# Appendix

## *Anomalies-by-Synthesis*: Anomaly Detection using Generative Diffusion Models for Off-Road Navigation

Sunshine Jiang[*1], Siddharth Ancha[*1], Travis Manderson[1], Laura Brandt[1],
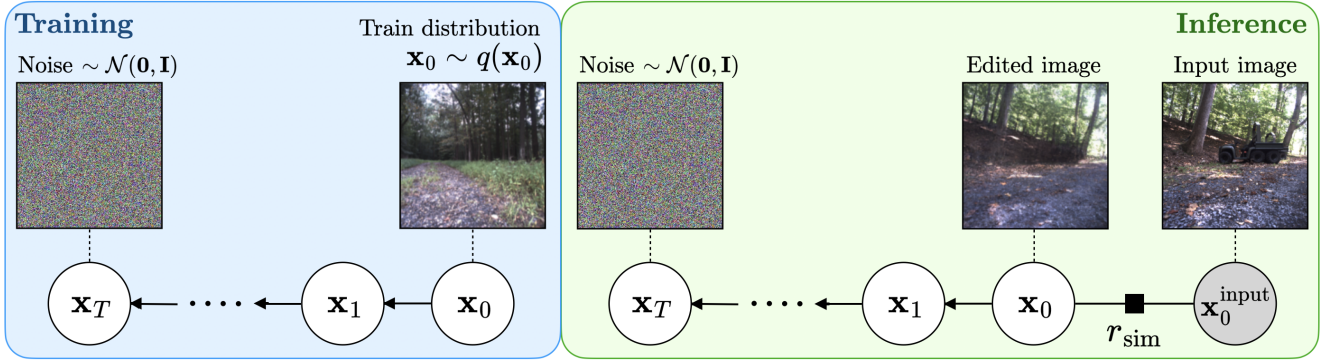Yilun Du[1], Philip R. Osteen[2] and Nicholas Roy[1]

**Fig. 4: Probabilistic graphical model of the *conditional* forward diffusion process.** We wish to sample the random variable $\mathbf{x}_0$ corresponding to the training data distribution $q(\mathbf{x}_0)$. The *directed* edges between $\mathbf{x}_{t-1}$ and $\mathbf{x}_t$ (for $t = 1, \ldots, T$) correspond to the vanilla forward diffusion process. Each directed edge denotes the sampling distribution $q(\mathbf{x}_t \mid \mathbf{x}_{t-1})$ which successively adds a small amounts of Gaussian noise: $q(\mathbf{x}_t \mid \mathbf{x}_t) = \mathcal{N}\left(\mathbf{x}_t ; \sqrt{1 - \beta_t}\,\mathbf{x}_{t-1}, \beta_t \mathbf{I}\right)$ [32, 38]. However, we are interested in sampling from $q(\mathbf{x}_0 \mid \mathbf{x}_0^{\text{input}}) = q(\mathbf{x}_0), r_{\text{sim}}(\mathbf{x}_0, \mathbf{x}_0^{\text{input}})$. This objective corresponds to adding an additional *undirected factor* $r_{\text{sim}}(\mathbf{x}_0, \mathbf{x}_0^{\text{input}})$ between $\mathbf{x}_0$ and $\mathbf{x}_0^{\text{input}}$; $\mathbf{x}_0^{\text{input}}$ is treated as constant. Our task is to perform inference over this graphical model and sample from $q(\mathbf{x}_0 \mid \mathbf{x}_0^{\text{input}})$ using a diffusion model that was trained to perform reverse diffusion in the absence of the $r_{\text{sim}}$ factor.

