

OpenGaussian++ - Refining OpenGaussian: The 3D Point-Level Understanding Framework with Open Vocabulary

Anonymous CVPR submission

Paper ID ****

Abstract

001 *OpenGaussian is limited by inconsistent per-view SAM*
002 *segmentations and a fragile, single-view semantic association.* We independently evaluate two solutions: multi-view
003 mask refinement and a more effective multi-view feature fusion framework. On ScanNet, our visibility-weighted fusion
004 model improves average mIoU from 0.299 to 0.318, producing cleaner object boundaries. This demonstrates that over-
005 coming single-view dependencies with a robust, visibility-aware multi-view consensus is critical for reliable 3D
006 semantic understanding.

1. Introduction

012 3D Gaussian Splatting (3DGS) [1] has become the leading
013 neural rendering method for its speed and quality. Open-
014 Gaussian [5] extends 3DGS for view-independent 3D point-
015 level understanding using cluster-based feature embedding
016 with SAM masks and CLIP features.

017 We present three modifications to enhance OpenGaus-
018 sian's semantic consistency and training efficiency: (1)
019 Multi-view SAM mask refinement enforcing consistency
020 across overlapping views, (2) Multi-view CLIP feature
021 fusion that averages features across views where each
022 Gaussian is visible, and (3) Visibility-weighted fusion using
023 camera angles to weight view contributions.

2. Related Works

2.1. Neural Rendering

026 3DGS achieves superior performance over Neural Radiance
027 Fields (NeRFs) [2] through rasterization-based rendering
028 rather than slow volume rendering.

2.2. 3D Point-Level Understanding

030 LangSplat [3] and LEGaussian [4] assign compressed
031 CLIP features directly to Gaussians, leading to information
032 loss and 2D-to-3D mapping challenges. OpenGaussian's

cluster-based approach with mask-level feature consistency
provides more stable learning signals but trades point-wise
granularity for distinct segmentation.

3. Method

OpenGaussian assigns semantic understanding to 3D Gaus-
sians through three steps: (1) learning 6-dimensional fea-
ture maps for each Gaussian using SAM mask consistency,
(2) clustering Gaussians by feature similarity, and (3) as-
signing CLIP values to clusters. We propose three modifi-
cations to enhance this framework.

3.0.1 Modification 1: Consistent SAM-masks

- **Problem:** SAM generates per-frame masks indepen-
dently, causing geometric and semantic inconsistencies
for the same object across different viewpoints. Object
boundaries are imprecise, and many pixels remain unclas-
sified.
- **Proposed solution:** We utilize multi-view consistency of
Gaussian splats to refine SAM masks through a two-stage
approach that synchronizes object IDs across views and
expands segment coverage:

1. Cross-view ID synchronization for consistent ob- ject identification:

We address the problem that identical objects receive
different segment IDs across camera views. For each
Gaussian splat g_i with opacity $\alpha_i > \tau_{\text{opacity}}$ (we use
 $\tau_{\text{opacity}} = 0.8$), we render the individual splat in each
camera view to obtain its 2D footprint and weight dis-
tribution. We apply depth-based visibility testing: a
splat is considered visible if its projection lies in im-
age bounds and $|d_{\text{projected}} - d_{\text{rendered}}| < \tau_{\text{depth}}$ where
 $d_{\text{projected}}$ is the Euclidean distance from the image
plane to the splat and d_{rendered} is the depth from the
rendered depth map.

For each visible footprint, we compute the most dom-
inant segment ID by aggregating pixel IDs inside the
footprint weighted by the opacity function of the pro-

069 jected splat: $ID_{\text{dominant}} = \arg \max_k \sum_{p \in \text{footprint}} w_p \cdot$
 070 $1[\text{SAM}(p) = k]$, where w_p is the value of the opacity
 071 function evaluated at pixel p . Finally, we constrain ob-
 072 ject IDs across multiple views to the same new global
 073 identifier by applying it to all pixels belonging to the
 074 object mask.

