# Dynam3D: Dynamic Layered 3D Tokens Empower VLM for Vision-and-Language Navigation

Bridging the gap between **Geometric Map** and **Semantic VLM** via dynamic hierarchical memory

**Zihan Wang, Seungjun Lee, Gim Hee Lee**

*School of Computing, National University of Singapore*

# The Navigation Dilemma

## 1. The "Video Tape" Approach

Standard Video-VLMs treat the world as a linear stream of frames.

⚠️ **Spatial Amnesia**: Video-based models rely on context windows. When an object leaves the frame, it is forgotten.

⚠️ **Geometry Blindness**: 2D video frames lack explicit 3D structure, leading to collisions and poor planning.



*"Please go to the kitchen and take the bread out of the microwave for me."*

Video-VLM → Action
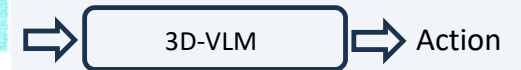
# The Navigation Dilemma

## 2. The "Frozen Map" Approach

Traditional 3D Maps assume a static world.

❄ **Static Assumptions**: Most mapping systems assume a static world, failing when object moves or environments change.

❄ **Granularity-Efficiency Conflict**: Dense representations (e.g., voxels) are computationally expensive for real-time reasoning, while sparse ones fail to capture fine-grained semantics for interaction.
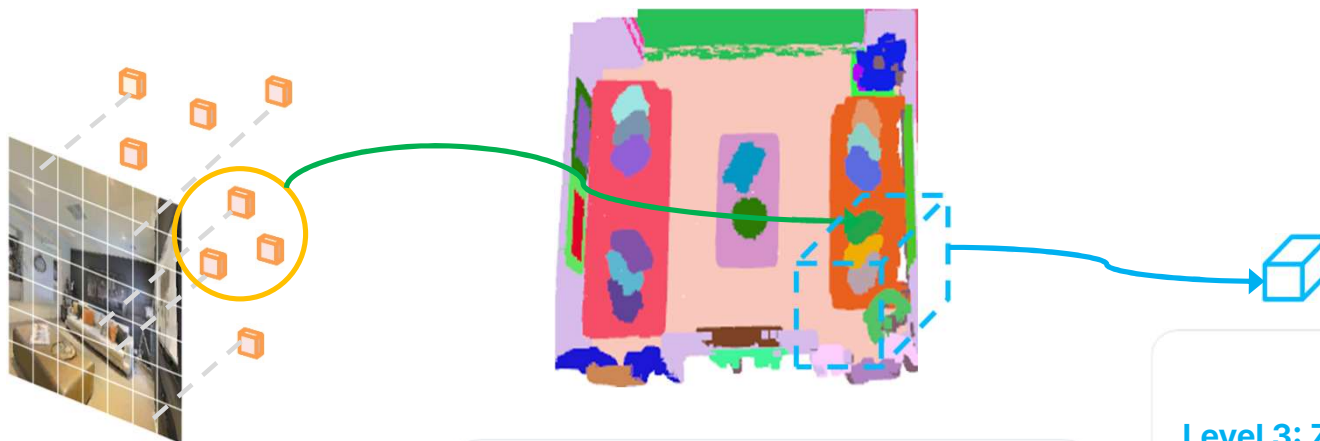


*"Please go to the kitchen and take the bread out of the microwave for me."*

3D-VLM → Action

# The Semantic Pyramid Tokenization

How do we compress a 1M-point world into a 1K-token VLM context window?



## Level 1: Patch

**Fine-grained Semantics**

Fine-grained semantics and geometry details from CLIP features.

