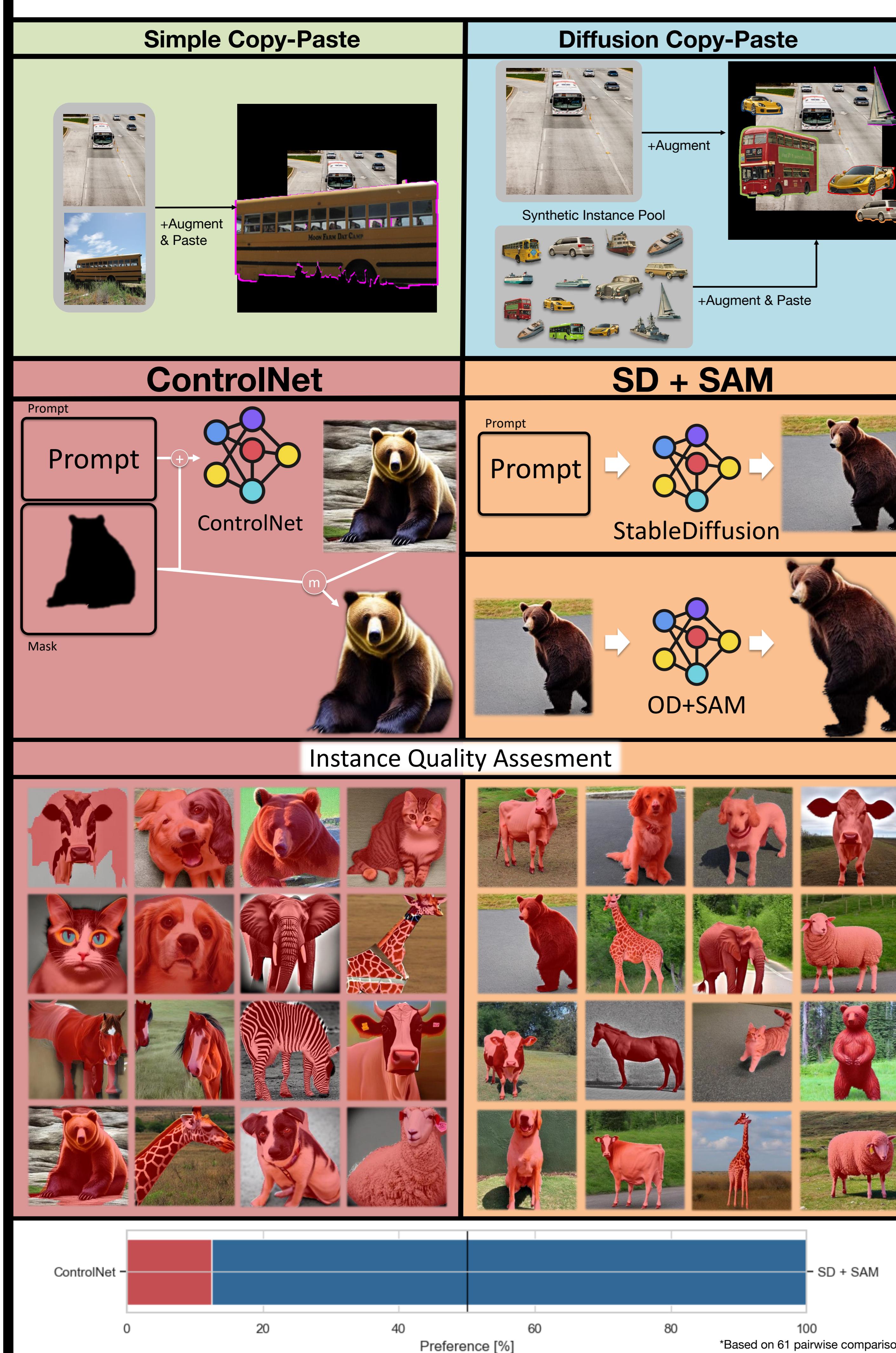
