| Method | RGB | Full mIOU | Invisible mIOU |
|---|---|---|---|
| Amodal-VAE | ✗ | **94.68** | **62.85** |
| Deocclusion | ✓ | 94.04 | 57.19 |
| ResNet-Amodal-VAE | ✓ | 94.53 | 61.97 |

Table 1: **Amodal Completion on KINS**. Checkmark on RGB column means the RGB image is used as input to the model.

| Method | Full mIOU | Inv. mIOU |
|---|---|---|
| GT Instance Mask | 87.03 | 0 |
| Pred. Mask | 80.83 | 0 |
| Amodal-VAE | **86.03** | **60.71** |
| Deocclusion | 84.94 | 52.93 |

Table 2: **Amodal Segmentation on KINS** Models take predicted segmentation mask as input.

| Method | Full mIOU | Invisible mIOU |
|---|---|---|
| Amodal-VAE | 94.68 | 62.85 |
| Search by visible area | 94.71 | 62.98 |
| Search by amodal GT | 95.82 | 69.96 |

Table 3: **Results for multiple predictions using samples.** We sample masks from approx. posteriors and search for best samples via visible area or full amodal GT area IOU. "Amodal-VAE" uses only the posterior mode.
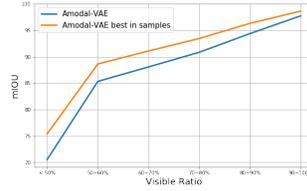


Figure 4: **Amodal Segmentation mIOU with different occlusion ratio on KINS.** X-axis represents visible area ratio. Better in zoomed-in view.

1 We would like to thank the reviewers for their constructive feedback.

2 **R1&R2: Novelty.** Although methodologically related approaches have appeared in the literature before, we focus
3 on amodal object mask completion where probabilistic methods have not been explored. For this task, we are the first to
4 propose a generative framework, which naturally captures the ambiguity inherent in the mask completion. Since naively
5 applying a VAE is not possible, we propose a novel training strategy. Compared to previous work that uses VAEs for
6 3D shape completion [30], we further perform occlusion simulation and bounding box prediction, add latent code
7 regularization during training. These contributions allow us to (1) achieve state-of-the-art results on the amodal mask
8 completion task without amodal ground truth supervision. (2) Our mask completions outperform the drawing skills of
9 human annotators. (3) We demonstrate the usefulness of our model on the downstream task of scene editing. (4) We
10 qualitatively and quantitatively (Tab. 3, Fig. 4 above) show the value of explicitly modeling the uncertainty in the task.

11 **R2&R4: Will RGB input make the task trivial?** We slightly modify the model architecture and provide additional
12 results. After the full-mask-only training stage, we use a ResNet-50 pretrained on image-net and taking RGB images
13 as input. We concatenate the ResNet's features and the mask encoder, add two convolution layers, and predict latent
14 code posterior distributions. As shown in Table 1, the additional RGB-based image features do not boost performance.
15 Hence, the mask together with its object class contains enough information and we discard RGB input for simplicity.

16 **R2: Motivation of amodal completion & Amodal instance segmentation is more interesting.** We agree that
17 amodal instance segmentation is also an interesting problem. However, amodal instance segmentation can be decom-
18 posed into instance segmentation and amodal completion. We provide additional experiments on amodal segmentation
19 (results in Tab. 2) and will update the paper. We crop images by GT bounding box and train a ResNet-PSP segmentation
20 model. We predict the instance masks on the test set, yielding 80.83% amodal mIOU. We use the predicted instance
21 masks as input to perform amodal completion using both our model and the baseline (Deocclusion [37]), which we
22 outperform. We also find that our model is robust to instance mask corruptions. Even though full mIOU drops by 8.65%
23 (94.68% vs. 86.03%), mIOU on invisible area only drops by just 2.14% (62.85% vs. 60.71%).

24 **R2: Lack of quantitative results for multiple posterior predictions.** We agree and quantitatively evaluate this
25 aspect (paper will be updated). For each instance, we sample 20 latent codes from the approx. posterior distribution.
26 We calculate mIOU using masks with the best visible area IOU or best amodal GT IOU. Results show that by sampling
27 we find masks that significantly better match amodal GT than using the approx. posterior mode (Tab. 3 above). Hence,
28 the approx. posteriors incorporate diverse plausible masks, correctly capturing the ambiguity. This suggests the use of
29 multiple posterior samples in downstream applications. Sampling is more beneficial for cases with significant occlusion
30 (Fig. 4 above) and we also show (Fig. 3 of supp. material) that more occlusion correctly results in wider posteriors.

31 **R1&R3&R4: Results are only shown on KINS dataset.** We primarily focus on driving scenes and exploit two
32 datasets with multiple classes, KINS and Cityscapes. KINS contains 7 categories (pedestrian, cyclist, person sitting, car,
33 van, tram, truck) and Cityscapes contains 8 categories (bicycle, bus, person, train, truck, motorcycle, car, rider). We
34 show quantitative results only on KINS, because it is the only available large-scale amodal dataset with accurate human
35 annotations. We plan to show quantitive results on more classes and non-rigid objects for the camera-ready version.

36 **R1 & R4: Comparisons for downstream tasks.** We appreciate the suggestions. However, our main contribution is
37 on the amodal completion task. Comparisons to other downstream task methods are beyond the scope of our paper.
38 Nevertheless, we compared FID scores with Deocclusion [37] on the image editing task.

39 **R1: Need to compare with more baselines & Only one category.** We apologize for the confusion and will clarify
40 this aspect in the camera-ready version. However, all experiments are conducted on datasets with multiple classes. We
41 use Deocclusion [37] as baseline, since it is the current state-of-the-art model on amodal completion with no amodal
42 ground truth supervision. Also, our user study demonstrates that we outperform the annotation skills of humans, which
43 further strengthens our conclusions. We will add additional details about the user study in the final version of the paper.

44 **R1: Transformation network.** The spatial transformation network is end-to-end differentiable and trained as part of
45 the pipeline in stages (2) and (3), using the mask reconstruction loss as supervision (we will clarify in the final version).

46 **R3: Accuracy versus occlusion percentage .** We agree. We show results in Fig. 4

47 **Details, etc:** We appreciate the suggestions and will incorporate them in the camera-ready version.