Count: High

## Level 2: Instance

**Object Entities**

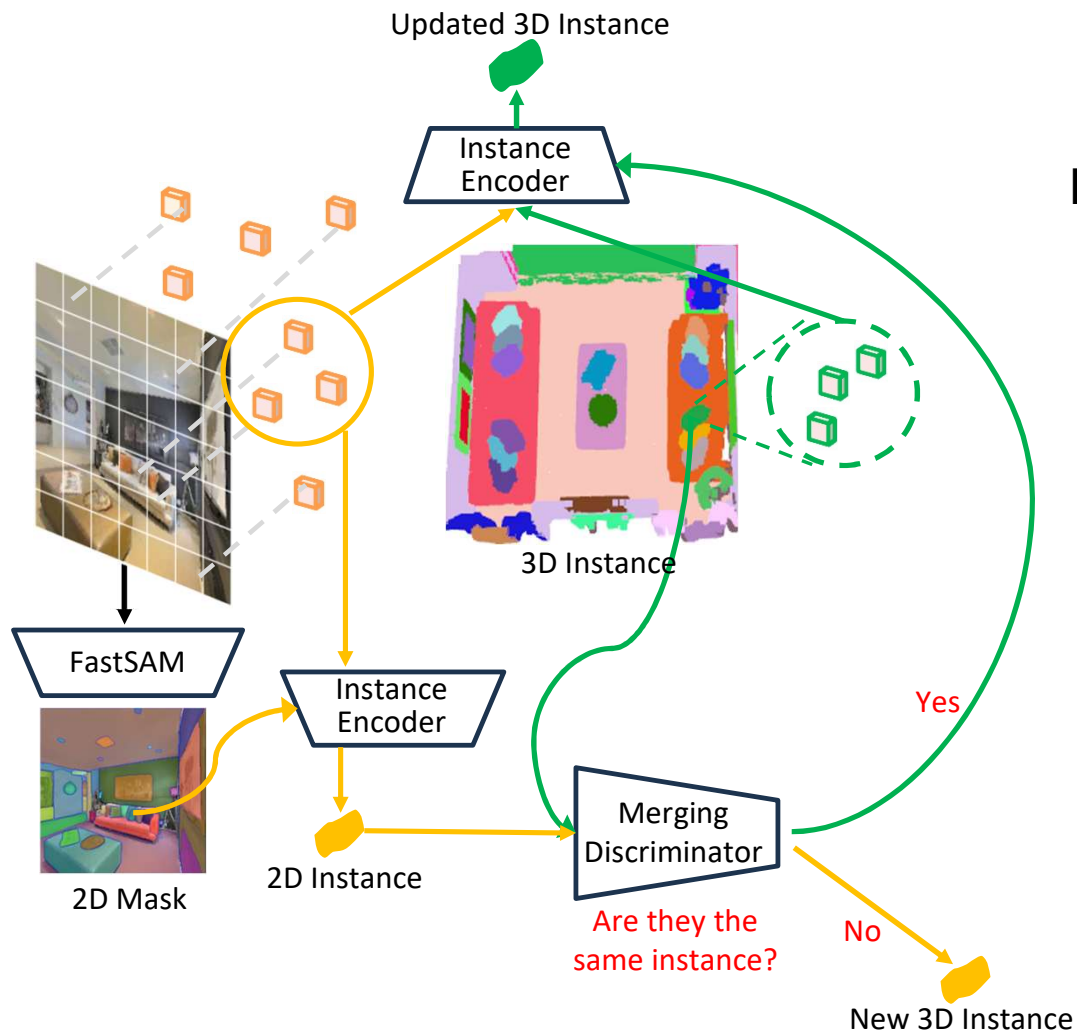3D objects (e.g., "Chair") aggregated from patches via FastSAM masks.

Count: Medium

## Level 3: Zone

**Spatial Layout**

High-level spatial regions (e.g., "Kitchen area") aggregated from objects for large-scale scene understanding.

Count: Very Low

# Online 3D Instance Construction



Updated 3D Instance

Instance Encoder

Instance Encoder

FastSAM

2D Mask

2D Instance

3D Instance

Merging Discriminator

Are they the same instance?

