

# Qwen-3D: A Generalist 3D Vision-Language Model for Spatial Understanding

Anonymous CVPR submission

Paper ID 00027

## Abstract

001 *Recent Large Multimodal Models (LMMs) achieve impres-*  
002 *sive performance on images and short videos, but long*  
003 *video reasoning remains computationally expensive and*  
004 *temporally inconsistent due to frame-level tokenization and*  
005 *limited context windows. We observe that 3D geometry pro-*  
006 *vides a natural compression mechanism for multi-view vi-*  
007 *sual streams. Geometric signals such as depth and camera*  
008 *pose allow RGB frames to be fused into persistent world-*  
009 *aligned representations, enabling efficient reasoning over*  
010 *space and time. Motivated by this insight, we introduce*  
011 *Qwen-3D, a geometry-aware LMM that leverages multi-*  
012 *view geometric signals to compress visual tokens within the*  
013 *Qwen backbone, enabling efficient processing of long and*  
014 *highly redundant video sequences. By modifying visual to-*  
015 *kens with 3D Rotary Positional Embedding, Qwen-3D per-*  
016 *forms attention in the world space rather than over inde-*  
017 *pendent image frames, enabling efficient cross-view reason-*  
018 *ing. Qwen-3D further integrates a query-based segmenta-*  
019 *tion decoder that grounds language tokens directly in vi-*  
020 *sual space, allowing unified reasoning across referential*  
021 *grounding, instance segmentation, and visual question an-*  
022 *swering for both images and videos. Across a broad suite of*  
023 *benchmarks, Qwen-3D outperforms large-scale proprietary*  
024 *2D models and existing 3D LMM approaches. Finally, we*  
025 *show that joint training on both 2D and 3D data preserves*  
026 *strong 2D vision-language capabilities while substantially*  
027 *improving 3D reasoning.*

## 028 1. Introduction

029 Current Vision-Language Models (VLMs) perform well on  
030 images and short video clips, but struggle with long multi-  
031 view video streams. Processing long sequences is computa-  
032 tionally expensive due to the quadratic cost of attention, and  
033 limited context windows prevent consistent spatio-temporal  
034 and 3D reasoning across frames. Multi-view 3D geometry,  
035 in the form of depth and camera poses, offers a more prin-  
036 ciple alternative, allowing video frames to be mapped into  
037 a shared 3D coordinate system. This enables compression

of multi-view streams into persistent scene representations 038  
where temporally distant frames correspond to nearby 3D 039  
locations. 040

Existing approaches to integrating 3D spatial awareness 041  
into VLMs generally follow one of two distinct paradigms. 042  
One line of work introduces 3D point clouds as auxiliary in- 043  
puts to the model [11, 19, 20], either alongside or in place of 044  
multi-view images. While these approaches expose explicit 045  
geometric structure, they treat point clouds as a modality 046  
separate from the visual tokens, overlooking the fact that 047  
point cloud features are inherently aligned with image fea- 048  
tures with corresponding depth. The second line of work 049  
integrates geometry directly into the visual token represen- 050  
tation. Methods such as LLaVA-3D[54] and Video-3D- 051  
LLM[53] modify positional encodings such that multi-view 052  
image tokens are embedded according to their 3D world co- 053  
ordinates rather than their 2D image-plane positions. This 054  
approach allows the model to reason over multi-view obser- 055  
vations in a shared spatial coordinate system while main- 056  
taining the strong visual representations learned by large 057  
VLM backbones. 058

Despite these advances, existing 3D large multimodal 059  
models (LMMs) still lag behind specialist 3D vision sys- 060  
tems. Dedicated models trained for tasks such as detection, 061  
segmentation, or grounding continue to outperform general- 062  
purpose 3D LMMs [22, 23, 55]. Moreover, most current 063  
3D LMMs do not attempt 3D object detection on standard 064  
benchmarks such as ScanNet [14, 37]. The only exception, 065  
Grounded-3D-LLM [12], achieves less than half the perfor- 066  
mance of state-of-the-art 3D detectors. These results high- 067  
light a key tradeoff: while general-purpose 3D LMM’s offer 068  
a broad range of reasoning capabilities, they currently lack 069  
the specialized grounding mechanisms of dedicated 3D sys- 070  
tems, preventing them from matching—or in many cases, 071  
even addressing—the high-level performance required for 072  
core 3D perception. 073

To bridge these performance gaps, we introduce Qwen- 074  
3D, a geometry-aware 3D LMM that extends the Qwen 075  
family of models [5] with explicit mechanisms for multi- 076  
view reasoning and object grounding. Built on the strong 077  
2D foundation of Qwen2.5-VL-3B, Qwen-3D integrates 078

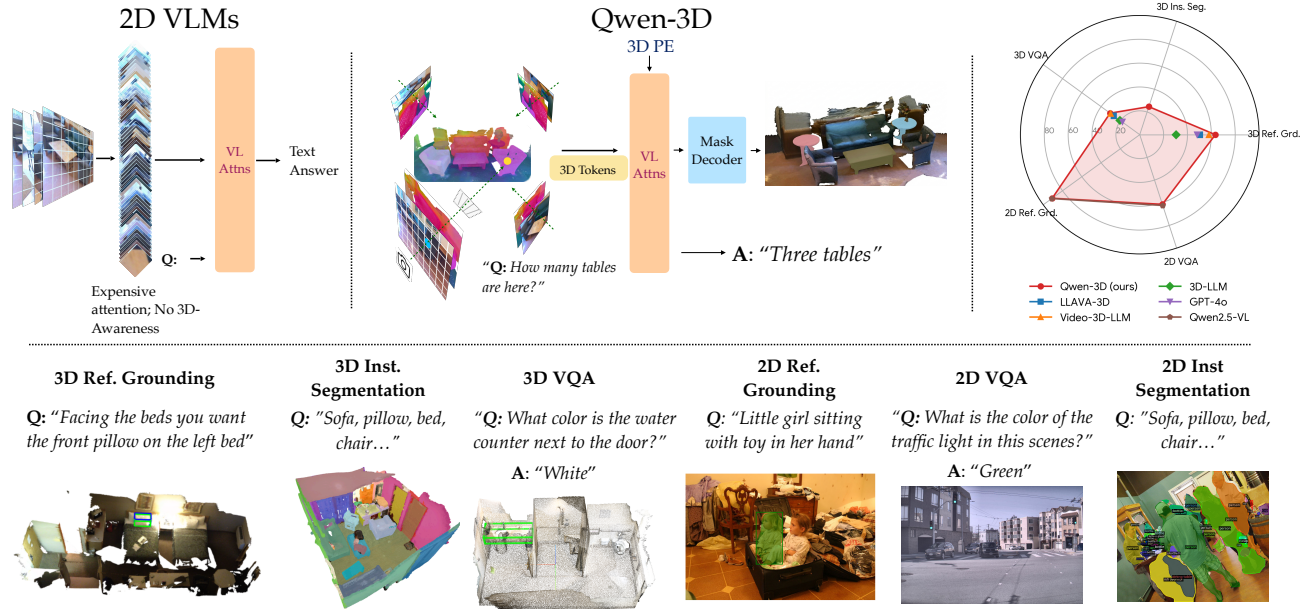


Figure 1. **Qwen-3D replaces frame-level attention with geometry-aligned world-space attention using 3D positional embeddings.** Unlike conventional 2D VLMs that rely on spatio-temporal attention over video frames, Qwen-3D builds a geometry-aware 3D representation from multi-view RGB, depth, and camera poses. This enables attention in the world space via 3D Rotary Positional Embeddings. Qwen-3D supports joint decoding of text responses and 2D/3D object masks. Qwen-3D achieves state-of-the-art performance across a multitude of vision–language tasks from images and multi-view videos.

079 3D structure directly into the vision–language backbone.  
 080 First, we leverage geometric information to compress vi-  
 081 sual tokens within the backbone, merging tokens origin-  
 082 ating from nearby 3D locations. This enables efficient rea-  
 083 soning over long multi-view videos while preserving spa-  
 084 tial consistency across views. Second, we introduce 3D Rotary  
 085 Positional Embeddings to the Qwen backbone, allowing at-  
 086 tention within the backbone to operate in a geometry-aware  
 087 coordinate system and improving cross-view spatial reason-  
 088 ing.

089 In addition, Qwen-3D incorporates a 3D query-based  
 090 segmentation decoder that grounds noun phrases from the  
 091 language stream directly in world space. Unlike prior ap-  
 092 proaches that communicate between the VLM backbone  
 093 and object decoder through a small set of query vectors,  
 094 our model shares full visual token representations between  
 095 the backbone and decoder, enabling richer visual–language  
 096 interactions and stronger grounding performance.

097 Across a wide range of 3D grounding benchmarks,  
 098 Qwen-3D outperforms both proprietary 2D VLMs and the  
 099 strongest existing 3D LMMs while maintaining strong per-  
 100 formance on 2D tasks. Compared to the previous 3D LMM  
 101 state-of-the-art, Qwen-3D improves 3D visual grounding by  
 102 4%, surpasses the best single-stage 3D LMMs by 12%, and  
 103 increases 3D instance segmentation accuracy by 12%. The  
 104 model also achieves state-of-the-art performance on 3D vi-  
 105 sual question answering while preserving the strong 2D ca-  
 106 pabilities of its backbone. Furthermore, Qwen-3D substan-

tially narrows the gap between general-purpose 3D LMMs  
 and specialist 3D grounding models on in-domain bench-  
 marks, while significantly outperforming specialist methods  
 on out-of-distribution 3D scenes and language instructions.

**Contributions.** Our contributions are as follows:

- A geometry-aware VLM backbone that integrates multi-  
view structure through 3D rotary positional embeddings  
and geometry-based token compression.
- A unified 3D grounding architecture with a query-based  
segmentation decoder that grounds language directly in  
world space and shares full visual token representations  
with the VLM backbone.
- A general-purpose 3D multimodal training framework  
that jointly learns from 2D and 3D data, preserving strong  
vision–language capabilities while improving 3D spatial  
reasoning and grounding.
- State-of-the-art performance among 3D-LMM methods  
across 3D benchmarks, improving 3D visual grounding  
by 4%, instance segmentation by 12%, while maintaining  
strong 2D multimodal performance.

We will open-source our model upon acceptance.