## APPENDIX I
### PROOF: GENERALIZED CONDITIONAL GUIDANCE GRADIENT

Classifier guidance is not only restricted to classifiers, it also requires training a classifier $p(y \mid \mathbf{x}_t)$ for each intermediate latent state [34, 39]. First, we extend classifier guidance to the more general setting where the conditioner is any non-negative function:

***Theorem 1:*** When a diffusion model $\epsilon_t^\theta$ is trained to sample from $q(\mathbf{x}_0)$, the conditional distribution $q(\mathbf{x}_0 \mid \mathbf{x}_0^{\text{input}}) \propto q(\mathbf{x}_0)\, r_{\text{sim}}(\mathbf{x}_0, \mathbf{x}_0^{\text{input}})$ can be sampled by using the following guidance gradient during reverse diffusion:

$$\mathbf{g}_t(\mathbf{x}_t) = \nabla_{\mathbf{x}_t} \log \int_{\mathbf{x}_0} q(\mathbf{x}_0 \mid \mathbf{x}_t)\, r_{\text{sim}}(\mathbf{x}_0, \mathbf{x}_0^{\text{input}})\, d\mathbf{x}_0$$

$$= \nabla_{\mathbf{x}_t} \log \mathbb{E}_{q(\mathbf{x}_0 \mid \mathbf{x}_t)} \left[ r_{\text{sim}}(\mathbf{x}_0, \mathbf{x}_0^{\text{input}}) \right] \quad (4)$$

### A. Proof using the variational inference perspective

$$q(\mathbf{x}_t \mid \mathbf{x}_{t+1}, \mathbf{x}_0^{\text{input}}) \propto q(\mathbf{x}_t, \mathbf{x}_{t+1} \mid \mathbf{x}_0^{\text{input}})$$

$$= \int_{\mathbf{x}_0} q(\mathbf{x}_t, \mathbf{x}_{t+1}, \mathbf{x}_0 \mid \mathbf{x}_0^{\text{input}})\, d\mathbf{x}_0$$

$$= \int_{\mathbf{x}_0} q(\mathbf{x}_0 \mid \mathbf{x}_0^{\text{input}})\, q(\mathbf{x}_t, \mathbf{x}_{t+1} \mid \mathbf{x}_0, \mathbf{x}_0^{\text{input}})\, d\mathbf{x}_0$$

$$= \int_{\mathbf{x}_0} q(\mathbf{x}_0 \mid \mathbf{x}_0^{\text{input}})\, \underbrace{q(\mathbf{x}_t, \mathbf{x}_{t+1} \mid \mathbf{x}_0)}\, d\mathbf{x}_0$$
*(since $\mathbf{x}_{1:T}$ is independent of $\mathbf{x}_0^{\text{input}}$ conditioned on $\mathbf{x}_0$)*

$$\propto \int_{\mathbf{x}_0} \underbrace{q(\mathbf{x}_0)\, r_{\text{sim}}(\mathbf{x}_0, \mathbf{x}_0^{\text{input}})}\, q(\mathbf{x}_t, \mathbf{x}_{t+1} \mid \mathbf{x}_0)\, d\mathbf{x}_0$$
*(by the definition of $q(\mathbf{x}_0 \mid \mathbf{x}_0^{\text{input}}) \propto q(\mathbf{x}_0)\, r_{\text{sim}}(\mathbf{x}_0, \mathbf{x}_0^{\text{input}})$)*

$$= \int_{\mathbf{x}_0} r_{\text{sim}}(\mathbf{x}_0, \mathbf{x}_0^{\text{input}})\, q(\mathbf{x}_t, \mathbf{x}_{t+1}, \mathbf{x}_0)\, d\mathbf{x}_0$$

$$= \int_{\mathbf{x}_0} r_{\text{sim}}(\mathbf{x}_0, \mathbf{x}_0^{\text{input}})\, q(\mathbf{x}_{t+1})\, q(\mathbf{x}_t \mid \mathbf{x}_{t+1})\, q(\mathbf{x}_0 \mid \mathbf{x}_t, \mathbf{x}_{t+1})\, d\mathbf{x}_0$$

$$= \int_{\mathbf{x}_0} r_{\text{sim}}(\mathbf{x}_0, \mathbf{x}_0^{\text{input}})\, q(\mathbf{x}_{t+1})\, q(\mathbf{x}_t \mid \mathbf{x}_{t+1})\, \underbrace{q(\mathbf{x}_0 \mid \mathbf{x}_t)}\, d\mathbf{x}_0$$
*(since $\mathbf{x}_0$ and $\mathbf{x}_{t+1}$ are independent conditioned on $\mathbf{x}_t$)*

