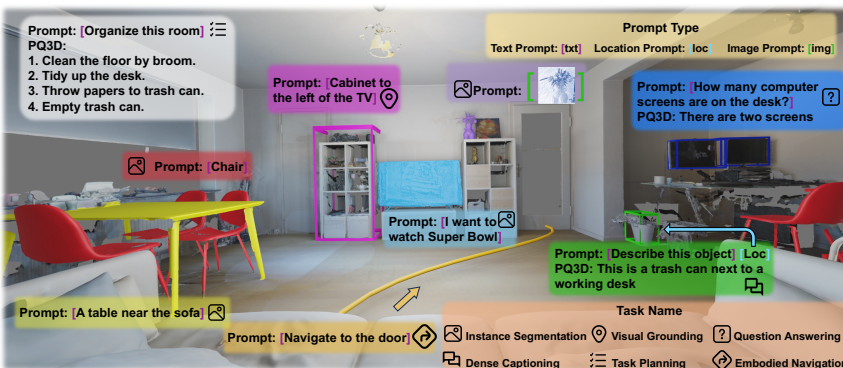


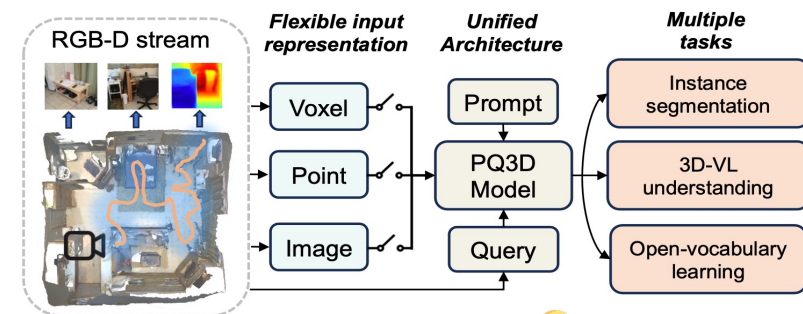
Contribution

- We introduce PQ3D, a **unified model** adept at solving a broad spectrum of 3D-VL tasks with **promptable queries**.
- PQ3D aligns *voxels*, *point clouds*, and *multi-view images* into a **shared 3D space** for joint training.
- PQ3D not only achieves **competitive results** but also sets new records across various 3D-VL tasks.



Highlight

Unifying multiple representations from RGB-D stream, including **Voxel**, **Point**, and **Multi-view images**.

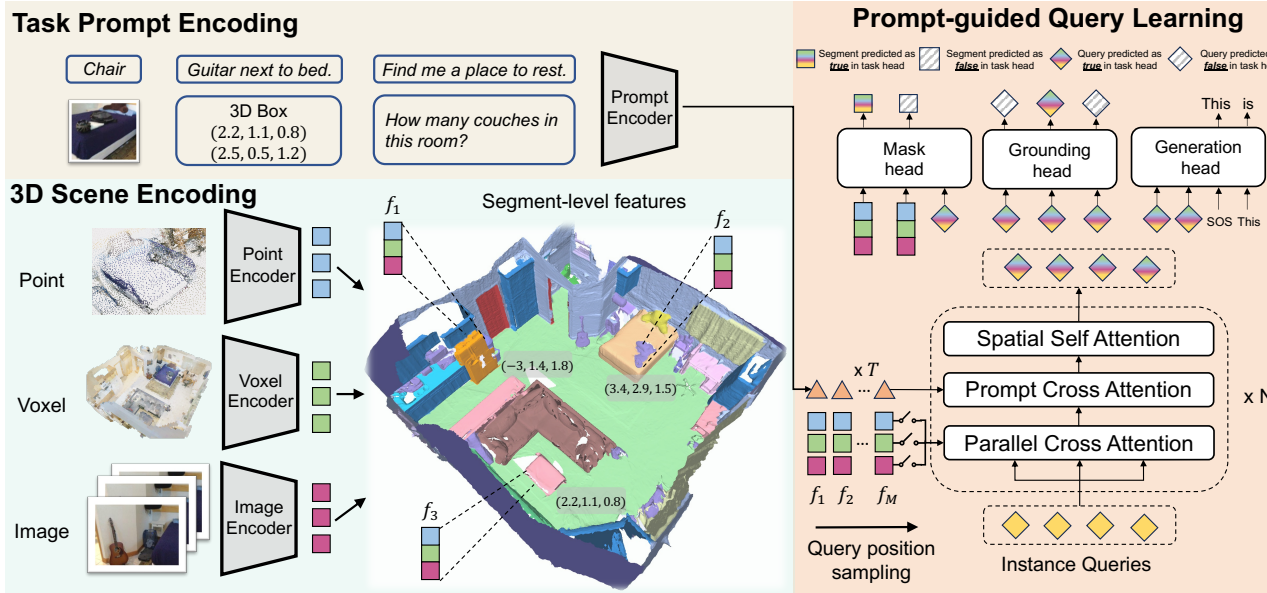


Unified training on **broad range** of tasks

Segmentation Refer
Captioning QA

PQ3D Model Design

- Diverse prompt and output format for multi-task learning**: Diverse prompts including text, image, and location are projected to a shared feature space. Instance queries can retrieve task-relevant information from prompts and be fed to different task heads for unified training.
- Leverage multiple representations in one model**: Point cloud, voxel grids, and multi-view images are aligned to segment level with shared coordinate space.



