1 We thank the four reviewers (**R1**-**R4**) for their constructive feedback and largely positive comments.

2 **Improving readability.** Several reviewers found the paper clear and "exceptionally well written" (**R3**), but they also
3 noted the density of the Methods and (esp. **R1**) felt it was overly formal. We agree that this section should better
4 get across the key ideas, so we will use the extra page allowed in revision to provide higher-level explanations and
5 background for each concept we introduce. To further improve readability, we will also make the following changes:

6 ★ Lines 108-120, the description of PSGs, will be made less formal. The key point is that PSGs are hierarchical graphs
7 with additional structures that connect to visual input. Within-level edges are meant to represent physical connections;
8 parent-to-child edges are meant to represent part-whole relationships; attribute vectors are meant to encode physical
9 properties of elements of the visual input; and spatial registrations map each vertex in the graph to a subset of pixels in
10 the input movie. Because a PSG is hierarchical, the registrations form a hierarchical segmentation of a scene.

11 ★ Lines 121-133 should better address **R1**'s questions about model training. In particular, PSGNets have two types
12 of losses that are *trained jointly*: (1) *rendering losses*, which train the parameters of Graph Vectorization and Feature
13 Extractor modules by backpropagation, because the rendered feature maps are differentiable functions of predicted
14 attribute vectors and image features, respectively; and (2) *perceptual grouping* losses, which train the affinity functions
15 in Graph Pooling modules. Rendering gradients *cannot* flow back into the affinity functions because Label Propagation
16 is not differentiable; perceptual grouping gradients *can* flow into the other modules. To clarify L258 (per **R2**), we
17 will emphasize that depths and normals can *supervise* rendering losses (1) but are never used as input. We describe
18 further experiments without any depths and normals supervision below. Perceptual grouping losses (2) never receive
19 supervision. We will implement **R3**'s idea to make Fig. 1 larger by splitting it into Architecture and Training figures.

20 ★ We can simplify the rest of the Methods as five shorter subsections: (1) ConvRNN Feature Extraction: explain that a
21 ConvRNN, unlike a deep CNN, generates features *from a single layer* with different useful properties after each recurrent
22 pass (see Supplement); (2) Graph Pooling: emphasize ideas of affinity functions and explain why Label Propagation
23 prevents rendering gradients from training this module; (3) Graph Vectorization: Move details of aggregation to
24 Supplement (e.g. L176-184, per **R1**), emphasize taking statistics over segment interiors and boundaries (to encode
25 shape information) and predicting attributes *via* Graph Convolution [Kipf & Welling, *ICLR* 2017]; (4) Rendering:
26 Combine the QTR, QSR, and Losses sections to explain why parameter-free decoding forces node attributes to encode
27 scene properties *explicitly*; discuss (per **R1**) why this leads to interpretable representations, as illustrated by graph
28 editing (Figure 4); (5) Perceptual Grouping: explain why pairwise inductive biases make sense for grouping and how
29 $\beta$-VAEs [Higgens *et al.*, *ICLR* 2016] can naturally encode the idea of node co-occurrence and motion-in-concert.

30 **Data to address Reviewer 1. R1** makes an excellent suggestion to compare the MultiDSprites (MDS) and CLEVR6
31 datasets with our custom datasets. We thus built PSGNetS-RGB, a compact version of the original model that does not
32 use depths/normals supervision. Trained/tested on the Playroom dataset, PSGNetS-RGB achieved scores of (**recall**
33 0.59, **mIoU** 0.55, **boundary F** 0.57), which are (as expected) somewhat, though not dramatically, lower than the
34 original PSGNetS scores. Without further hyperparameter tuning, PSGNetS-RGB trained/tested on MDS achieved
35 (**r**0.80, **m**0.70, **b**0.72); trained/tested on CLEVR6, it achieved (**r**0.73, **m**0.63, **b**0.67). The higher scores on MDS and
36 CLEVR6 suggest that the Playroom dataset is indeed harder to segment than those used in the literature. Moreover, in
37 running these experiments we found something else worth reporting: the Playroom-trained PSGNetS-RGB achieved
38 *zero-shot transfer performance* of (**r**0.75, **m**0.67, **b**0.71) on MDS and (**r**0.70, **m**0.60, **b**0.66) on CLEVR6, nearly as
39 high as the within-dataset scores. This further indicates that PSGNet inductive biases lead to fairly general object-
40 centric representation learning. We will include these points in the revision, along with MDS/CLEVR6 scores for our
41 implementations of the baseline models (which take longer to train than the allotted Author Response period.)

42 **Data to address Reviewers 2 & 3.** The Playroom-trained PSGNetS-RGB model above ablates the three $\delta x$ input
43 channels, as per **R2**. To measure the effect of ablation, we restored these channels to create PSGNetM-RGB and trained
44 five models with different random seeds. These achieved scores (mean +/- stdev) of (**r**0.65 +/ 0.01, **m**0.59 +/- 0.01,
45 **b**0.60 +/- 0.01), substantially higher than PSGNetS-RGB. This indicates that the $\delta x$ input channels are useful for
46 segmentation and supports our claim of low performance variability, noted by **R3**. Each model took ~24 hours to train
47 on one Titan Xp GPU; inference takes ~200 ms/image. PSGs are implemented on a GPU by choosing a *maximum* node
48 number per level, then masking out nodes that have no incoming child-to-parent edges (& thus represent nothing.)

49 **Response to Reviewer 4. R4** raised two concerns about (1) the types of edges in the PSG and (2) the inability to group
50 by real-world textures. We want to clarify that these are *not* inherent limitations of our approach, but rather choices
51 we made to limit the scope of the submission. As to (1), PSGs can naturally incorporate additional edge types, for
52 example to represent temporary occlusions, collisions, or support relationships between objects. We use these edge
53 types in ongoing work to predict physical dynamics with PSGs. As to (2), by self-supervising rendered feature maps on
54 higher-order statistics of the input image, such as the covariance of RGB pixels within each node segment, attribute
55 vectors can be trained to explictly encode texture. Perceptual grouping can take advantage of these texture components.