075 2. Mask expansion using weighted Gaussian foot- 076 prints:

077 To reduce regions classified as "void" and enhance the
 078 boundary quality between segments, we expand seg-
 079 mented regions with projected Gaussian splats. For
 080 each pixel in the image, we maintain a look-up table
 081 to accumulate values for different IDs. For each splat,
 082 we compute its underlying object ID and collect votes
 083 for the winning ID across multiple views. For those
 084 viewpoints where the ID of the projected splat matches
 085 the winning ID, we increase accumulated weights for
 086 this ID in all pixels belonging to the base segment. We
 087 then apply additional splat expansion for the splat pix-
 088 els that don't belong to the object segment -
 089 $\text{extension_mask} = (w_p > 0) \wedge \neg(\text{segment_mask}_k)$,
 090 updating $\mathbf{V}_{p,k}^{(c)} += w_p$. For the viewpoints where the
 091 computed ID of the projected splat doesn't match the
 092 winning ID, we introduce the multi-view constraint.
 093 Instead of accumulating the weights for the ID that be-
 094 longs to the splat in the viewpoint over the base object
 095 segment, we replace its ID with the winning ID and then
 096 accumulate the weights in the same way as above.
 097 By doing so, we indicate that the splat resides on the
 098 edge between different segments and can therefore be
 099 classified differently when rendered in different view-
 100 points. Therefore, we aim to invert the direction of
 101 expansion of the misclassified splat projection accord-
 102 ing to the direction seen in the majority of other view-
 103 points.

104 Due to efficient mask usage that is highly parallelizable
 105 on the GPU, the additional SAM masks refinement pipeline
 106 introduces relatively small overhead compared to the over-
 107 all time needed for OpenGaussian to process a scene.

108 3.1. Two-Level Code Book Discretization

109 In this step, the goal is to cluster Gaussians with similar
 110 features. Initially, the 6-dimensional feature embedding is
 111 concatenated with the x, y, and z coordinates to create a 9-
 112 dimensional feature embedding. The underlying rationale
 113 for this design choice is that Gaussians in proximity within
 114 Euclidean space are likely to exhibit semantic similarity.
 115 The codebook is then initialized by randomly selecting fea-
 116 ture embeddings from $k = 64$ distinct Gaussians, forming
 117 $\mathbf{C} \in \mathbb{R}^{k \times 9}$. Subsequently, each Gaussian is assigned to
 118 a cluster based on Euclidean distance. Following cluster
 119 assignment, the Gaussian-specific features undergo further
 120 refinement with a modified optimization objective. Specif-

ically, the optimization aims to minimize the intra-cluster
 variance by drawing feature embeddings within each clus-
 ter toward their respective cluster centroid. This is achieved
 by minimizing the following loss:

$$\mathcal{L}_p = \|M_p - M_c\|_1 \quad (1) \quad 125$$

In this case, the maps M_p and M_c do not contain the x,y,
 and z coordinates of the Gaussians. They are strictly used
 for the original clustering and are dropped once assigned.
 There now exists a one-level code book that discretizes the
 Gaussians. Now, each cluster in the code book is split up
 into $k = 10$ clusters, and the same process is repeated. For
 this fine level of clustering, the x,y, and z coordinate for
 each Gaussian are not used for the clustering.

134 3.2. Instance-Level 2D-3D Association without 135 Depth Test

The culmination of our 3D instance learning pipeline is the
 endowment of each geometric cluster with rich semantic
 meaning. We achieve this by associating each 3D cluster
 with a high-dimensional language feature vector from a pre-
 trained Vision-Language Model (VLM), specifically CLIP.
 This critical step bridges the gap between abstract 3D ge-
 ometry and open-vocabulary language understanding. The
 process requires establishing a robust and accurate mapping
 between the 3D world, as represented by our instance clus-
 ters, and the 2D observations from which semantic features
 are derived. In the following sections, we describe the base-
 line single-view association method and present our more
 robust multi-view fusion framework.

