



# Towards Visual Query Localization in the 3D World

Liang Peng<sup>1\*</sup>, Bohan Tan<sup>2\*</sup>, Zhipeng Zhang<sup>2,3 ‡</sup>, Haobo Li<sup>2</sup>, Yifan Jiao<sup>4,5</sup>, Xingping Dong<sup>1 †</sup>, Libo Zhang<sup>4,5</sup>

<sup>1</sup> Wuhan University <sup>2</sup> AutoLab, SAI, Shanghai Jiao Tong University <sup>3</sup> Anyverse Dynamics

<sup>4</sup> University of Chinese Academy of Sciences <sup>5</sup> Institute of Software, Chinese Academy of Sciences

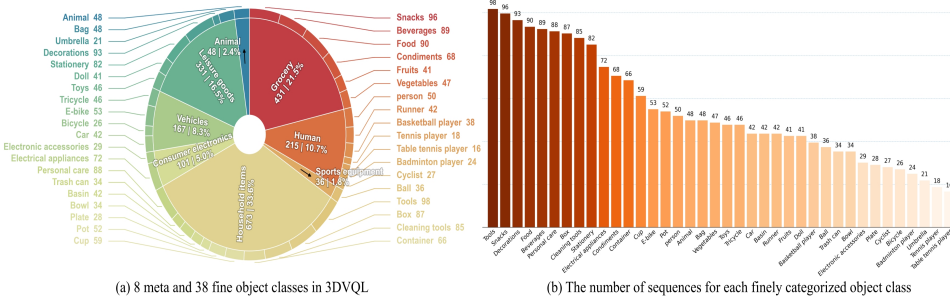
\*equal contributions; † first corresponding author; ‡ second corresponding author



## Contributions

1. A novel 3D visual query localization benchmark, 3DVQL, is introduced, covering RGB, depth, and point cloud data with fine-grained 9-DoF 3D bounding box annotations.
2. A query-guided 3D localization framework is proposed to promote future research on visual query localization in the 3D world.

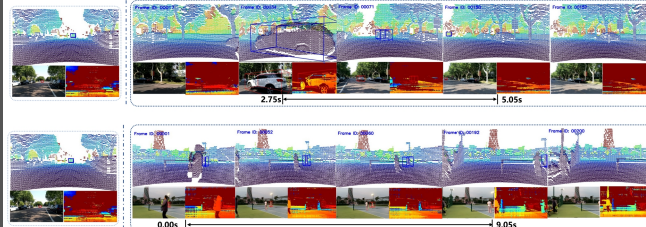
## Overview



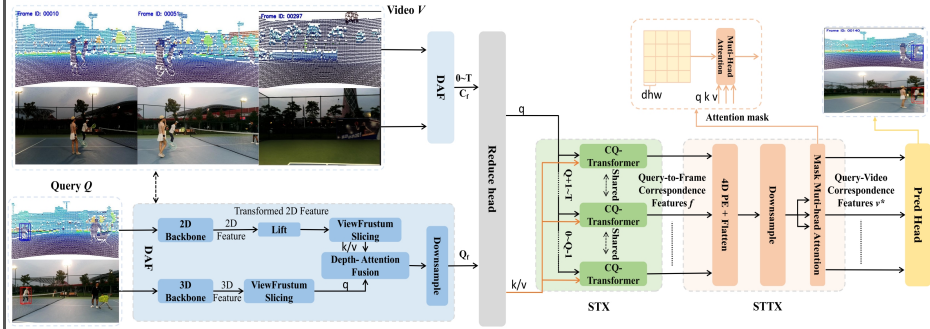
## 3DVQL Benchmark

We introduce **3DVQL**, the first benchmark for multimodal 3D visual query localization, with **2,002** sequences, **170K+** frames, **6.4K** response tracks, and **38** object categories.

3DVQL provides aligned point cloud, RGB, and depth modalities, supporting **PC-only**, **RGB-PC**, **RGB-D**, and other multimodal evaluation settings.