## 2. Related Work

**3D Large Vision-Language Models** Building upon the  
rapid progress of 2D Large Multimodal Models [5, 35, 42],  
recent work has focused on extending these models to un-  
derstand 3D scenes. Existing approaches can be grouped  
into four main categories:

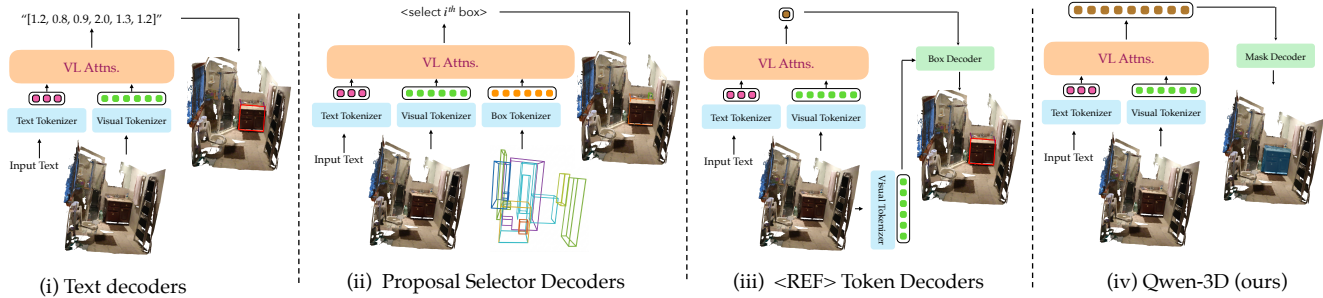


Figure 2. **Decoding object information in 3D LLMs.** From left to right: (i) Direct bounding box decoding in text space, as in 3D-LLM [17]; (ii) Two-stage grounding via proposal selection, as in Video-3D-LLM [19, 20, 53]; (iii) Special token decoding, as in LLaVA-3D [12, 54]; (iv) (Ours) Directly connecting contextualized vision–language features from the VLM backbone to a mask-based segmentation decoder

134 **(a) 3D point cloud encoders trained from scratch.** Methods  
 135 such as LL3DA [11], Scene-LLM [16], and Grounded-  
 136 3D-LLM [12] augment multi-view image streams with explicit  
 137 3D point cloud encoders. The resulting 3D features  
 138 are projected into an LLM backbone in addition to or as  
 139 a replacement for 2D image features. These approaches,  
 140 however, require large-scale point cloud–language datasets  
 141 for alignment—an acute limitation given the scarcity of 3D  
 142 data. In contrast, Qwen-3D builds upon powerful 2D pre-  
 143 trained features and augments them with 3D information via  
 144 positional encodings.  
 145 **(b) Learnable 3D feature compression.** 3D-LLM [17]  
 146 uses Q-Former layers [28] to compress large numbers of  
 147 2D foundation-model features into small sets of latent to-  
 148 kens. In contrast, Qwen-3D performs multi-view feature  
 149 compression in a parameter-free manner, directly guided by  
 150 the 3D spatial layout of the tokens.  
 151 **(c) Object-centric approaches.** Another line of work uses  
 152 object-level features as input to the language model. Ap-  
 153 proaches such as LEO [20] and ChatScene [19] first detect  
 154 objects with off-the-shelf 2D or 3D detectors, pool features  
 155 within detected regions, then feed these pooled features into  
 156 their VLMs. These methods may yield structured object  
 157 representations and improve grounding, but performance is  
 158 fundamentally constrained by the robustness of the detec-  
 159 tors themselves, which often struggle due to limited data  
 160 diversity. In contrast, Qwen-3D is a single-stage model that  
 161 directly grounds the language in the 3D visual stream.  
 162 **(d) Positional embedding adaptation.** Several methods  
 163 modify the positional embeddings of multi-view visual to-  
 164 kens to better encode 3D spatial relationships [22, 23, 53,  
 165 54]. Our model follows this general paradigm. Similar to  
 166 these models, we incorporate 3D information via positional  
 167 embeddings in the vision–language attention. We utilize 3D  
 168 Rotary Positional Encoding for this purpose.

169 Beyond spatial encoding, Qwen-3D also differs from  
 170 prior work in how grounded outputs are decoded. We dis-  
 171 cuss these grounding architectures in detail in the following

section.

**3D Visual Grounding.** Visual grounding—identifying  
 173 objects referred to by language—is a fundamental capabil-  
 174 ity for 3D vision–language systems. Early work [3, 21, 33]  
 175 achieved strong performance by designing specialized archi-  
 176 tectures tailored for 3D grounding. Subsequent meth-  
 177 ods [3, 23, 55, 56] unified grounding, question answer-  
 178 ing, and captioning within a single framework. More re-  
 179 cently, 3D LLMs have leveraged large-scale pretrained vi-  
 180 sion–language features to assist grounding in 3D scenes;  
 181 our method follows this paradigm. These newer approaches  
 182 typically adopt one of three designs ( Fig. 2):  
 183

**(a) Direct bounding-box decoding.** ( Fig. 2a) Models such  
 184 as 3D-LLM [17] directly decode 3D bounding boxes in the  
 185 text space. However, they achieve low performance on lo-  
 186 calization tasks, likely due to the scarcity of 3D-language  
 187 data and the unstructured nature of 3D scenes.  
 188

**(b) Two-stage grounding via proposal selection.**  
 189 ( Fig. 2b) Approaches such as Video-3D-LLM [53],  
 190 ChatScene [19], and LEO [20] first run an object detector  
 191 to generate candidate proposals, then select the object that  
 192 best matches the query. While this improves grounding  
 193 robustness, performance is bottle-necked by the quality of  
 194 the proposals.  
 195

**(c) Special token decoding.** ( Fig. 2c) Other methods de-  
 196 code special grounding tokens (e.g. < REF >) and local-  
 197 ize them using explicit decoder heads, such as in LLaVA-  
 198 3D [54] and Grounded-3D-LLM [12]. Essentially, the VLM  
 199 heads and the decoder heads are only connected via the gen-  
 200 erated < REF > tokens. Although these methods can, in  
 201 principle, ground multiple instances from a category, only  
 202 Grounded-3D-LLM has been applied to 3D object detec-  
 203 tion. We include ablations on the < REF > token design  
 204 with our architecture in our experiments.  
 205

Rather than relying on text-space bounding boxes, pro-  
 206 posal selection, or grounding tokens, **Qwen-3D directly**  
 207 **connects contextualized vision–language features from**  
 208

209 **the VLM backbone to a mask-based segmentation de-**  
 210 **coder** ( Fig. 2d). This avoids the bottleneck imposed by  
 211 the special token decoding and aligns more naturally with  
 212 contemporary multi-object detection and segmentation archi-  
 213 tectures. As a result, Qwen-3D achieves state-of-the-art  
 214 grounding and detection performance among 3D LMMs.

### 215 3. Method

216 We introduce Qwen-3D, a geometry-aware extension of  
 217 Qwen2.5-VL [5] that enables efficient reasoning over multi-  
 218 view 3D scenes. Given multi-view RGB-D frames, camera  
 219 poses, and a language query, Qwen-3D projects visual  
 220 features into a 3D feature cloud and compresses redundan-  
 221 cies via voxel-based token merging. A vision-language  
 222 backbone processes these geometry-aligned tokens using  
 223 3D Rotary Positional Embeddings for spatial reasoning in  
 224 world coordinates. Finally, a mask-based grounding de-  
 225 coder predicts 3D segmentation masks for queried objects  
 226 alongside text responses for question-answering. Notably,  
 227 Qwen-3D preserves the original 2D input–output interface  
 228 of Qwen2.5-VL while extending it to operate on multi-view  
 229 3D scenes. Figure 3 illustrates the architecture, detailed be-  
 230 low.

231 **Geometry-Aware Visual Encoding** We construct a 3D  
 232 feature cloud by projecting multi-view RGB features into  
 233 world coordinates using depth maps and camera poses. This  
 234 enables geometry-aware token compression based on spatial  
 235 proximity. We use the Qwen2.5-VL pre-trained ViT en-  
 236 coder to process RGB images into feature maps  $\mathcal{F}$ , down-  
 237 sampled by a factor of  $S$  to a size of  $N \times \frac{H}{S} \times \frac{W}{S} \times D$ ,  
 238 where  $D$  is the feature dimension. We then unproject the  
 239 pixel-aligned depth maps using camera intrinsics and poses,  
 240 nearest-downsampling them to match this spatial resolution  
 241 and yield 3D coordinates  $\mathcal{C}$  of size  $N \times \frac{H}{S} \times \frac{W}{S} \times 3$ .

242 **Geometry-based token compression** Multi-view obser-  
 243 vations inherently produce redundant tokens at overlapping  
 244 3D locations. To mitigate memory bottlenecks, we apply  
 245 voxel-based token merging (with a voxel size of 5cm) to the  
 246 feature-coordinate pairs, following [22, 23, 54]. This pro-  
 247 cess discretizes the space and mean-pools the features and  
 248 coordinates within each occupied voxel, yielding a compact,  
 249 unordered set of geometry-aligned tokens  $\mathcal{FC}' =$   
 250  $\{(f_j, p_j)\}_{j=1}^{M'}$ , where  $M' \leq M$ .

251 **Vision–Language Attention with 3D RoPE** We use the  
 252 Qwen2.5-VL language tokenizer [5] to embed the input  
 253 natural language query into a sequence of tokens  $\mathcal{T} =$   
 254  $\{t_k\}_{k=1}^L$ , where  $t_k \in \mathbb{R}^D$  and  $L$  is the number of tokens.  
 255 The concatenated sequence of voxelized 3D point features  
 256 and text tokens is then processed by  $N$  pre-trained multi-  
 257 modal attention layers from Qwen2.5-VL. To enable spatial  
 258 reasoning, we adapt Qwen’s Multimodal RoPE—originally

designed for 1D or 2D grids—to encode 3D world coordi-  
 nates  $(x, y, z)$ , resulting in four positional components

$$\text{PE}_{3D}(\mathbf{p}) = [\text{PE}(t); \text{PE}(x); \text{PE}(y); \text{PE}(z)].$$

where  $t$  represents a token’s temporal position in the se-  
 quence and  $x, y, z$  denote 3D spatial coordinates. For text,  
 the temporal ID increments sequentially (reducing to stan-  
 dard 1D RoPE), while point cloud tokens share a constant  
 temporal ID and use their world coordinates as spatial IDs.

This embedding defines a rotation matrix  $R(\mathbf{p})$  ap-  
 plied to query and key vectors before attention:  $\tilde{\mathbf{q}} =$   
 $R(\mathbf{p}_q)\mathbf{q}$ ,  $\mathbf{k} = R(\mathbf{p}_k)\mathbf{k}$

Because 3D point clouds are permutation-invariant, we  
 replace Qwen2.5-VL’s autoregressive causal masking with  
 full attention over visual tokens. This improves 3D perfor-  
 mance without degrading the backbone’s original capabili-  
 ties. Further attention masking details and RoPE ablations  
 are provided in the appendix and experiments section, re-  
 spectively.

U

**Grounding Decoder** Unlike prior VLM grounding ar-  
 chitectures that communicate with object decoders via a  
 small set of query tokens [12, 54], our approach directly at-  
 tends to the full set of contextualized visual tokens from the  
 backbone, enabling richer language-visual interactions. We  
 adapt a mask decoder design similar to UniVLG [23] [23].  
 We instantiate  $N$  learnable object queries that cross-attend  
 to the language and visual tokens from the VLM backbone,  
 followed by self-attention among the queries. For all at-  
 tention involving visual tokens, we 2D or 3D positional  
 embeddings corresponding to the input modality. Finally,  
 we decode segmentation masks via a dot product between  
 object queries and the updated visual tokens and identify  
 grounded object text spans via dot products with the lan-  
 guage tokens. For question answering, we use the Qwen2.5-  
 VL text-generation head to perform next-token prediction  
 over the contextualized scene-language features.

**Training Objectives** We supervise Qwen-3D on three  
 losses: (a) **Mask loss:** We assign predictions to ground-  
 truth instances via Hungarian matching [8] and supervise  
 matched masks with Binary Cross-Entropy (BCE) and Dice  
 loss, following Mask2Former [13].

(b) **Span alignment loss:** As in [21, 24, 29], we supervise  
 each predicted text span with the corresponding matched  
 ground-truth text span using Binary Cross Entropy loss.  
 The queries that remain unmatched are supervised to pre-  
 dict low probability over all text tokens.

(c) **Text generation loss:** For question-answering and cap-  
 tioning tasks, we use a token-level cross-entropy on the gen-  
 erated answer.