149 3.2.1 Original Method: Single-View Association

The baseline method employs a fragile "winner-takes-all"
 strategy. It associates each 3D cluster \mathcal{C}_i with the CLIP
 feature from a single best-matching 2D SAM mask found
 across all rendered views, discarding all other data. This
 reliance on a single observation makes the process highly
 susceptible to errors from occlusion or poor viewing angles.
 By ignoring the rich, consensus-building information avail-
 able in other views, this method leads to unstable and noisy
 semantic assignments.

159 3.2.2 Modification 2 & 3: Aggregation of CLIP- 160 features

To overcome the inherent fragility of the single-view ap-
 proach, we introduce a multi-view fusion framework de-
 signed for robustness and stability. Our paradigm shifts
 from finding the single best view to building a reliable se-
 mantic consensus from multiple high-quality observations.
 The framework consists of a shared initial association stage
 followed by two distinct feature aggregation strategies.

168 **View-wise Association and Filtering.** The foundational
 169 step of our framework is to first evaluate the quality of asso-
 170 ciation for a 3D cluster \mathcal{C}_i in *every* visible view v . For each
 171 view, we find the best possible match among the local 2D
 172 SAM masks $M_{v,j}$. The quality of this match is quantified
 173 by a holistic joint score $S_{i,v,j}$ that considers both geometric
 174 and feature-space alignment:

$$S_{i,v,j} = \text{IoU}(\pi(\mathcal{C}_i, v), M_{v,j}) \times \underbrace{(1 - d_{\text{L1}}(\mathbf{f}_{\mathcal{C}_i, v}, \mathbf{f}_{M_{v,j}}))}_{\text{Feature Similarity}}, \quad (2)$$

175 where:

- $\pi(\mathcal{C}_i, v)$ is the 2D silhouette obtained by rendering the 3D cluster \mathcal{C}_i from the perspective of camera v .
- $\mathbf{f}_{\mathcal{C}_i, v}$ and $\mathbf{f}_{M_{v,j}}$ are the mean instance features of the rendered cluster and the 2D mask, respectively.
- The IoU term ensures strong geometric alignment, while the feature similarity term (scaled to $[0, 1]$) penalizes feature-space dissimilarity. A perfect feature match yields a similarity of 1, while diverging features drive it towards 0.

186 For each view v , we determine the best local match by finding
 187 $S_{i,v} = \max_j S_{i,v,j}$.

188 This process yields a set of candidate associations, one
 189 for each view. As a crucial quality control step, we filter
 190 this set by retaining only the associations whose score
 191 $S_{i,v}$ exceeds a confidence threshold τ (e.g., $\tau = 0.2$). This
 192 filtering discards low-confidence matches that could arise
 193 from spurious overlaps or feature mismatches. The result
 194 of this stage is a curated set of high-quality views \mathcal{V}_i^* for
 195 each 3D cluster, along with their corresponding best-match
 196 CLIP features $\{\mathbf{f}_{i,v}^{\text{CLIP}}\}_{v \in \mathcal{V}_i^*}$.

197 **Feature Aggregation Strategies.** Having curated a set of
 198 high-confidence 2D semantic observations, we aggregate
 199 them to produce a single, robust feature vector for the 3D
 200 cluster. We propose two strategies for this final step.

201 • **Modification 2: Democratic (Unweighted) Aggrega-
 202 tion.** This strategy is founded on the principle that every
 203 validated viewpoint should have an equal say in determin-
 204 ing the final semantic identity. The definitive feature \mathbf{F}_i
 205 is computed as the arithmetic mean of the CLIP features
 206 from all views in the high-quality set \mathcal{V}_i^* :

$$\mathbf{F}_i = \frac{1}{|\mathcal{V}_i^*|} \sum_{v \in \mathcal{V}_i^*} \mathbf{f}_{i,v}^{\text{CLIP}}. \quad (3)$$

208 By averaging, this method effectively mitigates the im-
 209 pact of random, uncorrelated noise present in individual
 210 views. The normalization by $|\mathcal{V}_i^*|$ ensures that the mag-
 211 nitude of the final feature vector is not dependent on the
 212 number of available high-quality views, promoting stabil-
 213 ity across different instances. Its principal strength lies in
 214 its simplicity and effectiveness at noise reduction.

