# 3DIS: DEPTH-DRIVEN DECOUPLED IMAGE SYNTHESIS FOR UNIVERSAL MULTI-INSTANCE GENERATION

Dewei Zhou[1*], Ji Xie[1*], Zongxin Yang[2*], Yi Yang[1]

[1] ReLER, CCAI, Zhejiang University  [2] Harvard University  * Equal Contribution

**Code**

## Introduction



TL;DR: 3DIS allows users to perform MIG using various foundational models (including SD1.5, SD2, SDXL, FLUX), where they can specify the position and attributes of each instance in one image.
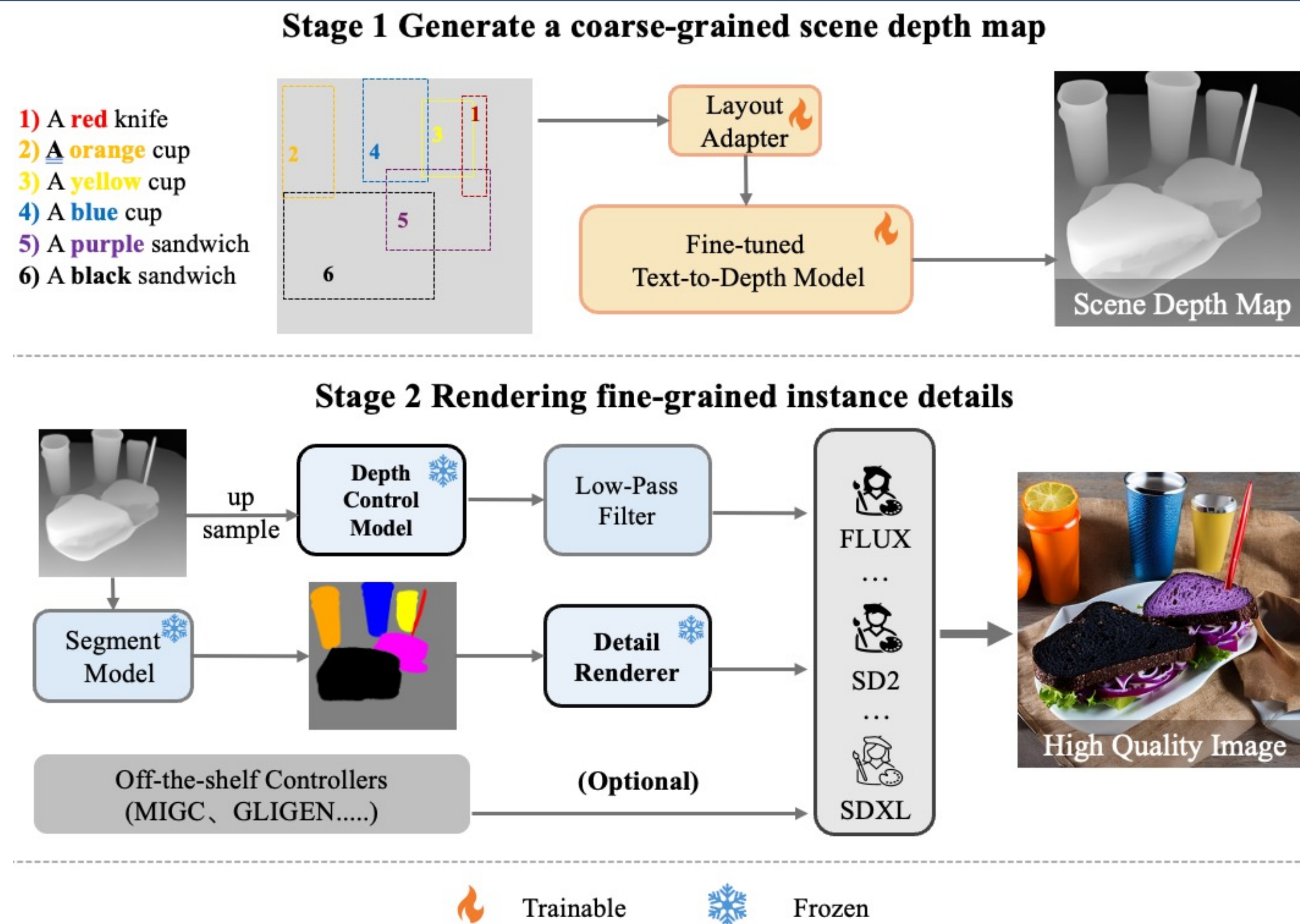
### Background:
- Multi-Instance Generation (MIG) allows users to define the locations and attribute of multiple instances in the generated image.
- The mainstream MIG methods involve training an adapter directly on the generative model to control both the position and attributes of each instance. Whenever a more powerful base model emerges, the adapter needs to be retrained, which is resource-intensive.

