# LaRender: Training-Free Occlusion Control in Image Generation via Latent Rendering - Supplementary Materials

Xiaohang Zhan
Tencent
xiaohangzhan@outlook.com

Dingming Liu
Tencent
dingmingliu3722@gmail.com

## 1. The RealOCC Dataset

**Statistics**. As shown in Figure 1, the RealOCC dataset maintains a balanced distribution of images with 2 to 5 objects, capturing diverse occlusion scenarios. Images with more objects usually contain fewer occlusion pairs, which reflects real-world patterns. On average, each image includes 3.56 objects and 3.73 occlusion pairs, offering varied and realistic examples for evaluating occlusion relationships. This distribution ensures the dataset is both diverse and practical for testing occlusion handling methods.
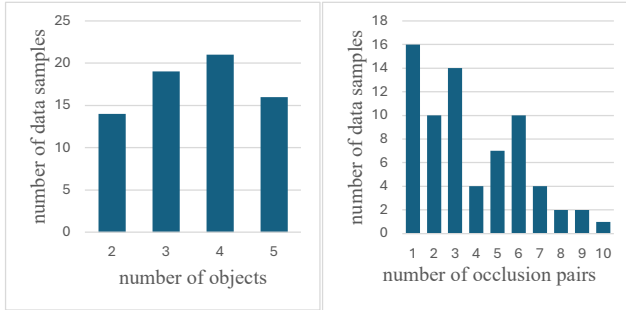


Figure 1. The statistics of the RealOcc dataset.

**Examples**. Figure 3 shows RealOCC examples with 2 to 4 objects, covering both simple and complex occlusions.

## 2. More Visualizations.

### 2.1. More results

Figure 6 compares LaRender with SDXL and FLUX. Despite style differences, both handle occlusion well. Failure cases are also shown. Figure 8 highlights a subject occluding various concepts as an interesting example.

### 2.2. Visualization of Ablation Study

Figure 6 illustrates the effects of different $\sigma_i(t)$ schedules and the use of cross-attention maps. A fixed opaque $\sigma_i(t)$ causes excessive opacity and object disappearance, making it only suitable in early steps. Fixing $\sigma_i(t)$ at a nor-
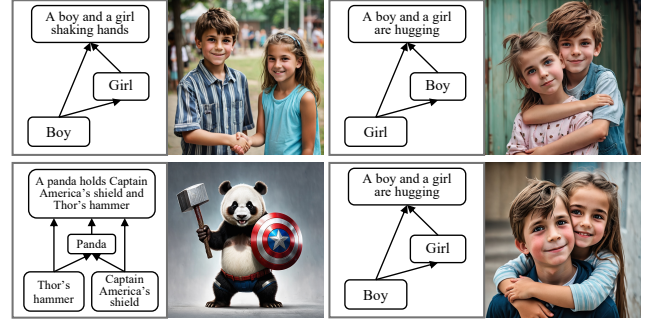


Figure 2. Examples of generating inter-object interactions via background prompts.

mal value throughout may lead to concept mixing and misaligned occlusions. In contrast, our inverse-proportional schedule generates clearer objects and occlusion relationships. Additionally, cross-attention maps help refine object contours better than bounding box masks alone.

### 2.3. Visualization of interactions.

Results of inter-object interactions is shown in Figure 2.

### 2.4. Visualization of 2-element effects.

To showcase LaRender's semantic opacity control, we visualize two elements in Figure 7, showing it can independently control multiple effects and generate high-quality results.