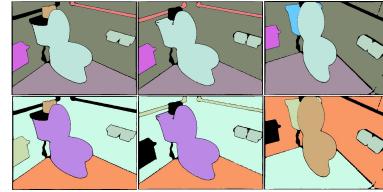


Figure 1. SAM mask consistency before (bottom) and after (top) our cross-view ID synchronization.

- **Modification 3: Confidence-and-Visibility-Weighted Aggregation.** This more refined strategy is founded on the intuition that not all "good" views are equally "great". A view where an object is large, central, and unoccluded should contribute more to the final identity than a view where it is small, peripheral, or partially occluded. We therefore introduce a weighting scheme that prioritizes more salient and confident observations. The weight $w_{i,v}$ for each view $v \in \mathcal{V}_i^*$ is defined as:

$$w_{i,v} = \underbrace{S_{i,v}}_{\text{Match Confidence}} \times \underbrace{V_{i,v}}_{\text{Visibility/Salience}}, \quad (4)$$

where $S_{i,v}$ is the association score representing the *confidence* of the 2D-3D match, and $V_{i,v}$ is a *visibility* score. As implemented in our code, $V_{i,v}$ is operationally defined as the total pixel count of the rendered silhouette, naturally giving more weight to views where the object appears larger. The final feature \mathbf{F}_i is the weighted average:

$$\mathbf{F}_i = \frac{\sum_{v \in \mathcal{V}_i^*} w_{i,v} \cdot \mathbf{f}_{i,v}^{\text{CLIP}}}{\sum_{v \in \mathcal{V}_i^*} w_{i,v}}. \quad (5)$$

This ensures that the final semantic embedding is anchored by the most informative and unambiguous perspectives available in the dataset, leading to a more accurate and nuanced representation that reflects the quality of the underlying observations.

Both modifications provide a significant improvement in robustness over the baseline, with the weighted approach offering a more sophisticated and potentially more accurate fusion of multi-view semantic information.

4. Results

We evaluate our methods on 9 challenging scenes from the ScanNet dataset presented in Table 1.

Regarding *modification 1*, due to computational constraints, Table 1 presents results from a preliminary implementation using point approximations for Gaussian projections. These results fall within the margin of error relative to baseline. Evaluation of our current implementation on scene0062_00 yielded mIoU: 0.3217 and mAcc.: 0.5224, showing no significant improvement over baseline.

Scene	Baseline		M1 (Ours)		M2 (Ours)		M3 (Ours)	
	mIoU	mAcc	mIoU	mAcc	mIoU	mAcc	mIoU	mAcc
scene0000_00	0.318	0.458	0.3285	0.4855	0.3192	0.4609	0.3304	0.4941
scene0062_00	0.345	0.555	0.3972	0.6591	0.3300	0.5682	0.3667	0.5767
scene0070_00	0.228	0.339	0.2163	0.3110	0.2367	0.3343	0.2131	0.3033
scene0097_00	0.361	0.521	0.3792	0.5496	0.3041	0.4636	0.3593	0.5453
scene0140_00	0.227	0.401	0.2497	0.4336	0.2241	0.4004	0.2562	0.4447
scene0200_00	0.402	0.550	0.3879	0.5235	0.4164	0.5595	0.4384	0.5802
scene0347_00	0.291	0.406	0.3218	0.4437	0.2948	0.4118	0.2934	0.4213
scene0590_00	0.298	0.432	0.2991	0.4048	0.3028	0.4185	0.3332	0.4471
scene0645_00	0.299	0.488	0.2919	0.4703	0.3140	0.4910	0.2851	0.4712
Average	0.299	0.451	0.312	0.466	0.305	0.456	0.318	0.477

Table 1. ScanNet per-scene semantic segmentation results.