### Motivations:
- Decoupling MIG into scene building and detail rendering. **We only train one Layout-to-Depth model** to control instance positions, and then use a training-free approach to **render details using multiple models**.
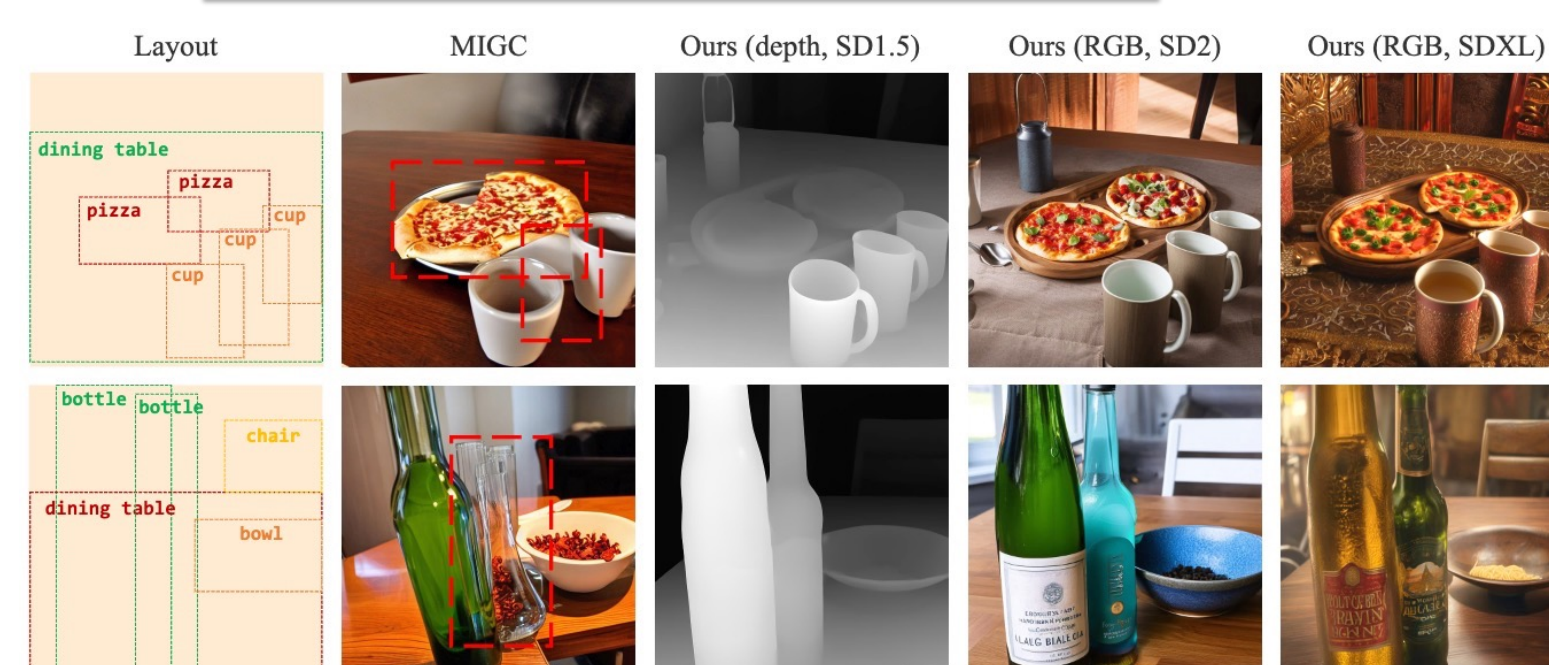
## Methodology: 3DIS

### Stage 1 Generate a coarse-grained scene depth map



1) A **red** knife
2) **A orange** cup
3) A **yellow** cup
4) A **blue** cup
5) A **purple** sandwich
6) A **black** sandwich

### Stage 2 Rendering fine-grained instance details



🔥 Trainable    ❄ Frozen

- This method divides the multi-instance generation process into two stages.
- First, a Layout-to-depth network is trained to generate a scene depth map.
- Secondly, we use Depth-ControlNet to precisely position each instance and a Detail Renderer to accurately render the attributes of each instance without any training.

