 11 views, while the CLEVR-based datasets
 39 have 10. This sampling procedure is introduced as a *randomized partition of T observations into two subsets \mathcal{T} and*
 40 *\mathcal{Q}* in lines 191-192 in the paper, but we will revise the text to further clarify this. **‘A video that shows segmentations
 41 when moving the camera...’** As we cannot include external links here, we will create one for the final version.

42 **R3: ‘include a score that ignores background, or ... discuss this difference’.** We thank R3 for pointing this out, and
 43 we agree that IODINE’s poor mIoU is mostly due to its handling of the background. We will discuss this issue in the
 44 paper, and add a score (e.g. Adjusted Rand Index) to our comparisons that ignores the background. However, good
 45 models shouldn’t split up such simple backgrounds—see Fig. 3 in our paper and Fig. 1 in IODINE paper. **‘In what
 46 sense is the set $\mathbf{z} = \{z_k\}$ grouped into T groups based on the views?’** R3 is correct in that this design is just to show
 47 that we update \mathbf{z} one piece (i.e. z^t) at a time w.r.t. its corresponding x^t . We will add some discussions to clarify this.