$$\propto q(\mathbf{x}_t \mid \mathbf{x}_{t+1}) \int_{\mathbf{x}_0} q(\mathbf{x}_0 \mid \mathbf{x}_t)\, r_{\text{sim}}(\mathbf{x}_0, \mathbf{x}_0^{\text{input}})\, d\mathbf{x}_0$$

$$= q(\mathbf{x}_t \mid \mathbf{x}_{t+1})\, \mathbb{E}_{q(\mathbf{x}_0 \mid \mathbf{x}_t)} \left[ r_{\text{sim}}(\mathbf{x}_0, \mathbf{x}_0^{\text{input}}) \right]$$

Therefore, following the same analysis as the first-order Gaussian approximation, the guidance gradient is $\nabla_{\mathbf{x}_t} \log \mathbb{E}_{q(\mathbf{x}_0 \mid \mathbf{x}_t)}\big[r_{\text{sim}}(\mathbf{x}_0, \mathbf{x}_0^{\text{input}})\big]$ evaluated at $\mathbf{x}_t = \boldsymbol{\mu}_{t+1}(\mathbf{x}_{t+1})$.

### B. Proof using the score functions perspective

The score function of intermediate states $\mathbf{x}_t$ of the vanilla forward diffusion process is defined as $\mathbf{s}(\mathbf{x}_t) = \nabla_{\mathbf{x}_t} \log q(\mathbf{x}_t)$. However, we're interested in the score function of the *conditional* forward diffusion process $\mathbf{s}(\mathbf{x}_t \mid \mathbf{x}_0^{\text{input}}) = \nabla_{\mathbf{x}_t} \log q(\mathbf{x}_t \mid \mathbf{x}_0^{\text{input}})$. The additional term that needs to be added to $\mathbf{s}(\mathbf{x}_t)$ to obtain $\mathbf{s}(\mathbf{x}_t \mid \mathbf{x}_0^{\text{input}})$ is the guidance gradient. Therefore, we now derive $\mathbf{s}(\mathbf{x}_t \mid \mathbf{x}_0^{\text{input}})$ in terms of $\mathbf{s}(\mathbf{x}_t)$:

$$\underbrace{\mathbf{s}(\mathbf{x}_t \mid \mathbf{x}_0^{\text{input}})}_{\text{conditional score function}} = \nabla_{\mathbf{x}_t} \log q(\mathbf{x}_t \mid \mathbf{x}_0^{\text{input}})$$

$$= \nabla_{\mathbf{x}_t} \log \int_{\mathbf{x}_0} q(\mathbf{x}_t, \mathbf{x}_0 \mid \mathbf{x}_0^{\text{input}}) \, \mathrm{d}\mathbf{x}_0$$

$$= \nabla_{\mathbf{x}_t} \log \int_{\mathbf{x}_0} q(\mathbf{x}_0 \mid \mathbf{x}_0^{\text{input}}) \, q(\mathbf{x}_t \mid \mathbf{x}_0, \mathbf{x}_0^{\text{input}}) \, \mathrm{d}\mathbf{x}_0$$

$$= \nabla_{\mathbf{x}_t} \log \int_{\mathbf{x}_0} q(\mathbf{x}_0 \mid \mathbf{x}_0^{\text{input}}) \, q(\mathbf{x}_t \mid \mathbf{x}_0) \, \mathrm{d}\mathbf{x}_0$$

*(since $\mathbf{x}_{1:T}$ is independent of $\mathbf{x}_0^{\text{input}}$ conditioned on $\mathbf{x}_0$)*

$$= \nabla_{\mathbf{x}_t} \log \int_{\mathbf{x}_0} \underbrace{q(\mathbf{x}_0) \, r_{\text{sim}}(\mathbf{x}_0, \mathbf{x}_0^{\text{input}})}_{} \, q(\mathbf{x}_t \mid \mathbf{x}_0) \, \mathrm{d}\mathbf{x}_0$$