Our complete loss is formulated as:

$$\mathcal{L} = \alpha_{\text{mask}} \mathcal{L}_{\text{mask}} + \alpha_{\text{span}} \mathcal{L}_{\text{span}} + \alpha_{\text{gen}} \mathcal{L}_{\text{gen}} \quad (1)$$

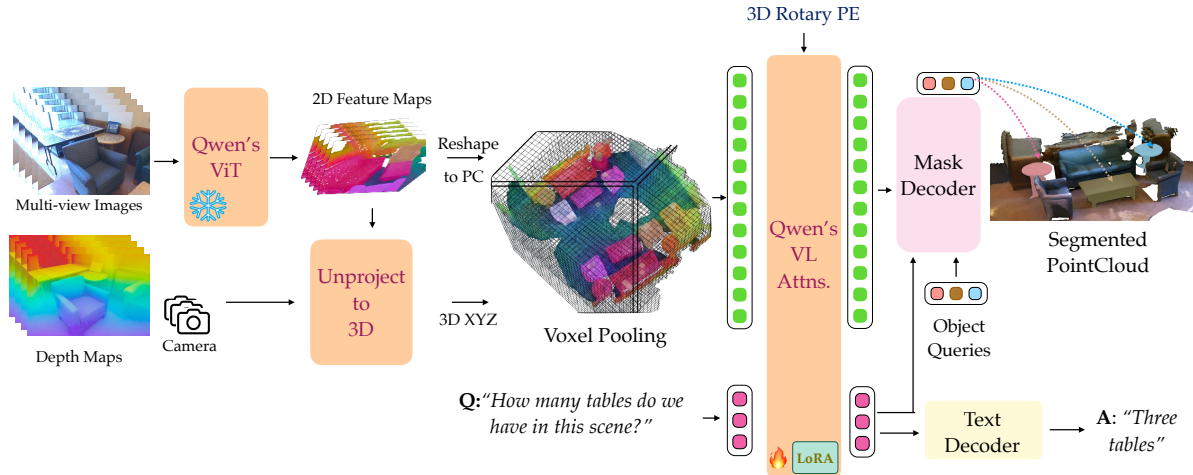


Figure 3. **Qwen-3D architecture.** Given a natural language query and multi-view RGB-D inputs, the Qwen2.5-VL vision encoder extracts multi-view 2D features, unprojects them into world-space XYZ coordinates, and voxel-pools them to reduce spatial redundancy. The resulting tokens are fused with text via Qwen vision–language attention layers augmented with 3D Rotary Positional Embeddings and LoRA adaptation. Two task heads operate on the shared tokens: a 3D mask decoder for referential grounding and instance segmentation and a text decoder for open-ended question answering.

where  $\mathcal{L}_{\text{mask}}$  is the mask loss,  $\mathcal{L}_{\text{span}}$  is the span alignment loss,  $\mathcal{L}_{\text{gen}}$  is the text generation loss, and the  $\alpha$ 's are the loss weight terms.

**Joint 2D-3D Training** Qwen-3D shares parameters across 2D and 3D modalities, natively processing single or multiple RGB images as well as multi-view posed RGB-D frames. Within both the VLM backbone and the grounding decoder, 2D visual tokens are structured on a regular grid with 2D positional embeddings, whereas 3D inputs are represented as unordered point tokens with 3D embeddings. To improve alignment between 2D and 3D representations during training, we lift 2D data to 3D with probability  $p$  using reconstruction models such as MoGE [43].

**Implementation Details** Qwen-3D introduces only  $\sim 50\text{M}$  trainable parameters. We freeze the Qwen2.5-VL ViT backbone, fine-tune the vision–language attention layers via LoRA [18], and train the mask decoder from scratch. We train jointly on 2D and 3D datasets runs for 200k iterations (learning rate  $10^{-4}$ ) on eight 48GB L40S GPUs with an effective batch size of 8, which takes approximately three days. Text-generation loss is applied exclusively to captioning and question-answering tasks; for detection, we construct prompts [21] by concatenating object class names (e.g., “find chair. table. sofa.”).

Following [22, 23], we subsample 15 frames per scene during training. At inference time, we feed all posed RGB-D frames to our model ( $\sim 90$  on average for ScanNet), which takes about 2 seconds per scene end-to-end. Our voxel-pooling strategy is critical for enabling this scalability - without it, the model runs out of memory even with sub-

stantially fewer input frames. Finally, due to computational constraints, we adopt the 3B VLM backbone variant rather than the 7B models common in prior 3D-LMMs [53, 54]. Additional hyperparameters are detailed in the appendix.

## 4. Experiments

We evaluate Qwen-3D against existing LMMs and specialized 3D vision models on visual grounding in both in-domain (Sec. 4.1) and out-of-domain settings (Sec. 4.2), along with 3D instance segmentation (Sec. 4.3) and 3D VQA (Sec. 4.4). We also assess how well Qwen-3D retains its 2D multimodal capabilities (Sec. 4.5) and analyze which design choices most significantly impact performance (Sec. 4.6). Qualitative results, failure mode analysis, robustness to depth and camera pose noise, and performance on text-only tasks are provided in the appendix.

**Training Datasets** We train jointly on a mixture of 3D and 2D datasets to enable 3D comprehension while preserving the base model’s pre-trained capabilities. The 3D datasets include referential grounding (SR3D, NR3D [1], ScanRefer [10]), instance segmentation (ScanNet200 [37], Matterport [9]), and question answering (ScanQA [4], SQA3D [34]). To mitigate catastrophic forgetting of Qwen’s original capabilities, we co-train on 2D datasets: referential grounding (RefCOCO, RefCOCO+, RefCOCOG [25]), instance segmentation (COCO [30]), captioning and QA (LLaVA-Instruct-150k [31]), and instruction fine-tuning (Alpaca [41]).

Table 1. Results on 3D Visual Grounding and QA for both experts and LMMs. \*2D VLM numbers obtained from prior works [3, 40]

Methods		ScanRefer		ScanQA (Val)				SQA3D (Test)
		Acc @25 (Det)	Acc @50 (Det)	EM@1	C	M	R	EM@1
Expert Models	BUTD-DETR [21]	52.2	39.8	-	-	-	-	-
	ScanQA [4]	-	-	23.5	67.3	13.6	34.3	-
	PQ3D [56]	56.7	51.8	20.0	65.2	13.9	-	47.1
	3D-VisTA [55]	51.0	46.2	22.4	69.6	13.9	35.7	48.5
	ODIN [22]	43.1	33.4	-	-	-	-	-
	Locate-3D [3]	61.1	50.9	-	-	-	-	-
	UniVLG [23]	63.5	56.4	25.7	78.5	15.2	40.0	50.2
<b>2D VLMs*</b>								
	LLaMA [2]	37.3	24.2	-	-	-	-	-
	GPT-4o [35]	48.2	32.6	18.0	58.3	14.2	33.4	-
<b>Two-Stage</b>								
	ChatScene [19]	55.5	50.2	21.6	87.7	18.0	41.6	54.6
	LEO [20]	-	-	24.5	101.4	20.0	49.2	50.0
	Video-3D LLM - 7B [53]	58.1	51.7	30.1	102.1	-	-	<b>58.6</b>
<b>Single-Stage</b>								
	NaviLLM [52]	-	-	23.0	75.9	15.4	38.4	-
	LL3DA [11]	-	-	-	76.8	15.9	37.3	-
	SceneLLM [16]	-	-	27.2	80.0	16.6	40.0	54.2
	3D-LLM [17]	30.3	-	20.5	69.4	14.5	35.7	-
	Grounded 3D-LLM [12]	48.6	44.0	-	75.9	-	-	-
	LLaVA-3D - 7B [54]	50.1	42.7	27.0	<b>103.1</b>	20.8	<b>49.6</b>	55.6
	Qwen-3D - 3B (Ours)	<b>65.6</b>	<b>57.7</b>	<b>30.5</b>	95.2	<b>38.3</b>	46.9	56.0

#### 367 4.1. Evaluation on 3D Referential Grounding

368 **Datasets.** We evaluate on the validation sets of three  
369 ScanNet-based [14] 3D referential grounding benchmarks:  
370 SR3D, NR3D [1], and ScanRefer [10]. While SR3D com-  
371 prises 88k synthetic utterances, NR3D (41k) and ScanRefer  
372 (51k) feature complex, human-annotated queries. Follow-  
373 ing recent work [3, 7, 22, 23], we operate directly on noisy,  
374 raw sensor RGB-D point clouds rather than clean, post-  
375 processed meshes. Although this setup introduces sensor-  
376 mesh misalignments that can degrade performance [22, 23],  
377 it better reflects practical embodied learning scenarios.

378 **Evaluation Metrics.** We report standard Top-1 accuracy,  
379 where a prediction is correct if the highest-confidence pre-  
380 dicted bounding box achieves an Intersection over Union  
381 (IoU) with the ground-truth box above a threshold (0.25,  
382 0.5). Since our model predicts segmentation masks, we  
383 convert masks to bounding boxes by thresholding at their  
384 extreme corners.

385 **Baselines.** Following prior work [53, 54], we compare  
386 Qwen-3D against state-of-the-art expert (non-LMM) and  
387 LLM-based approaches. Expert baselines include two-  
388 stage [55, 56] and single-stage [3, 23] methods. LLM-  
389 based baselines include: (i) two-stage models relying on  
390 detector proposals (LEO [20], Chat-Scene [19], Video-

391 3D-LLM [53]); (ii) single-stage text-space decoders (3D-  
392 LLM [17]); and (iii) single-stage  $\langle REF \rangle$  token decod-  
393 ers (Grounded 3D-LLM [17], LLaVA-3D [54]). Unlike  
394 these, our single-stage method directly decodes segmenta-  
395 tion masks by routing LMM backbone features to an ob-  
396 ject mask decoder. We also evaluate against proprietary 2D  
397 VLMs (GPT-4o [35], LLaMA [2], Qwen2-VL [32]). Tab. 1  
398 presents quantitative results, and full results on Referit3D  
399 are available in the appendix.

#### 400 **Qwen-3D outperforms all prior 3D LMM-based models.**

401 It surpasses the text-decoding single-stage model of 3D-  
402 LLM [17] by over 30%, the recent single-stage state-of-the-  
403 art LLaVA-3D [54] by 12%, and the two-stage Video-3D-  
404 LLM [53] by 4%. This establishes Qwen-3D as the new  
405 state-of-the-art for 3D referential grounding among LMM-  
406 based models. Notably, our strongest baseline, Video-3D-  
407 LLM, uses a larger backbone (Qwen2-VL-7B), whereas  
408 Qwen-3D uses Qwen2.5-VL-3B (e.g., 70.1 vs. 65.3 on  
409 RealWorldQA [45]), indicating that the improvements of  
410 Qwen-3D are not solely explained by a stronger VLM back-  
411 bone.

#### 412 **Qwen-3D closes the gap with expert 3D grounding mod-**

413 **els.** On ScanRefer [10], Qwen-3D outperforms state-of-the-  
414 art expert model UniVLG and all other models and sub-  
415 stantially narrows the gap between specialist 3D models

Table 2. Evaluation on Locate3D ScanNet++.

Model	Acc@25	Acc@50
UniVLG [23]	32.3	24.6
Video-3D LLM [53]	33.2	27.6
Qwen-3D (Ours)	<b>50.5</b>	<b>28.5</b>

Table 3. Evaluation on ScanNet200 Instance Segmentation.