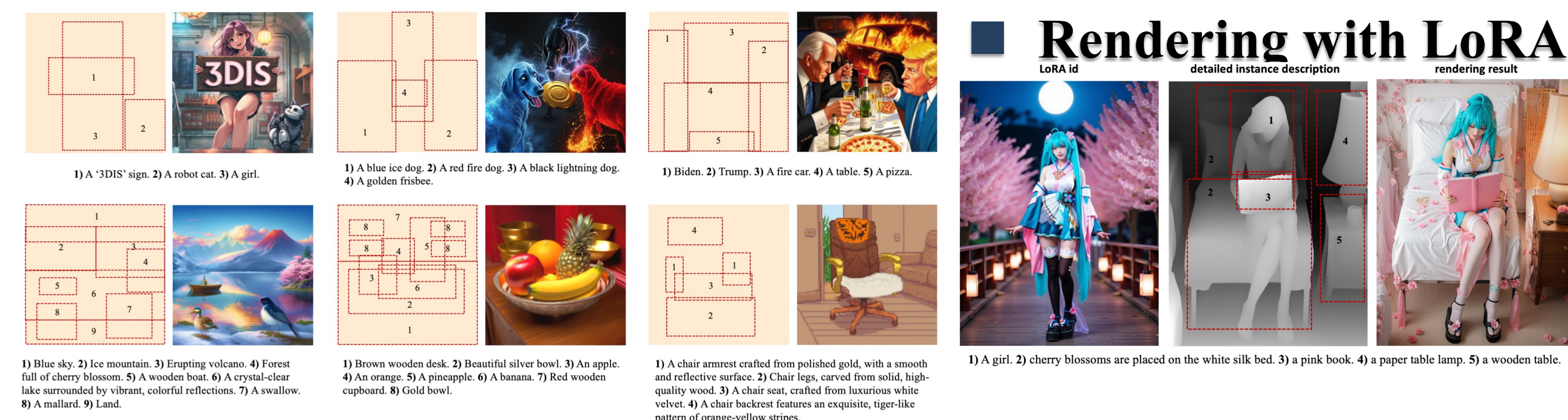## Experiments & Application

### Results on COCO-POS



| Method | Layout Accuracy | | | Instance Accuracy | | | Image Quality | |
|---|---|---|---|---|---|---|---|---|
| | $AP\uparrow$ | $AP_{75}\uparrow$ | $AP_{50}\uparrow$ | $SR_{inst}\uparrow$ | MIoU | $CLIP\uparrow$ | $SR_{img}\uparrow$ | $FID\downarrow$ |
| BoxDiff [ICCV23] | 3.15 | 2.12 | 10.92 | 22.74 | 27.28 | 18.82 | 0.53 | 25.15 |
| MultiDiff [ICML23] | 6.37 | 4.24 | 13.22 | 28.75 | 34.17 | 20.12 | 0.80 | 33.20 |
| GLIGEN [CVPR23] | 38.49 | 40.75 | 63.79 | 83.31 | 70.14 | 19.61 | 40.13 | 26.80 |
| MIGC [CVPR24] | 45.03 | 46.15 | 80.09 | 83.37 | 71.92 | 20.07 | 43.25 | 24.52 |
| 3DIS | 56.83 | 62.40 | 82.29 | 84.71 | 73.32 | 20.84 | 46.50 | 23.24 |
| vs. prev. SoTA | +11.8 | +16.3 | +2.2 | +1.3 | +1.4 | +0.8 | +3.3 | +1.3 |

## Results on COCO-MIG



### MIG using various models（e.g., FLUX and SDXL）



### Rendering with LoRA



### Rendering Real Scene Depth Maps



### Control Depth Order



### Complex Attribute Rendering