

PlanGen: Towards Unified Layout Planning and Image Generation in Auto-Regressive Vision Language Models

Supplementary Material

A. Inference Details

We used classifier-free guidance during the inference process of image generation, and the guidance scale is 5. We only apply negative layout guidance when performing layout-guided image manipulation including object deletion. To speed up the inference, we use kv-cache. It takes about 17 seconds to generate 8 images in one batch on a single A100 GPU.

B. Experimental Details

In-context Prompt for Baselines in Layout Planning. Qwen-2.5-7b-instruct and Llama-3.1-8b-instruct themselves do not naturally support the layout planning task. We leverage LLMs’ in-context learning capabilities to generate layouts from global captions following Layout-GPT, and the in-context prompts are as follows:

You are an intelligent bounding box generator. I will provide you with a caption for a photo, image, or painting. Your task is to generate the bounding boxes for the objects mentioned in the caption, along with a background prompt describing the scene. The images are of size 512 x 512. The top-left corner has coordinate [0, 0]. The bottom-right corner has coordinate [512, 512]. The bounding boxes should not overlap or go beyond the image boundaries. Each bounding box should be in the format of (object name, [top-left x coordinate, top-left y coordinate, box width, box height]) and should not include more than one object. Do not put objects that are already provided in the bounding boxes into the background prompt. Do not include non-existing or excluded objects in the background prompt. Use "A realistic scene" as the background prompt if no background is given in the prompt. If needed, you can make reasonable guesses. Please refer to the example

below for the desired format.

```
input: A realistic image of landscape scene depicting a green car parking on the left of a blue truck, with a red air balloon and a bird in the sky
you need output: [('a green car', [21, 281, 211, 159]), ('a blue truck', [269, 283, 209, 160]), ('a red air balloon', [66, 8, 145, 135]), ('a bird', [296, 42, 143, 100])]
```

```
input: A realistic top-down view of a wooden table with two apples on it
you need output: [('a wooden table', [20, 148, 472, 216]), ('an apple', [150, 226, 100, 100]), ('an apple', [280, 226, 100, 100])]
```

```
input: An oil painting of a pink dolphin jumping on the left of a steam boat on the sea
you need output: [('a steam boat', [232, 225, 257, 149]), ('a jumping pink dolphin', [21, 249, 189, 123])]
```

The input for you to process is: {}

Baselines Details on Image Layout Understanding. For Grounding-DINO, we perform grounding detection by giving the image and the global caption of the image. For two other LLM-based baselines, i.e. CogVLM-grounding and Qwen-VL-Chat, we take prompts suitable for the corresponding models to help the model to output accurate detection results. Specifically, for CogVLM-grounding, we give the image and ask “*Can you provide a description of the image and include the coordinates [[x0,y0,x1,y1]] for each mentioned object?*” following its formal demo. For Qwen-VL-Chat, we apply two rounds of questions. First, we give the image and ask the model “*What objects are in the image?*”. Then, after the model answers this question, we ask the model “*Box out the positions of these objects in the figure*”. We find that for Qwen-VL-Chat, the effect of such two-round Q&A will be much better than a single question.