	Model	mAP	mAP25
Closed Vocabulary	Mask3D [39]	27.4	42.3
	PQ3D [56]	27.0	46.3
	MAFT [27]	29.2	43.3
	ODIN [22]	31.5	53.1
Language-Prompted	PQ3D [56]	20.2	32.5
	UniVLG [23]	27.9	46.1
LLM-Based	Grounded-3D-LLM [12]	12.1	16.8
	Qwen-3D (Ours)	<b>27.7</b>	<b>46.2</b>

416 and LMM-based approaches on other 3D grounding bench-  
417 marks.

## 418 4.2. Out-of-Domain 3D Referential Grounding

419 While LMMs often trail specialists in-domain, they typ-  
420 ically excel at out-of-domain (OOD) generalization. We  
421 evaluate Qwen-3D on Locate-3D [3], which provides hu-  
422 man instructions for ScanNet++ [49] scenes. ScanNet++  
423 introduces a distinct domain shift from our fine-tuning data  
424 as it is captured via iPhone LiDAR rather than ScanNet’s  
425 iPad Structure sensor.

426 We compare against public checkpoints of UniVLG and  
427 Video-3D-LLM (supplied with state-of-the-art ODIN [22]  
428 box proposals). We omit LLaVA-3D [54] as its ground-  
429 ing model weights and code are not publicly released. As  
430 shown in Tab. 2, Qwen-3D significantly outperforms both  
431 baselines on these OOD tasks. We attribute our better  
432 generalization over UniVLG to Qwen-3D’s stronger pre-  
433 training. Furthermore, while Video-3D-LLM uses a large  
434 VLM backbone, its reliance on off-the-shelf 3D detectors  
435 bottlenecks OOD robustness. Conversely, Qwen-3D di-  
436 rectly decodes boxes from VLM features, allowing it to bet-  
437 ter exploit the underlying representation for superior gener-  
438 alization.

## 439 4.3. 3D Instance Segmentation

440 We evaluate Qwen-3D on the ScanNet200 [37] instance  
441 segmentation benchmark. While traditional methods [22,  
442 39] assume a closed vocabulary setup, recent models [23,  
443 56]—like ours—adopt a language-prompted paradigm  
444 (e.g., “find chairs. tables. sofa.”). Furthermore, because  
445 most 3D grounding VLMs predict only a few bounding  
446 boxes, they fail as full scene detectors, with Grounded-

Table 4. 2D Ref. grounding datasets and RealWorldQA

	RefCOCO	RefCOCO+	RefCOCOg	RealWorldQA
LAVT [48] (B)	72.7	62.4	61.2	-
ReSTR [26]	67.2	55.7	54.5	-
X-Decoder (L) [57]	-	-	64.6	-
UniVLG [23]	69.2	61.3	64.1	-
Qwen-2.5 3B [5]	89.1	82.4	85.2	<b>62.6</b>
Qwen-3D (Ours)	<b>90.7</b>	85.3	<b>86.0</b>	61.4

3D-LLM [12] being the only exception to our knowledge. 447  
As shown in Tab. 3, Qwen-3D outperforms Grounded-3D- 448  
LLM by 15% AP and 30% AP25, effectively bridging the 449  
performance gap between generalist VLMs and matching 450  
the performance of language-prompted specialist models. 451

## 452 4.4. 3D Visual Question Answering

453 We evaluate Qwen-3D on two 3D question answering 454  
benchmarks: ScanQA [4] and SQA3D [34]. Both datasets 455  
use visual scenes from ScanNet [14], with ScanQA focus- 456  
ing on spatial-relation questions and SQA3D emphasizing 457  
situational reasoning.

458 Following prior work, we report Exact Match (EM@1), 459  
ROUGE (R), CIDEr (C), and METEOR (M). As shown 460  
in Tab. 1, Qwen-3D outperforms state-of-the-art expert 461  
model UniVLG [23], and achieves comparable perfor- 462  
mance to the single-stage LMM state-of-the-art LLaVA- 463  
3D-7B [54].

## 464 4.5. 2D Vision-Language tasks

465 To prevent degradation of the base model’s strong 2D ca- 466  
pabilities, we co-train Qwen-3D on 2D datasets, including 467  
the RefCOCO family [25] and LLaVA-Instruct-150k [31]. 468  
Tab. 4 evaluates Qwen-3D against its pre-trained base 469  
model (Qwen2.5-VL 3B [5]) on RefCOCO 2D grounding 470  
benchmarks and the held-out RealWorldQA [45] dataset. 471  
As in our 3D setup, Qwen-3D predicts segmentation masks 472  
which we convert to bounding boxes for evaluation. Results 473  
show that Qwen-3D performs comparably to the 2D-only 474  
base model across both tasks, successfully maintaining its 475  
2D strengths while extending effectively to 3D.

476 For additional visualizations of Qwen-3D, ablations on 477  
pre-training retention, and results on text-only tasks, please 478  
refer to the appendix.

## 479 4.6. Additional Analysis and Ablations

480 We ablate several of our key design choices on ScanRefer 481  
(Top-1@0.5).

482 **Full Text Tokens vs. Grounding Token** We ablate the 483  
full token interface between the vision-language attention 484  
module and the object mask decoder against a variant that 485  
replaces the full text token injection with a special ground- 486  
ing token < REF > into the mask decoder, similar to 487  
LLAVA-3D [54] and Grounded-3D-LLM [12]. As shown in

Table 5. Ablations of Qwen-3D

(a) Text vs. < REF >		(b) Pos. Embed.		(c) Attn. Mask		(d) VLM Tuning	
Model	Top1 Acc.	Pos. Embed.	Top1 Acc	Mask Type	Top1 Acc.	Finetune	Top1 Acc.
Grounded 3D-LLM	44.0	2D	53.2	Causal	36.4	Frozen	37.0
LLAVA-3D - 7B	42.7	Naive 3D mRoPE	53.5	Full	<b>57.7</b>	Finetune	<b>57.7</b>
Qwen-3D < REF >	43.0	3D RoPE (ours)	<b>57.7</b>				
Qwen-3D (ours)	<b>57.7</b>						

488 Tab. 5a, with the < REF > token, Qwen-3D achieves com- 528  
 489 parable performance to both LLAVA-3D and Grounded-3D- 529  
 490 LLM, however, Qwen-3D with the full text and visual token 530  
 491 interface greatly outperforms all < REF > token baselines. 531

492 **Positional Embedding:** We compare the 3D positional 532  
 493 embeddings against the 2D MRoPE utilized in Qwen2.5- 533  
 494 VL and a naive implementation of 3D RoPE in Qwen’s 534  
 495 MRoPE on 3D grounding. As shown in Tab. 5b, the 3D 535  
 496 RoPE outperforms all other RoPE variants. We hypothesize 536  
 497 this is because 3D RoPE improves the model’s cross-view 537  
 498 spatial reasoning over the 2D variant, while the distribution 538  
 499 of the frequencies in Qwen’s original MRoPE hinders the 539  
 500 model’s capability to discern spatial relationships along the 540  
 501 point cloud’s axes. We provide further details of the alter- 541  
 502 native RoPE designs in the appendix. 542

503 **Causal Mask vs. Full Attention Mask:** Qwen2.5-VL 543  
 504 applies causal masking in its vision–language attention lay- 544  
 505 ers, even for non-autoregressive tasks such as grounding. 545  
 506 We compare this original design—which stays closer to 546  
 507 the model’s pre-training distribution—against a full atten- 547  
 508 tion mask where all vision and language tokens attend to 548  
 509 each other except during autoregressive text generation. As 549  
 510 shown in Tab. 5c, the full-attention variant significantly out- 550  
 511 performs the causal-mask variant. 551

512 **Freezing vs. Fine-tuning the Base VLM:** As shown in 552  
 513 Tab. 5d, fine-tuning the underlying Qwen VLM is essen- 553  
 514 tial to achieve strong grounding performance. Keeping the 554  
 515 VLM frozen leads to a substantial drop in accuracy. 555

## 516 5. Limitations and Future Directions 556

517 In this work, we focus on visual grounding from multi-view 557  
 518 RGB videos that depict largely *static scenes*. Extending 558  
 519 Qwen-3D to handle dynamic environments is an important 559  
 520 direction for future work. This will likely require integrat- 560  
 521 ing temporal motion representations, dynamic scene recon- 561  
 522 struction, or object-centric tracking into the model’s geo- 562  
 523 metric reasoning pipeline. Additionally, as with prior 3D 563  
 524 vision-language models, the performance of Qwen-3D de- 564  
 525 pends on the quality of the input 3D geometry. Although 565  
 526 our experiments already operate on real-world noisy 3D re- 566  
 527 constructions and demonstrate strong performance, jointly 567

refining geometry during inference remains an interesting 528  
 direction for future work. 529

Another promising direction is enabling *multi-step* 530  
*grounded reasoning*. In the current formulation, Qwen- 531  
 3D primarily performs single-step grounding by directly lo- 532  
 calizing language-referred entities in the scene. However, 533  
 many real-world tasks require sequential reasoning, such as 534  
 composing multiple grounding operations, performing spa- 535  
 tial comparisons, or executing multi-stage instructions. In- 536  
 corporating multi-step reasoning mechanisms, for example 537  
 through iterative grounding or reasoning frameworks such 538  
 as [38], could enable richer interaction between language 539  
 and spatial representations. We expect that combining 540  
 geometry-aware perception with structured reasoning will 541  
 further improve performance on complex visual–language 542  
 tasks in 3D environments. 543

## 544 6. Conclusion 544

545 We introduced Qwen-3D, a geometry-aware 3D vi- 545  
 546 sion–language model that integrates explicit multi-view ge- 546  
 547 ometric reasoning into a strong multimodal backbone. By 547  
 548 leveraging depth and camera pose to guide token compres- 548  
 549 sion and introducing scale-aligned 3D Rotary Positional 549  
 550 Embeddings, Qwen-3D efficiently processes long multi- 550  
 551 view video streams while maintaining consistent cross-view 551  
 552 spatial understanding. We further proposed a tightly in- 552  
 553 tegrated grounding architecture that directly connects con- 553  
 554 textualized vision–language features from the backbone to 554  
 555 a 3D query-based segmentation decoder, enabling direct 555  
 556 language-to-3D alignment. Extensive experiments show 556  
 557 that Qwen-3D significantly advances the state of the art 557  
 558 among general-purpose 3D LMMs, achieving strong im- 558  
 559 provements in 3D grounding and instance segmentation 559  
 560 while preserving competitive 2D multimodal performance. 560  
 561 Our results substantially narrow the gap between generalist 561  
 562 multimodal models and specialist 3D grounding systems, 562  
 563 and demonstrate superior robustness in out-of-domain set- 563  
 564 tings. Overall, our findings highlight the importance of 564  
 565 tightly integrating geometric structure with pretrained vi- 565  
 566 sion–language representations, suggesting a promising path 566  
 567 toward unified models capable of reasoning across both 2D 567  
 and 3D visual environments. 568

