

001 **Appendix**

002 This appendix provides additional details on implementa-
003 tion, evaluation metrics, and qualitative results that could
004 not be included in the main document due to space con-
005 straints. The document is organized as follows:

- 006 • **Section A** presents extended experiments in which
007 Comp3D- \mathcal{G} is trained on multiple diverse datasets to im-
008 prove its scalability and generalizability.
- 009 • **Section B** provides additional implementation details re-
010 garding the architecture of the Gaussian Decoder and the
011 training details for Any-Feature 3D Lifting.
- 012 • **Section C** details the experimental settings, including
013 dataset descriptions, baseline configurations, and evalua-
014 tion protocols for efficient novel view synthesis, 3D scene
015 understanding, and multi-view feature encoding.
- 016 • **Section D** presents additional ablation studies comparing
017 compact reconstruction strategies and analyzing compo-
018 nents of Comp3D- \mathcal{F} , followed by discussions on the au-
019 toencoder design, FPS analysis, and limitations.
- 020 • **Section E** provides additional quantitative experimental
021 results, including latent decoding or two-view novel view
022 synthesis.
- 023 • **Section F** presents additional qualitative results, includ-
024 ing visualizations of attention maps, 3D scene under-
025 standing, feature PCA, and renderings from 24-view in-
026 puts.

027 **A. Multi Dataset Training**

028 We extend Comp3D- \mathcal{G} to improve its scalability and gen-
029 eralizability by incorporating diverse training datasets, in-
030 creasing the number of input views, and scaling up the
031 number of Gaussians. However, since each dataset ex-
032 hibits different camera parameters and scene scales, naively
033 combining them makes training C3G unstable and difficult
034 to converge. To address this, we first describe the archi-
035 tectural modifications made to accommodate these varia-
036 tions (§ A.1), followed by implementation details covering
037 dataset configurations and training procedures (§ A.2). Fi-
038 nally, we present experimental results demonstrating the ef-
039 fectiveness of this extension (§ A.3).

040 **A.1. Methodology**

041 **Architecture modification.** To support training on diverse
042 datasets with varying numbers of input views and an in-
043 creased number of Gaussians, we introduce minimal mod-
044 ifications to the original Comp3D- \mathcal{G} architecture. To fully
045 exploit the pretrained weights of Comp3D- \mathcal{G} , which were
046 already well-trained on the RealEstate10K dataset [52], we
047 keep the architectural changes as minimal as possible to pre-
048 serve the original structure.

049 First, to handle differences in the intrinsics of the cam-
050 era and the scale of the scene between training datasets, we

incorporate the intrinsic embeddings as follows [19, 48].
Specifically, we additionally take camera intrinsics as in-
put and pass them through a projection layer consisting of
a 2-layer MLP initialized with zero weights. The projected
intrinsic features are then added to the patch-embedded fea-
tures within VGGT [42].

Next, we increase the number of Gaussians to better
cover larger scenes and capture finer details. To achieve
this, we replicate the weights of the Gaussian head. In
the original architecture, the Gaussian head generates one
Gaussian per token. By replicating its weights N_G times,
we instead produce N_G Gaussians per token, resulting in a
total of $N \times N_G$ Gaussians. Since the replicated heads share
the same initialization, Gaussians originating from the same
token are localized in the same spatial region, encouraging
them to focus on capturing finer local details.

Additional loss. To improve training stability, we intro-
duce additional geometry supervision in the form of depth
and normal losses, encouraging the network to develop a
stronger understanding of scene geometry. We first extract
pseudo ground-truth normal maps N_t and depth maps D_t
from the point maps produced by VGGT [42] for the corre-
sponding view I_t . We then render the depth maps \hat{D}_t and
normal maps \hat{N}_t using the same rasterizer as 3DGS [20],
replacing the color attributes with z -value and surface nor-
mals. For supervision, we apply scale-shift invariant loss
term for depth and MSE loss for normals,

$$\mathcal{L}_{\text{depth}} = \|(\alpha \cdot D_t + \beta) - \hat{D}_t\|, \quad (1)$$

$$\mathcal{L}_{\text{normal}} = \|N_t - \hat{N}_t\|, \quad (2)$$

where α and β is the scale and shift parameters obtained via
least-squares alignment between rendered depth and pseudo
ground-truth depth, following [43].

The overall loss term is followed as:

$$\mathcal{L} = \mathcal{L}_{\text{novel}} + \lambda_{\text{depth}} \mathcal{L}_{\text{depth}} + \lambda_{\text{normal}} \mathcal{L}_{\text{normal}}, \quad (3)$$

where λ_{depth} and λ_{normal} denotes the loss weight of normal
and depth.

058 **A.2. Implementation details**

059 For efficient training, we initialize from the pretrained
060 weights of Comp3D- \mathcal{G} which were trained solely on the
061 RealEstate10K dataset [52]. We set the number of Gaus-
062 sians per query token to $N_G = 16$. During training, we ran-
063 domly sample 2 to 24 images per scene, with a maximum of
064 192 images per epoch, and train at a resolution of 224×224 .
065 We use the AdamW optimizer [21] with a learning rate of
066 $1e-6$. The geometry loss weights are set to $\lambda_{\text{depth}} = 0.01$
067 and $\lambda_{\text{normal}} = 0.001$. We train for 150K steps in total, with
068 all other hyperparameters following Comp3D- \mathcal{G} .

069 For training, we incorporate 8 datasets spanning di-
070 verse domains. Specifically, we use RealEstate10K [52],

Table 10. **Comparison of novel view synthesis with multi-view input images on DL3DV [27].** Our method generates fewer Gaussians while achieving competitive or superior quality.

Methods	12 view				24 view				36 view			
	PSNR↑	SSIM↑	LPIPS↓	#G↓	PSNR↑	SSIM↑	LPIPS↓	#G↓	PSNR↑	SSIM↑	LPIPS↓	#G↓
AnySplat [19]	17.004	0.439	0.431	570K	17.263	0.445	0.439	1,140K	18.286	0.525	0.286	1,576K
VGGT+NoPo [42, 48]	11.714	0.320	0.550	602K	10.285	0.262	0.607	1,204K	9.65	0.236	0.628	1,806K
C3G (Ours)	18.768	0.505	0.408	2K	18.843	0.511	0.404	2K	18.675	0.504	0.412	2K
C3G++ (Ours)	19.368	0.512	0.431	32K	19.570	0.521	0.425	32K	19.736	0.523	0.424	32K

101 DL3DV [27], HyperSim [36], WildRGBD [46], Tar-
 102 tanAir [44], ARKitScenes [3], BlendedMVS [47], and
 103 MapFree [1]. These datasets collectively cover a wide range
 104 of scenarios, including indoor, outdoor, and object-centric
 105 scenes across both real and synthetic domains.

106 A.3. Experiments

107 We evaluate novel-view synthesis performance on the test
 108 split of the DL3DV dataset [27]. As shown in Tab. 10, our
 109 model outperforms the baselines, benefiting from its im-
 110 proved generalizability and the increased number of Gaus-
 111 sians, which together enable better coverage of large scenes
 112 and finer detail capture. This is particularly evident on
 113 DL3DV, which predominantly consists of outdoor scenes
 114 with large spatial extents that are difficult to represent with
 115 only 2K Gaussians. By scaling up the number of Gaussians,
 116 our model can cover a significantly larger portion of each
 117 scene, leading to improved reconstruction quality. Further-
 118 more, the results demonstrate that our model is robust to
 119 varying numbers of input views.

120 B. Additional implementation details

121 B.1. Details of Gaussian decoder

122 **Architecture details.** We first extract feature maps
 123 $\mathbf{F}_v = \mathcal{E}(I_v)$ from the visual encoder $\mathcal{E}(\cdot)$ given V
 124 input images $\{I_v\}_{v=1}^V$. For VGGT [42] which supports
 125 multi-view inputs, we directly obtain V feature maps
 126 $\{\mathbf{F}_v\}_{v=1}^V = \mathcal{E}(\{I_v\}_{v=1}^V)$ in a single feed-forward pass. For
 127 DINOv3 [39], we obtain V feature maps by passing each
 128 image to the encoder independently. We then concatenate
 129 the feature maps with the learnable tokens \mathbf{Q} , which are
 130 randomly initialized. This concatenated representation is
 131 processed through L transformer layers $\mathcal{T}_G(\cdot)$, which con-
 132 sists of self-attention layers and MLP layers with ReLU ac-
 133 tivation functions and layer normalization [2]:

$$134 \quad [\bar{\mathbf{Q}}; \bar{\mathbf{F}}] = \mathcal{T}_G([\mathbf{Q}; \mathbf{F}]), \quad (4)$$

135 where $[\cdot, \cdot]$ denotes concatenation of dimension axis, $\bar{\mathbf{Q}}$
 136 denotes refined learnable tokens and $\bar{\mathbf{F}}$ denotes refined fea-
 137 tures. In the self-attention layer, learnable tokens and the
 138 feature tokens are each projected to query, key, and value
 139 features to process the attention calculation. The query, key

and value features of learnable tokens \mathbf{Q} and visual encoder
 features \mathbf{F} is calculated as follows:

$$142 \quad Q = [Q_G; Q_I] = \mathcal{P}_Q([\mathbf{Q}; \mathbf{F}]), \quad (5)$$

$$143 \quad K = [K_G; K_I] = \mathcal{P}_K([\mathbf{Q}; \mathbf{F}]), \quad (6)$$

$$144 \quad V = [V_G; V_I] = \mathcal{P}_V([\mathbf{Q}; \mathbf{F}]), \quad (7)$$

145 where Q_G, Q_I denotes query features, K_G, K_I denotes key
 146 features, and V_G, V_I denotes value features of learnable to-
 147 ken and visual encoder features, respectively. And $\mathcal{P}_Q, \mathcal{P}_K$,
 148 and \mathcal{P}_V denote the projection layer of query, key, and value,
 149 respectively.

150 Then, the resulting output of self-attention layer in \mathcal{T}_G
 151 is computed as:

$$152 \quad \text{Attn}(Q, K, V) = \text{Softmax}\left(\frac{QK^\top}{\sqrt{d_Q}}\right)V, \quad (8)$$

153 where d_Q denotes the channel dimension of query features.

