

000 *Where Am I and What Will I See:*
 001 **AN AUTO-REGRESSIVE MODEL FOR SPATIAL LOCAL-**
 002 **IZATION AND VIEW PREDICTION**
 003
 004
 005

006 **Anonymous authors**

007 Paper under double-blind review
 008
 009
 010
 011

012 **A ADDITIONAL EXPERIMENTAL DETAILS**
 013

014 **A.1 PRE-PROCESSING**
 015

016 **Camera System Unification.** In our experiments, due to the utilization of diverse camera systems
 017 across different datasets, we initially unify the camera coordinate systems of all datasets into a
 018 common reference frame, the RUB coordinate system.

019 **Standardization of Camera Distribution.** Since the camera positions obtained by COLMAP
 020 (Schönberger & Frahm, 2016) lack scale information and the datasets used in training stage en-
 021 compass both synthetic object datasets and real-world scene datasets, significant variations in cam-
 022 era position scales exist across different datasets and scenes. To mitigate the scale discrepancies
 023 among different scenes, we first employ a fixed intrinsic camera matrix. Although this approach
 024 may introduce certain perspective issues, it does not impact the final results of camera positions and
 025 orientations.

026 Subsequently, we computed the variance of camera positions across different scenes and scaled the
 027 camera positions of scenes within the same dataset by a common factor β such that the variance
 028 of camera positions in the dataset is standardized to 1. We present the β corresponding to different
 029 datasets in the table 1. This standardization process compacts the relative camera positions within
 030 the training set, facilitating the modeling of the overall camera distribution of the dataset.

031 Finally, as our cameras are randomly sampled, cases where two distant cameras capture images with
 032 little to no overlap are prevalent. To address this, during the training of the camera tokenizer and
 033 auto-regressive model, we filter out such instances by setting a distance threshold $\delta = 5$ to restrict
 034 excessive distances between two cameras.
 035

036 **Table 1: Dataset scaling factor.**
 037

Dataset	Scaling Factor β
Objaverse	1.0
CO3Dv2	0.1
MVImgNet	0.5
RealEstate10K	10.0

044
 045
 046 **A.2 DATASET SELECTION**
 047

048 **Objaverse** (Deitke et al., 2023). We incorporated a multi-view subset of Objaverse rendered by
 049 Zero-1-to-3 (Liu et al., 2023) as part of our training set. However, this dataset significantly out-
 050 weighed the remaining real-world image dataset. To enable the model to better comprehend the
 051 spatial distribution of real-world scenes, we opted to utilize only a subset of Objaverse. Upon in-
 052 vestigation, we observed substantial variance in the rendering quality of different objects within
 053 Zero-1-to-3. Thus, we selected a subset of higher quality renderings from Tang et al. (2024) to
 enhance the efficiency and efficacy of the training process.

054 **CO3Dv2** (Reizenstein et al., 2021). We selected a subset of categories from CO3Dv2 as our train-
 055 ing categories (seen categories), while another subset was designated as unseen categories. The
 056 categorization of classes was guided by Ray Diffusion (Zhang et al., 2024), as shown in the table 2.

057 **MVImgNet** (Yu et al., 2023) and **RealEstate10K** (Zhou et al., 2018). Both datasets consist of
 058 real-world scene data, encompassing object-centric scenes and authentic indoor environments. We
 059 integrated the complete data from these two datasets into our training process.

061 Table 2: Partition of CO3Dv2 (Reizenstein et al., 2021) .

062

Seen Categories						Unseen Categories	
apple	backpack	banana	baseballbat	baseballglove	bench	ball	book
bicycle	bottle	bowl	broccoli	cake	car	couch	frisbee
carrot	cellphone	chair	cup	donut	hairdryer	hotdog	kite
handbag	hydrant	keyboard	laptop	microwave	motorcycle	remote	sandwich
mouse	orange	parkingmeter	pizza	plant	stopsign	skateboard	suitcase
teddybear	toaster	toilet	toybus	toyplane	toytrain		
toytruck	tv	umbrella	vase	wineglass			