569

## References

570

571

572

573

574

575

576

577

578

579

580

581

582

583

584

585

586

587

588

589

590

591

592

593

594

595

596

597

598

599

600

601

602

603

604

605

606

607

608

609

610

611

612

613

614

615

616

617

618

619

620

621

622

623

624

625

- [1] Panos Achlioptas, Ahmed Abdelreheem, Fei Xia, Mohamed Elhoseiny, and Leonidas Guibas. ReferIt3D: Neural Listeners for Fine-Grained 3D Object Identification in Real-World Scenes. In *Proc. ECCV*, 2020. 5, 6, 12, 13, 14
- [2] Meta AI. Llama 3: The llama-3 herd of models. <https://ai.meta.com/llama/>, 2024. Large language model. 6, 14
- [3] Sergio Arnaud, Paul McVay, Ada Martin, Arjun Majumdar, Krishna Murthy Jatavallabhula, Phillip Thomas, Ruslan Partsey, Daniel Dugas, Abha Gejji, Alexander Sax, Vincent-Pierre Berges, Mikael Henaff, Ayush Jain, Ang Cao, Ishita Prasad, Mrinal Kalakrishnan, Michael Rabbat, Nicolas Ballas, Mido Assran, Oleksandr Maksymets, Aravind Rajeswaran, and Franziska Meier. Locate 3d: Real-world object localization via self-supervised learning in 3d, 2025. 3, 6, 7, 14
- [4] Daichi Azuma, Taiki Miyanishi, Shuhei Kurita, and Motoaki Kawanabe. Scanqa: 3d question answering for spatial scene understanding. In *proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 19129–19139, 2022. 5, 6, 7, 12
- [5] Shuai Bai, Keqin Chen, Xuejing Liu, and Jialin Wang et. al. Qwen2.5-vl technical report, 2025. 1, 2, 4, 7, 12, 13
- [6] Daniel Bolya, Cheng-Yang Fu, Xiaoliang Dai, Peizhao Zhang, Christoph Feichtenhofer, and Judy Hoffman. Token merging: Your ViT but faster. In *International Conference on Learning Representations*, 2023. 13
- [7] Ang Cao, Sergio Arnaud, Oleksandr Maksymets, Jianing Yang, Ayush Jain, Sriram Yenamandra, Ada Martin, Vincent-Pierre Berges, Paul McVay, Ruslan Partsey, Aravind Rajeswaran, Franziska Meier, Justin Johnson, Jeong Joon Park, and Alexander Sax. From thousands to billions: 3d visual language grounding via render-supervised distillation from 2d vlms, 2025. 6
- [8] Nicolas Carion, Francisco Massa, Gabriel Synnaeve, Nicolas Usunier, Alexander Kirillov, and Sergey Zagoruyko. End-to-End Object Detection with Transformers. In *Proc. ECCV*, 2020. 4
- [9] Angel Chang, Angela Dai, Thomas Funkhouser, Maciej Halber, Matthias Niessner, Manolis Savva, Shuran Song, Andy Zeng, and Yinda Zhang. Matterport3d: Learning from rgb-d data in indoor environments. *arXiv preprint arXiv:1709.06158*, 2017. 5
- [10] Dave Zhenyu Chen, Angel Chang, and Matthias Nießner. ScanRefer: 3D Object Localization in RGB-D Scans using Natural Language. In *Proc. ECCV*, 2020. 5, 6, 12, 14
- [11] Sijin Chen, Xin Chen, Chi Zhang, Mingsheng Li, Gang Yu, Hao Fei, Hongyuan Zhu, Jiayuan Fan, and Tao Chen. Ll3da: Visual interactive instruction tuning for omni-3d understanding reasoning and planning. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 26428–26438, 2024. 1, 3, 6, 14
- [12] Yilun Chen, Shuai Yang, Haifeng Huang, Tai Wang, Runsen Xu, Ruiyuan Lyu, Dahua Lin, and Jiangmiao Pang. Grounded 3d-llm with referent tokens. *arXiv preprint arXiv:2405.10370*, 2024. 1, 3, 4, 6, 7, 14
- [13] Bowen Cheng, Ishan Misra, Alexander G. Schwing, Alexander Kirillov, and Rohit Girdhar. Masked-attention mask transformer for universal image segmentation. 2022. 4, 13
- [14] 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. 1, 6, 7
- [15] Angela Dai, Matthias Nießner, Michael Zollöfer, Shahram Izadi, and Christian Theobalt. Bundlefusion: Real-time globally consistent 3d reconstruction using on-the-fly surface re-integration. *ACM Transactions on Graphics 2017 (TOG)*, 2017. 12
- [16] Rao Fu, Jingyu Liu, Xilun Chen, Yixin Nie, and Wenhan Xiong. Scene-llm: Extending language model for 3d visual understanding and reasoning. *arXiv preprint arXiv:2403.11401*, 2024. 3, 6, 14
- [17] Yining Hong, Haoyu Zhen, Peihao Chen, Shuhong Zheng, Yilun Du, Zhenfang Chen, and Chuang Gan. 3d-llm: Injecting the 3d world into large language models. *Advances in Neural Information Processing Systems*, 36:20482–20494, 2023. 3, 6, 14
- [18] Edward J. Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. Lora: Low-rank adaptation of large language models, 2021. 5
- [19] Haifeng Huang, Yilun Chen, Zehan Wang, Rongjie Huang, Runsen Xu, Tai Wang, Luping Liu, Xize Cheng, Yang Zhao, Jiangmiao Pang, et al. Chat-scene: Bridging 3d scene and large language models with object identifiers. *Advances in Neural Information Processing Systems*, 37:113991–114017, 2024. 1, 3, 6, 14
- [20] Jianguo Huang, Silong Yong, Xiaojian Ma, Xiongkun Linghu, Puhao Li, Yan Wang, Qing Li, Song-Chun Zhu, Baoxiang Jia, and Siyuan Huang. An embodied generalist agent in 3d world. *arXiv preprint arXiv:2311.12871*, 2023. 1, 3, 6, 14
- [21] Ayush Jain, Nikolaos Gkanatsios, Ishita Mediratta, and Katerina Fragkiadaki. Bottom up top down detection transformers for language grounding in images and point clouds. In *European Conference on Computer Vision*, pages 417–433. Springer, 2022. 3, 4, 5, 6, 14
- [22] Ayush Jain, Pushkal Katara, Nikolaos Gkanatsios, Adam W Harley, Gabriel Sarch, Kriti Aggarwal, Vishrav Chaudhary, and Katerina Fragkiadaki. Odin: A single model for 2d and 3d segmentation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 3564–3574, 2024. 1, 3, 4, 5, 6, 7, 14
- [23] Ayush Jain, Alexander Szwedlow, Yuzhou Wang, Sergio Arnaud, Ada Martin, Alexander Sax, Franziska Meier, and Katerina Fragkiadaki. Unifying 2d and 3d vision-language understanding, 2025. 1, 3, 4, 5, 6, 7, 12, 13, 14, 15
- [24] Aishwarya Kamath, Mannat Singh, Yann André LeCun, Ishan Misra, Gabriel Synnaeve, and Nicolas Carion. MDETR - Modulated Detection for End-to-End Multi-Modal Understanding. In *Proc. ICCV*, 2021. 4
- [25] Sahar Kazemzadeh, Vicente Ordonez, Mark Matten, and Tamara Berg. ReferItGame: Referring to objects in pho-

- 684 tographs of natural scenes. In *Proceedings of the 2014 Con-*  
685 *ference on Empirical Methods in Natural Language Process-*  
686 *ing (EMNLP)*, pages 787–798, Doha, Qatar, 2014. Associa-  
687 tion for Computational Linguistics. 5, 7, 12
- [26] Namyup Kim, Dongwon Kim, Cuiling Lan, Wenjun Zeng,  
688 and Suha Kwak. Restr: Convolution-free referring image  
689 segmentation using transformers, 2022. 7
- [27] Xin Lai, Yuhui Yuan, Ruihang Chu, Yukang Chen, Han Hu,  
690 and Jiaya Jia. Mask-attention-free transformer for 3d in-  
691 stance segmentation. In *Proceedings of the IEEE/CVF Inter-*  
692 *national Conference on Computer Vision*, pages 3693–3703,  
693 2023. 7
- [28] Junnan Li, Dongxu Li, Silvio Savarese, and Steven Hoi.  
694 Blip-2: Bootstrapping language-image pre-training with  
695 frozen image encoders and large language models. In *In-*  
696 *ternational conference on machine learning*, pages 19730–  
697 19742. PMLR, 2023. 3
- [29] Liunian Harold Li, Pengchuan Zhang, Haotian Zhang, Jian-  
700 wei Yang, Chunyuan Li, Yiwu Zhong, Lijuan Wang, Lu  
701 Yuan, Lei Zhang, Jenq-Neng Hwang, et al. Grounded  
702 language-image pre-training. In *Proceedings of the*  
703 *IEEE/CVF Conference on Computer Vision and Pattern*  
704 *Recognition*, pages 10965–10975, 2022. 4
- [30] Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays,  
705 Pietro Perona, Deva Ramanan, Piotr Dollár, and C Lawrence  
706 Zitnick. Microsoft coco: Common objects in context. In  
707 *Computer Vision–ECCV 2014: 13th European Conference,*  
708 *Zurich, Switzerland, September 6–12, 2014, Proceedings,*  
709 *Part V 13*, pages 740–755. Springer, 2014. 5, 12
- [31] Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee.  
710 Visual instruction tuning. *Advances in neural information*  
711 *processing systems*, 36:34892–34916, 2023. 5, 7, 12
- [32] Yuanyuan Liu, Haiyang Mei, Li Zhang, et al. View-on-  
712 graph: Zero-shot 3d visual grounding via vision-language  
713 reasoning on scene graphs. In *Proceedings of the AAAI Con-*  
714 *ference on Artificial Intelligence (AAAI)*, 2026. 6
- [33] Junyu Luo, Jiahui Fu, Xianghao Kong, Chen Gao, Haibing  
715 Ren, Hao Shen, Huaxia Xia, and Si Liu. 3d-sps: Single-stage  
716 3d visual grounding via referred point progressive selection.  
717 In *2022 IEEE/CVF Conference on Computer Vision and Pat-*  
718 *tern Recognition (CVPR)*. IEEE, 2022. 3
- [34] Xiaojian Ma, Silong Yong, Zilong Zheng, Qing Li, Yi-  
719 tao Liang, Song-Chun Zhu, and Siyuan Huang. Sqa3d:  
720 Situated question answering in 3d scenes. *arXiv preprint*  
721 *arXiv:2210.07474*, 2022. 5, 7, 12
- [35] OpenAI, Aaron Hurst, Adam Lerer, Adam P. Goucher, Adam  
722 Perelman, Aditya Ramesh, and Aidan Clark et al. Gpt-4o  
723 system card, 2024. 2, 6, 14
- [36] Junha Roh, Karthik Desingh, Ali Farhadi, and Dieter  
724 Fox. Languagerefer: Spatial-language model for 3d visual  
725 grounding. In *Conference on Robot Learning*, pages 1046–  
726 1056. PMLR, 2022. 14
- [37] David Rozenberszki, Or Litany, and Angela Dai. Language-  
727 grounded indoor 3d semantic segmentation in the wild. In  
728 *European Conference on Computer Vision*, pages 125–141.  
729 Springer, 2022. 1, 5, 7, 12
- [38] Gabriel Sarch, Snigdha Saha, Naitik Khandelwal, Ayush  
730 Jain, Michael J. Tarr, Aviral Kumar, and Katerina Fragki-  
731 adaki. Grounded reinforcement learning for visual reason-  
732 ing, 2025. 8
- [39] Jonas Schult, Francis Engelmann, Alexander Hermans, Or  
733 Litany, Siyu Tang, and Bastian Leibe. Mask3d: Mask trans-  
734 former for 3d semantic instance segmentation. In *2023*  
735 *IEEE International Conference on Robotics and Automation*  
736 *(ICRA)*, pages 8216–8223. IEEE, 2023. 7, 15
- [40] Simranjit Singh, Georgios Pavlakos, and Dimitrios Stam-  
737 moulis. Evaluating zero-shot gpt-4v performance on 3d vi-  
738 sual question answering benchmarks, 2024. 6
- [41] Rohan Taori, Ishaan Gulrajani, Tianyi Zhang, Yann Dubois,  
739 Xuechen Li, Carlos Guestrin, Percy Liang, and Tatsunori B.  
740 Hashimoto. Stanford alpaca: An instruction-following llama  
741 model. [https://crfm.stanford.edu/2023/03/](https://crfm.stanford.edu/2023/03/13/alpaca.html)  
742 [13/alpaca.html](https://crfm.stanford.edu/2023/03/13/alpaca.html), 2023. Dataset and model release. 5,  
743 12
- [42] Gemini Team, Rohan Anil, Sebastian Borgeaud, Jean-  
744 Baptiste Alayrac, Jiahui Yu, and Radu Soricut et al. Gemini:  
745 A family of highly capable multimodal models, 2025. 2
- [43] Ruicheng Wang, Sicheng Xu, Cassie Dai, Jianfeng Xiang,  
746 Yu Deng, Xin Tong, and Jiaolong Yang. Moge: Unlocking  
747 accurate monocular geometry estimation for open-domain  
748 images with optimal training supervision. *arXiv preprint*  
749 *arXiv:2410.19115*, 2024. 5
- [44] Yubo Wang, Xueguang Ma, Ge Zhang, Yuansheng Ni,  
750 Abhranil Chandra, Shiguang Guo, Weiming Ren, Aaran  
751 Arulraj, Xuan He, Ziyang Jiang, et al. Mmlu-pro: A more  
752 robust and challenging multi-task language understanding  
753 benchmark. *arXiv preprint arXiv:2406.01574*, 2024. 12, 14
- [45] X.AI. Realworldqa. Blog post, 2024. Accessed: 2025-05-12.  
754 6, 7
- [46] An Yang, Baosong Yang, Beichen Zhang, Binyuan Hui, and  
755 Bo Zheng et al. Qwen2.5 technical report. *arXiv preprint*  
756 *arXiv:2412.15115*, 2024. 14
- [47] Zhengyuan Yang, Songyang Zhang, Liwei Wang, and Jiebo  
757 Luo. Sat: 2d semantics assisted training for 3d visual  
758 grounding. In *Proceedings of the IEEE/CVF International*  
759 *Conference on Computer Vision*, pages 1856–1866, 2021. 14
- [48] Zhao Yang, Jiaqi Wang, Yansong Tang, Kai Chen, Heng-  
760 shuang Zhao, and Philip H. S. Torr. Lavt: Language-aware  
761 vision transformer for referring image segmentation, 2022. 7
- [49] Chandan Yeshwanth, Yueh-Cheng Liu, Matthias Nießner,  
762 and Angela Dai. Scannet++: A high-fidelity dataset of 3d in-  
763 door scenes. In *Proceedings of the IEEE/CVF International*  
764 *Conference on Computer Vision*, pages 12–22, 2023. 7
- [50] Zhihao Yuan, Xu Yan, Yinghong Liao, Ruimao Zhang,  
765 Sheng Wang, Zhen Li, and Shuguang Cui. Instancerefer:  
766 Cooperative holistic understanding for visual grounding on  
767 point clouds through instance multi-level contextual refer-  
768 ring. In *Proceedings of the IEEE/CVF International Confer-*  
769 *ence on Computer Vision*, pages 1791–1800, 2021. 14
- [51] Yiming Zhang, ZeMing Gong, and Angel X Chang.  
770 Multi3drefer: Grounding text description to multiple 3d ob-  
771 jects. In *Proceedings of the IEEE/CVF International Con-*  
772 *ference on Computer Vision*, pages 15225–15236, 2023. 13