154 Each of the refined learnable tokens $\bar{\mathbf{Q}}_i$ are decoded as
 155 a single Gaussian \mathbf{G}_i by passing the Gaussian head $H_G(\cdot)$,
 156 which consists of a single linear layer:

$$157 \quad \mathbf{G}_i = H_G(\bar{\mathbf{Q}}_i). \quad (9)$$

158 **Attention visualization.** To visualize the attention map in
 159 Fig. 3-(b), we select one Gaussian \mathbf{G}_i and its correspond-
 160 ing learnable tokens $\bar{\mathbf{Q}}_i$. We then calculate the attention
 161 weights between the query features of the learnable token
 162 Q_{G_i} and the key features K_I as follows:

$$163 \quad A_i = Q_{G_i} \cdot K_I^\top, A_i \in \mathbb{R}^{N_{\text{head}} \times h \times w}, \quad (10)$$

164 where N_{head} denotes number of heads. We average the A_i
 165 at the head dimension, then apply min-max normalization
 166 for better visualization.

167 B.2. Details of feature decoder

168 **Architecture details.** We initialize Comp3D- \mathcal{F} by copying
 169 the Comp3D- \mathcal{G} 's architecture and parameters, except for the
 170 learnable tokens and the Gaussian head. We first extract
 171 feature maps $\mathbf{F}'_v = \mathcal{E}'(I_v)$ using the user-desired visual en-
 172 coder $\mathcal{E}'(\cdot)$. We then introduce new learnable feature tokens
 173 \mathbf{Q}' and concatenate them with \mathbf{F}'_v , which are then processed
 174 through $\mathcal{T}_{\mathcal{F}}(\cdot)$ with the same architecture as $\mathcal{T}_G(\cdot)$:

$$175 \quad [\bar{\mathbf{Q}}'; \bar{\mathbf{F}}'] = \mathcal{T}_{\mathcal{F}}([\mathbf{Q}'; \mathbf{F}']), \quad (11)$$

176 where $\bar{\mathbf{Q}}'$ denotes the refined learnable tokens and $\bar{\mathbf{F}}'$ denotes the refined features. In the self-attention layers of $\mathcal{T}_{\mathcal{F}}$,
177 we reuse query Q and key K features from $\mathcal{T}_{\mathcal{G}}$ with stop-
178 gradients, while only the value projection layer is newly
179 trained:
180

$$181 \quad V' = [V'_G; V'_I] = \mathcal{P}'_V([\mathbf{Q}'; \mathbf{F}']), \quad (12)$$

182 where V'_G, V'_I denote the value features of the learnable fea-
183 ture tokens and features from \mathcal{E}' , respectively, and \mathcal{P}'_V de-
184 notes the projection layer of value in $\mathcal{T}_{\mathcal{F}}$.

185 Then, the resulting output of self-attention layer in $\mathcal{T}_{\mathcal{F}}$ is
186 computed as:

$$187 \quad \text{Attn}(Q, K, V') = \text{Softmax}\left(\frac{QK^\top}{\sqrt{d_Q}}\right)V'. \quad (13)$$

188 These refined learnable feature tokens $\bar{\mathbf{Q}}'$ are converted
189 to multi-view aggregated features \mathbf{F}''_i for each Gaussian \mathbf{G}_i
190 by passing the feature head $H_{\mathcal{F}}$, which consists of a single
191 linear layer:

$$192 \quad \mathbf{F}''_i = H_{\mathcal{F}}(\bar{\mathbf{Q}}'_i). \quad (14)$$

193 Note that the introduction of new learnable queries \mathbf{Q}'
194 is to enable Comp3D- \mathcal{F} to take features extracted from
195 any user-desired visual encoder $\mathcal{E}'(\cdot)$ as input, considering
196 that different features ($\mathcal{E}(\cdot), \mathcal{E}'(\cdot)$) will have different fea-
197 ture dimensions. When the two visual encoders are identi-
198 cal ($\mathcal{E}(\cdot) = \mathcal{E}'(\cdot)$), we can additionally copy the learned
199 queries \mathbf{Q} from Comp3D- \mathcal{G} , only training the value projec-
200 tion layer $\mathcal{P}'_V(\cdot)$ and the feature head $H_{\mathcal{F}}(\cdot)$.

201 C. Experimental settings

202 C.1. Novel view synthesis

203 **Datasets.** To evaluate novel view synthesis, we use the
204 RealEstate10K [52] dataset. We adopt the same train-test
205 split as prior work [5, 6, 48]. RealEstate10K primarily con-
206 tains indoor real estate videos with camera poses computed
207 using COLMAP [38]. For multi-view evaluation, we ran-
208 domly sample 1,000 image sets from the test split, each con-
209 sisting of multiple context views and 3 target views that are
210 disjoint from the context views.

211 **Baselines.** For 2-view input settings, we compare against
212 SOTA generalizable feed-forward methods on novel view
213 synthesis. We compare PixelSplat [5] and MVSplat [6]
214 as *Pose-dependent* categories, which require ground-truth
215 pose information for input. We also compare our model
216 with *Pose-free* categories which use only rgb images as in-
217 puts such as CoPoNeRF [17], Splatt3R [40], PF3plat [16],
218 SPFSplat [18], NoPoSplat [48]. We also reimplement the
219 VGGT+NoPo, which replaces NoPoSplat’s MAST3R [24]
220 backbone with VGGT [42] while maintaining NoPoSplat’s
221 pipeline to estimate per-pixel Gaussians. Ours also fall into
222 *Pose-free* settings.

For the multi-view ($V > 2$) settings, we compare against
state-of-the-art generalizable methods that support arbitrary
numbers of input views, such as AnySplat [19]. We evaluate
using various numbers of input views: 12, 24, and 36. We
also evaluate VGGT+NoPo under the same settings. For
AnySplat, we provide input images at 448×448 resolutions,
which perform better than 256×256 , and render target view
images at 256×256 resolution.

Evaluation protocol. Given unposed images as inputs,
our method reconstructs and represents 3D scenes using
3D Gaussians. We rasterize the 3D Gaussians at ground-
truth camera poses and compare the rendered images with
ground-truth target views. We report standard novel view
synthesis metrics such as PSNR, SSIM, and LPIPS.

However, as discussed in NoPoSplat [48], reconstructing
3D scenes from sparse unposed views is inherently scale-
ambiguous. Although our method successfully generates
3D Gaussians, they may not fully align with the ground-
truth scene scale in the validation dataset. To ensure fair
comparison with other baselines, we follow pose-free meth-
ods by optimizing the target view camera pose while freez-
ing all other parameters. Specifically, we first reconstruct
3D Gaussians, then freeze them and optimize only the target
camera pose such that the rendered image closely matches
the ground-truth target view. Finally, we compute metrics
using the optimized target view camera poses. Note that this
optimization scheme is only needed for evaluation purposes
and is not required in real-world scenarios.

For the multi-view settings, we additionally perform
short test-time optimization following 3DGS [20], denoted
as C3G w/ TTO, NoPo+VGGT w/ TTO and AnySplat w/
TTO. We use the predicted Gaussians as initialization and
conduct per-scene optimization for a limited number of
steps. Specifically, we set the optimization steps to 1,000
steps with densification intervals of 100 steps. We set the
learning rate as follows: means at $1.6e-4$, scales at $3e-4$, ro-
tations at $1e-3$, harmonics at $2.5e-3$, and densities at $5e-3$,
using Adam optimizer [21]. Our loss function is formulated
as $\mathcal{L} = 0.8\mathcal{L}_{\text{MSE}} + 0.2\mathcal{L}_{\text{SSIM}}$, where \mathcal{L}_{MSE} denotes the MSE
loss and $\mathcal{L}_{\text{SSIM}}$ denotes the SSIM loss between the rendered
image and ground-truth images. Other settings follow the
original 3DGS. For AnySplat, we follow their default test-
time optimization and do not conduct densification to avoid
the Out-Of-Memory due to the number of Gaussians.

267 C.2. 3D scene understanding

Datasets. To evaluate 3D scene understanding capabilities,
we follow previous approaches [4, 11] and use the Scan-
Net [9] and Replica [41] datasets.

For ScanNet, we use 40 scenes following LSM [11], se-
lected based on valid pose and depth data where Feature-
3DGS [51] performs well. We evaluate 8 categories: wall,
floor, ceiling, chair, table, sofa, bed, and others. We select

275 30 images with a stride of 10, and target views are chosen
276 as the 1st and 4th images within every 8-image interval.

277 For Replica, we select 80 images with a stride of 3, and
278 the target view is the 2nd image within every 8-image in-
279 terval. We use camera poses obtained from COLMAP [38].
280 We evaluate 3 scenes where LSeg [25] performs well: of-
281 fice3, office4, and room1. We select 5-6 categories per
282 scene. Specifically, for office3, we use wall, ceiling, floor,
283 chair, and table; for office4, we use wall, ceiling, floor,
284 chair, tv-screen, and table; for room1, we use wall, ceiling,
285 floor, bed, and blinds.

286 **Baselines.** To validate our effectiveness, we compare our
287 model with per-scene optimized feature novel-view synthe-
288 sis tasks and feed-forward feature novel-view synthesis.

289 For per-scene optimization methods (Feature-3DGS [51]
290 and CF³ [23]), we optimize the per-scene 3D Gaussians
291 using all the posed inputs except for target view images.
292 We then render the 3D Gaussians to target view images to
293 calculate the metrics. For the ScanNet dataset, there are
294 no readily available sparse initial point clouds, so we train
295 with randomly initialized point clouds following the previ-
296 ous works [26, 34]. For feature-3DGS [51], we optimize
297 for 5,000 steps to avoid overfitting. For CF³ [23], we first
298 optimize 3DGS with 30,000 steps and additionally optimize
299 the CF³ with 3,000 steps.

300 For feature novel-view synthesis (LSM [11] and Ours),
301 we first select two input views to generate 3D Gaussians,
302 which are then projected to target views disjoint from the
303 input views. Following the same protocol as novel view
304 synthesis evaluation, these methods require target pose op-
305 timization to resolve scale ambiguity.

306 For LSeg [25] and MaskCLIP [50], we directly extract
307 features from the target viewpoint images, since these meth-
308 ods cannot render features at novel viewpoints like 3D-
309 based approaches.

