

GP-NeRF: Generalized Perception NeRF for Context-Aware 3D Scene Understanding

Supplementary Material

7.1. Implementation of our Transformers

We provide a simple and efficient pytorch pseudo-code to implement the attention operations in the field-aggregation, ray-aggregation transformer blocks in Alg. 1, 2. We use ray features to generate attention maps A_{field} and A_{ray} , and reuse them to construct semantic-embedding field as well as semantic features rendering.

Algorithm 1: Field-Aggregation Transformer

Input:

X_0 \rightarrow coordinate aligned features ($N_{rays}, N_{pts}, D_{rgb}$)
 X_{rgb} \rightarrow epipolar view Ray feats ($N_{rays}, N_{pts}, N_{views}, D_{rgb}$)
 X_{sem} \rightarrow epipolar view Sem feats ($N_{rays}, N_{pts}, N_{views}, D_{sem}$)

Δd \rightarrow relative directions ($N_{rays}, N_{pts}, N_{views}, 3$)

Network: $f_Q, f_K, f_V, f_P, f_A, f_{rgb}$ \rightarrow MLP layers

Output: $S_{rgb}^{3D}, S_{sem}^{3D}$

Forward: **Red** for semantic-embedding field aggregation

- 1 $Q = f_Q(X_0), K = f_K(X_{rgb}), V = f_V(X_{rgb})$
 - 2 $P_{field} = f_P(\Delta d)$
 - 3 $A_{field} = K - Q[:, :, None, :] + P$
 - 4 $A_{field} = \text{softmax}(A, \text{dim} = -2)$
 - 5 $A'_{field} = A_{field} \cdot \text{repeat_interleave}(4)$
 - 6 $P'_{field} = P \cdot \text{repeat_interleave}(4)$
 - 7 $S_{rgb}^{3D} = ((V + P) \cdot A) \cdot \text{sum}(\text{dim} = 2)$
 - 8 $S_{rgb}^{3D} = f_{rgb}(S_{rgb}^{3D})$
 - 9 $S_{sem}^{3D} = ((X_{sem} + P'_{field}) \cdot A'_{field}) \cdot \text{sum}(\text{dim} = 2)$
-

7.2. Reconstruction results in instance setting

During the novel view instance segmentation task, we evaluate our reconstruction results and compare them with SOTA method DM-NeRF[38]. As shown in Table 5, our approach surpasses DM-NeRF in terms of SSIM and LPIPS metrics by 0.02% and 0.065%, respectively. It demonstrates that contextual information from semantic features can enhance the geometry reconstruction in our jointly optimized field and rendering framework.

7.3. Few-step Finetuning Comparison

Tab. 6 presents a comparison of different models, showcasing their mIoU and finetuning times on the ScanNet [9] dataset, along with the AP75 metric in Replica [33]. We observe that by finetuning with limited time, our model is able to achieve a better perception accuracy than a well-trained per-scene optimized method, such as 3.45% in

Algorithm 2: Ray-Aggregation Transformer

Input:

X_0^{rgb} \rightarrow coordinate aligned rgb features ($N_{rays}, N_{pts}, D_{rgb}$)

X_0^{sem} \rightarrow coordinate aligned sem features ($N_{rays}, N_{pts}, D_{sem}$)

x \rightarrow point coordinates (after PE) ($N_{rays}, N_{pts}, D_{rgb}$)

d \rightarrow target view direction (after PE) ($N_{rays}, N_{pts}, D_{rgb}$)

Network: $f_Q, f_K, f_V, f_P, f_A, f_{rgb}, f_{sem}$ \rightarrow MLP layers

Output: $S_{rgb}^{2D}, S_{sem}^{2D}$

Forward: **Red** for semantic-embedding field aggregation

- 1 $X_0^{rgb} = f_P(\text{concat}(X_0^{rgb}, d, x))$
 - 2 $Q = f_Q(X_0^{rgb}), K = f_K(X_0^{rgb}), V = f_V(X_0^{rgb})$
 - 3 $A_{ray} = \text{matmul}(Q, K^T) / \sqrt{D}$
 - 4 $A_{ray} = \text{softmax}(A_{ray}, \text{dim} = -1)$
 - 5 $A'_{ray} = A_{ray} \cdot \text{repeat_interleave}(4)$
 - 6 $S_{rgb}^{2D} = \text{matmul}(V, A_{ray})$
 - 7 $S_{rgb}^{2D} = f_{rgb}(S_{rgb}^{2D})$
 - 8 $S_{sem}^{2D} = \text{matmul}(X_0^{sem}, A'_{ray})$
-

Table 5. Quantitative results of reconstruction task in Replica[33] during instance segmentation setting.

| Scene | DM-NeRF | | | Ours | | |
|----------|-----------------|-----------------|--------------------|--|--|---|
| | PSNR \uparrow | SSIM \uparrow | LPIPS \downarrow | PSNR \uparrow | SSIM \uparrow | LPIPS \downarrow |
| Office.0 | 40.66 | 0.972 | 0.07 | 39.25 | 0.984 | 0.027 |
| Office.2 | 36.98 | 0.964 | 0.115 | 36.01 | 0.974 | 0.042 |
| Office.3 | 35.34 | 0.955 | 0.078 | 36.02 | 0.982 | 0.027 |
| Office.4 | 32.95 | 0.921 | 0.172 | 32.75 | 0.94 | 0.085 |
| Room.0 | 34.97 | 0.94 | 0.127 | 34.29 | 0.972 | 0.049 |
| Room.1 | 34.72 | 0.931 | 0.134 | 36.45 | 0.968 | 0.043 |
| Room.2 | 37.32 | 0.963 | 0.115 | 34.75 | 0.960 | 0.085 |
| Average | 36.13 | 0.949 | 0.116 | 35.64 ^{0.49\downarrow} | 0.969 ^{0.02\uparrow} | 0.051 ^{0.065\downarrow} |

mIoU with Semantic-NeRF [52] and 3.7% in AP75 with DM-NeRF [38]. Specifically, we observe that our method surpasses Semantic-Ray, requiring only half as many finetuning steps, and improves the mIoU by 0.74%, which further demonstrates that our semantic embedding field with more discrimination successfully improves the generalized ability.