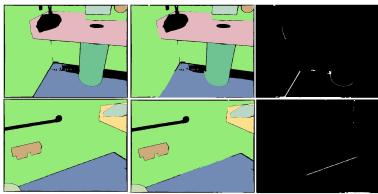


Figure 2. Mask expansion stages: original SAM mask, after ID synchronization, after mask expansion, and difference map.

Qualitative analysis demonstrates enhanced mask consistency across viewpoints in the initial stage, as shown in Figure 1. Toilet objects and structural elements maintain consistent IDs across different viewing angles after cross-view synchronization.

Figure 2 shows the secondary refinement stage reduces unclassified regions and improves boundary definition. However, segmentation quality correlates inversely with splat sampling density, as discussed in Section 5.

The baseline’s single-view approach is prone to artifacts like semantic bleeding (Fig. 4b). In contrast, our multi-view fusion framework improves robustness. While unweighted averaging (M2) mitigates noise through consensus (Fig. 4c), our visibility-weighted strategy (M3) achieves optimal performance. By prioritizing high-quality, unoccluded views (Fig. 3), M3 yields the most accurate segmentation with the sharpest boundaries (Fig. 4d). This validates our hypothesis that a weighted multi-view consensus is critical for robust 3D semantic understanding.

5. Conclusion

Several technical challenges merit further investigation. The initial challenge concerns the observed phenomenon wherein larger Gaussian splats, particularly those rep-

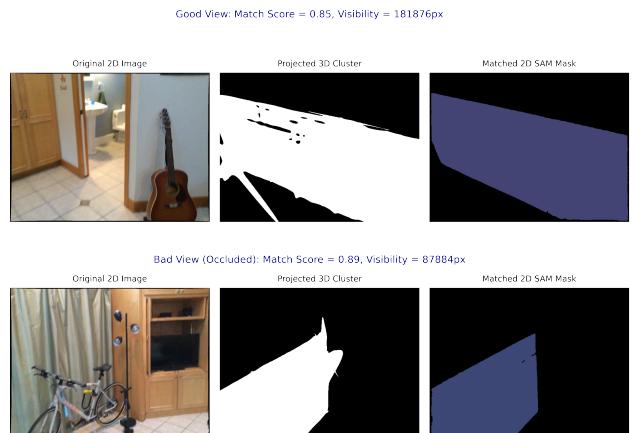


Figure 3. Mechanism of the Visibility-Weighted Fusion. Our method determines each view’s contribution to the final semantic feature by calculating a weight, which is a product of its Match Score and Visibility. This figure contrasts two views of the same 3D cluster. (Top) The "Good View" combines a strong Match Score (0.85) with high Visibility (181,876px), resulting in a large overall weight. (Bottom) The partially occluded "Bad View" has a high Match Score (0.89), but its Visibility is drastically reduced by occlusion (87,884px).

resenting substantial structural elements such as walls, demonstrate dominant characteristics that effectively occlude smaller splats corresponding to less volumetrically significant objects. This occlusion mechanism results in suboptimal representation of smaller entities within the refined segmentation masks. One possible solution could be a splat sampling strategy that is less penalizing for smaller objects. The second challenge pertains to the cross-view constraint methodology illustrated in Figure 1. While this approach demonstrates efficacy in enhancing segmentation