- 797 [52] Duo Zheng, Shijia Huang, Lin Zhao, Yiwu Zhong, and Li-  
798 wei Wang. Towards learning a generalist model for embod-  
799 ied navigation. In *Proceedings of the IEEE/CVF Conference*  
800 *on Computer Vision and Pattern Recognition*, pages 13624–  
801 13634, 2024. 6, 14
- 802 [53] Duo Zheng, Shijia Huang, and Liwei Wang. Video-3d llm:  
803 Learning position-aware video representation for 3d scene  
804 understanding. In *Proceedings of the IEEE/CVF Conference*  
805 *on Computer Vision and Pattern Recognition (CVPR)*, pages  
806 8995–9006, 2025. 1, 3, 5, 6, 7, 12, 13, 14
- 807 [54] Chenming Zhu, Tai Wang, Wenwei Zhang, Jiangmiao Pang,  
808 and Xihui Liu. Llava-3d: A simple yet effective pathway to  
809 empowering llms with 3d-awareness, 2024. 1, 3, 4, 5, 6, 7,  
810 12, 14
- 811 [55] Ziyu Zhu, Xiaojian Ma, Yixin Chen, Zhidong Deng, Siyuan  
812 Huang, and Qing Li. 3d-vista: Pre-trained transformer for 3d  
813 vision and text alignment. In *Proceedings of the IEEE/CVF*  
814 *International Conference on Computer Vision*, pages 2911–  
815 2921, 2023. 1, 3, 6, 14
- 816 [56] Ziyu Zhu, Zhuofan Zhang, Xiaojian Ma, Xuesong Niu, Yixin  
817 Chen, Baoxiong Jia, Zhidong Deng, Siyuan Huang, and Qing  
818 Li. Unifying 3d vision-language understanding via prompt-  
819 able queries. *arXiv preprint arXiv:2405.11442*, 2024. 3, 6,  
820 7, 14
- 821 [57] Xueyan Zou, Zi-Yi Dou, Jianwei Yang, Zhe Gan, Linjie Li,  
822 Chunyuan Li, Xiyang Dai, Harkirat Behl, Jianfeng Wang, Lu  
823 Yuan, et al. Generalized decoding for pixel, image, and lan-  
824 guage. In *Proceedings of the IEEE/CVF Conference on Com-*  
825 *puter Vision and Pattern Recognition*, pages 15116–15127,  
826 2023. 7

827  
828  
829  
830  
831  
832  
833  
834  
835  
836  
837  
838  
839  
840  
841  
842  
843  
844  
845  
846  
847  
848  
849  
850  
851  
852  
853  
854  
855  
856  
857  
858  
859  
860  
861  
862  
863  
864  
865  
866  
867  
868  
869  
870  
871  
872  
873  
874

## 7. Appendix

### 7.1. Sensitivity to Pose and Depth Noise

While Qwen-3D already operates on real-world sensor noise in all experiments reported in the paper, we further stress-test the model under controlled settings with substantial depth and camera pose noise.

Following UniVLG [23], we simulate depth noise by injecting Gaussian noise with increasing variance into the raw depth maps before unprojection. To simulate additional pose noise, we add Gaussian noise with increasing variance to both the translation and rotation components of the provided camera poses.

We evaluate on ScanRefer (Top-1@0.25) against UniVLG, the strongest reported baseline. We choose UniVLG as the primary baseline because many strong 3D LMMs, such as Video-3D-LLM [53], adopt two-stage pipelines that rely on external object detection models and operate on pre-processed detections at inference time. Fairly evaluating such methods under additional noise would require these external detectors to be re-run on the perturbed point clouds, making comparisons difficult to standardize. In contrast, existing single-stage 3D LMMs either substantially underperform our model or do not release the necessary visual grounding code (e.g., LLaVA-3D [54]). Nevertheless, UniVLG remains the strongest reported baseline on ScanRefer, and we therefore report comparisons against it.

As shown in Fig. 4, Qwen-3D avoids catastrophic failure even under significant noise. Specifically, Qwen-3D demonstrates high robustness to depth noise, showing no significant performance drop even under high variance noise, matching UniVLG’s resilience to spurious points. Qwen-3D is similarly robust to pose noise, degrading gracefully and even outperforming the state-of-the-art model under extreme misalignment. We attribute this robustness to the strong integration of 2D pretrained vision-language priors from the Qwen backbone into explicit 3D structures, allowing the model to project reasoning over 2D inputs into accurate 3D masks even when the underlying 3D structures are corrupted.

### 7.2. Additional 3D Grounding Results

We report detailed 3D grounding evaluations on the SR3D, NR3D [1], and ScanRefer [10] benchmarks in Tab. 6. To our knowledge, no prior 3D LLM-based models report on the SR3D and NR3D benchmarks. Qwen-3D outperforms all prior LLM methods on ScanRefer using only a 3 billion parameter backbone, despite many of the strongest prior methods utilizing more powerful backbones of 7 billion parameters or more.

### 7.3. Text-only Results

We jointly train our model on visual reasoning tasks alongside the pure text instruction fine-tuning dataset, Alpaca [41], and evaluate Qwen-3D’s instruction-following capabilities on the MMLU-Pro [44] benchmark. As shown in Tab. 7, our model performs competitively with Qwen2.5-VL-3B [5] on general instruction-following tasks. This demonstrates no significant degradation of the original backbone’s text generation capabilities, despite our adaptations from causal attention masking to full attention over the inputs.

### 7.4. End-to-end Latency

We estimate the end-to-end latency of Qwen-3D on 3D tasks, tracking from raw scene inputs to the final 3D grounding outputs. As noted in Section 3 of the main text, doing an inference on the full set of multi-view RGB-D frames through Qwen-3D takes approximately 2 seconds. This includes the time required to unproject RGB-D frames (20 ms) and voxelize the 3D point cloud (1 ms). Standard RGB-D SLAM reconstruction methods [15] used for ScanNet scenes add about 1 to 5 seconds, yielding a total end-to-end latency of roughly 3 to 7 seconds. Notably, SLAM reconstruction costs can be amortized per scene as the reconstruction only needs to be run once; subsequent queries to the same environment only require the 2-second model forwarding time. Additionally, we note that all 3D LMMs need to run SLAM to obtain their point clouds; Qwen-3D does not introduce any new dependencies over prior methods.

### 7.5. Detailed Comparison of Training Datasets

Qwen-3D is trained on a mix of 2D and 3D datasets. The 3D datasets comprises approximately 255k samples spanning 3D question answering [4, 34], referential grounding [1, 10], and instance segmentation [37]. The 2D datasets include referential grounding (RefCOCO, RefCOCO+, RefCOCOg [25]), instance segmentation (COCO [30]), captioning and question answering (LLaVA-Instruct-150k [31]), and instruction fine-tuning (Alpaca [41]). Notably, all of these 2D datasets are included in the pre-training dataset of many modern 2D VLMs, including Qwen2.5VL, the backbone we use. We include them to avoid catastrophic forgetting of 2D capabilities as our model’s VLM backbone as it trains on new 3D data.