Yes

No

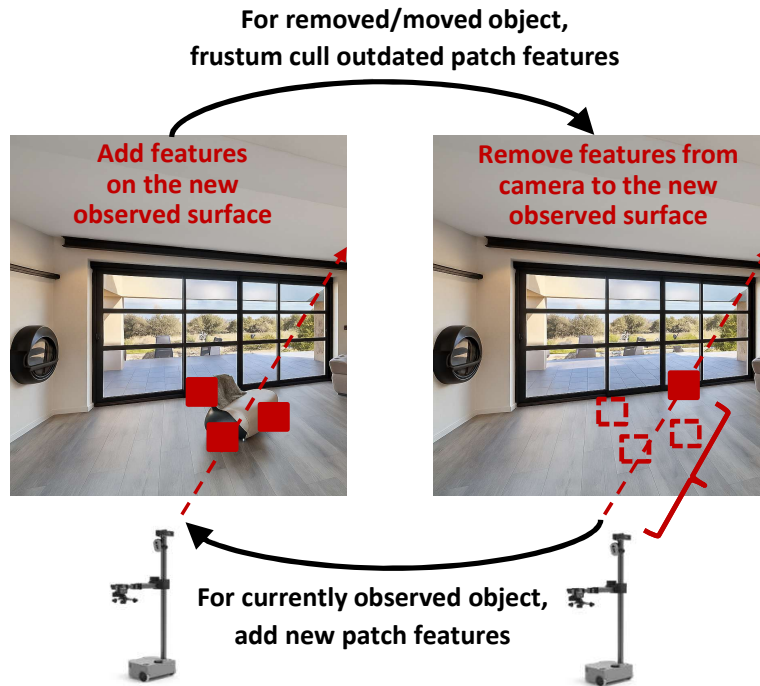New 3D Instance

## FastSAM + Merging Discriminator

1. Aggregate patches via 2D mask for 2D instance

2. Retrieve Top-K nearest existing 3D instances

3. A learned Merging Discriminator predicts if a 2D-3D pair is the same instance based on:

- Feature Similarity (Semantic)
- Euclidean Distance (Geometric)

4. Concatenate their patch features and update the 3D instance representation

# Adapt to the Dynamic World



Video or Static maps always fail here.
We need a map that "breathes".

# Dynamic Frustum Culling (Forget outdated information)

**For removed/moved object, frustum cull outdated patch features**

**Add features on the new observed surface**

**Remove features from camera to the new observed surface**

**For currently observed object, add new patch features**
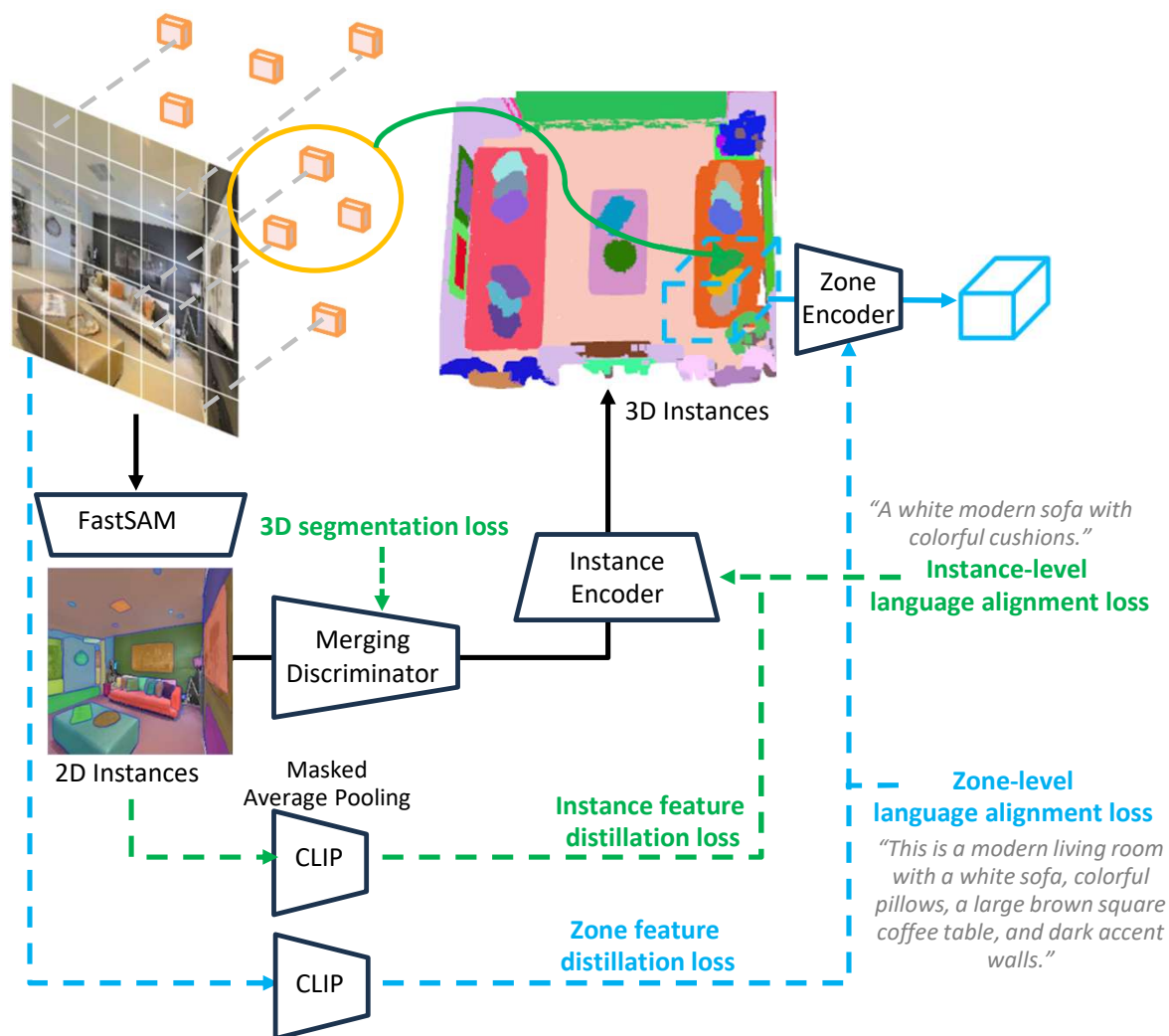
$$P_c^\top = \begin{bmatrix} x_c \\ y_c \\ z_c \end{bmatrix} = \mathbf{R}P_w^\top + \mathbf{T}, \quad \begin{bmatrix} u \\ v \\ 1 \end{bmatrix} = \frac{1}{z_c}\mathbf{K}\begin{bmatrix} x_c \\ y_c \\ z_c \end{bmatrix},$$