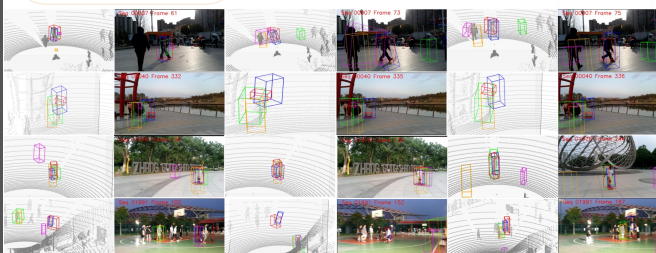
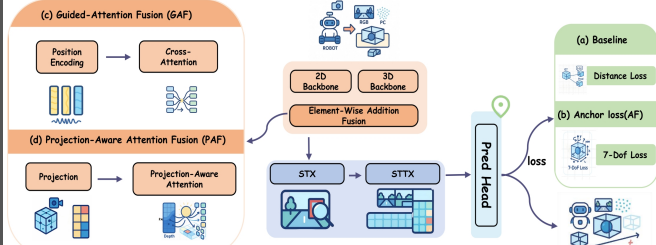
## LAF: Lift and Attention Fusion



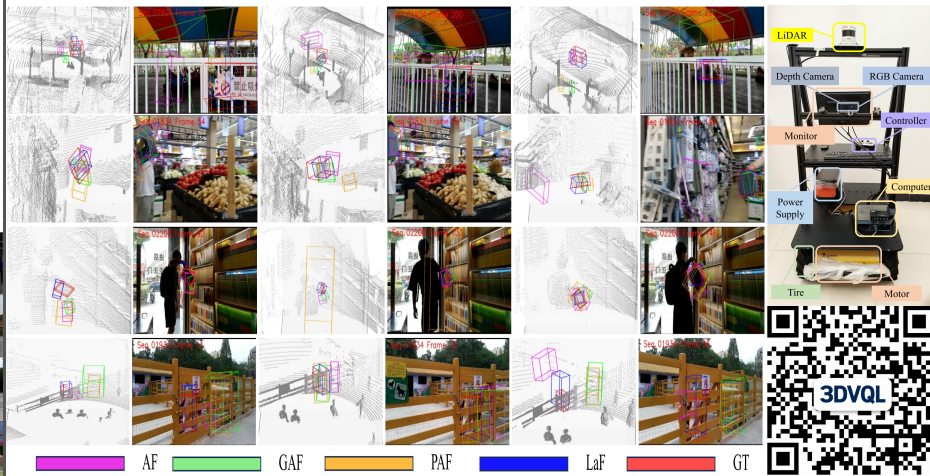
Method	tAP	tAP <sub>0.25</sub>	stAP	stAP <sub>0.05</sub>	rec.%	Succ.
AF	0.181	0.442	0.003	0.015	0.093	11.693
GAF	0.291	0.597	0.015	0.075	0.049	26.309
PAF	0.224	0.577	0.021	0.104	0.115	32.156
<b>LaF</b>	<b>0.293</b>	<b>0.607</b>	<b>0.044</b>	<b>0.222</b>	<b>0.264</b>	<b>46.041</b>

## Baselines

We establish three RGB-PC baselines for 3DVQL, including **AF** with anchor-based 3D box prediction, **GAF** with depth-guided attention fusion, and **PAF** with projection-aware **RGB-PC** feature fusion.



## Qualitative Results



Benchmark	Tot. Seq.	Ann. fr.	Avg trk.	Res trk.	Obj. cls.	Mod.		
						RGB	PC	D
2DVQL [10]	2,538	-	90	3.2k	-	✓	✗	✗
<b>3DVQL</b>	<b>2,002</b>	<b>170k</b>	<b>84</b>	<b>6.4k</b>	<b>38</b>	✓	✓	✓

