Figure 1: **To Reviewer vV5n.** Detailed modifications of Figure 1.



Figure 2: **To Reviewer 16K4.** The visualization of reasoning failure cases. In the bottom right corner of the image, we re-select the qualitative results from our real-world demonstration. Additionally, we replace the red dot and virtual end-effector with a physical Franka Panda Robot.

Table 1: **To Reviewer 16K4.** The Scalability Exploration of RoboMamba. For parameter scalability, since the official 7B model is not public, we compare the reasoning abilities of Mamba 1.4B and 2.7B on the same training dataset in the first two rows. The results indicate that increasing the model parameters effectively enhances reasoning ability under efficient LLM settings. For training data scalability, we introduce more training data during the co-training stage. Although there is a slight improvement on some MLLM benchmarks, the average reasoning did not improve.

| Parameters | ShareGPT4V-SFT 665K | LVIS-Instruct-4V 220K | OKVQA | GQA | POPE |
|---|---|---|---|---|---|
| 1.4B | | | 28.5 | 40.8 | 66.8 |
| 2.7B | | | 62.3 | 63.8 | 86.9 |
| 2.7B | + | | 56.8 | 57.5 | 85.9 |
| 2.7B | | + | 62.4 | 64.4 | 86.0 |