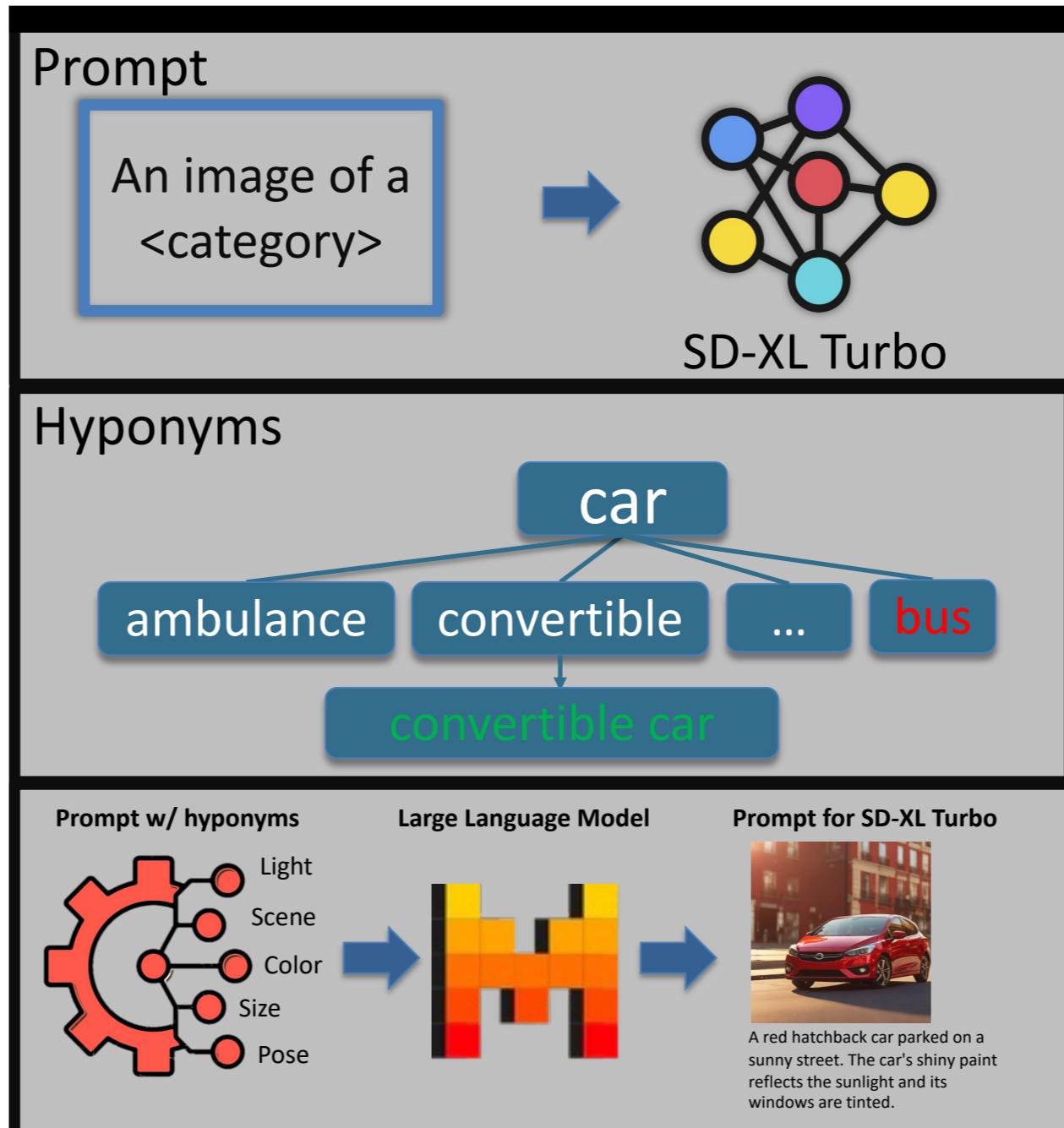