$$\text{FrustumCulling}(P_w), \text{ if } 0 < z_c < \min(d_{u,v} + \delta, \Delta), \ 0 < u < H, \text{ and } 0 < v < W.$$

- *Where $d_{u,v}$ is the observed depth. If patch $z_c$ is closer than the current observed surface, it will be removed.*

- *$\delta$ is a noise threshold and $\Delta$ is the farthest culling distance.*
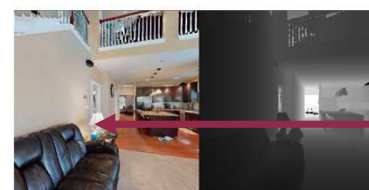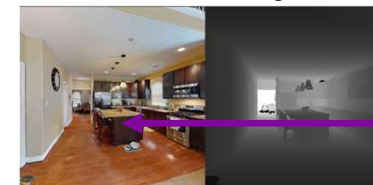
# Contrastive Learning for Semantic Alignment



FastSAM

2D Instances

Masked Average Pooling

CLIP

CLIP

Merging Discriminator

**3D segmentation loss**

Instance Encoder

Zone Encoder

3D Instances

*"A white modern sofa with colorful cushions."*

**Instance-level language alignment loss**

**Instance feature distillation loss**

**Zone-level language alignment loss**

*"This is a modern living room with a white sofa, colorful pillows, a large brown square coffee table, and dark accent walls."*