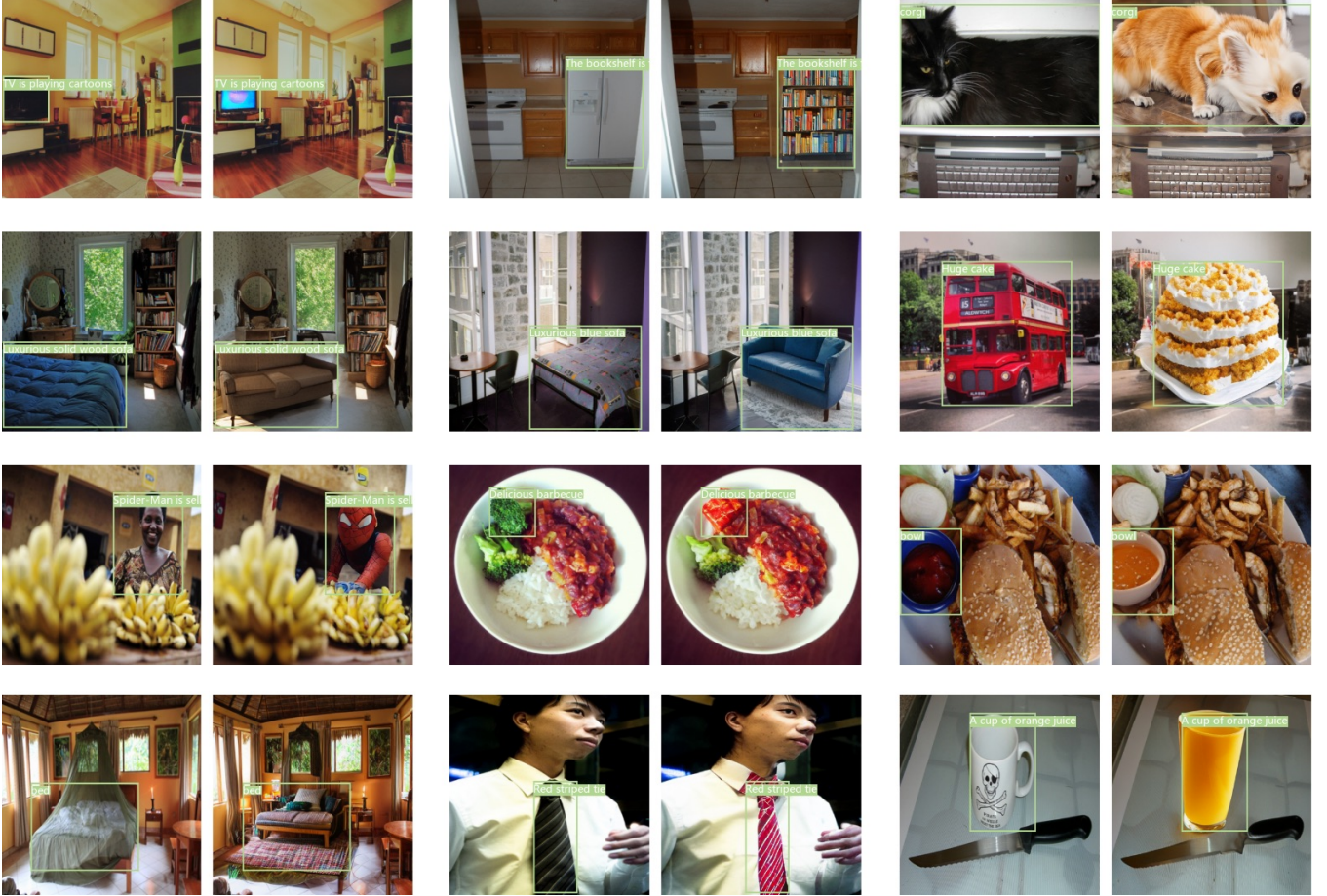


Figure 1. More examples for Layout-guided Image Manipulation. The contents to be edited are drawn in the form of bounding boxes on the original images and the edited images for easy comparisons.

Evaluation of Image Layout Understanding. Similar to HiCo, we calculate the maximum IoU between the boxes predicted by different models and the ground truth boxes. If the maximum IoU is higher than the threshold 0.5, we calculate the clip text similarity between corresponding local descriptions. If the CLIP score is higher than 0.2, we mark it as a correct prediction. We use AR, AP, AP50 and AP75 to evaluate the performance of image layout understanding.

C. Additional Results

We show more examples of layout-guided image manipulation in Figure 1, more examples of layout-image joint generation in Figure 3, more examples of layout-to-image in Figure 4, and more results of image layout understanding in Figure 5. These rich examples show PlanGen’s excellent performance on multiple related tasks.

Comparisons with Janus-Pro Baseline. We supplement comparisons with the Janus-Pro (1B) baseline across multiple tasks. PlanGen significantly outperforms Janus-Pro in

Method	Spatial \uparrow	Color. \uparrow	Textual. \uparrow	Shape \uparrow	FID \downarrow
Janus-Pro	72.91	59.67	62.85	61.03	18.52
PlanGen	92.21 (+19.30)	82.69 (+23.02)	86.53 (+23.68)	85.36 (+24.33)	13.91 (-4.61)

Table 1. Layout-to-image comparison with Janus-Pro baseline.