## Introduction

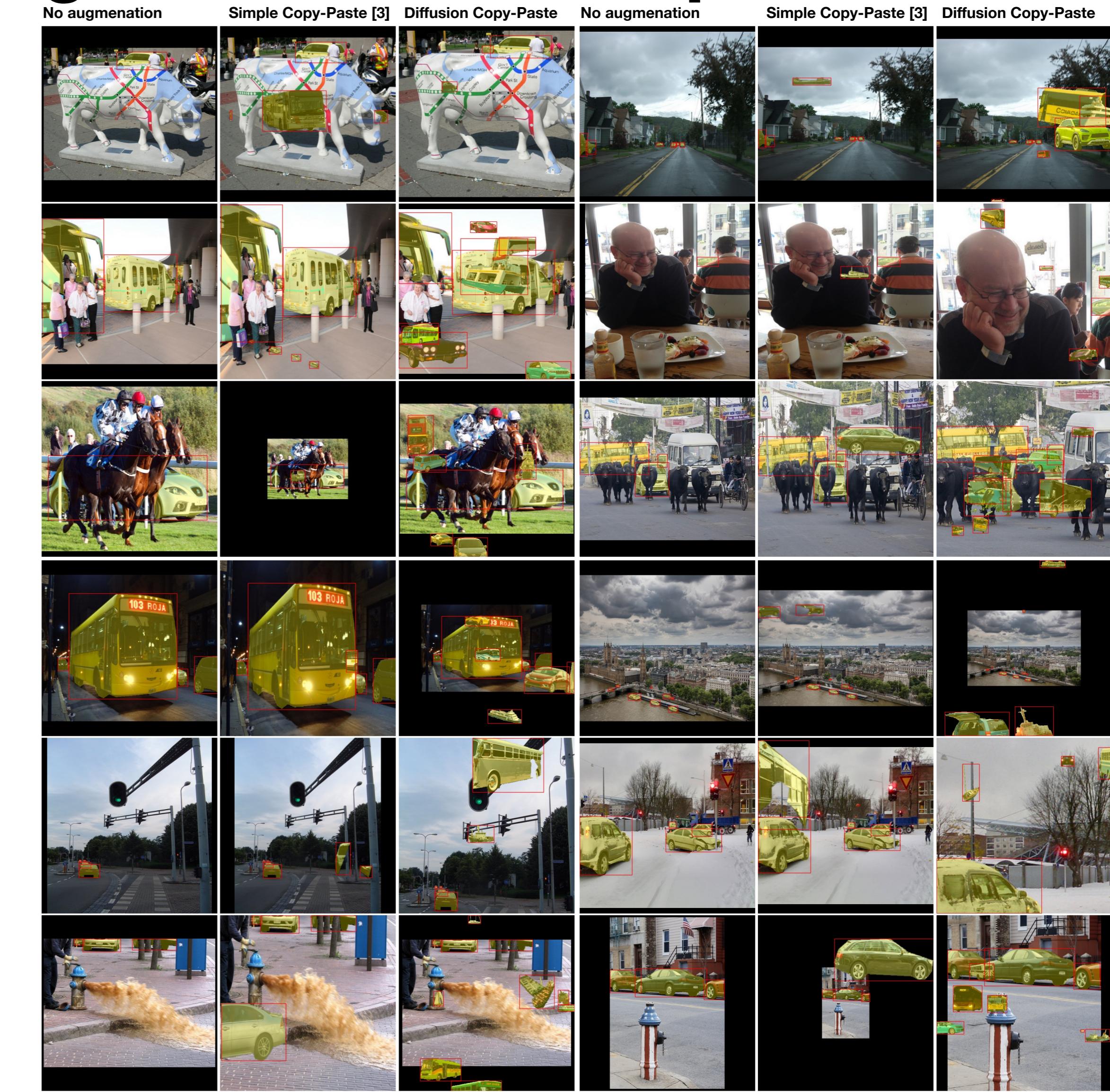


## Generation Ablations

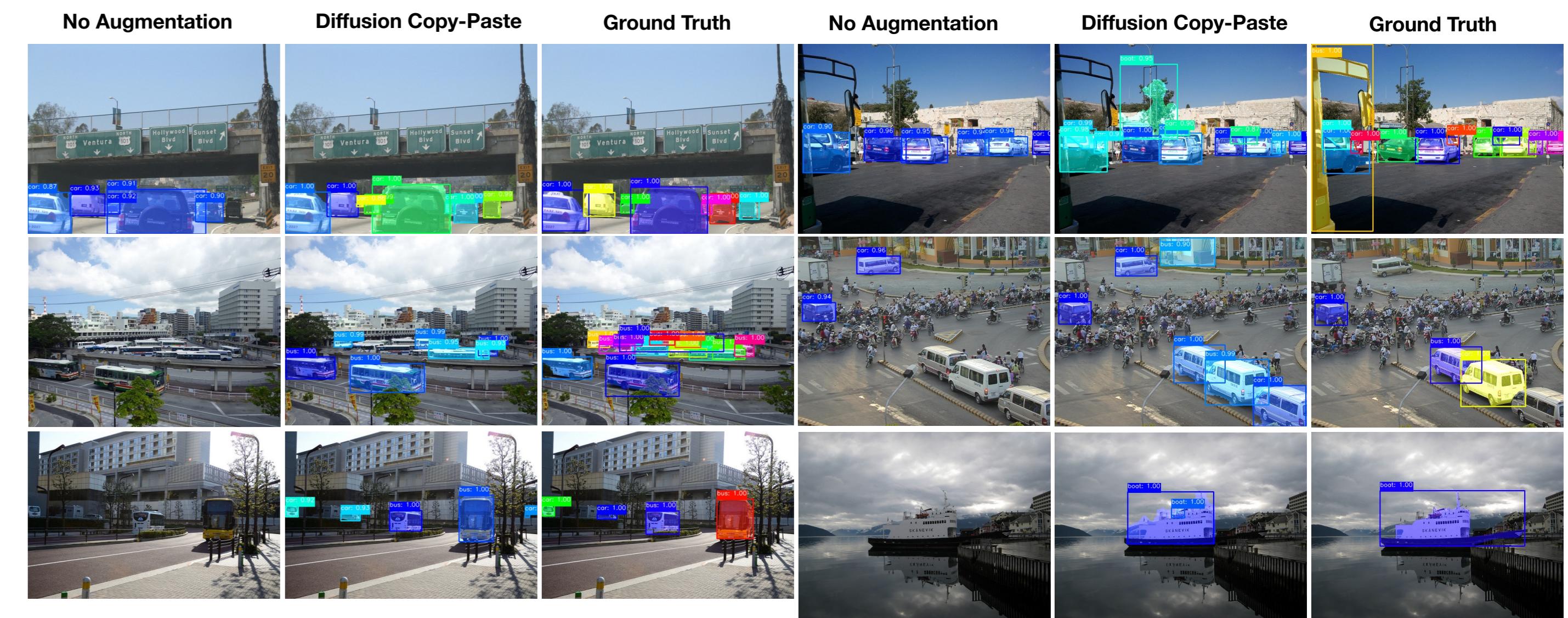
- Simple text template
- Hyponyms and synonyms from WordNet
- Manually remove **misleading** and modify **ambiguous**
- We compare this method with using the class name
- We compare 1 vs 4 steps for SD-XL Turbo
- Experiment with LLM for added diversity
- Calculate CMMD and DINOv2 feature standard deviation



