

ULIP-2: Towards Scalable Multimodal Pre-training for 3D Understanding

Supplementary Material

A. Appendix

A.1. Ablation on 3D Input

In order to fairly compare to OpenShape on Objaverse-LVIS benchmark, which utilizes 10k colored point clouds as the 3D input, we adopt the same 3D input preprocessing as in OpenShape. We conduct this ablation study to assess how color information influences zero-shot classification results on Objaverse-LVIS. To do this, we evaluated the ULIP-2 pre-trained Point-BERT model, pre-trained on Objaverse + ShapeNet. Our findings from Table 10 show that ULIP-2 maintains strong performance on Objaverse-LVIS zero-shot classification tasks, even without using color information.

3D Encoder Input	Objaverse-LVIS	
	top-1	top-5
8k xyz	48.9	77.1
10k xyzrgb	50.6	79.1

Table 10. Point-BERT w/ ULIP-2 zero-shot 3D classification on Objaverse-LVIS, pre-trained on Objaverse and ShapeNet jointly with OpenCLIP ViT-G encoders.

A.2. Different Kinds of 3D Backbones

To verify ULIP-2’s improvement is agnostic to 3D backbones, we conduct experiments on the PointNeXt backbone. Results in Table 11 show that, with another kind of intrinsically different 3D backbone, ULIP-2 can still improve the performance significantly. Given that Point-BERT is a scale-up-friendly transformer-based architecture and has better zero-shot classification results, we mainly conduct experiments on Point-BERT for all our experiments.

Model	Pre-train method	ModelNet40	
		top-1	top-5
PointNeXt [36]	ULIP [52]	56.2	77.0
	ULIP-2	72.8	95.7
Point-BERT [55]	ULIP [52]	60.4	84.0
	ULIP-2	75.2	95.0

Table 11. Zero-shot 3D classification on ModelNet40 with different 3D backbones, pre-trained on ShapeNet with SLIP ViT-B encoders.

A.3. More Details of the Generated Triplets

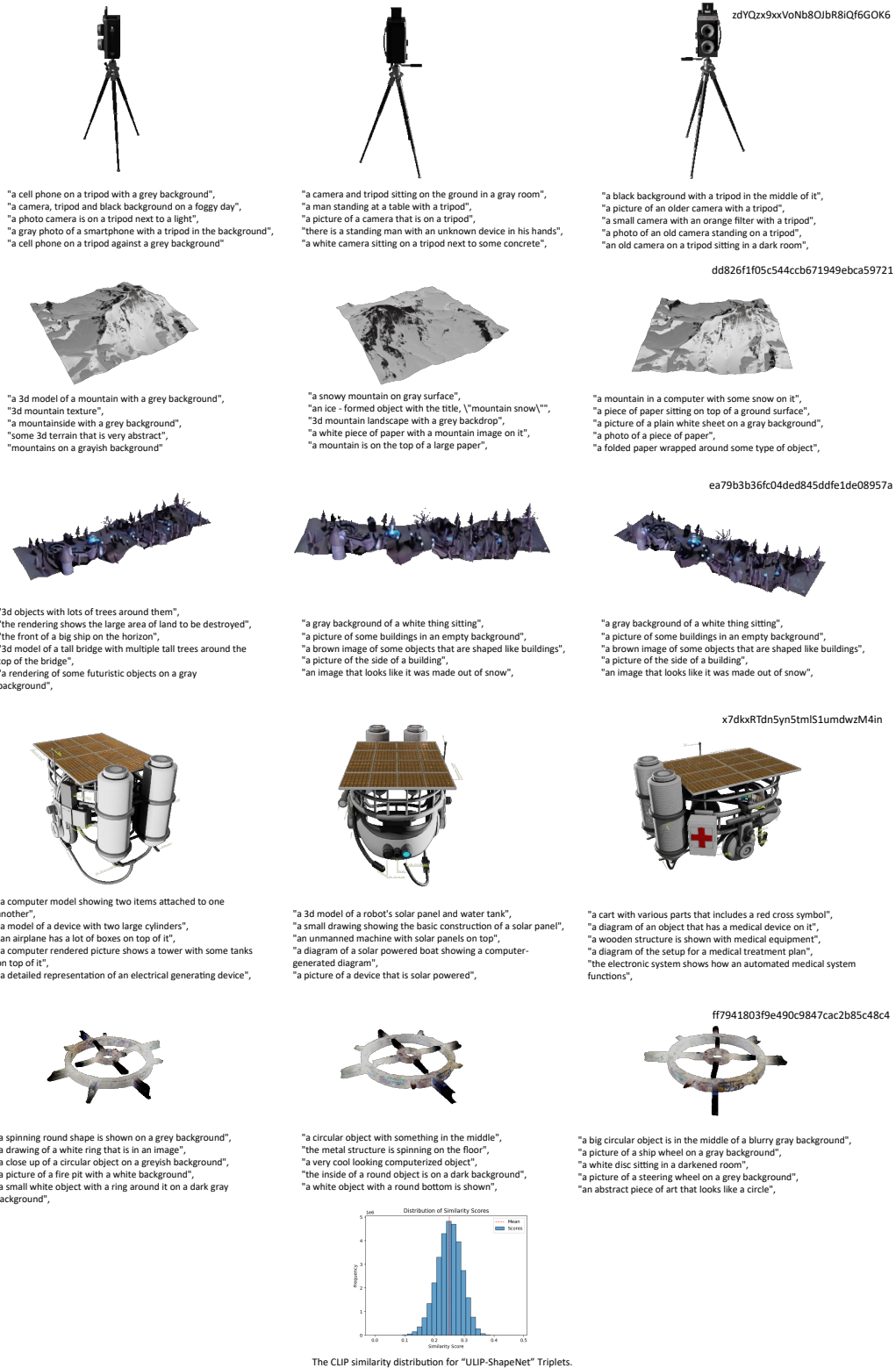


Figure 4. More qualitative samples from ULIP-Objaverse Triplets.