063

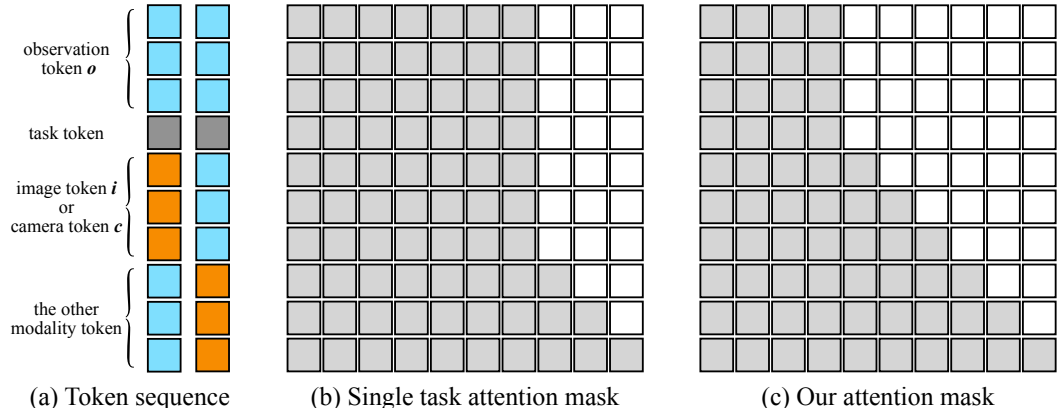
070

071 A.3 AUTO-REGRESSIVE MODEL TRAINING

072 **Task Tokens.** In this study, we focus solely on two tasks: novel view synthesis and camera pose
 073 estimation. Therefore, our task tokens are limited to these two, placed at the end of the codebook
 074 for easier future expansion with additional tasks.

075 **Training Process.** At the initial phase, we evenly allocate training resources among four conditional
 076 distributions. Subsequently, we observe that reducing the occurrences of $p(i|o)$ and $p(c|o)$ during
 077 training gradually enhances the numerical outcomes. However, this approach also entails a trade-off
 078 by compromising a portion of training stability.

079 **Attention Mechanisms.** In our training stage, we tokenize the initial observed image o along with
 080 the image i and camera c corresponding to random sampled viewpoints, resulting in three token
 081 sequences of equal length. These sequences are concatenated to form a fixed-length token sequence
 082 denoted as s . The concatenation order of the three tokens plays a crucial role in how the model
 083 interprets and integrates the information. Two concatenation orders, were employed, as illustrated
 084 in Figure 1, denoted as (t_o, t_i, t_c) and (t_o, t_c, t_i) .



095

096 **Figure 1: The comparison of Attention Mechanisms.** (a) We standardize the token sequences s
 097 of two tasks to the same length. (b) In the conventional alternating training scheme, target modality
 098 tokens can attend to all preceding conditional tokens. (c) We propose to model the joint distribution
 099 such that, in (a), the token sequences for the two tasks have visibility only to the currently observa-
 100 tion tokens o .

101

102 The attention mask in figure 1 illustrates the distinction between our training approach and alternat-
 103 ing training of two targets.

108 **Training Resources.** Our camera tokenizer was trained for approximately 2 days on 4 NVIDIA
109 A100 GPUs, while our autoregressive model was trained for about 3 weeks on 16 NVIDIA A100
110 GPUs.

112 B VISUALIZATION RESULTS

114 B.1 NOVEL VIEW SYNTHESIS

116 We selected a number of representative images, including those from the training dataset, virtu-
117 ally synthesized images, real-world images, and stylistic images, as initial observations. Due to the
118 uncertainty regarding the scale of the scenes, we first employed GST sampling to determine reason-
119 able camera positions. These positions were then used as conditions in conjunction with the initial
120 observations to generate new perspective images, as illustrated in the figure 2.