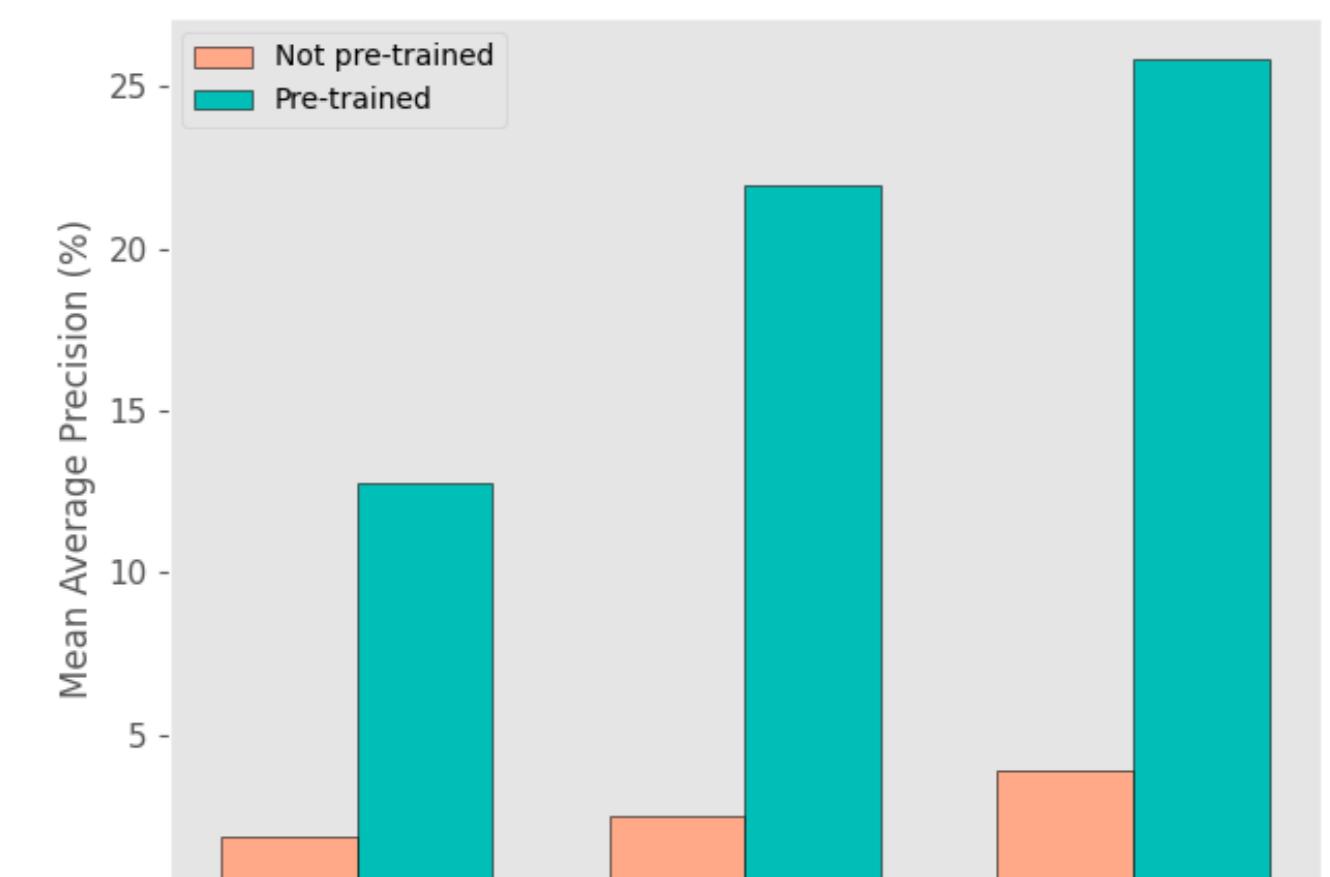
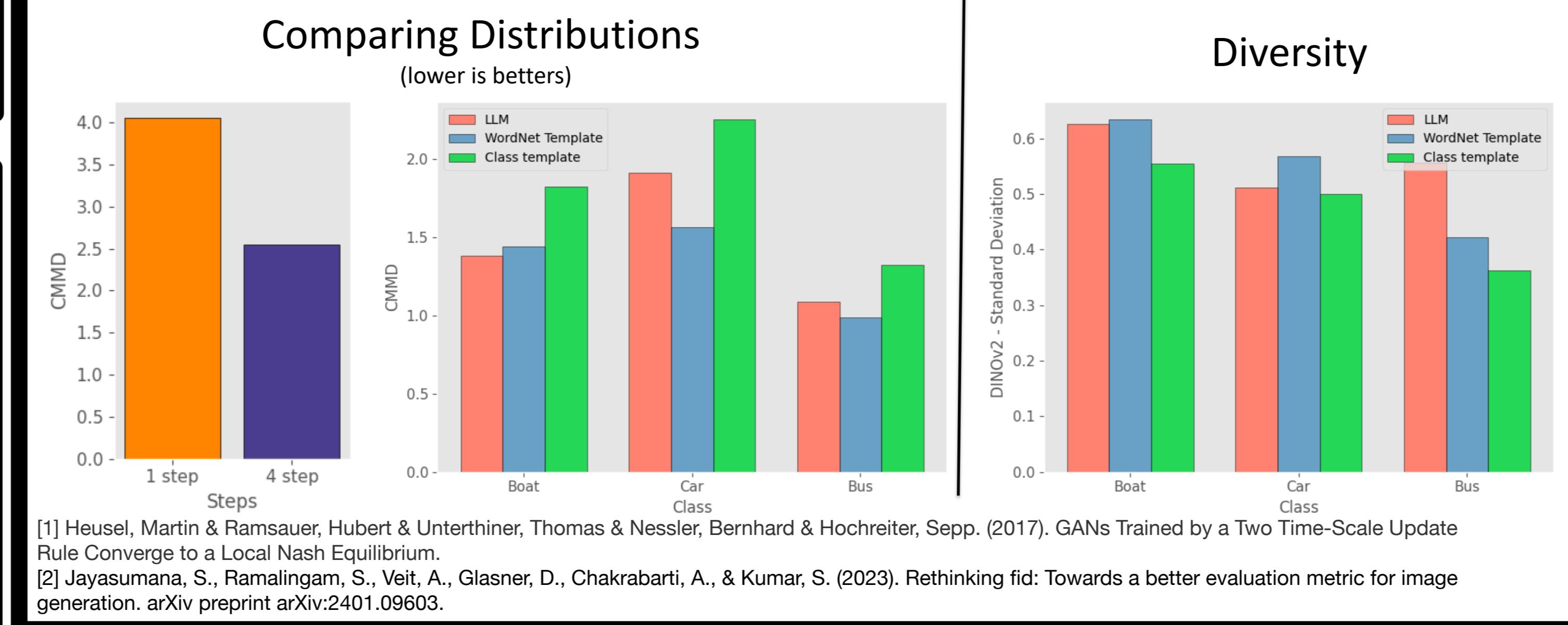
## Augmentation Examples