**Zone feature distillation loss**

**Instance ID**: 132
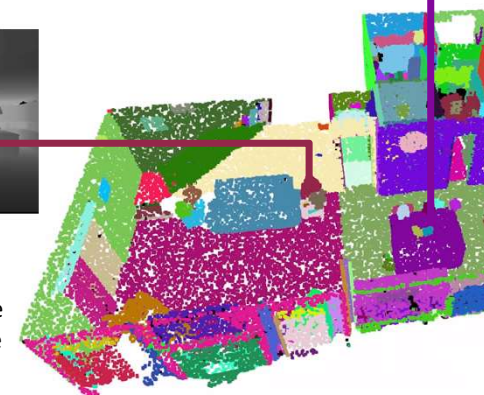**Object category**: dining table
**Language description**：  The dining table is in the kitchen, close to the refrigerator and sink.
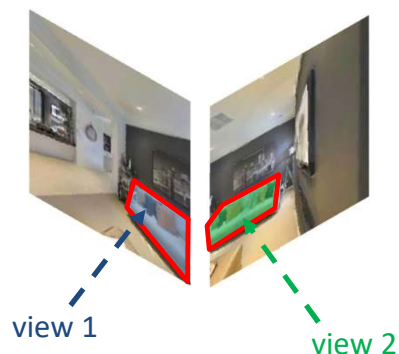
**Instance ID**: 568
**Object category**: table lamp
**Language description**：  A white table lamp sits on the side table next to the leather sofa.

1,883 Object categories, 5K+ 3D scenes, 2M+ language descriptions
*from SceneVerse, ScanNet, HM3D, Matterport3D, 3RScan, ARKitScenes, and Structured3D.*

# Subspace Contrastive Learning for 3D Consistency



view 1

view 2

📉 **The Challenge: View Inconsistency**
- Naive feature distillation is interfered by **background noise**
- Results in significant feature gaps for the same instance $O$ across different views
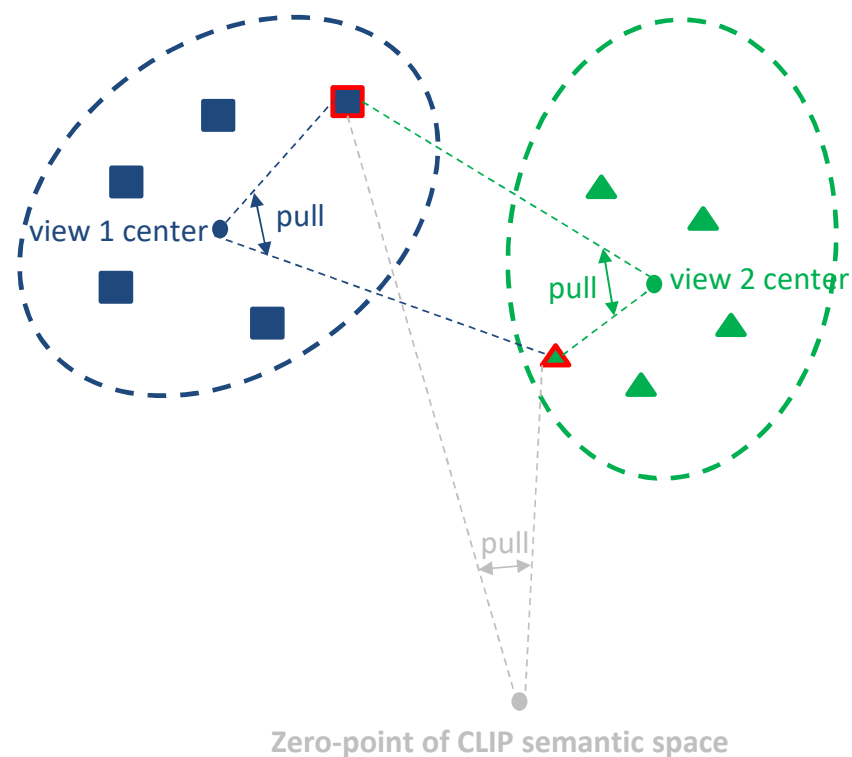
🛠 **The Method: Shift the Anchor**

$$\mathcal{L}_{subspace\_distillation} = \frac{1}{I} \sum_{i=1}^{I} \mathrm{CrossEntropy}(\{\mathrm{CosSim}(\,(\mathcal{O}_i - \mathcal{V}_j), (\mathcal{O}_j^{gt} - \mathcal{V}_j)\,)/\tau\}_{j=1}^{J}, i)$$

- **Calculate $V_j$:** The "Semantic Center" of the current view (average of all patches)
- **Subspace Alignment:** Optimization targets $(O - V_j)$ instead of absolute $O$

