# Map2Thought: Explicit 3D Spatial Reasoning via Metric Cognitive Maps

Xiangjun Gao[1]    Zhensong Zhang[2]    Dave Zhenyu Chen[2]    Songcen Xu[2]

Long Quan[1]    Eduardo Pérez-Pellitero[2]    Youngkyoon Jang[2]

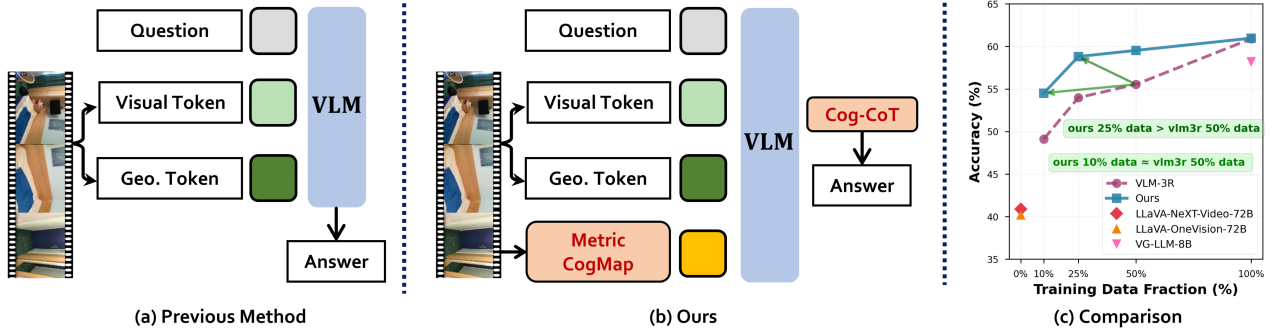[1]The Hong Kong University of Science and Technology    [2]Huawei Noah's Ark Lab

Figure 1. **Comparison with prior approaches and efficiency gains.** (a) Previous 3D VLMs fuse visual and geometric tokens but through implicit latent reasoning, limiting spatial interpretability. (b) Our approach introduces a metrically grounded cognitive map (*Metric-CogMap*) and a chain-of-thought-style reasoning process (*Cog-CoT*), enabling explicit and interpretable spatial reasoning from the same multimodal inputs. (c) This design yields substantial data efficiency: with just 10% or 25% of the training data, our model is comparable to or surpasses the performance of the baseline model trained with 50% of the data.

## Abstract

*We propose Map2Thought, a framework that enables explicit and interpretable spatial reasoning for 3D VLMs. The framework is grounded in two key components: Metric Cognitive Map (Metric-CogMap) and Cognitive Chain-of-Thought (Cog-CoT). Metric-CogMap provides a unified spatial representation by integrating a discrete grid for relational reasoning with a continuous, metric-scale representation for precise geometric understanding. Building upon the Metric-CogMap, Cog-CoT performs explicit geometric reasoning through deterministic operations (e.g., vector operations, bounding-box distances, and occlusion-aware appearance order cues) producing interpretable inference traces grounded in 3D structure. Experimental results show that Map2Thought enables explainable 3D understanding, achieving 59.9% accuracy using only half the supervision—closely matching the 60.9% baseline trained with full dataset. It consistently outperforms state-of-the-art methods by 5.3%, 4.8%, and 4.0% under 10%, 25%, and 50% training subsets, respectively, on the VSI-Bench.*

## 1. Introduction

Recent advances in 3D Vision-Language Models (3D-VLMs) have increasingly explored bridging multi-modal input signals to harness the rich knowledge priors of large language models (LLMs) for 3D spatial understanding. Beyond projecting 3D geometric tokens from point clouds [20, 49] or RGB-D sequences [54, 56] into the latent space of large language models, recent advances have leveraged pretrained 3D vision foundation models—such as VGGT [38] and CUT3R [39]—which extract 3D structure directly from monocular RGB videos to support more grounded and scalable spatial reasoning. This integration allows 3D-VLMs to access geometry-aware representations without relying on additional modalities. While these methods represent a step forward in generalizability compared to earlier spatial reasoning techniques [2, 5, 6, 9, 10, 15, 46], several core challenges remain unresolved:

- **Implicit fusion without geometric grounding:** Multimodal representations are fused implicitly without enforcing alignment with physical constraints, limiting transparency and verifiability in spatial reasoning.