274
275
276
277
278
279
280
281
282
283

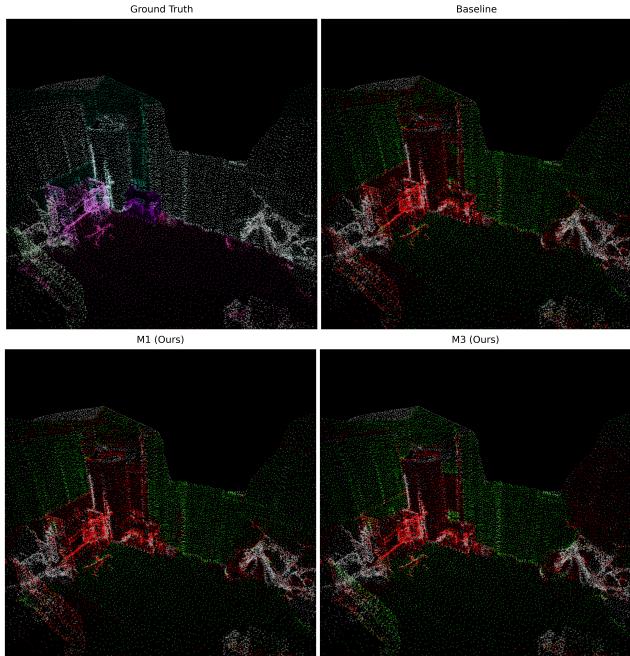


Figure 4. Qualitative Comparison of Semantic Segmentation. This figure illustrates the semantic segmentation of a complex scene. (a) Ground Truth: Shows distinct, clean segments for different objects. (b) Baseline: The "winner-takes-all" approach suffers from significant noise and semantic bleeding, misclassifying the floor (green) as part of the central object (red). (c) M2 (Ours): Our unweighted averaging method significantly cleans up the segmentation, improving the distinction between objects. (d) M3 (Ours): Our visibility-weighted averaging method provides the most accurate result, yielding sharper boundaries and a segmentation that most closely resembles the ground truth.

284 consistency across multiple perspectives, implementation
 285 difficulties were encountered when scaling the algorithm
 286 to accommodate both comprehensive scene-wide application
 287 and diminutive object instances. These implementation
 288 constraints may potentially be mitigated through para-
 289 metric optimization of point sampling density and opacity
 290 threshold configurations. An additional algorithmic con-
 291 straint worth consideration involves the selective exclusion
 292 of splats exhibiting elevated variance ratios along their prin-
 293 cipal axes, as these elements frequently generate projections
 294 spanning multiple segmentation boundaries. These identi-
 295 fied limitations constitute priority areas for subsequent re-
 296 search iterations.

297 References

- 298 [1] Bernhard Kerbl, Georgios Kopanas, Thomas Leimkühler, and
 299 George Drettakis. 3d gaussian splatting for real-time radiance
 300 field rendering. *ACM Transactions on Graphics*, 42(4), 2023.
 301 1
- 302 [2] Ben Mildenhall, Pratul P. Srinivasan, Matthew Tancik,
 303 Jonathan T. Barron, Ravi Ramamoorthi, and Ren Ng. Nerf:

- | | |
|---|--|
| Representing scenes as neural radiance fields for view synthesis, 2020. 1
[3] Minghan Qin, Wanhua Li, Jiawei Zhou, Haoqian Wang, and Hanspeter Pfister. Langsplat: 3d language gaussian splatting. <i>arXiv preprint arXiv:2312.16084</i> , 2023. 1
[4] Jin-Chuan Shi, Miao Wang, Hao-Bin Duan, and Shao-Hua Guan. Language embedded 3d gaussians for open-vocabulary scene understanding. <i>arXiv preprint arXiv:2311.18482</i> , 2023. 1
[5] Yanmin Wu, Jiarui Meng, Haijie Li, Chenming Wu, Yahao Shi, Xinhua Cheng, Chen Zhao, Haocheng Feng, Errui Ding, Jingdong Wang, and Jian Zhang. Opengaussian: Towards point-level 3d gaussian-based open vocabulary understanding, 2024. 1 | 304
305
306
307
308
309
310
311
312
313
314
315
316
317 |
|---|--|