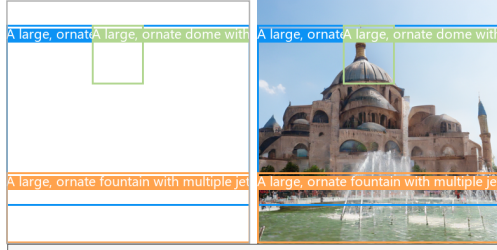
layout-to-image as shown in Table 1, with superior region-wise scores and a lower FID. We observe that even with in-context examples, Janus-Pro (1B & 7B) struggles to generate reasonable layouts from captions, potentially due to its limited training. Similarly, Janus-Pro faces challenges in understanding image layouts. All these experiments show that PlanGen holds significant advantages over the Janus-Pro baseline.

Failure cases. We also show several failure cases on the task of layout-image joint generation in Figure 2. We observed that PlanGen may experience distortion when generating human bodies, as shown in Figure 2, which is a common challenge faced by some previous autoregressive image generation models. When multiple identical objects



Figure 2. Failure cases.

are generated, additional ones may appear, which is caused by incomplete object annotations in the training data. In the 2nd row, the geometry of generated mask straps is slightly inappropriate. When generating dogs in smaller size, the quality declines also occur. More efficient image modeling or wider training should alleviate these problems.



This is a photo showcasing a grand mosque, with its iconic dome and minaret standing out against the blue sky. The mosque is located in a spacious square, surrounded by a fountain with water splashing, and the square is filled with people. The surrounding trees and the blue sky add a touch of natural beauty to this religious building.



This is a photo showcasing fashionable clothing in a store. The focus is on a mannequin dressed in a floral-patterned jacket and a blue scarf, with a plaid-patterned hat on its head. The background is blurred, but other mannequins and clothing can be seen. The entire scene is illuminated by soft indoor lighting, creating a warm shopping atmosphere.



This is a photograph showcasing a city port landscape. In the foreground, a red boat is docked at the pier, with a few people walking on the pier. The midground features a busy harbor with numerous sailboats and yachts moored, and a modern building with a blue roof. The background is a cityscape, with buildings of various styles and sizes, as well as distant hills. The sky is a light blue, with clouds scattered across it.



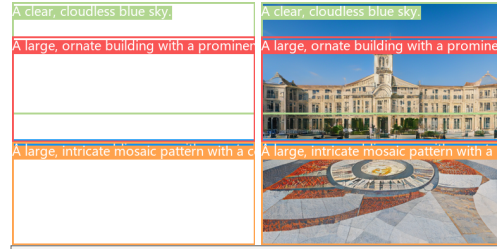
This is a photo showcasing a modern shopping center, capturing the bustling scene of shoppers and the surrounding environment. In the foreground of the photo, there are several tall palm trees, adding a touch of tropical flair to the scene. The center of the photo is a pedestrian walkway, with people walking on it, some carrying shopping bags, others chatting. The shops on both sides of the walkway are brightly lit, displaying various fashionable clothes and accessories. The signboards of the shops are clearly visible, such as 'ALBERTO' and 'MELISSA ODINGA'. The entire scene is illuminated by the soft sunlight, creating a warm shopping atmosphere.



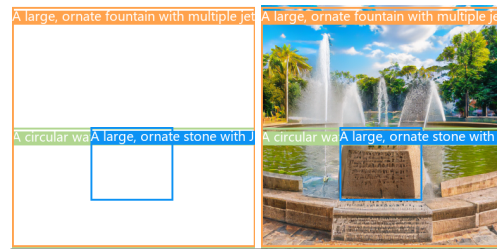
This is a photo showcasing a natural landscape, with a majestic mountain range in the center of the frame, its surface covered with rocks and sparse vegetation. The mountain range is illuminated by sunlight, with the shadows of the peaks clearly visible. In the foreground of the image, there is a small building with a traditional style, its color contrasting with the surrounding environment. The entire scene is captured under clear weather, with a blue sky and a few white clouds, adding a sense of tranquility and vastness to the image.



This is a photo showcasing three stone tablets with Chinese inscriptions. The tablets are arranged in a row, with the one in the middle being the largest, and the two on either side being smaller. Each tablet has a unique design, with the middle tablet having a red background and black Chinese characters, while the two on either side have a yellow background and black Chinese characters. Below the middle tablet is a stone slab with red Chinese characters, and the three tablets are placed on a set of stone steps. The background features some trees and a building, but they are not the focus.