## 3. LLM instructions

We design the prompt templates in Figure 8 with instructions, in-context examples, and a test case from the user's input. The LLM follows the instruction to identify occlusion pairs and generate object layouts based on occlusion order.
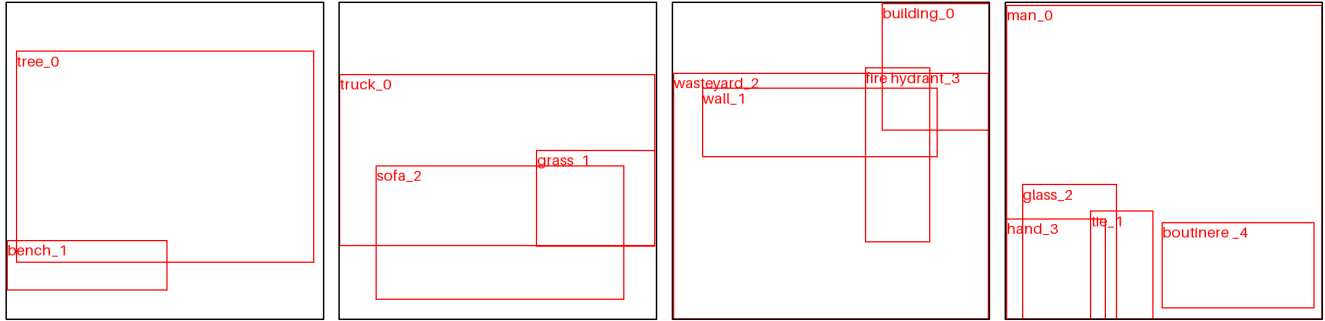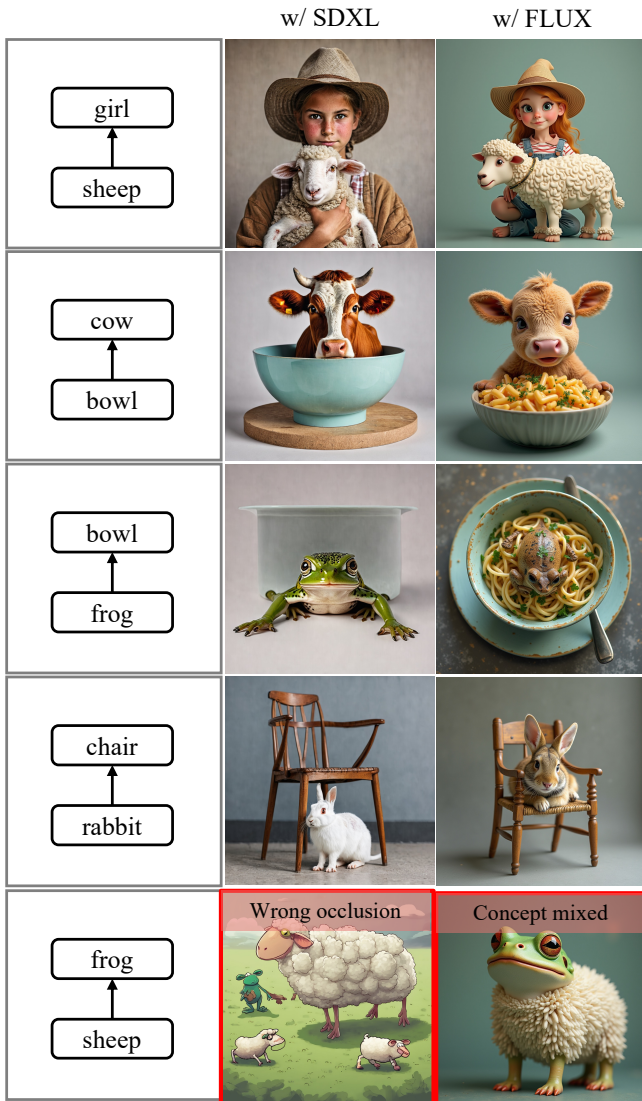
Figure 3. Examples with different number of objects in RealOcc dataset. The bounding boxes are inferred from the amodal masks in the COCOA dataset. The index after the name at the upper-left corner means the ordering from bottom to top. For example, the second example demonstrates grass occludes a truck and a sofa occludes both the track and the grass.
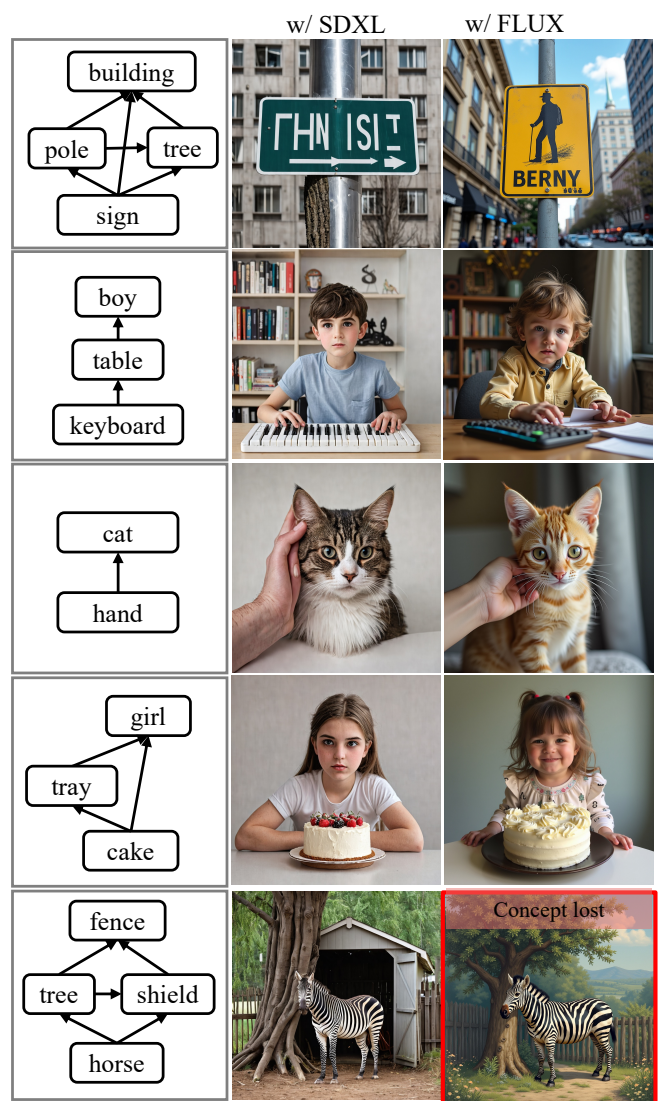


Figure 4. More results by our method, based on SDXL and FLUX. Failure cases are marked in red boxes.

*A boy in front of a [cat, turtle, pig, house, mouse]*



*A girl hidden by a [bag, sheep, key, chicken, bird]*



*A vase hidden by a [television, airplane, book, candle, clock]*



Figure 5. This figure show results of a subject occluding with other objects.



Figure 6. The visual results of ablation study, including the impact of different density schedules (left) and the presence of cross-attention maps in computing transmittance maps (right).

Figure 7. Grid figure of dual element changes. Please zoom in to see the changes.

**Output**:
prompt = "On a street, a man and a motorcycle are each obstructing one side of a car"
background = "a street"
objects = ["a car", "a motorcycle", "a man"]
locations = [[0.3, 0.3, 0.7, 0.7], [0.4, 0.1, 0.75, 0.4], [0.2, 0.5, 0.9, 0.75]]
occlusion_pairs = [("a motorcycle", "a car"), ("a man ", "a car")]

DeepSeek

LaRender

Figure 8. The instruction for DeepSeek to parse the occlusion graph.