*(by the definition of $q(\mathbf{x}_0 \mid \mathbf{x}_0^{\text{input}}) \propto q(\mathbf{x}_0) \, r_{\text{sim}}(\mathbf{x}_0, \mathbf{x}_0^{\text{input}})$)*

$$= \nabla_{\mathbf{x}_t} \log \int_{\mathbf{x}_0} r_{\text{sim}}(\mathbf{x}_0, \mathbf{x}_0^{\text{input}}) \, q(\mathbf{x}_t, \mathbf{x}_0) \, \mathrm{d}\mathbf{x}_0$$

$$= \nabla_{\mathbf{x}_t} \log \int_{\mathbf{x}_0} r_{\text{sim}}(\mathbf{x}_0, \mathbf{x}_0^{\text{input}}) \, q(\mathbf{x}_t) \, q(\mathbf{x}_0 \mid \mathbf{x}_t) \, \mathrm{d}\mathbf{x}_0$$

$$= \nabla_{\mathbf{x}_t} \log q(\mathbf{x}_t) + \nabla_{\mathbf{x}_t} \log \int_{\mathbf{x}_0} q(\mathbf{x}_0 \mid \mathbf{x}_t) \, r_{\text{sim}}(\mathbf{x}_0, \mathbf{x}_0^{\text{input}}) \, \mathrm{d}\mathbf{x}_0$$

$$= \nabla_{\mathbf{x}_t} \log q(\mathbf{x}_t) + \nabla_{\mathbf{x}_t} \log \mathbb{E}_{q(\mathbf{x}_0 \mid \mathbf{x}_t)}\big[r_{\text{sim}}(\mathbf{x}_0, \mathbf{x}_0^{\text{input}})\big]$$

$$= \mathbf{s}(\mathbf{x}_t) + \underbrace{\nabla_{\mathbf{x}_t} \log \mathbb{E}_{q(\mathbf{x}_0 \mid \mathbf{x}_t)}\big[r_{\text{sim}}(\mathbf{x}_0, \mathbf{x}_0^{\text{input}})\big]}_{\text{guidance gradient}}$$

### APPENDIX II
### PROOF: APPROXIMATING EXPECTED DENOISED IMAGE $\mathbf{x}_0$ GIVEN $\mathbf{x}_t$ USING THE DIFFUSION MODEL

### APPENDIX III
### TRAINING ON THE RUGD DATASET

### A. Information about the RUGD dataset

The RUGD dataset (Fig. 5, [15]) is an off-road dataset of video sequences captured from a small, unmanned mobile robot traversing in unstructured environments. It contains over 7,453 labeled images from 17 scenes, annotated with pixel-level segmentation over 24 semantic classes. The annotated frames are spaced five frames apart.

We split the 24 semantic categories as 16 in-distribution labels: $\mathbb{C}_{\text{ID}}$ = {dirt, sand, grass, tree, pole, sky, asphalt, gravel, mulch, rock-bed, log, fence, bush, sign, rock, concrete}, and 8 OOD labels

corresponding to "obstacle" classes: $\mathbb{C}_{\text{OOD}}$ = {vehicle, container/generic-object, building, bicycle, person, bridge, picnic-table, water}

### B. Training a diffusion model on the RUGD Dataset

Our diffusion model is trained on samples from the RUGD train split that does not contain humans and artificial constructs (Fig. 6 *(left)*). The out-of-distribution and anomalous images are 'held out' for evaluation (Fig. 6 *(right)*). As can be seen from Fig. 7, our trained diffusion model successfully generates realistic images containing only in-distribution classes such as trees, grass, and the occasional footpath.

### C. Performance on RUGD Dataset

We evaluate our method on the RUGD [15] dataset and share qualitative results in Fig. 8. The diffusion model is trained on RUGD data without artificial constructs. At test-time, we present our method with OOD images containing anomalies from held-out classes. Fig. 9 shows the impact of each component of our *analysis* pipeline.