Our most competent expert model baseline, UniVLG [23], is trained on a nearly identical dataset mixture, excluding only LLaVA-Instruct-150k and the text-only Alpaca dataset. In contrast, the state-of-the-art single-stage 3D LMM, LLaVA-3D [54] utilizes a much larger 3D fine-tuning dataset, compiling over 860K 3D visual-reasoning samples together from various benchmarks alongside the 2D LLaVA-Instruct-150k dataset. The leading two-stage

875  
876  
877  
878  
879  
880  
881  
882  
883  
884  
885  
886  
887  
888  
889  
890  
891  
892  
893  
894  
895  
896  
897  
898  
899  
900  
901  
902  
903  
904  
905  
906  
907  
908  
909  
910  
911  
912  
913  
914  
915  
916  
917  
918  
919  
920  
921  
922  
923  
924  
925

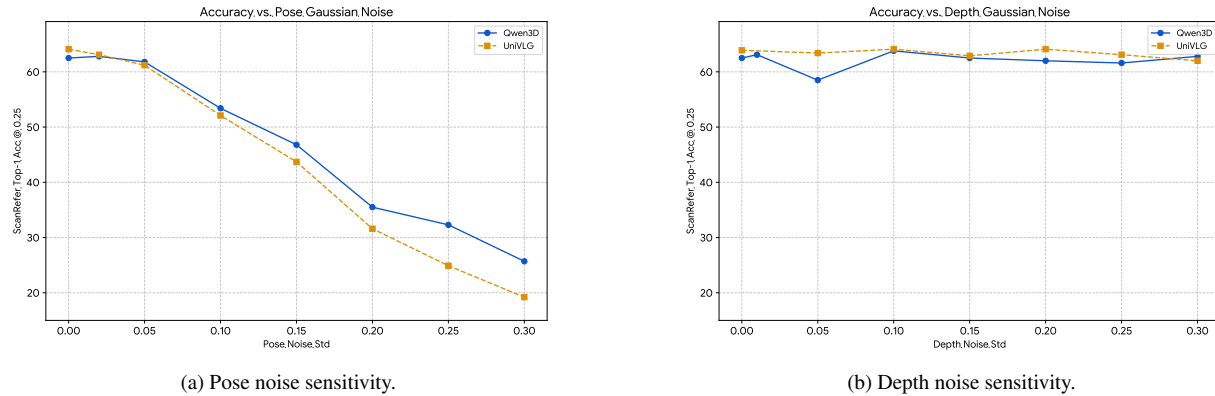


Figure 4. Sensitivity analysis to camera pose noise (left) and depth noise (right) on ScanRefer (Top-1@0.25). We compare Qwen-3D against state-of-the-art expert model UniVLG [23]

926 3D LMM, Video-3D-LLM [53], uses a 3D training mixture similar to ours - substituting Referit3D [1] with Multi3DRefer [51]. However, it also relies heavily on external object detectors, which are typically trained on Scan-Net200 instance segmentation datasets.

## 931 7.6. Rotary Position Embedding Implementation

932 To process interleaved visual and textual data in the vision-attention layers, Qwen2.5-VL [5] utilizes a Multimodal Rotary Positional Embedding (mRoPE). This method decomposes positional encodings into temporal, height, and width components. Each token  $i$  in a sequence of length  $K$  is assigned a position ID in each dimension, and the final positional embedding is the concatenation of the sinusoidal encodings of these components:

$$940 \quad \text{PE}(\mathbf{p}) = \text{PE}(t_i; h_i; w_i)_{i=1}^K,$$

941 For text tokens, the temporal ID  $t$  increments linearly, while  $h$  and  $w$  are assigned the same values as  $t$ , effectively reducing mRoPE to 1D-RoPE. For image tokens,  $t$  is held constant while  $(h, w)$  are assigned to corresponding pixel coordinates. To maintain temporal consistency across all modalities and dimensions, the position IDs of each modality are incremented by the maximum position ID of the preceding modality, ensuring visual features and text instructions occupy distinct, non-overlapping positions in the positional embeddings.

951 **2D RoPE Projection:** This extension aims to stay as close as possible to the Qwen2.5-VL pre-training distribution. Temporal IDs follow the original scheme, while the height and width IDs for voxelized point-cloud tokens are obtained by sampling corresponding 2D pixel coordinates from one of the views where the voxel is visible.

957 **Naive 3D MRoPE:** This extension aims to inject absolute 3D spatial information by replacing the height and width

components with embeddings of absolute XYZ coordinates: 959

$$960 \quad \text{PE}(\mathbf{p}) = \text{PE}(t_i; x_i; y_i; z_i)_{i=1}^K, \quad 960$$

961 This method is similar to our final 3D RoPE, except it utilizes Qwen2.5-VL’s partitioning of the embedding channels among the four dimensions, resulting in an imbalanced frequency spectrum across the four dimensions, with high-frequency channels being assigned to the temporal axis and low-frequency channels being assigned to the z-axis. We find that this degrades the model’s 3D grounding capabilities in our ablations.

969 Our 3D RoPE instead initializes the same range of frequencies for each dimension, resulting in a uniform distribution of high and low frequency channels across all components. We find that this method significantly outperforms all other variants on 3D grounding tasks.

## 974 7.7. Additional Hyperparameters

975 Hungarian Matching cost weights during training have a sizable impact on our model’s performance on 3D grounding tasks. Our final configuration utilizes the standard cost coefficients from Mask2Former [13], as we found that alternative weighting schemes proposed in other 3D grounding works [23] can hinder performance.

## 981 7.8. Comparison with Token Merging

982 While Qwen-3D utilizes voxelization pooling to reduce multi-view redundancy, a popular alternative for reducing redundancy in video model architectures is token merging [6], which iteratively fuses visually similar tokens within ViT attention layers. We attempted to compare these methods to our voxelization pooling by replacing the voxelization with a plug-and-play token merging strategy within our vision-language attention layers [6]. This implementation proved infeasible for dense 3D scenes, as

Table 6. **Results on 3D language grounding for both experts and LMMs.** We evaluate top-1 accuracy on the official validation set. \* UniVLG numbers reproduced using the official code with a single 40GB 8 GPU node instead of their 32 GPU setup, and verified with the authors.

Methods	SR3D				NR3D				ScanRefer		
	Acc @25 (Det)	Acc @50 (Det)	Acc @75 (Det)	Acc (GT)	Acc @25 (Det)	Acc @50 (Det)	Acc @75 (Det)	Acc (GT)	Acc @25 (Det)	Acc @50 (Det)	Acc @75 (Det)
ReferIt3DNet [1]	27.7	-	-	39.8	24.0	-	-	-	26.4	16.9	-
ScanRefer [10]	-	-	-	-	-	-	-	-	35.5	22.4	-
InstanceRefer [50]	31.5	-	-	48.0	29.9	-	-	-	40.2	32.9	-
LanguageRefer [36]	39.5	-	-	56.0	28.6	-	-	-	-	-	-
SAT-2D [47]	35.4	-	-	57.9	31.7	-	-	-	44.5	30.1	-
BUTD-DETR [21]	52.1	-	-	67.0	43.3	-	-	54.6	52.2	39.8	-
<b>Experts</b> PQ3D [56]	62.0	55.9	46.2	79.7	52.2	45.0	37.6	66.7	56.7	51.8	43.3
3D-VisTA [55]	56.5	51.5	42.8	76.4	47.7	42.2	35.5	65.1	51.0	46.2	36.7
ODIN [22]	38.1	29.3	23.1	-	31.6	20.8	15.8	-	43.1	33.4	26.2
Locate-3D [3]	65.8	52.9	-	-	53.7	40.5	-	-	59.9	49.6	-
Locate-3D+ [3]	68.2	54.8	-	-	56.1	43.2	-	-	61.1	50.9	-
UniVLG (1 GPU node) [23]*	71.7	63.6	51.2	-	52.8	42.0	33.0	-	64.1	52.7	43.5
UniVLG [23]	<b>73.0</b>	<b>64.8</b>	<b>51.8</b>	-	<b>58.3</b>	<b>49.8</b>	<b>39.1</b>	-	<b>63.5</b>	<b>56.4</b>	<b>46.0</b>
<b>2D VLMs*</b>											
LLaMA [2]	21.3	13.9	-	-	28.0	16.9	-	-	37.3	24.2	-
Qwen2-VL-2B [46]	-	-	-	-	31.1	-	-	-	35.9	31.9	-
Qwen2-VL-72B [46]	-	-	-	-	47.6	-	-	-	44.8	40.3	-
GPT-4o [35]	29.2	18.9	-	-	38.2	25.1	-	-	48.2	32.5	-
<b>Two-Stage</b>											
ChatScene [19]	-	-	-	-	-	-	-	-	55.5	50.2	-
LEO [20]	-	-	-	-	-	-	-	-	-	-	-
<b>LMMs</b> Video-3D LLM - 7B [53]	-	-	-	-	-	-	-	-	58.1	51.7	-
<b>Single-Stage</b>											
NaviLLM [52]	-	-	-	-	-	-	-	-	-	-	-
LL3DA [11]	-	-	-	-	-	-	-	-	-	-	-
SceneLLM [16]	-	-	-	-	-	-	-	-	-	-	-
3D-LLM [17]	-	-	-	-	-	-	-	-	30.3	-	-
Grounded 3D-LLM [12]	-	-	-	-	-	-	-	-	47.9	44.1	-
LLaVA-3D - 7B [54]	-	-	-	-	-	-	-	-	50.1	42.7	-
Qwen-3D(Ours)	<b>61.8</b>	<b>55.8</b>	<b>45.6</b>	-	<b>55.1</b>	<b>48.7</b>	<b>39.5</b>	-	<b>65.6</b>	<b>57.7</b>	<b>45.0</b>

Table 7. **Qwen-3D on instruction-following benchmark MMLU-Pro [44]**

Model	EM@1
Qwen2.5-VL-72B	71.2
Qwen2.5-VL-3B	57.1
Qwen-3D (Ours)	55.6

991 with token merging, the first vision-language attention layer  
 992 must still process the point cloud at full resolution (aver-  
 993 aging around 33K points for ScanNet scenes). This con-  
 994 sistentlly resulted in GPU out-of-memory errors during in-  
 995 ference, regardless of how aggressively the token merging  
 996 would downsample in subsequent layers. Conversely, our

voxelization pooling preemptively downsamples the point 997  
 cloud, reducing the number of points to an average of about 998  
 14K points prior to the attention layers. This acts as a nat- 999  
 ural, 3D-aware compression mechanism that bypasses the 1000  
 memory overhead of full-resolution multi-view attention. 1001

Table 8. **Mask Decoder Text-Query Cross Attention Ablations**

Mask Decoder	Top1 Acc
Q->V cross-attn	<b>52.5</b>
Q->V cross-attn	44.3

## 7.9. Additional Implementation Details 1002

For our model, we use Qwen2.5-VL as the backbone due 1003  
 to its open-source implementation and its tight integration 1004

1005 of visual and textual features in its visual-language atten-  
1006 tion layers. Within our codebase, we also introduce several  
1007 engineering modifications to the original Qwen2.5-VL ar-  
1008 chitecture.

1009 The Qwen2.5-VL vision encoder handles the batching  
1010 of multiple images by sequentially concatenating images  
1011 along the token dimension and using attention masks to  
1012 prevent cross-image interaction, resulting in  $O(N^2)$  com-  
1013 plexity in the vision encoder with respect to the number of  
1014 images  $N$ . We optimize this by implementing batched for-  
1015 warding in the attention mechanism to perform per-image  
1016 batched attention. This modification reduces the complex-  
1017 ity from  $O(N^2)$  to  $O(N)$  while preserving the original be-  
1018 havior of the vision encoder. This also substantially reduces  
1019 GPU memory usage, which is essential for processing scene  
1020 videos with large numbers of input RGB images.