We further evaluate the above experiments in instance segmentation setting, shown in the bottom column in Tab. 6. Not surprising, compared with SOTA method DM-NeRF[38], we achieve better performance with only 4k training steps, by 3.7% in AP75.

7.4. Additional Visualization Results

Fig. 9 shows the additional qualitative results of semantic prediction and reconstruction.

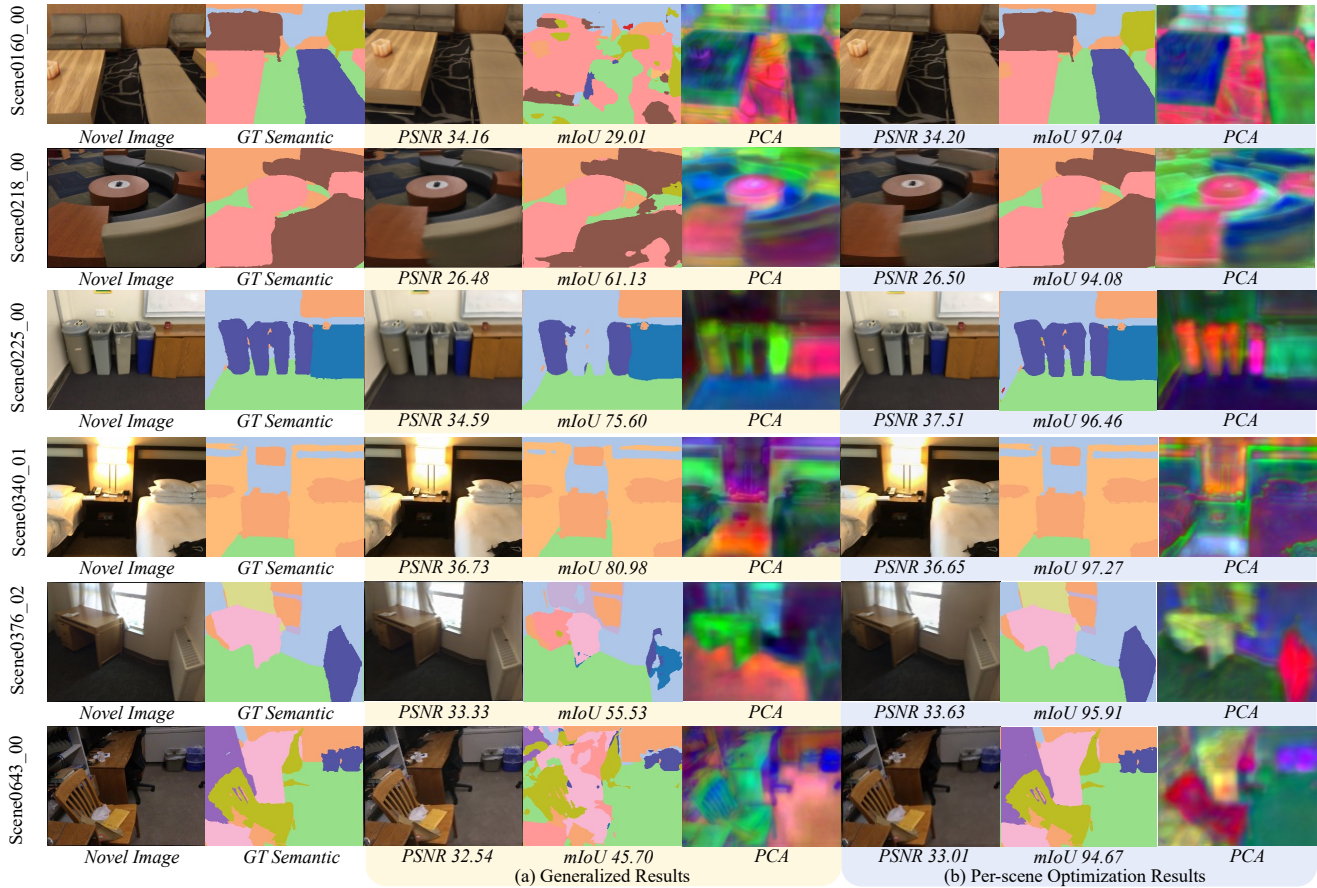


Figure 9. The visualization results in ScanNet[9]. Here we visualize the semantic as well as reconstruction results in both generalized and finetuning settings.

| Method | Train Step | Train Time | mIoU/AP75 |
|----------------------|------------|------------|------------------------|
| Semantic-NeRF [52] | 50k | ~2h | 89.33 |
| MVSNeRF w/s-Ft | 5k | ~20min | 52.02 |
| NeuRay [25] w/s-Ft | 5k | ~32min | 79.23 |
| Semantic-Ray [20]-Ft | 5k | ~20min | 92.04 |
| Ours-Ft | 2.5k | ~20min | 92.78 ^{0.74↑} |
| DM-NeRF | 200k | ~2h | 81.03 |
| Ours-Ft | 4k | ~30min | 84.73 ^{3.7↑} |

Table 6. mIoU and training steps/time on ScanNet [9]. "w/ s" means adding a semantic head on the baseline architectures.