### D. Tuning the Guidance Strength

The strength of the guidance term in our diffusion model can be tuned to enforce a variable level of consistency between the image being generated $\mathbf{x}'$ and the target image $\mathbf{x}$. Fig. 10 *(left)* shows the impact of the guidance term on the image generated, for a sample anomalous image $\mathbf{x}$ from the RUGD dataset.

### E. Limitations

Our method's performance appears to degrade on samples with a large number of anomalies (Fig. 10 *(right)*), compared to input images with a fewer number of anomalies. This can be largely attributed to under-segmentation by SAM, which uses an (input-independent) hyperparameter that dictates the number of segments to output. Future work could aim to develop a criterion by which SAM can vary its segmentation resolution, without requiring human intervention.

### APPENDIX IV
### TOY EXPERIMENTS ON THE CLEVR DATASET

### A. Validating our improved diffusion guidance

We show that the SoftRect energy function improves our ability to remove anomalies via guided diffusion from the CLEVR dataset [79] in Section IV-B . In Section III-D, we also show what happens as we tune the guidance strength hyperparameter.

### B. Validating the SoftRect Guidance Function

To validate our choice of SoftRect over L2 guidance, we train two diffusion models — one for each guidance method — on examples from the CLEVR data [79] containing no reds, yellows, or browns (Fig. 11 *(left)*). We then present the trained model with anomalous images containing those held out colours, and compare their ability to remove anomalies without modifying non-anomalous parts of the image (Fig. 11 *(right)*). We find that SoftRect indeed outperforms L2 guidance when it comes to diffusion model-based anomaly removal.
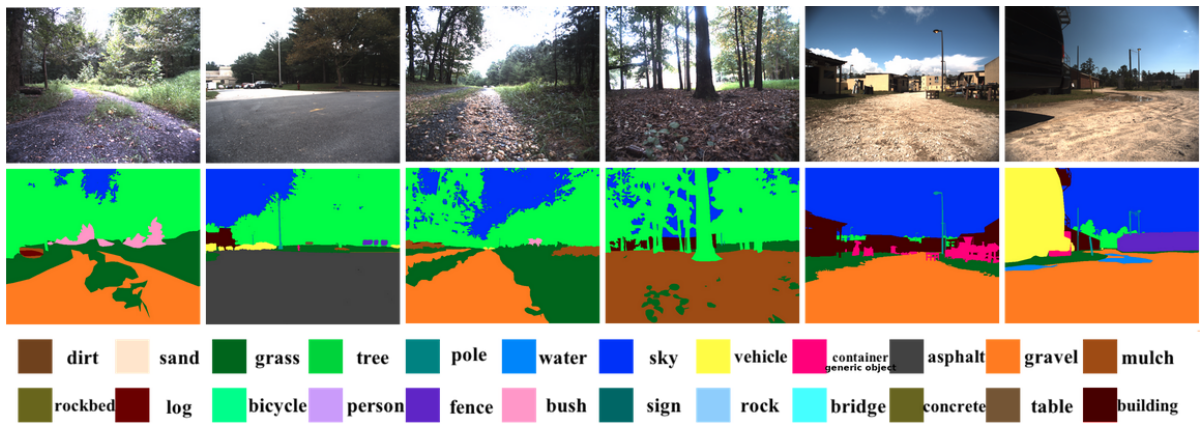
**Fig. 5:** Examples of video frames, annotations and semantic classes from the full RUGD dataset [15].



**Fig. 6: In-distribution and out-of-distribution images from the RUGD dataset.** *Left:* Examples of the in-distribution images on which our RUGD diffusion model was trained. In general, these images contain natural, off-road vegetation — a mixture of forest, meadow, mulch, and paths, without any humans or artificial constructions like buildings or vehicles. *Right:* Examples of held-out, out-of-distribution RUGD images the robot might encounter. These contains anomaly objects like buildings and vehicles. The diffusion model trained on the images on the left must remove anomalies from the images on the right.