310 **Evaluation protocol.** Given RGB images as input, we
311 first extract language-aware features using models such as
312 LSeg [25] or MaskCLIP [50]. 3D Gaussians with lifted
313 features are generated following each method’s procedure,
314 then projected to the target and source views. The target
315 view denotes an unseen viewpoint not present in the in-
316 put, while the source view denotes a seen viewpoint that
317 every method has observed at least once. We render both
318 RGB images and features at each viewpoint. The rendered
319 features are then processed into segmentation maps by se-
320 lecting the most relevant CLIP [35] text embeddings corre-
321 sponding to each label. We report open-vocab segmentation
322 metrics (mIOU and Accuracy), and novel view synthesis
323 metrics (PSNR, SSIM, and LPIPS).

324 C.3. Multi-view feature encoding

325 Obtaining multi-view invariant features, also termed as 3D
326 aware features, has been a long-standing goal in computer

vision and graphics. Probe3D [10] defines a set of tasks, in-
327 cluding two-view correspondences [7, 8, 12–15] and single-
328 view depth estimation [33], to probe the 3D awareness of
329 the features. FiT3D proposes a two-stage fine-tuning task
330 to build view-invariant features, where they first train more
331 than 1000 per-scene feature-3DGS with the desired fea-
332 tures. With the pre-trained 3DGS over multiple scenes, they
333 finetune the original vision encoder to follow the rendered
334 features from the 3D Gaussians. They show that after train-
335 ing, they achieve higher correspondence scores in the task
336 defined in Probe3D. As our Comp3D- \mathcal{F} can aggregate the
337 input features, we analyze the effectiveness of Comp3D- \mathcal{F}
338 as a view-invariant feature decoder. 339

340 **Datasets.** To validate the effectiveness of our Comp3D-
341 \mathcal{F} as a multi-view invariant feature decoder, we use the
342 ScanNet [9] dataset following Probe3D [10]. ScanNet is
343 a large-scale dataset of indoor scenes with RGB images,
344 depth maps, and camera poses. We evaluate on 1,500 image
345 pairs from the test split following SuperGlue [37].

346 **Baselines.** We compare DINOv2 [31], DINOv3 [39], and
347 VGGT tracking features [42], which are known to perform
348 well on zero-shot correspondence estimation.

349 **Evaluation protocols.** Given two images, we first extract
350 a feature map for each image. For ours, we first predict
351 3D Gaussians with features and project the 3D Gaussians to
352 each input-view pose. Then, we estimate correspondences
353 between the two images using nearest neighbors. Following
354 Probe3D [10], we filter the correspondence using Lowe’s
355 ratio test [30] to find the strong unique matches to reduce the
356 noisy correspondence. We rank the correspondences using
357 the ratio test and keep the top 1,000 correspondences.

358 We evaluate correspondence quality using PCK @ 10px,
359 which measures the accuracy of correspondence within 10
360 pixels between estimated and ground-truth matches. We di-
361 vide the image pairs into 4 groups based on viewing angle
362 differences: $0\sim 15^\circ$ (θ_0^{15}), $15\sim 30^\circ$ (θ_{15}^{30}), $30\sim 60^\circ$ (θ_{30}^{60}), and
363 $60\sim 180^\circ$ (θ_{60}^{180}). The results can be seen in Tab. 4.

364 C.4. Multi-view feature upsampling

365 Feature upsampling is the task of upsampling the extracted
366 features from pre-trained vision encoders to a higher resolu-
367 tion for more fine-grained downstream tasks. Although our
368 Comp3D- \mathcal{F} also takes the low-resolution extracted features
369 from $\mathcal{E}'(\cdot)$ as input, when combined with the 3D Gaussians
370 and rendered to a specific viewpoint, since Gaussians can
371 be rendered at any camera pose, our model can generate
372 features at arbitrary resolutions by setting the desired reso-
373 lution and intrinsics in the CUDA rasterizer. In this section,
374 we analyze whether the rendered features from Comp3D- \mathcal{F}
375 can be used as an effective multi-view feature upsampler.

376 **Datasets.** We use the ScanNet [9] dataset, following the
377 same setup as multi-view feature encoding, to evaluate the
378 probing capability of the upsampled features.

379 **Baselines.** We compare DINOv2 [31] and DINOv3 [39]
380 features as baselines against their upsampled versions,
381 which match the input image resolution using AnyUp [45],
382 a generalizable feature upsampling module.

383 **Evaluation protocols.** Following the same protocol
384 as multi-view feature encoding, we first extract high-
385 resolution feature maps, then estimate correspondences be-
386 tween two images using nearest neighbors. We evaluate
387 correspondence quality using PCK@10px with image pairs
388 divided by viewing angle differences. We also present qual-
389 itative results using PCA visualization to demonstrate how
390 upsampled features capture finer details. The results can be
391 seen in Tab. 4.

Table 11. Comparison of novel view synthesis with compact Gaussians generation strategies.

	Methods	Average			
		PSNR \uparrow	SSIM \uparrow	LPIPS \downarrow	#G \downarrow
2-view image input	Sampling	21.340	0.665	0.272	2K
	Voxelize	19.957	0.609	0.403	4K
	Ours	22.387	0.713	0.259	2K
24-view image input	Sampling	22.446	0.714	0.219	25K
	Sampling w/ TTO	25.487	0.832	0.158	28K
	Voxelize	19.052	0.558	0.446	81K
	Voxelize w/ TTO	27.751	0.885	0.165	90K
	NoPo+VGGT	21.244	0.664	0.200	1,204K
	NoPo+VGGT w/ TTO	28.463	0.902	0.135	1,204K
	Ours	23.797	0.747	0.198	2K
	Ours w/ TTO	29.987	0.916	0.136	27K

Table 12. Analysis of tradeoff between rendering performance and number of Gaussians.

# of Gaussians	PSNR \uparrow	SSIM \uparrow	LPIPS \downarrow	FPS \uparrow
2048 \times 1	20.625	0.623	0.321	3083.50
2048 \times 2	20.844	0.639	0.303	3015.03
2048 \times 4	21.080	0.644	0.297	2828.14
2048 \times 8	21.218	0.653	0.284	2476.01

392 D. Additional ablation and discussion

393 D.1. Comparison of compact 3D reconstruction 394 strategies

395 To validate the effectiveness of our compact reconstruction
396 strategy, we conduct ablation studies comparing different
397 approaches to obtain compact representations. For the sam-
398 pling baseline, we first estimate per-pixel Gaussian centers
399 and their corresponding Gaussians, then downsample by $8\times$
400 to predict only 2,048 Gaussians, matching the number of
401 Gaussians with our method. For the voxelization baseline,
402 we follow AnySplat’s [19] strategy: we first estimate per-
403 pixel Gaussians, then voxelize them with additional layers.
404 Since AnySplat’s original voxel size cannot sufficiently re-
405 duce the number of Gaussians, we increase the voxel size
406 to 0.2, which results in approximately 4K Gaussians to best
407 match the 2,048 Gaussians produced by our method.

408 As shown in Tab. 11, our strategy demonstrates superior
409 performance compared to sampling or voxelization methods
410 while using fewer Gaussians than voxelization. Both sam-
411 pling and voxelization rely on per-pixel estimation, which
412 restricts them to pixel locations and prevents accurate es-
413 timation of necessary regions. Additionally, voxelization
414 requires heuristic determination of voxel size and can re-
415 duce details in high-frequency regions. In contrast, our ap-
416 proach does not depend on pixel locations as learnable to-
417 kens can flexibly estimate Gaussians at appropriate loca-
418 tions, enabling more effective compact Gaussian field re-
419 construction.

Table 13. Additional ablation studies for Comp3D- \mathcal{F} .

Detach attn.	Detach $\mathcal{L}_{\text{feat}}$ to G_i	Lseg [25]				
		mIOU \uparrow	Acc. \uparrow	PSNR \uparrow	SSIM \uparrow	LPIPS \downarrow
\times	\checkmark	0.490	0.761	23.078	0.750	0.293
\checkmark	\times	0.490	0.757	23.243	0.754	0.286
\checkmark	\checkmark	0.513	0.783	23.886	0.770	0.285

D.2. Additional analysis of number of Gaussians

420 We first analyze the tradeoff between rendering perfor-
421 mance and the number of Gaussians in Tab. 6. Notably, this
422 experiment reveals a slight performance drop when increas-
423 ing the number of learnable tokens from 2,048 to 4,096.
424 We attribute this to our framework’s design of using a fixed
425 number of learnable queries, which encourages each query
426 to learn cross-view correspondences. Increasing the num-
427 ber of tokens appears to weaken this inductive bias toward
428 coherent multi-view aggregation, making attention scores
429 less localized and yielding no additional benefit. In con-
430 trast, our experiments in Tab. 12 demonstrate that increas-
431 ing the number of decoded Gaussians *per* token consistently
432 improves performance, enabling the model to capture fine-
433 grained details from scenes coarsely represented by 2K to-
434 kens. We believe this analysis highlights a promising direc-
435 tion for applications that prioritize feed-forward NVS per-
436 formance. Note that the results in Tab. 12 are reported after
437 75K training steps (out of 450K) to assess feasibility.
438

D.3. Additional ablation for feature decoder

439 We additionally ablate the components of Comp3D- \mathcal{F} . In
440 Tab. 13, we conduct experiments on (1) detaching the
441 copied attention weights from Comp3D- \mathcal{G} , and (2) propa-
442 gating feature loss to Gaussian geometry attributes (means,
443 covariances, spherical harmonics, and opacities). When
444 feature loss propagates to the copied attention weights,
445 Comp3D- \mathcal{G} receives ambiguous gradients because the fea-
446 tures are not perfectly multi-view consistent. Consequently,
447 the attention mechanism cannot correctly identify corre-
448 spondences between learnable query tokens and the encoder
449 features $\mathcal{E}(\cdot)$ (e.g., VGGT) for Comp3D- \mathcal{G} . To propagate
450 feature loss to Gaussian attributes, we modify the CUDA
451 rasterizer from Feature-3DGS [51], which originally prop-
452 agates feature loss only to Gaussian feature attributes. We
453 extend the CUDA rasterizer to propagate feature loss to all
454 Gaussian attributes. However, this also degrades geome-
455 try estimation results because foundation model features are
456 not perfectly multi-view consistent. We hypothesize that
457 with perfectly multi-view invariant features, feature loss
458 could improve reconstruction quality, especially since pho-
459 tometric loss is also imperfect, as real-world RGB images
460 contain visual artifacts such as appearance variations, light-
461 ing changes, and noise.
462

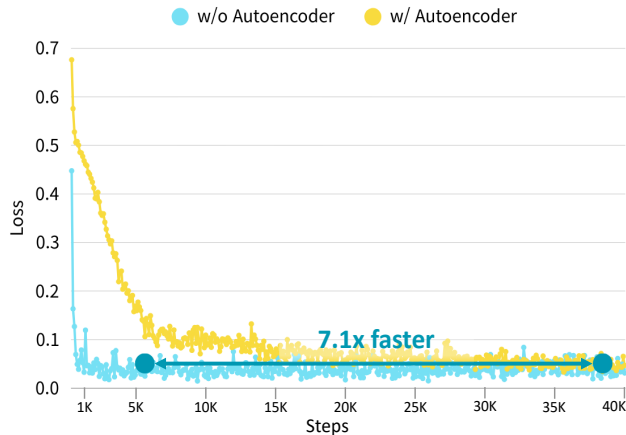


Figure 6. Convergence speed improvement by eliminating autoencoder in our framework.