This is a photo showcasing a grand building with a classic European style. The building is located in the center of the picture, with a prominent clock tower at the top, and the facade is decorated with intricate carvings and windows. In front of the building is a spacious square, with a large mosaic pattern in the middle, composed of red, black, and white tiles, and the word 'EUROPA' is written in the middle. The square is paved with white marble, and there are some people and vehicles around. The background is a clear blue sky, with a few white clouds leisurely drifting by.



This is a photograph showcasing a fountain in a park. The fountain is located in the center of the image, with water shooting up from the center, forming a series of elegant arcs. The fountain is surrounded by a circular water pool, with a commemorative stone in the middle. The stone is engraved with Japanese text, which may be the name of the fountain or a commemorative text. The background of the fountain is a lush green park, with a few people walking in the distance. The sky is clear, with a few white clouds scattered across the blue sky.

Figure 3. More examples for Layout-Image Joint Generation. Global captions for layout-image joint generation are attached below the images.



Figure 4. More examples for Layout-to-Image Generation.

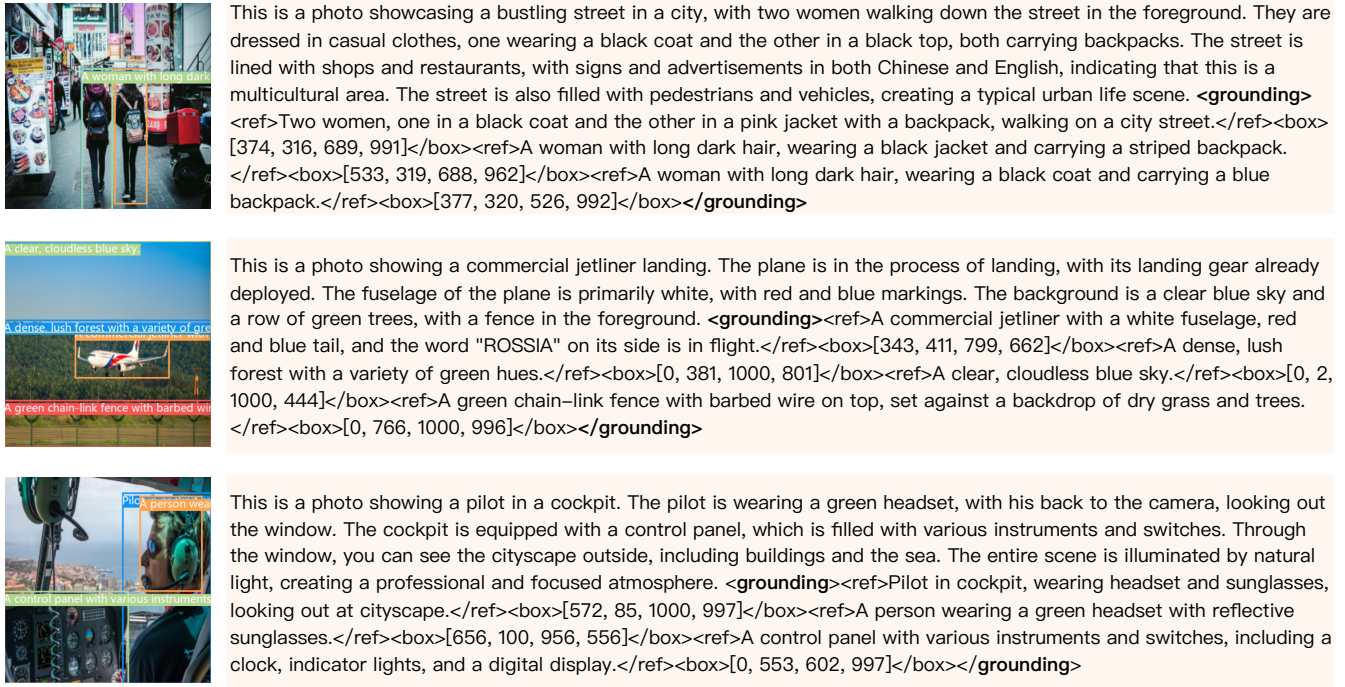


Figure 5. More examples for Image Layout Understanding.