🚀 **The Effect: Bias Mitigation**
- Moves the anchor from **CLIP Origin** to the **View Center** $V_j$
- Effectively removes view-specific bias, enforcing stronger multi-view consistency

view 1 center ● pull

pull ● view 2 center

pull

Zero-point of CLIP semantic space

# The Brain: 3D-VLM Architecture



**Depth image**

**RGB image**

**3D Project**

**CLIP Visual Encoder**

**FastSAM**

**Patch Feature Points in 3D Space**

**Render Panoramic Rays within Generalizable Feature Fields**

**Panoramic 3D Patch Tokens**

Aggregate patch features within 2D mask

**Instance Encoder**

2D Instance Masks

Merge 2D instances into existing 3D instances

**Merging Discriminator**

2D Instance Features

Aggregate patch features within 3D mask

**3D Instance Masks**

**Instance Encoder**

**3D Instance Tokens**

**Zone Encoder**

**3D Zone Tokens**

Aggregate 3D instance features within cube zone

**3D Vision-Language Large Model**

Turn right θ degree

Turn left θ degree

Forward d cm

Stop

Instruction: "Please go to the kitchen and take the bread out of the microwave for me."

**Instruction Tokens**

Turn left 30 degree. Forward 50 cm. Turn left 45 degree. Forward 75 cm.

**History Actions**

**LLaVA-Phi-3-mini**

A lightweight (3.8B) Multimodal LLM.

**INPUT:**
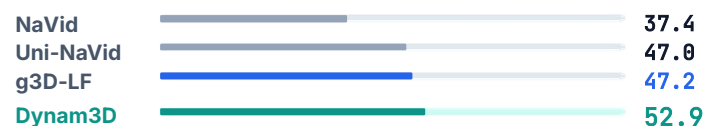{patch_tokens}{instance_tokens}{zone_tokens}{instruction_tokens} {history_action_tokens}

**OUTPUT:**
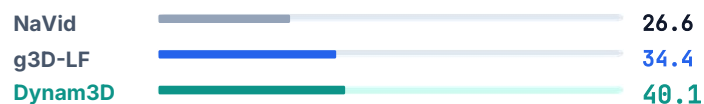1) Turn left θ degree. 2) Turn right θ degree.
3) Forward d cm. 4) Stop.

# Navigation Performance

We outperform both video-based and map-based baselines, specifically in **Success Rate (SR)** and **Path Efficiency (SPL)**.

**R2R-CE**    **Step-by-step following, *e.g.,*** *"Walk through the bedroom around the bed. Walk out of the door into the hallway. Walk towards the closet area in the hallway."*
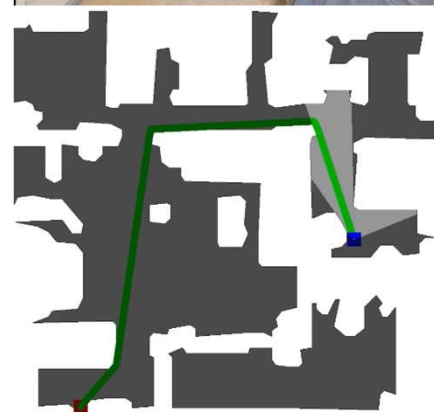
| | |
|---|---|
| NaVid | 37.4 |
| Uni-NaVid | 47.0 |
| g3D-LF | 47.2 |
| Dynam3D | 52.9 |

**REVERIE-CE**    **High-level instruction, *e.g.,*** *"Go to the familyroom and bring me the pillow from the couch closest to the entrance."*

| | |
|---|---|
| NaVid | 26.6 |
| g3D-LF | 34.4 |
| Dynam3D | 40.1 |

**NavRAG-CE**    **User-demand instruction, *e.g.,*** *"Walk to the warm hall featuring elegant wooden accents and set the large wooden table with candles and napkins for a lovely dinner ambiance."*

| | |
|---|---|
| NaVid | 19.4 |
| g3D-LF | 21.4 |
| Dynam3D | 24.7 |

*"After exiting the bedroom, walk straight along the hallway, then turn left at the end of the hallway to enter the kitchen, and walk to the stove."*