122 B.2 RELATIVE CAMERA POSE ESTIMATION

124 We selected several highly challenging examples to test the spatial localization capabilities of GST.
125 As illustrated in the figure 3, the selected image pairs include real-world images, images of the same
126 subject taken under different shooting conditions, and images of the same object depicted under
127 various artistic styles. GST demonstrated outstanding performance across all these examples.

129 C LIMITATIONS AND FUTURE WORKS

131 The scarcity of multi-view datasets with precise camera annotations poses a significant barrier to
132 scaling up GST. In the current work, we only explored the most fundamental scenario involving
133 a single observation image and one novel perspective. Consequently, when sampling multiple im-
134 ages and camera positions simultaneously, issues of consistency may arise, although this problem
135 decreases as the number of training viewpoints increases. In the future, we will aim to collect more
136 multi-view data with precise poses and will explore extending GST to accommodate an arbitrary
137 number of views as conditions, thereby broadening its applicability. Additionally, we trained on
138 datasets without scale, and the potential for extension to scenes with real-world scale remains to be
139 investigated.

141 REFERENCES

- 142 Matt Deitke, Dustin Schwenk, Jordi Salvador, Luca Weihs, Oscar Michel, Eli VanderBilt, Ludwig
143 Schmidt, Kiana Ehsani, Aniruddha Kembhavi, and Ali Farhadi. Objaverse: A universe of anno-
144 tated 3d objects. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern
145 Recognition*, pp. 13142–13153, 2023.
- 146 Ruoshi Liu, Rundi Wu, Basile Van Hoorick, Pavel Tokmakov, Sergey Zakharov, and Carl Vondrick.
147 Zero-1-to-3: Zero-shot one image to 3d object. In *Proceedings of the IEEE/CVF international
148 conference on computer vision*, pp. 9298–9309, 2023.
- 149 Jeremy Reizenstein, Roman Shapovalov, Philipp Henzler, Luca Sbordone, Patrick Labatut, and
150 David Novotny. Common objects in 3d: Large-scale learning and evaluation of real-life 3d cat-
151 egory reconstruction. In *Proceedings of the IEEE/CVF international conference on computer
152 vision*, pp. 10901–10911, 2021.
- 153 Johannes Lutz Schönberger and Jan-Michael Frahm. Structure-from-motion revisited. In *Confer-
154 ence on Computer Vision and Pattern Recognition (CVPR)*, 2016.
- 155 Jiaxiang Tang, Zhaoxi Chen, Xiaokang Chen, Tengfei Wang, Gang Zeng, and Ziwei Liu. Lgm:
156 Large multi-view gaussian model for high-resolution 3d content creation. *arXiv preprint
157 arXiv:2402.05054*, 2024.
- 158 Xianggang Yu, Mutian Xu, Yidan Zhang, Haolin Liu, Chongjie Ye, Yushuang Wu, Zizheng Yan,
159 Chenming Zhu, Zhangyang Xiong, Tianyou Liang, et al. Mvimngnet: A large-scale dataset of
160
161

162 multi-view images. In *Proceedings of the IEEE/CVF conference on computer vision and pattern*
163 *recognition*, pp. 9150–9161, 2023.

164

165 Jason Y Zhang, Amy Lin, Moneish Kumar, Tzu-Hsuan Yang, Deva Ramanan, and Shubham Tul-
166 siani. Cameras as rays: Pose estimation via ray diffusion. *arXiv preprint arXiv:2402.14817*,
167 2024.

168 Tinghui Zhou, Richard Tucker, John Flynn, Graham Fyffe, and Noah Snavely. Stereo magnification:
169 Learning view synthesis using multiplane images. *arXiv preprint arXiv:1805.09817*, 2018.