Table 14. Computational Cost.

	Methods	Inference time (s) ↓	GPU Mem. (MiB) ↓	FPS ↑
2-view image input	NoPoSplat [48]	0.092	6,063	898.28
	NoPo+VGGT	0.093	7,495	1566.28
	Ours	0.091	4,150	2742.68
24-view image input	AnySplat [15]	1.449	6,722	225.50
	AnySplat w/ TTO [15]	35.907	6,722	225.50
	Ours	0.466	4,372	3083.50
	Ours w/ TTO	4.273	4,372	2511.49
Feature lifting	LSM [9]	0.165	5,333	625.12
	Ours	0.093	5,896	785.35

D.4. Discussion of autoencoder in existing methods

To reduce the memory consumption of Gaussians, prior works [11, 23] commonly adopt autoencoders to compress Gaussian feature dimensions and restore them via a decoder. However, autoencoder-based compression inevitably introduces information loss. Moreover, as shown in Fig. 6, removing the autoencoder allows our model to converge $7.1\times$ faster. With an autoencoder framework, additional training time is required for the encoder–decoder architecture. In contrast, without an autoencoder, we directly leverage the attention mechanism of Comp3D- \mathcal{G} to aggregate and store features. Since the attention maps in Comp3D- \mathcal{G} already identify the necessary features for each learnable token, each feature token naturally attends to the same multi-view regions as its corresponding Gaussian token, effectively reusing learned correspondences for feature aggregation. This emergent property enables highly efficient feature lifting with minimal additional training overhead.

D.5. Discussion of Computational Costs

We include a detailed analysis of computational costs in Tab. 14. In *novel-view synthesis* settings, compared to prior works [19, 48] that employ a DPT head — which incurs significant computational overhead due to the large feature res-

olution space — our shallow transformer decoder achieves faster inference with lower peak GPU memory usage. For *feature lifting*, our model achieves higher rendering speed and FPS compared to LSM [11]. While LSM must compress features to remain tractable due to its large number of Gaussians, our compactness enables feature lifting and rendering without any compression. Although increasing the feature dimension slightly raises peak GPU memory, the overhead is modest and predictable, enabling our method to deliver higher PSNR and mIoU — validating that compact Gaussians provide an efficient and effective substrate for feature lifting.

D.6. Discussion of FPS

As shown in Tab. 3 and Tab. 4, our method reduces the number of Gaussians by $65\times$ (to only **2K**) compared to per-pixel approaches [11, 48], yet FPS does not scale linearly with this reduction. We further analyze FPS in Tab 14. In the main table, FPS measurements include camera matrix multiplication and tensor operations for pre-processing prior to Gaussian rasterization. We additionally report FPS measured over the rendering step alone, which demonstrates that our compactness becomes more valuable when isolating the rendering stage, yielding substantially higher FPS.

D.7. Limitations

Although our view-invariant feature decoder, Comp3D- \mathcal{F} , allows lifting arbitrary 2D features into 3D, we have not evaluated all recent foundation models. For instance, we believe that leveraging features from the Segment Anything Model (SAM) [22] could enable robust multi-view consistent segmentation if the features are aggregated within our framework before being decoded by a pre-trained SAM decoder; however, we did not analyze these specific features in this work.

In addition, following prior works, our experimental validation primarily focuses on multi-view open-vocabulary segmentation. To achieve holistic 3D scene understanding, future work could explore integrating our feature fields with recent Multimodal Large Language Models (MLLMs) [28, 29] to enable 3D Scene Question Answering. Furthermore, since our framework enables novel view feature rendering without information loss or compression artifacts, it holds significant promise for integration with Vision-Language-Action (VLA) models or robotics applications. However, such scenarios are frequently dynamic, whereas our current framework is limited to static scene reconstruction. Extending our compact representation to dynamic scenes remains an exciting avenue that would unlock potential across various autonomous fields.

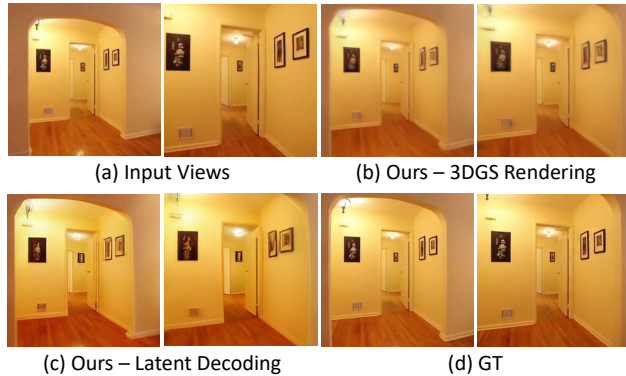


Figure 7. **Novel view synthesis via latent decoding.** We explore the potential of combining our view-invariant feature decoder, $\text{Comp3D-}\mathcal{F}$, with generative models. We lift DINOv2-base features (which serve as latents for a Representation Autoencoder (RAE) [49]) extracted from the input views ((a)) to 3D Gaussians and render them at novel viewpoints. (b) **Ours – 3DGS Rendering:** Standard RGB rendering from the estimated Gaussians provides faithful reconstruction but exhibits some blurriness in high-frequency regions. (c) **Ours – Latent Decoding:** By decoding the rendered feature maps through the RAE’s diffusion transformer with a 1-step denoising process, we recover significantly sharper textures and edges. Although the generative nature introduces slight variations in fine details compared to the Ground Truth (d), this validates our model’s ability to provide consistent 3D-aware latents for diffusion-based pipelines.

534 E. Additional experiment results

535 E.1. Novel view synthesis via latent decoding

536 As our multi-view invariant feature decoder $\text{Comp3D-}\mathcal{F}$
 537 can take any desired feature as input, we conduct an additional
 538 experiment of leveraging the latents of the representation
 539 autoencoder (RAE) [49] as input. Specifically,
 540 we leverage the open-sourced RAE, which learns a diffusion
 541 transformer [32] and a decoder based on the frozen
 542 DINOv2-base [31] as the latent. By lifting the DINOv2-
 543 base features and rendering novel viewpoints, we can de-
 544 code the rendered features using the pre-trained RAE de-
 545 coder or DiT after adding a small noise to the rendered fea-
 546 tures. In Fig. 7, we show the results of novel view syn-
 547 thesis obtained from the estimated 3D Gaussians (Fig. 7-(b))
 548 and the novel view synthesis obtained by decoding the ren-
 549 dered features through the DiT (Fig. 7-(c)). Specifically,
 550 for this experiment, we add a 1-step noise from the original
 551 50-step denoising schedule of RAE-DiT [49], and denoise
 552 with 1-step. As shown in Fig. 7, renderings from 3D Gaus-
 553 sians already provide sufficiently good results but still con-
 554 tain some blurry regions, whereas the latent decoding yields
 555 significantly sharper results. However, since the pipeline in-
 556 volves a generative process, even with a single step, we ob-
 557 serve slight deformations in fine details. Nevertheless, this
 558 experiment demonstrates the potential of combining novel

Table 15. **Comparison of novel view synthesis on RealEstate10K [52].** Our method maintains competitive results while using far fewer Gaussians.

Pose-free	Methods	Average					
		#G ↓	Memories ↓	FPS ↑	PSNR ↑	SSIM ↑	LPIPS ↓
✗	PixelSplat [5]	131K	33.6MB	388.03	23.848	0.806	0.185
	MVSplat [6]	131K	33.6MB	392.6	23.977	0.811	0.176
✓	CoPoNeRF [17]	-	-	0.4	18.938	0.619	0.388
	Splatt3R [40]	131K	33.6MB	393.1	15.318	0.490	0.436
	PF3splat [16]	131K	33.6MB	397.1	21.042	0.739	0.233
	SPFSplat [18]	131K	33.6MB	397.3	25.845	0.852	0.152
	NoPoSplat [48]	131K	33.6MB	369.8	25.033	0.838	0.160
	VGGT+NoPo [42, 48]	100K	25.6MB	419.8	23.015	0.762	0.187
	AnySplat [19]	320K	81.9MB	319.2	18.828	0.656	0.358
	C3G (Ours)	2K	0.1MB	451.7	22.387	0.713	0.259

view synthesis with diffusion processes by leveraging la- 559
 tentals as input to our $\text{Comp3D-}\mathcal{F}$. 560

E.2. Novel view synthesis in two-view setting 561

In this section, we evaluate novel view synthesis perfor- 562
 mance on the RealEstate10K dataset [52], following No- 563
 PoSplat’s [48] protocol where two input images are used 564
 to estimate 3D Gaussians and render a target view. We 565
 compare with both *pose-dependent* models [5, 6] and *pose- 566*
free models [16–18, 40, 48]. Since we use VGGT [42] as 567
 our visual encoder for $\text{Comp3D-}\mathcal{G}$, we additionally train 568
 VGGT+NoPo, which replaces NoPoSplat’s MAST3R [24] 569
 backbone with VGGT [42] while maintaining NoPoS- 570
 plat’s pipeline to estimate per-pixel Gaussians. Note that 571
 $\text{Comp3D-}\mathcal{G}$ directly estimates Gaussians from unposed im- 572
 ages, falling into the *pose-free* category. For NoPoSplat, 573
 VGGT+NoPo, and ours, we follow NoPoSplat’s test-time 574
 camera pose optimization, which is only necessary for eval- 575
 uation. As shown in Tab. 15, despite estimating $65\times$ fewer 576
 Gaussians than per-pixel methods [16–18, 40, 48], our ap- 577
 proach achieves comparable rendering quality with much 578
 faster speeds, validating that our compact Gaussians is suf- 579
 ficient for 3D scene reconstruction. 580