**Fig. 7:** Samples generated from the trained diffusion model (without conditioning). The training data is shown in **??** *(left)*. The generated samples are photorealistic, and appear very similar to the training images.
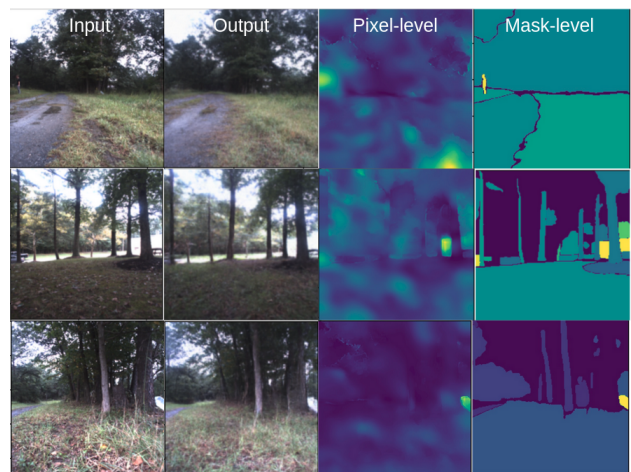


**Fig. 8:** Qualitative results on small, anomalous examples from the RUGD dataset. Our full `DiffUnc` pipeline does particularly well at detecting small and camouflaged/human-imperceptible anomalies.
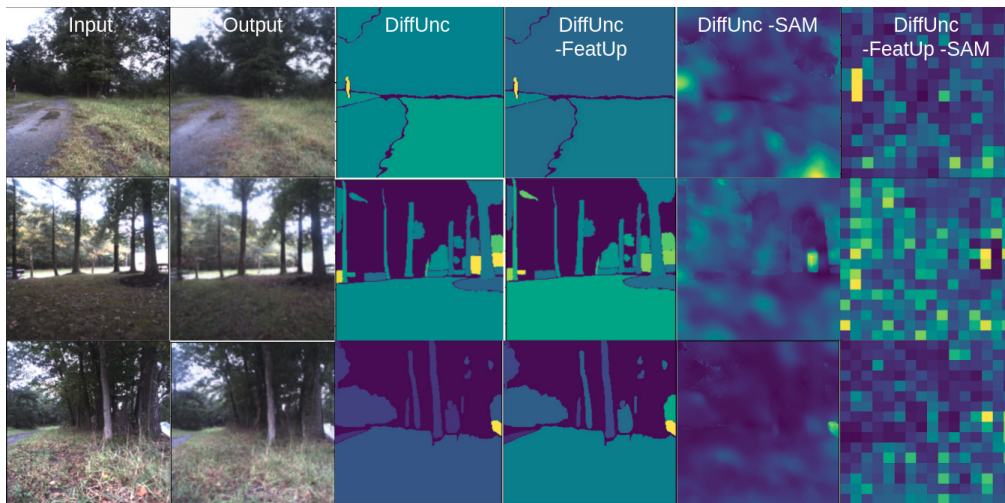
**Fig. 9:** Impact of SAM and FeatUp on `DiffUnc` performance. Removing FeatUp lowers the contrast between in- and out-of- distribution segments. Removing SAM particularly degrades performance on small anomalies.
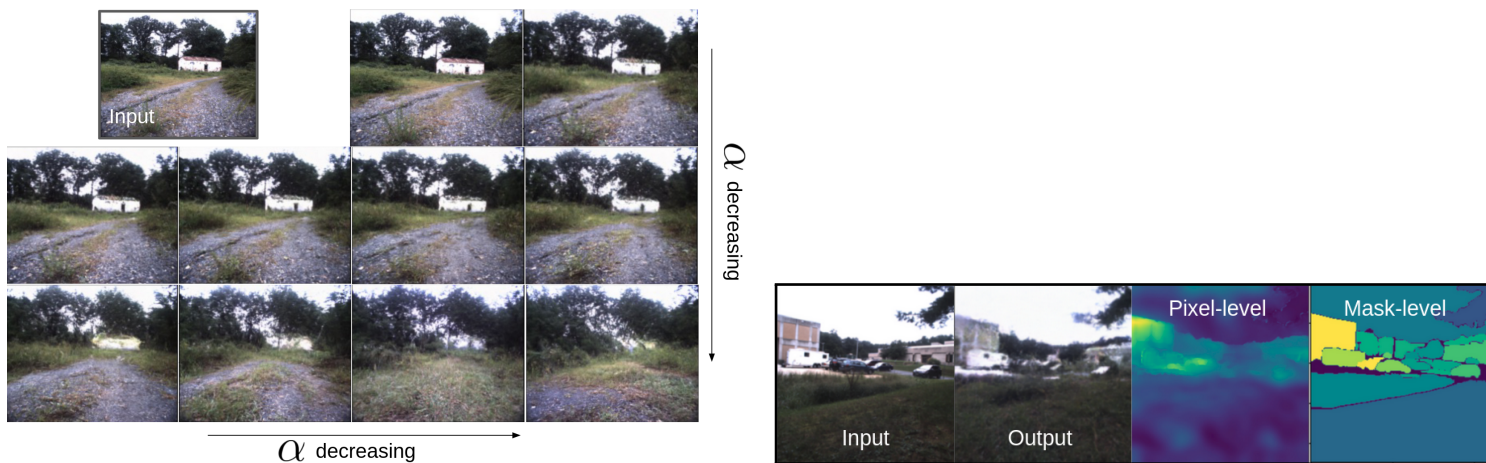