170

171

172

173

174

175

176

177

178

179

180

181

182

183

184

185

186

187

188

189

190

191

192

193

194

195

196

197

198

199

200

201

202

203

204

205

206

207

208

209

210

211

212

213

214

215

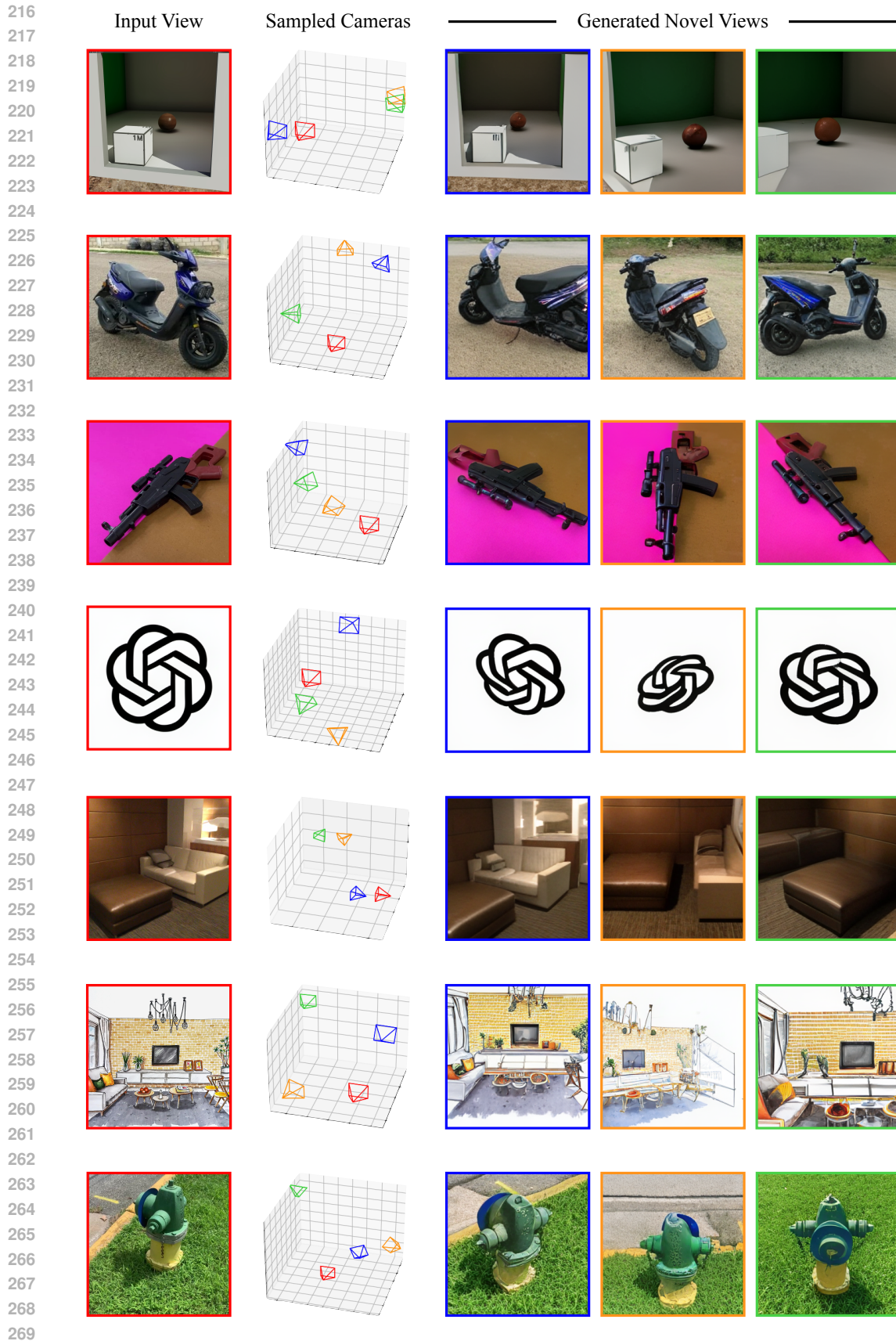


Figure 2: Visualization results of novel view synthesis.