1021 Additionally, Qwen2.5-VL utilizes causal masking in its  
1022 visual-language attention layers, allowing each vision and  
1023 text token to attend only to itself and to earlier tokens in the  
1024 sequence. This is problematic for our point cloud inputs be-  
1025 cause this operation is not permutation-invariant and would  
1026 cause the model to learn artifacts of the arbitrary ordering  
1027 of the point cloud sequences rather than the underlying ge-  
1028 ometry. We rectify this by replacing this causal-masked at-  
1029 tention with all-to-all attention for all text and visual inputs  
1030 and finetuning our backbone with causal masking only over  
1031 the answer tokens. We ablate these changes in the main text.

### 1032 7.10. Visualizations of Failure cases

1033 In Fig. 5, we outline three common failure modes for Qwen-  
1034 3D on 3D referential grounding tasks:

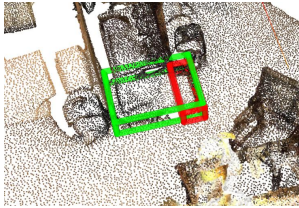
1035 **Incomplete object selection (Left):** The model occasion-  
1036 ally masks only a partial section of the target. This typically  
1037 occurs when the underlying point cloud contains significant  
1038 holes or artifacts, skewing the geometric understanding of  
1039 the complete object shape.

1040 **Confusion between instances (Middle):** The model cor-  
1041 rectly identifies the object class but grounds the wrong in-  
1042 stance. Confusion between multiple instances is a com-  
1043 mon failure mode seen in mask-decoding architectures, and  
1044 similar issues have been noted in previous works, such as  
1045 Mask3D [39] and UniVLG [23].

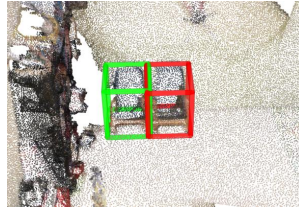
1046 **Language ambiguity (Right):** The model occasionally  
1047 confuses the target object with other reference objects men-  
1048 tioned in more ambiguous querying expressions.

### 1049 7.11. Visualizations of Qwen-3D Results

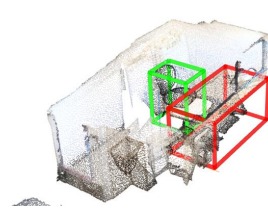
1050 We include visualizations of Qwen-3D predictions in 3D  
1051 referential grounding tasks in Fig. 6, 2D referential ground-  
1052 ing in Fig. 7, instance segmentation in Fig. 8, and visual  
1053 question-answering tasks in Fig. 9.



select the suitcase that is in the center of the bathtub and the desk



chair with back against the wall directly next to book shelf



there is a black desk chair in the center of the room. it is in front of the desk with a phone sitting on it

Figure 5. Failure cases of Qwen-3D on 3D grounding tasks. The red segmentation masks and boxes refer to Qwen-3D's prediction and the green masks and boxes indicate the ground truth.

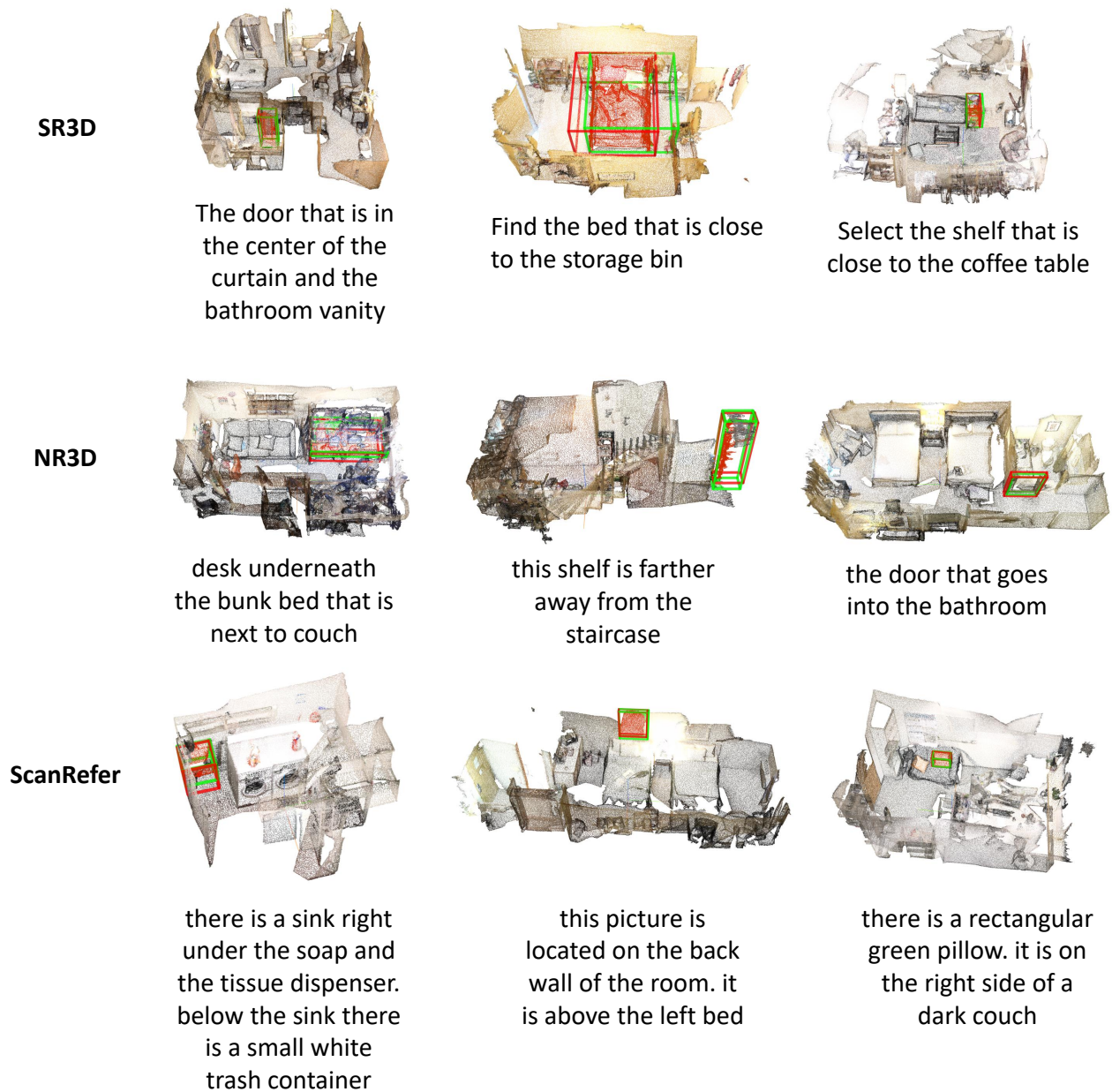


Figure 6. Visualizations of Qwen-3D’s predictions on **3D Referential Grounding Datasets SR3D, NR3D, and Scanrefer**. The red segmentation masks and boxes refer to Qwen-3D’s prediction and the green masks and boxes indicate the ground truth.



Figure 7. Visualizations of Qwen-3D’s predictions on 2D Referential Grounding Datasets RefCOCOg, RefCOCO+, and RefCOCO. The model’s predictions are indicated by the green mask and bounding box.



Figure 8. Visualizations of Qwen-3D’s instance segmentation predictions on COCO and ScanNet200.

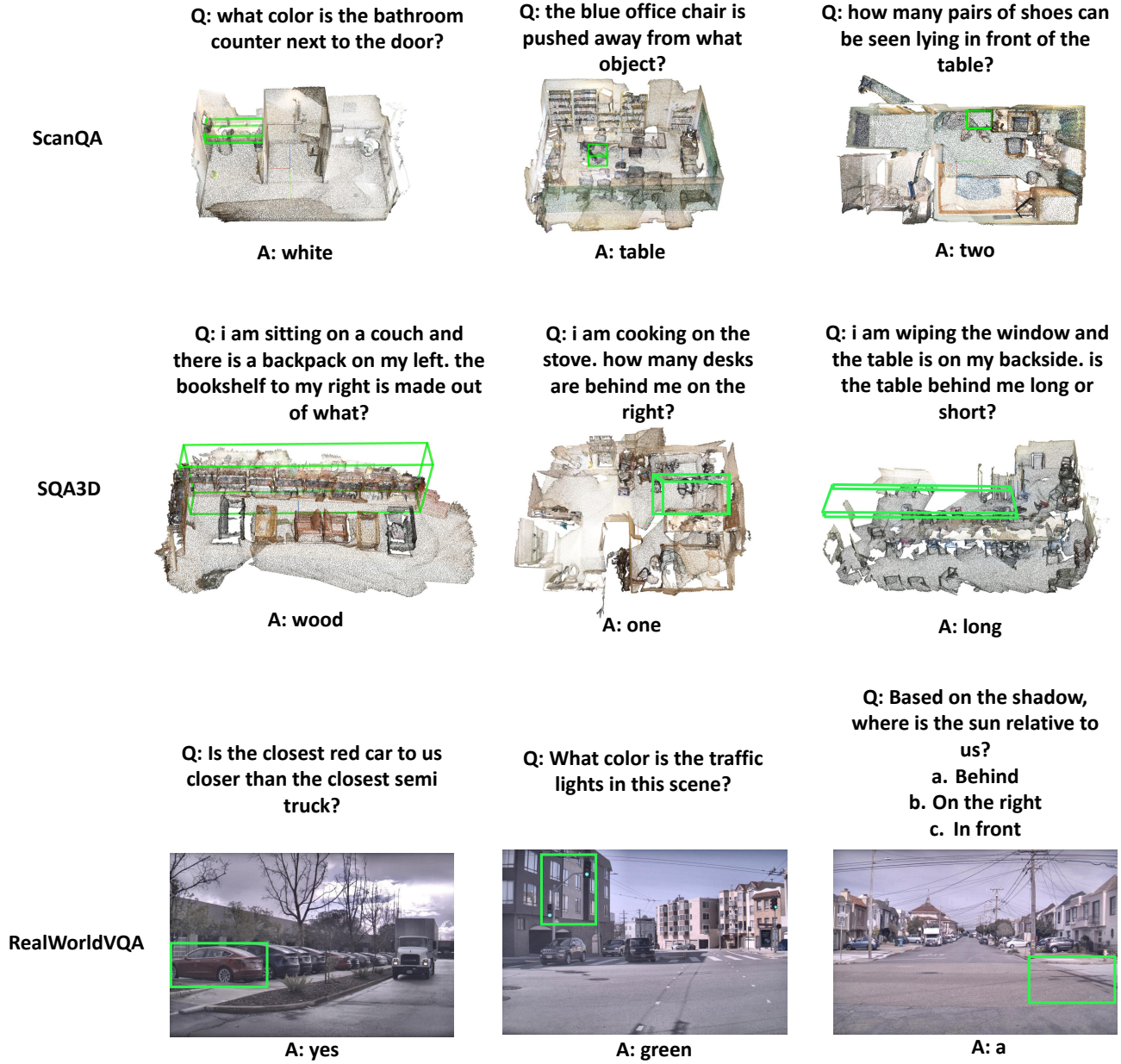


Figure 9. Visualizations of Qwen-3D’s responses to visual question-answering tasks on 3D benchmarks SQA3D and ScanRefer and the 2D benchmark RealWorldVQA. Green boxes denote objects relevant to the posed question.