F. Additional qualitative results 581

F.1. Additional results of attention map visualization 582

We additionally visualize the attention maps between learn- 584
 able tokens \mathbf{Q} and features \mathbf{F} in $\mathcal{T}_{\mathcal{G}}$, extending Fig. 3-(b). 585
 As illustrated in Fig. 8, $\text{Comp3D-}\mathcal{G}$ exhibits sharp, focused 586
 attention patterns on spatially coherent regions across mul- 587
 tiple views for all Gaussian tokens and input images. These 588
 results demonstrate an emergent behavior: to accurately re- 589
 construct novel views with a limited number of N Gaus- 590
 sians, the model learns to position 3D Gaussians at geomet- 591
 rically coherent regions. 592

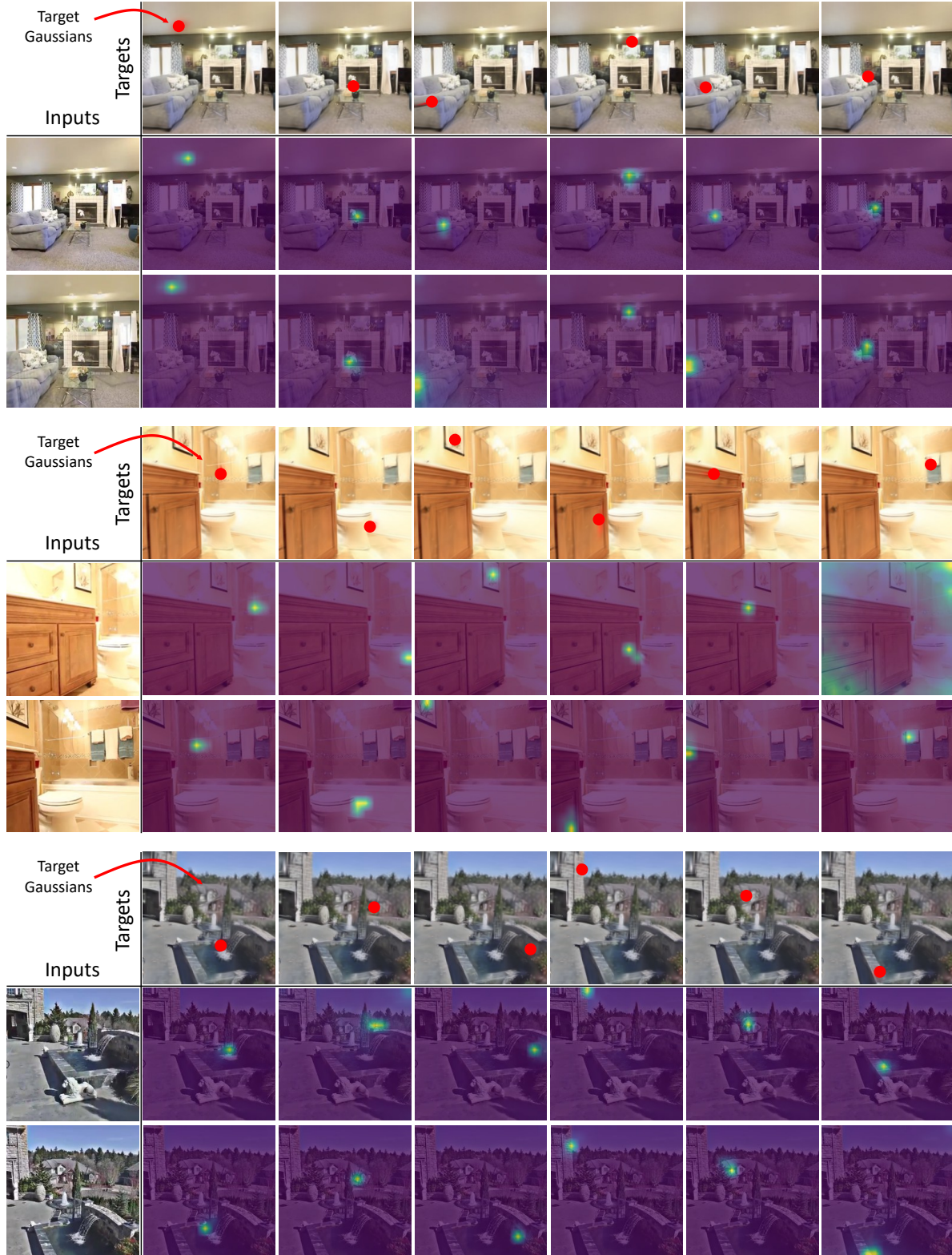


Figure 8. **Additional visualization of learned attention patterns between a target Gaussian and input images.** Without explicit supervision, each query token (red dots) learns to attend to spatially coherent regions across multiple views, naturally discovering corresponding regions.

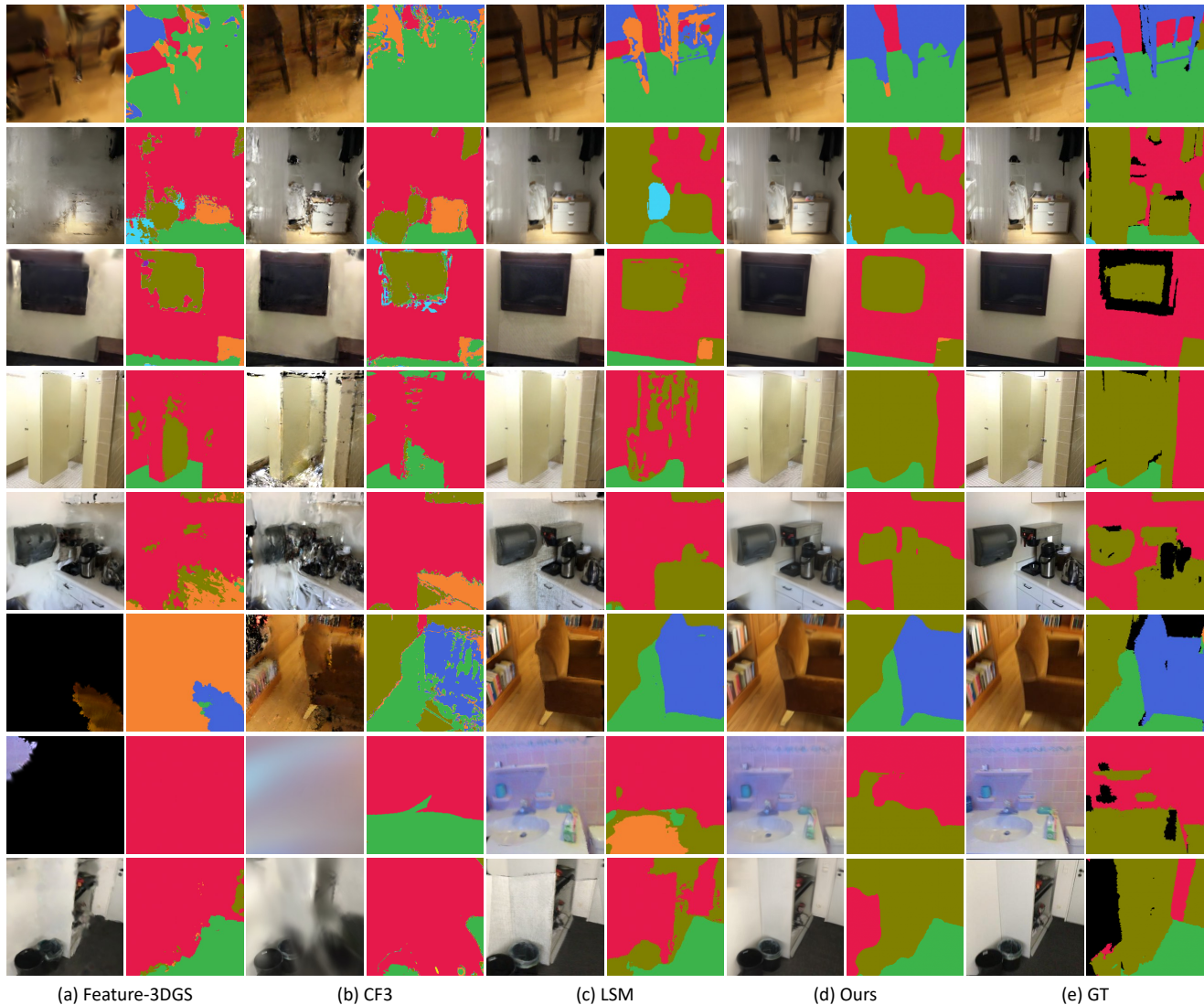


Figure 9. **Additional qualitative results of 3D scene understanding on ScanNet [9].** We conduct qualitative comparison for 3D scene understanding via novel view synthesis and open-vocabulary segmentation. When compared to both per-scene optimization ((a), (b)) and feed-forward ((c), (d)) methods, ours show the most high-fidelity renderings and accurate segmentation maps compared to the ground-truth.

593
594

F.2. Additional qualitative results of 3D scene understanding

595
596
597
598
599
600
601

We additionally present qualitative results for 3D scene understanding on the ScanNet dataset. In Fig. 9, we visualize novel view synthesis results and open-vocabulary segmentation results from feature-lifted 3D Gaussians. Our method generates more geometrically accurate 3D Gaussians and more effectively aggregates multi-view features than competing methods.

F.3. Qualitative results of multi-view feature encoding

602
603

In Fig. 10, we additionally visualize the comparison between pre-lifting features and features aggregated by Comp3D- \mathcal{F} using PCA. As shown, our model effectively aggregates multi-view features before lifting, producing representations that are both view-invariant and more semantically discriminative.