# Sim-to-Real Deployment

Deploy Dynam3D on **Hello Robot Stretch 3** in **NUS Robotics Living Studio**.
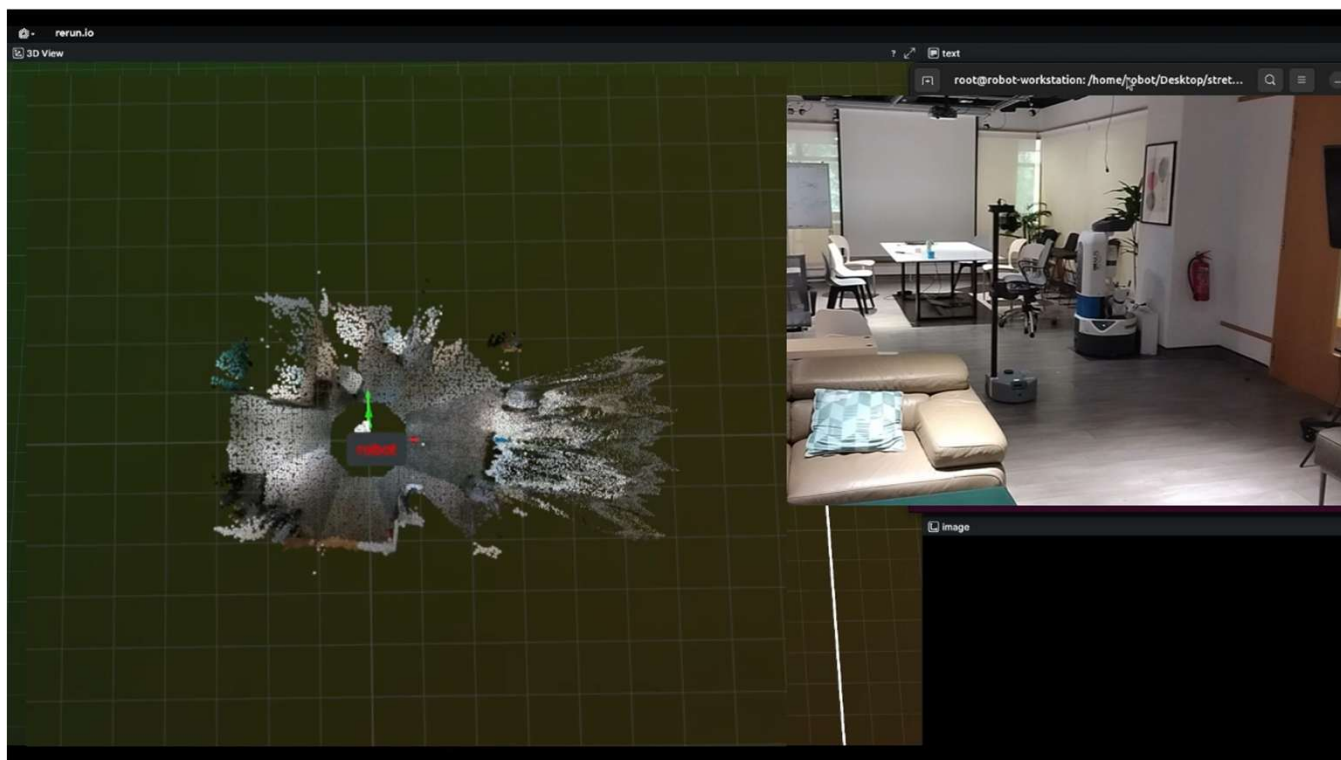
**70%**
Static Success Rate

**45%**
Dynamic Success Rate

The dynamic layered 3D tokens effectively handles moved objects.

*"Please pick up the blue cup on the table and place it in the kitchen sink."*

# Conclusion: Towards Dynamic Embodied Memory

## Hierarchical

**Patch → Instance → Zone**

Bridging the gap between fine-grained geometric details and high-level VLM reasoning.

## Dynamic

**Active Update**

Frustum Culling enables the map to "breathe" and adapt to changes with **83ms** latency.

## Aligned

**3D Consistency**

Shifting anchors to local view centers effectively denoises 2D-to-3D feature distillation.

**Code Available**
github.com/MrZihan/Dynam3D