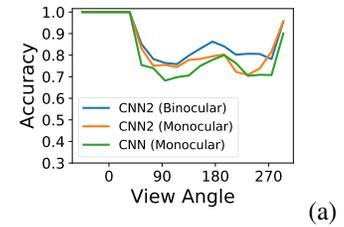
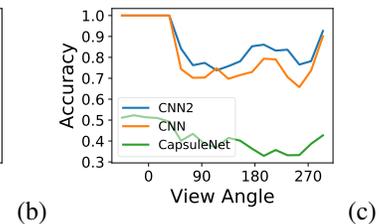
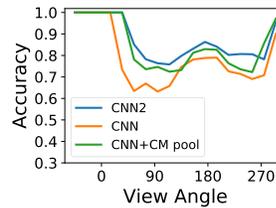


1 **To Reviewer 1.** Thanks for your positive comments. **Q1:** How the proposed archite-
 2 cture would fair on ... single-image object classification? **A1:** Good question! With
 3 monocular images, the parallax channels contain all zeros, therefore the CNN² de-
 4 generates into a conventional CNN gracefully. Fig. (a) shows the performance of
 5 degenerated CNN² with the single-eye images from the RGB-D Object dataset. **Q2:**
 6 Does the accuracy of CNN² improve as more filters are added? **A2:** As shown in Table
 7 3 in the supplementary file, increasing the number of filters in CNN² does not guarantee
 8 performance gain. **Q3:** Does the size of the features remain the same as the input as
 9 one moves up the layers due to CM pooling? Is it necessary? **A3:** The size does not
 10 need to be the same as the input nor across layers, but the size of feature maps for the
 11 same filter must be of the same size in order to allow the filter to detect stereoscopic patterns.

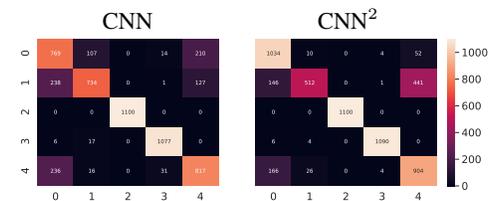


12 **To Reviewer 2.** Thanks for your constructive comments.

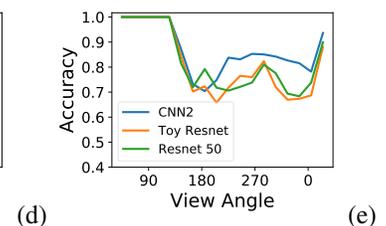
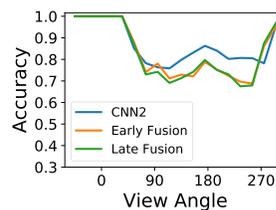
13 **Q1:** Ablation study of CM pooling on vanilla CNN. **A1:**
 14 Fig. (b) shows the performance of vanilla CNN with
 15 CM pooling over the RGB-D Object dataset. The CM
 16 pooling can indeed help the vanilla CNN detect useful
 17 features. **Q2:** Vanilla CNN tuning details. **A2:** The
 18 table below summarizes our model candidates for the
 19 RGB-D Object dataset. The vanilla CNN gives relatively
 20 stable performance during the hyperparameter search.
 21 **Q3:** Scale up the number of instances in each class ...
 22 using ShapeNet. **A3:** As suggested, we conduct new
 23 experiments using the ShapeNet dataset following the settings for ModelNet2D. Now, each class (airplanes, cars,
 24 cameras, lamps, and chairs) has at least 100 instances. Fig. (c) shows that CNN² still outperforms CNN and CapsuleNet.
 25 **Q4:** Confusion matrices for classification. **A4:** Below please see the confusion matrices of the predictions made by
 26 CNN and CNN² at unseen view angles on the RGBD-Object dataset. The CNN² outperforms CNN in most cases,
 27 except when classifying the classes 1 (flashlight) and 4 (stapler) that are similar in shape but different in texture at
 28 certain view angles. This suggests that the CNN² relies more on shapes than textures to generalize, a bias that humans
 29 have been shown to possess (Matthias Bethge et al., "ImageNet-trained CNNs are biased towards texture; increasing
 30 shape bias improves accuracy and robustness," ICLR'19).



# Channels					Task Model		Unseen Avg.
16(5)	16(5)	32(3)	32(3)	32(3)	256(2)	256(1)	0.776
16(5)	16(5)	32(3)	32(3)	32(3)	512(2)	512(1)	0.781
16(5)	16(5)	32(3)	48(3)	48(3)	512(2)	512(1)	0.788
16(5)	16(5)	48(3)	48(3)	48(3)	512(2)	512(1)	0.795
16(5)	16(5)	32(3)	48(3)	64(3)	512(2)	512(1)	0.783



33 **To Reviewer 3.** Thanks for your comments. **Q1:** How
 34 about the late fusion networks...? **A1:** As suggested, we
 35 compare CNN² with two new baselines that perform early
 36 and late "fusion" (i.e., dual parallax augmentation) in only
 37 the first and last layer, respectively. Fig. (d) shows the
 38 results on the RGB-D Object dataset. CNN² outperforms
 39 other baselines because it has fusion at all layers, which
 40 allows small differences between the feature maps in two
 41 paths to add up to a big difference at a deeper layer. **Q2:**
 42 The experimental settings are toy-like. **A2:** Please note
 43 that our classification tasks are tested at view angles that
 44 are *unseen* during the training time. Comparing to traditional image classification, these tasks are very challenging,
 45 and our settings are already more complex than the ones used by Hinton et al. in their CapsuleNet work published in
 46 ICLR'18, which considered only grayscale images. The CNN² has advanced the state-of-the-art performance on the
 47 grayscale ModelNet2D and SmallNORB datasets and, for the first time, gives improved 3D viewpoint generalizability
 48 on the colored RGB-D Object dataset. **Q3:** The neural network backbone is weak. **A3:** A backbone, e.g. ResNet,
 49 that is strong to make predictions at seen angles does *not* imply that it is strong at unseen angles. To show this, we
 50 compare the performance of CNN² with ResNet-50 and a toy ResNet having a similar number of parameters as CNN²
 51 on the SmallNORB dataset. The results are shown in Fig. (e). Although the SmallNORB dataset contains only
 52 grayscale images and looks "easy," neither of the ResNet variants generalizes better than CNN². We will add the above
 53 experiments to the paper.



54 We hope our above explanation relieves your concerns, and if so, please consider raising your score.