604
605
606
607
608
609

F.4. Qualitative results of multi-view feature upsampling

610
611

In Fig. 11, we present qualitative results of multi-view feature upsampling results.

612
613

614 **F.5. Qualitative results of novel view synthesis with**
615 **multi-view inputs**

616 In Fig. 12, we provide qualitative comparisons of novel
617 view synthesis using 24 input images. We compare
618 our method against AnySplat [19], a state-of-the-art feed-
619 forward approach that supports arbitrary numbers of input
620 views. Despite using significantly fewer Gaussians (approx.
621 2K) compared to AnySplat (approx. 2.6M), our method
622 produces more geometrically consistent renderings with
623 higher quality renderings ((b) vs. (d)). Furthermore, when
624 applying short test-time optimization (denoted as Ours w/
625 TTO and AnySplat w/ TTO), our method significantly out-
626 performs AnySplat in recovering high-frequency details,
627 demonstrating that our compact representation is sufficient
628 for reconstruction and novel view synthesis while also serv-
629 ing as a robust initialization for per-scene optimization.

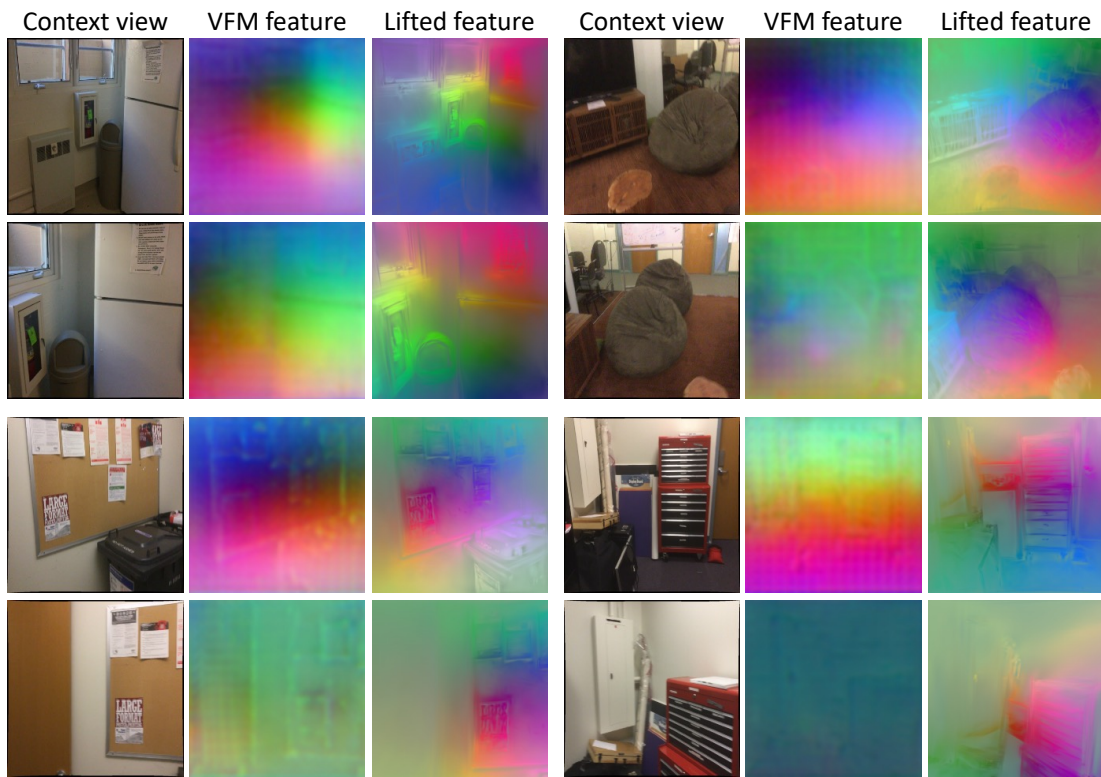
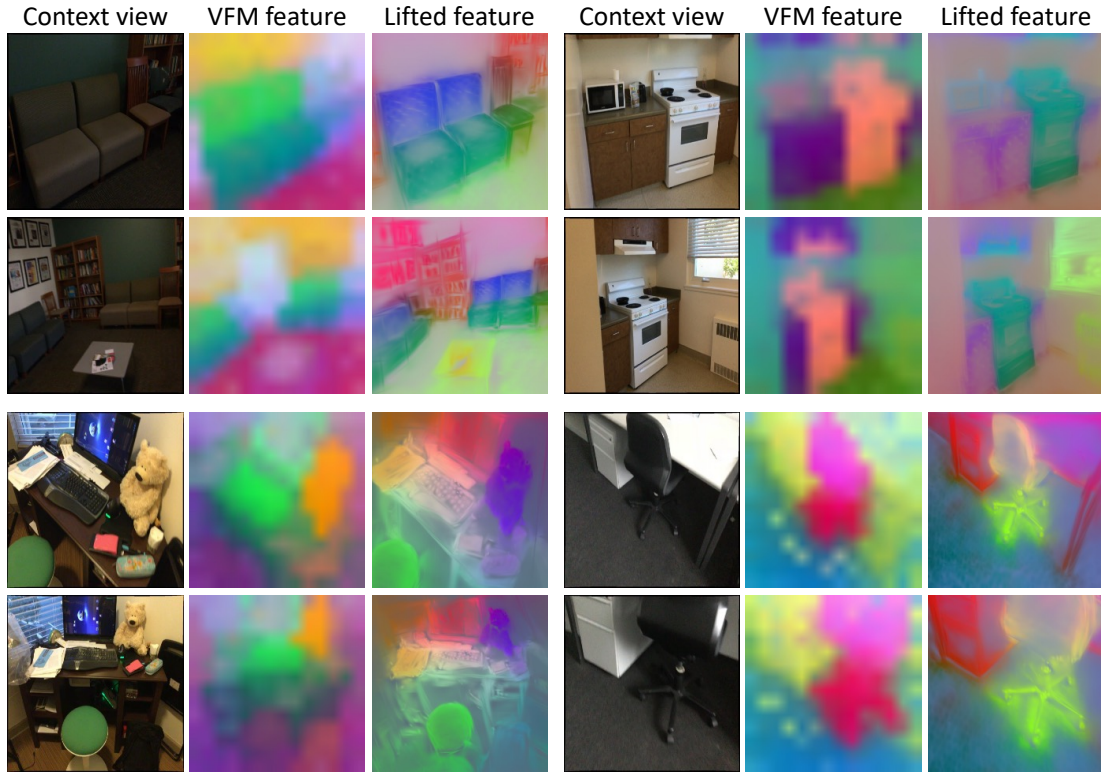


Figure 10. **Additional PCA visualization of multi-view features on ScanNet [9].** We visualize the PCA results of encoded multi-view features. Our method improves multi-view consistency compared to the original visual features.

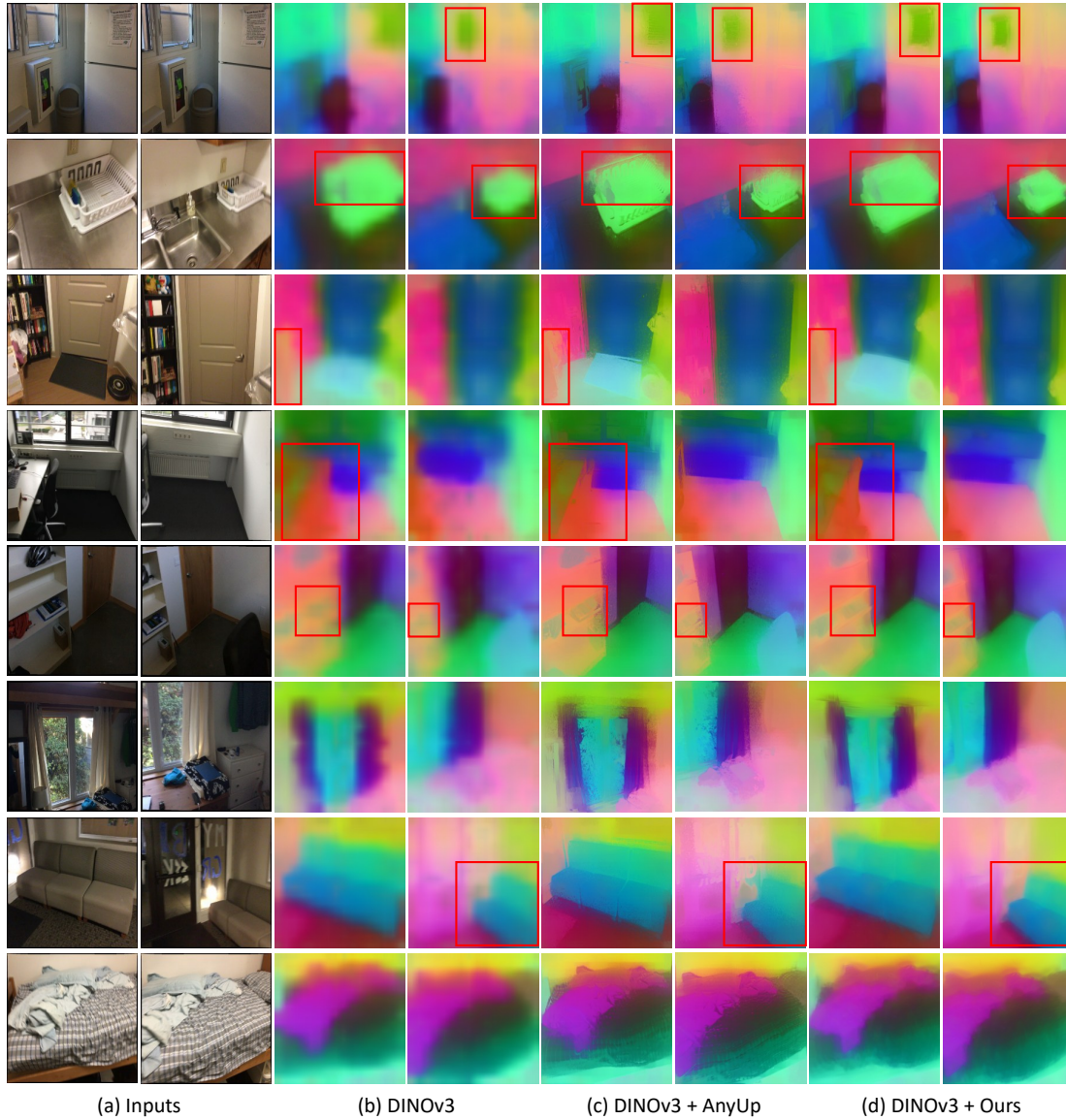


Figure 11. **Additional PCA visualization of upsampled feature on ScanNet [9].** We visualize the PCA results of the upsampled feature. Our model upsamples the features while maintaining the multi-view consistencies compared to other baselines.