270
271
272
273
274
275
276
277
278
279
280
281
282
283
284
285
286
287
288
289
290
291
292
293
294
295
296
297
298
299
300
301
302
303
304
305
306
307
308
309
310
311
312
313
314
315
316
317
318
319
320
321
322
323

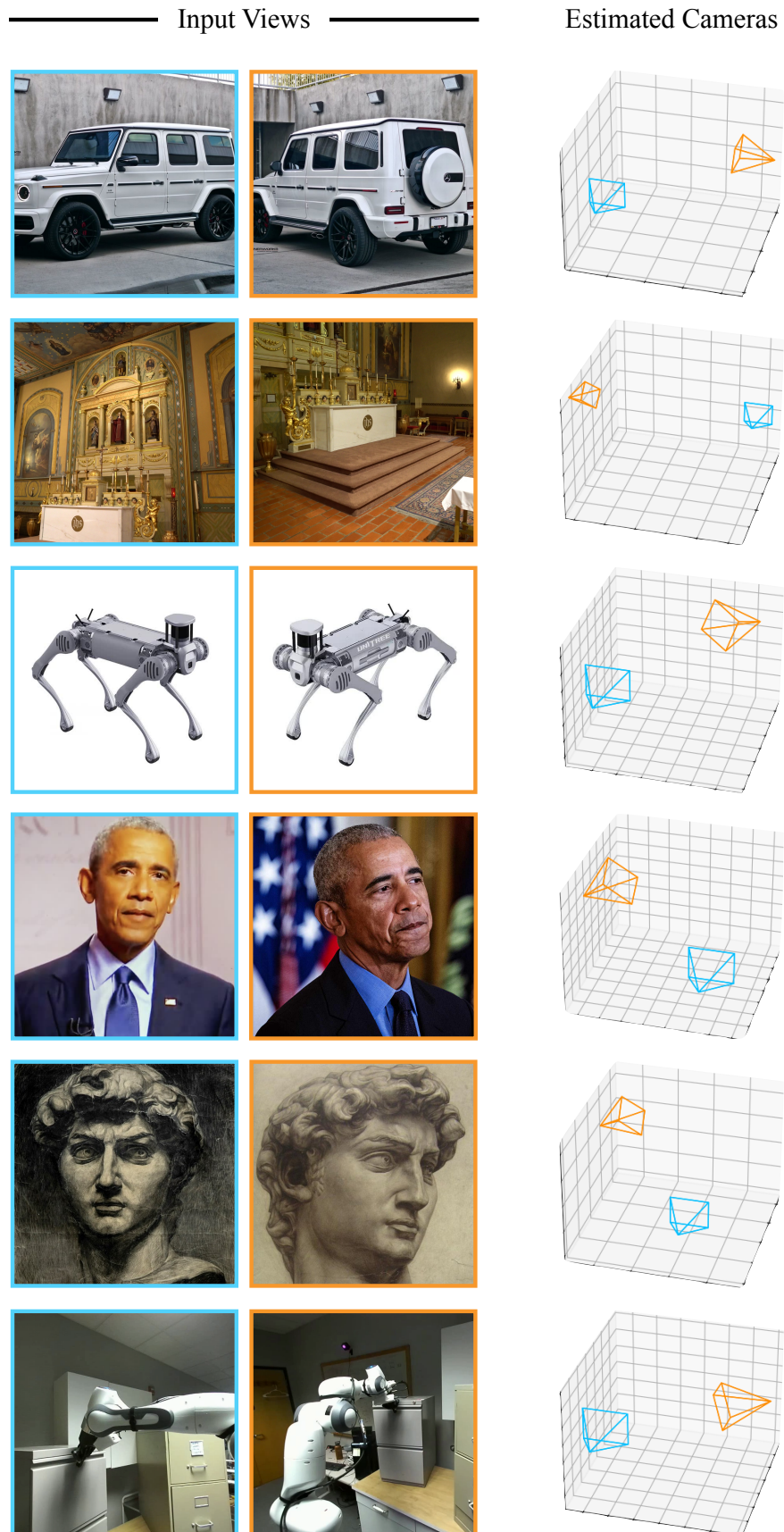


Figure 3: Visualization results of relative camera pose estimation.