## Experiments with Mask-RCNN [4]

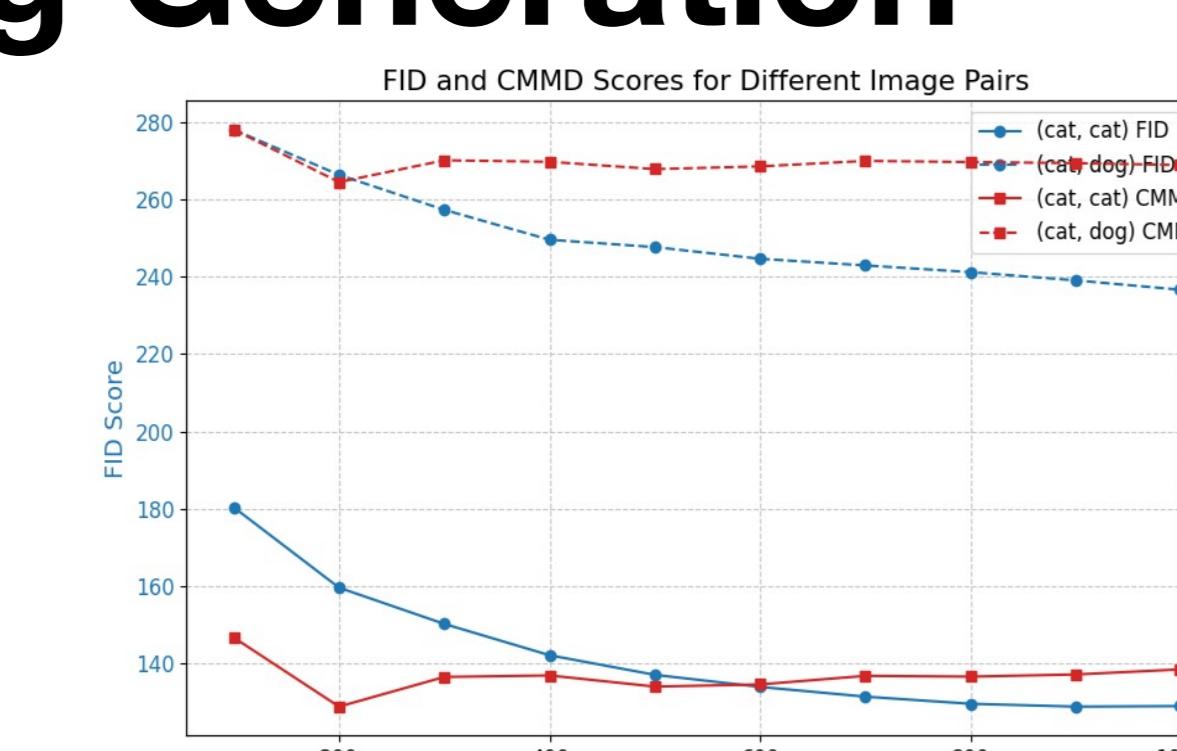


- Mask-RCNN [4]
- 420 training images with three classes
- Trained for 100 epochs, AdamW optimizer, 5 epoch linear lr warmup to 0.0003, followed by cosine annealing.
- /w & /wo pretrained backbone on ImageNet1k
- Box mAP evaluated on COCO validation set



## Metric: Assessing Generation

- We compare the generations using FID [1] and CMMD [2].
- FID and **CMMD** is calculated by masking out everything but the instance.
- FID and CMMD are calculated for multiple numbers (100, ..., 1000) of images to check the stability of the metrics with limited data.



[1] Heusel, Martin & Ramsauer, Hubert & Unterthiner, Thomas & Nessler, Bernhard & Hochreiter, Sepp. (2017). GANs Trained by a Two Time-Scale Update Rule Converge to a Local Nash Equilibrium.  
[2] Jayasumana, S., Ramalingam, S., Veit, A., Glasner, D., Chakrabarti, A., & Kumar, S. (2023). Rethinking fid: Towards a better evaluation metric for image generation. arXiv preprint arXiv:2401.09603.