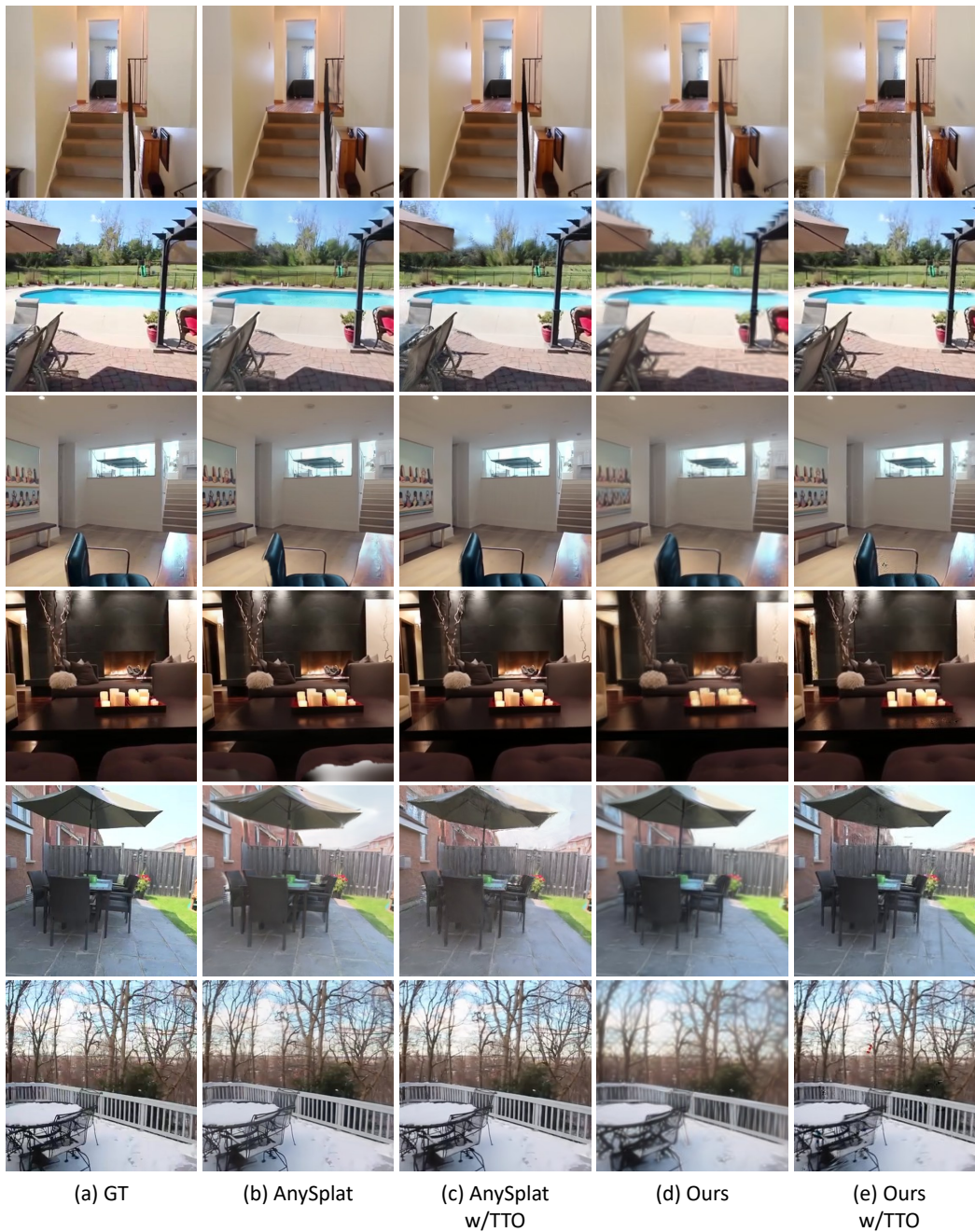


Figure 12. **Additional qualitative results of novel view synthesis on RealEstate10K [52].** We conduct a qualitative comparison for novel view synthesis with available multi-view images. Our method produces the highest quality rendering results, with or without test-time Gaussian optimization. TTO denotes that test-time optimization is applied to the Gaussians.

630
631
632
633
634
635
636
637
638
639
640
641
642
643
644
645
646
647
648
649
650
651
652
653
654
655
656
657
658
659
660
661
662
663
664
665
666
667
668
669
670
671
672
673
674
675
676
677
678
679
680
681
682
683
684
685
686

References

- [1] Eduardo Arnold, Jamie Wynn, Sara Vicente, Guillermo Garcia-Hernando, Aron Monszpart, Victor Prisacariu, Daniyar Turmukhambetov, and Eric Brachmann. Map-free visual relocalization: Metric pose relative to a single image. In *European Conference on Computer Vision*, pages 690–708. Springer, 2022. 2
- [2] Jimmy Lei Ba, Jamie Ryan Kiros, and Geoffrey E Hinton. Layer normalization. *arXiv preprint arXiv:1607.06450*, 2016. 2
- [3] Gilad Baruch, Zhuoyuan Chen, Afshin Dehghan, Tal Dimry, Yuri Feigin, Peter Fu, Thomas Gebauer, Brandon Joffe, Daniel Kurz, Arik Schwartz, et al. Arkitscenes: A diverse real-world dataset for 3d indoor scene understanding using mobile rgb-d data. *arXiv preprint arXiv:2111.08897*, 2021. 2
- [4] Neil Burgess. Spatial memory: how egocentric and allocentric combine. *Trends in cognitive sciences*, 10(12):551–557, 2006. 3
- [5] David Charatan, Sizhe Lester Li, Andrea Tagliasacchi, and Vincent Sitzmann. pixelsplat: 3d gaussian splats from image pairs for scalable generalizable 3d reconstruction. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 19457–19467, 2024. 3, 8
- [6] Yuedong Chen, Haofei Xu, Chuanxia Zheng, Bohan Zhuang, Marc Pollefeys, Andreas Geiger, Tat-Jen Cham, and Jianfei Cai. Mvsplat: Efficient 3d gaussian splatting from sparse multi-view images. In *European Conference on Computer Vision*, pages 370–386. Springer, 2024. 3, 8
- [7] Seokju Cho, Sunghwan Hong, Sangryul Jeon, Yunsung Lee, Kwanghoon Sohn, and Seungryong Kim. Cats: Cost aggregation transformers for visual correspondence. *Advances in Neural Information Processing Systems*, 34:9011–9023, 2021. 4
- [8] Seokju Cho, Sunghwan Hong, and Seungryong Kim. Cats++: Boosting cost aggregation with convolutions and transformers. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 45(6):7174–7194, 2022. 4
- [9] Angela Dai, Angel X Chang, Manolis Savva, Maciej Halber, Thomas Funkhouser, and Matthias Nießner. Scannet: Richly-annotated 3d reconstructions of indoor scenes. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 5828–5839, 2017. 3, 4, 10, 12, 13
- [10] Mohamed El Banani, Amit Raj, Kevis-Kokitsi Maninis, Abhishek Kar, Yuanzhen Li, Michael Rubinstein, Deqing Sun, Leonidas Guibas, Justin Johnson, and Varun Jampani. Probing the 3d awareness of visual foundation models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 21795–21806, 2024. 4
- [11] Zhiwen Fan, Jian Zhang, Wenyan Cong, Peihao Wang, Renjie Li, Kairun Wen, Shijie Zhou, Achuta Kadambi, Zhangyang Wang, Danfei Xu, et al. Large spatial model: End-to-end unposed images to semantic 3d. *Advances in neural information processing systems*, 37:40212–40229, 2024. 3, 4, 7
- [12] Sunghwan Hong and Seungryong Kim. Deep matching prior: Test-time optimization for dense correspondence. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 9907–9917, 2021. 4
- [13] Sunghwan Hong, Seokju Cho, Jisu Nam, Stephen Lin, and Seungryong Kim. Cost aggregation with 4d convolutional swin transformer for few-shot segmentation. In *European Conference on Computer Vision*, pages 108–126. Springer, 2022.
- [14] Sunghwan Hong, Jisu Nam, Seokju Cho, Susung Hong, Sangryul Jeon, Dongbo Min, and Seungryong Kim. Neural matching fields: Implicit representation of matching fields for visual correspondence. *Advances in Neural Information Processing Systems*, 35:13512–13526, 2022.
- [15] Sunghwan Hong, Seokju Cho, Seungryong Kim, and Stephen Lin. Unifying feature and cost aggregation with transformers for semantic and visual correspondence. *arXiv preprint arXiv:2403.11120*, 2024. 4
- [16] Sunghwan Hong, Jaewoo Jung, Heeseong Shin, Jisang Han, Jiaolong Yang, Chong Luo, and Seungryong Kim. Pf3plat: Pose-free feed-forward 3d gaussian splatting. *arXiv preprint arXiv:2410.22128*, 2024. 3, 8
- [17] Sunghwan Hong, Jaewoo Jung, Heeseong Shin, Jiaolong Yang, Seungryong Kim, and Chong Luo. Unifying correspondence pose and nerf for generalized pose-free novel view synthesis. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 20196–20206, 2024. 3, 8
- [18] Ranran Huang and Krystian Mikolajczyk. No pose at all: Self-supervised pose-free 3d gaussian splatting from sparse views. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 27947–27957, 2025. 3, 8
- [19] Lihan Jiang, Yucheng Mao, Linning Xu, Tao Lu, Kerui Ren, Yichen Jin, Xudong Xu, Mulin Yu, Jiangmiao Pang, Feng Zhao, et al. Anysplat: Feed-forward 3d gaussian splatting from unconstrained views. *arXiv preprint arXiv:2505.23716*, 2025. 1, 2, 3, 6, 7, 8, 11
- [20] Bernhard Kerbl, Georgios Kopanas, Thomas Leimkühler, and George Drettakis. 3d gaussian splatting for real-time radiance field rendering. *ACM Trans. Graph.*, 42(4):139–1, 2023. 1, 3
- [21] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014. 1, 3
- [22] Alexander Kirillov, Eric Mintun, Nikhila Ravi, Hanzi Mao, Chloe Rolland, Laura Gustafson, Tete Xiao, Spencer Whitehead, Alexander C Berg, Wan-Yen Lo, et al. Segment anything. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 4015–4026, 2023. 7
- [23] Hyunjoon Lee, Joonkyu Min, and Jaesik Park. Cf3: Compact and fast 3d feature fields. *arXiv preprint arXiv:2508.05254*, 2025. 4, 7
- [24] Vincent Leroy, Yohann Cabon, and Jérôme Revaud. Grounding image matching in 3d with mast3r. In *European Conference on Computer Vision*, pages 71–91. Springer, 2024. 3, 8
- [25] Boyi Li, Kilian Q Weinberger, Serge Belongie, Vladlen Koltun, and René Ranftl. Language-driven semantic segmentation. *arXiv preprint arXiv:2201.03546*, 2022. 4, 6