a) Task Prompt Encoding

CLIP text, image encoder for VL prompts
MLP for encoding location prompt

b) 3D Scene Encoding

Point cloud: PointNet++ feature for each point
Voxel: Apply SparseConvUNet to voxel grid.
Multi-view image: OpenSeg pixel feature

c) Prompt-guided Query Learning

Parallel cross attention between queries and features.
Cross attention between queries and prompts.
Different task heads for unified learning.

$$\text{Mask } p_{\text{mask}} = \sigma(f_s(\mathbf{V} + \mathbf{I} + \mathbf{P}) \cdot f_a(\mathbf{Q})^T)$$

$$\text{Grounding } p_{\text{grd}} = \sigma(f_g(\mathbf{Q}))$$

$$\text{Generation } T5 \text{ small autoregressive}$$

$$\mathbf{Q}'_i = \text{FFN}(\text{Norm}(\mathbf{Q}_i + \sum_{\mathbf{F} \in \{\mathbf{V}, \mathbf{I}, \mathbf{P}\}} \text{MaskedCrossAttn}(\mathbf{Q}_i, \mathbf{F})))$$

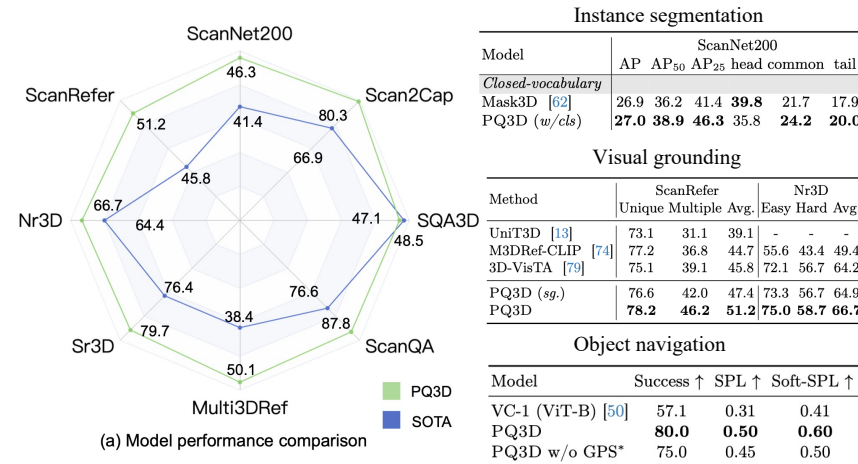
$$\mathbf{Q}''_i = \text{FFN}(\text{Norm}(\mathbf{Q}'_i + \text{CrossAttn}(\mathbf{Q}'_i, \mathbf{t})))$$

$$\mathbf{Q}_{i+1} = \text{FFN}(\text{Norm}(\text{SpatialSelfAttn}(\mathbf{Q}''_i)))$$

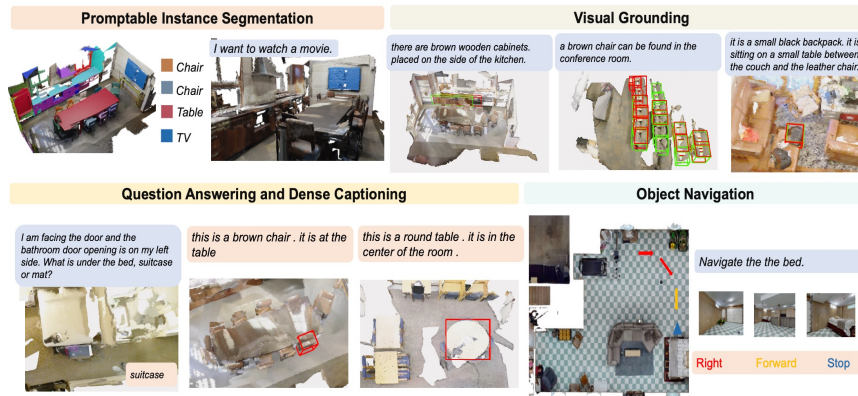
$$\text{Attn}(\mathbf{Q}, \mathbf{K}, \mathbf{V}, \mathbf{S}) = \text{softmax}\left(\frac{\mathbf{Q}\mathbf{K}^T}{\sqrt{d_h}} + \log \sigma(\mathbf{S}\mathbf{w})\right) \mathbf{V}$$

Results and Insights

a) PQ3D demonstrates **superior performance on most tasks, from low-level segmentation to high-level reasoning**.



b) PQ3D is a **versatile model** can take various prompts and solve a broad range of tasks in a flexible way.



c) **Unifying multiple representations and tasks benefits 3D vision language understanding.**

V	P	I	Refer	QA	Caption	Task Data	Refer	QA	Caption
✓			46.1 / 47.1	43.7 / 44.2	67.8 / 68.1	+Refer	-	1.8 ↑	2.2 ↑
✓	✓		49.2 / 49.4	45.4 / 45.8	74.6 / 74.7	+QA	0.0 ↑	-	1.2 ↑
✓	✓	✓	51.2	47.1	80.3	+Caption	0.6 ↑	0.7 ↑	-