**Fig. 10:** *Left:* When there are many anomalies in a single image (in this case all vehicles and buildings are anomalies), the performance of our model appears to degrade. This is likely because SAM is under-segmenting an image with a high number of objects. The number of masks SAM generates is a hyperparameter in our pipeline and so could be adjusted, but the question of how to handle images with a varying number of objects is a non-trivial one. *Right:* The strength of our diffusion model's guidance term can be tuned by a hyperparameter $\alpha$. As $\alpha$ decreases, the guidance enforcing consistency between the target image and the diffusion process' output weakens, and an increasing number of changes are allowed.
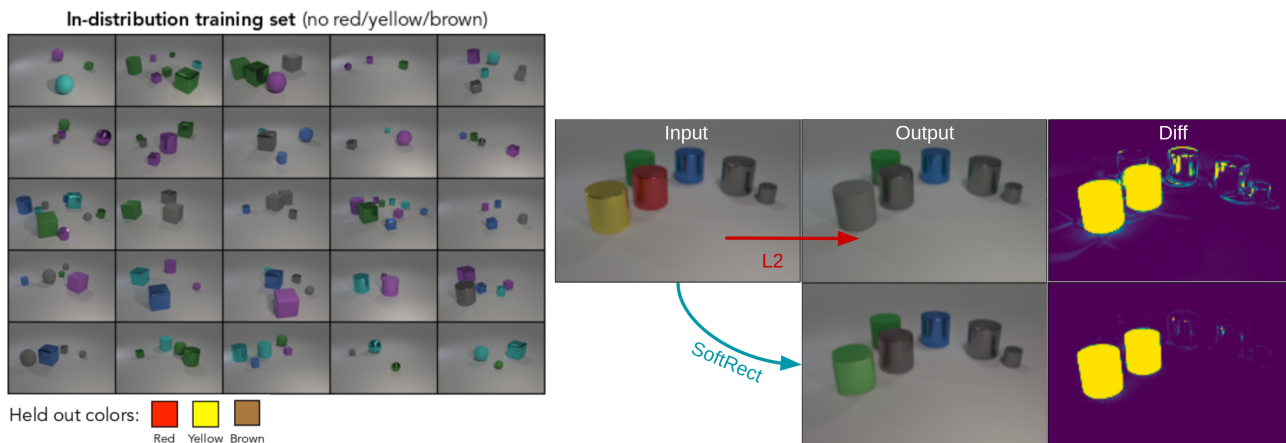


**Fig. 11:** *Left:* Examples of CLEVR images without reds, yellows, or browns, upon which our tabletop diffusion models were trained. *Right:* Comparing the performance of SoftRect- and L2- guided diffusion models shows that SoftRect successfully guides the diffusion process toward fewer unecessary changes on in-distribution objects. As a result, SoftRect-guided diffusion yields an uncertainty map with few false-positive pixels.