687
688
689
690
691
692
693
694
695
696
697
698
699
700
701
702
703
704
705
706
707
708
709
710
711
712
713
714
715
716
717
718
719
720
721
722
723
724
725
726
727
728
729
730
731
732
733
734
735
736
737
738
739
740
741
742
743
744

- 745 [26] Wanhua Li, Yujie Zhao, Minghan Qin, Yang Liu, Yuanhao
746 Cai, Chuang Gan, and Hanspeter Pfister. Langsplatv2: High-
747 dimensional 3d language gaussian splatting with 450+ fps.
748 *arXiv preprint arXiv:2507.07136*, 2025. 4
- 749 [27] Lu Ling, Yichen Sheng, Zhi Tu, Wentian Zhao, Cheng Xin,
750 Kun Wan, Lantao Yu, Qianyu Guo, Zixun Yu, Yawen Lu,
751 et al. DI3dv-10k: A large-scale scene dataset for deep
752 learning-based 3d vision. In *Proceedings of the IEEE/CVF*
753 *Conference on Computer Vision and Pattern Recognition*,
754 pages 22160–22169, 2024. 2
- 755 [28] Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee.
756 Visual instruction tuning. *Advances in neural information*
757 *processing systems*, 36:34892–34916, 2023. 7
- 758 [29] Haotian Liu, Chunyuan Li, Yuheng Li, and Yong Jae Lee.
759 Improved baselines with visual instruction tuning. In *Pro-*
760 *ceedings of the IEEE/CVF conference on computer vision*
761 *and pattern recognition*, pages 26296–26306, 2024. 7
- 762 [30] David G Lowe. Distinctive image features from scale-
763 invariant keypoints. *International journal of computer vi-*
764 *sion*, 60(2):91–110, 2004. 4
- 765 [31] Maxime Oquab, Timothée Darcet, Théo Moutakanni, Huy
766 Vo, Marc Szafraniec, Vasil Khalidov, Pierre Fernandez,
767 Daniel Haziza, Francisco Massa, Alaaeldin El-Nouby, et al.
768 Dinov2: Learning robust visual features without supervision.
769 *arXiv preprint arXiv:2304.07193*, 2023. 4, 5, 8
- 770 [32] William Peebles and Saining Xie. Scalable diffusion models
771 with transformers. In *Proceedings of the IEEE/CVF inter-*
772 *national conference on computer vision*, pages 4195–4205,
773 2023. 8
- 774 [33] Luigi Piccinelli, Yung-Hsu Yang, Christos Sakaridis, Mattia
775 Segu, Siyuan Li, Luc Van Gool, and Fisher Yu. Unidepth:
776 Universal monocular metric depth estimation. In *Proceed-*
777 *ings of the IEEE/CVF Conference on Computer Vision and*
778 *Pattern Recognition*, pages 10106–10116, 2024. 4
- 779 [34] Minghan Qin, Wanhua Li, Jiawei Zhou, Haoqian Wang, and
780 Hanspeter Pfister. Langsplat: 3d language gaussian splatting.
781 In *Proceedings of the IEEE/CVF Conference on Computer*
782 *Vision and Pattern Recognition*, pages 20051–20060, 2024.
783 4
- 784 [35] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya
785 Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry,
786 Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning
787 transferable visual models from natural language supervi-
788 sion. In *International conference on machine learning*, pages
789 8748–8763. PmLR, 2021. 4
- 790 [36] Mike Roberts, Jason Ramapuram, Anurag Ranjan, Atulit
791 Kumar, Miguel Angel Bautista, Nathan Paczan, Russ Webb,
792 and Joshua M Susskind. Hypersim: A photorealistic syn-
793 thetic dataset for holistic indoor scene understanding. In
794 *Proceedings of the IEEE/CVF international conference on*
795 *computer vision*, pages 10912–10922, 2021. 2
- 796 [37] Paul-Edouard Sarlin, Daniel DeTone, Tomasz Malisiewicz,
797 and Andrew Rabinovich. Superglue: Learning feature
798 matching with graph neural networks. In *Proceedings of*
799 *the IEEE/CVF conference on computer vision and pattern*
800 *recognition*, pages 4938–4947, 2020. 4
- 801 [38] Johannes L Schonberger and Jan-Michael Frahm. Structure-
802 from-motion revisited. In *Proceedings of the IEEE con-*
ference on computer vision and pattern recognition, pages
4104–4113, 2016. 3, 4
- [39] Oriane Siméoni, Huy V Vo, Maximilian Seitzer, Federico
Baldassarre, Maxime Oquab, Cijo Jose, Vasil Khalidov,
Marc Szafraniec, Seungeun Yi, Michaël Ramamonjisoa,
et al. Dinov3. *arXiv preprint arXiv:2508.10104*, 2025. 2,
4, 5
- [40] Brandon Smart, Chuanxia Zheng, Iro Laina, and Vic-
tor Adrian Prisacariu. Splatt3r: Zero-shot gaussian
splatting from uncalibrated image pairs. *arXiv preprint*
arXiv:2408.13912, 2024. 3, 8
- [41] Julian Straub, Thomas Whelan, Lingni Ma, Yufan Chen, Erik
Wijmans, Simon Green, Jakob J Engel, Raul Mur-Artal, Carl
Ren, Shobhit Verma, et al. The replica dataset: A digital
replica of indoor spaces. *arXiv preprint arXiv:1906.05797*,
2019. 3
- [42] Jianyuan Wang, Minghao Chen, Nikita Karaev, Andrea
Vedaldi, Christian Rupprecht, and David Novotny. Vggt: Vi-
sual geometry grounded transformer. In *Proceedings of the*
Computer Vision and Pattern Recognition Conference, pages
5294–5306, 2025. 1, 2, 3, 4, 8
- [43] Ruicheng Wang, Sicheng Xu, Cassie Dai, Jianfeng Xiang,
Yu Deng, Xin Tong, and Jiaolong Yang. Moge: Unlocking
accurate monocular geometry estimation for open-domain
images with optimal training supervision. In *Proceedings of*
the IEEE/CVF Conference on Computer Vision and Pattern
Recognition, pages 5261–5271, 2025. 1
- [44] Wenshan Wang, Delong Zhu, Xiangwei Wang, Yaoyu Hu,
Yuheng Qiu, Chen Wang, Yafei Hu, Ashish Kapoor, and Se-
bastian Scherer. Tartanair: A dataset to push the limits of
visual slam. In *2020 IEEE/RSJ International Conference*
on Intelligent Robots and Systems (IROS), pages 4909–4916.
IEEE, 2020. 2
- [45] Thomas Wimmer, Prune Truong, Marie-Julie Rakotosaona,
Michael Oechsle, Federico Tombari, Bernt Schiele, and
Jan Eric Lenssen. Anyup: Universal feature upsampling.
arXiv preprint arXiv:2510.12764, 2025. 5
- [46] Hongchi Xia, Yang Fu, Sifei Liu, and Xiaolong Wang. Rgb-
d objects in the wild: Scaling real-world 3d object learning
from rgb-d videos. In *Proceedings of the IEEE/CVF Con-*
ference on Computer Vision and Pattern Recognition, pages
22378–22389, 2024. 2
- [47] Yao Yao, Zixin Luo, Shiwei Li, Jingyang Zhang, Yufan Ren,
Lei Zhou, Tian Fang, and Long Quan. Blendedmvs: A large-
scale dataset for generalized multi-view stereo networks. In
Proceedings of the IEEE/CVF conference on computer vi-
sion and pattern recognition, pages 1790–1799, 2020. 2
- [48] Botao Ye, Sifei Liu, Haofei Xu, Xueting Li, Marc Pollefeys,
Ming-Hsuan Yang, and Songyou Peng. No pose, no problem:
Surprisingly simple 3d gaussian splats from sparse unposed
images. *arXiv preprint arXiv:2410.24207*, 2024. 1, 2, 3, 7, 8
- [49] Boyang Zheng, Nanye Ma, Shengbang Tong, and Saining
Xie. Diffusion transformers with representation autoen-
coders. *arXiv preprint arXiv:2510.11690*, 2025. 8
- [50] Chong Zhou, Chen Change Loy, and Bo Dai. Extract free
dense labels from clip. In *European conference on computer*
vision, pages 696–712. Springer, 2022. 4

- 860 [51] Shijie Zhou, Haoran Chang, Sicheng Jiang, Zhiwen Fan, Ze-
861 hao Zhu, Dejie Xu, Pradyumna Chari, Suyu You, Zhangyang
862 Wang, and Achuta Kadambi. Feature 3dgs: Supercharging
863 3d gaussian splatting to enable distilled feature fields. In *Pro-
864 ceedings of the IEEE/CVF Conference on Computer Vision
865 and Pattern Recognition*, pages 21676–21685, 2024. 3, 4, 6
- 866 [52] Tinghui Zhou, Richard Tucker, John Flynn, Graham Fyffe,
867 and Noah Snavely. Stereo magnification: Learning
868 view synthesis using multiplane images. *arXiv preprint
869 arXiv:1805.09817*, 2018. 1, 3, 8, 14