- **Lack of accumulated spatial context:** Temporal and scene-level spatial cues remain weakly supervised, limiting the model's ability to reason about appearance order and global layout.

- **Biased supervision and poor generalization:** Training on subsampled or biased visual inputs leads to overfitting, impairing reasoning about object scale, type, and spatial anchoring across diverse scene types.

To address these challenges, we propose Map2Thought,

an explicit 3D spatial reasoning framework grounded in Metric Cognitive Maps that enables reliable and interpretable 3D understanding, as shown in Fig. 2. Map2Thought is built upon two core components: 1) *Metric-CogMap*, a unified spatial representation that combines a discrete grid for symbolic relational reasoning with a continuous metric-scale grid for precise geometric perception; and 2) *Cog-CoT*, an explicit chain-of-thought reasoning module that performs interpretable geometric computations over the *Metric-CogMap*. Unlike prior cognitive maps used in [45, 48], our *Metric-CogMap* encodes richer spatial detail, including object occupancy, metric-scale positions, and real-world scale bounding boxes—enabling more fine-grained spatial inference. Our proposed *Metric-CogMap* is constructed through a robust video-to-map pipeline that extends state-of-the-art 2D object detection and segmentation models—originally designed for frame-based input and output—by integrating 3D vision foundation models and a novel covisibility map–based geometric validation step to ensure accurate and consistent 3D spatial grounding. *Cog-CoT* complements this representation by performing transparent, modular geometric reasoning (e.g., distance, direction) through deterministic computations, enabling verifiable inference traces and making our framework easily extensible to other 3D-VLMs without retraining. With this design, Map2Thought attains 59.9% accuracy on VSI-Bench using only 50% of the training data—nearly matching the 60.9% full-data baseline—while still exceeding the baseline by a consistent 4.0% margin under equal 50% supervision. The main contributions of Map2Thought are threefold:

- **Explicit Metric-CogMap representation:** We introduce *Metric-CogMap*, a unified spatial representation that integrates discrete symbolic grids with continuous metric-scale geometry, enabling LLMs to interpret 3D scenes with structured, interpretable reasoning.
- **Cog-CoT reasoning paradigm:** We introduce *Cog-CoT*, an interpretable CoT that operates over the accumulated spatial context of the Metric-CogMap to perform explicit and verifiable 3D spatial reasoning.
- **Data-efficient 3D-VLMs:** Our framework exhibits strong data efficiency, achieving better accuracy than comparable baselines under limited training data, while improving generalization by mitigating data-dependent overfitting.

## 2. Related Work

Large Language Models (LLMs) have progressed rapidly in recent years, with active research expanding their capabilities across diverse domains such as code generation, commonsense reasoning, and mathematics [36]. A prominent direction in this evolution involves extending LLMs to handle multiple modalities—enabling them to develop visual [3] and spatial understanding [32]. Early

Visual-Language Models (VLMs), such as CLIP [33] and ALIGN [24], learn joint image-text embeddings, while later works [1, 27] decouple vision and language modules to better support cross-modal reasoning. While these models perform well on single-image tasks, they remain limited in reasoning about 3D spatial context intuitive to humans—such as distance, relative positioning, and object counting.

**Point cloud encoding in 3D VLMs.** Recent 3D-VLMs integrate point clouds with RGB and text inputs to provide geometric grounding, but their fusion in latent space often limits its spatial precision and interpretability [26]. Models such as [14, 20] use spatial transformers [7] to encode object-centric 3D tokens and capture inter-object geometry, while [8, 18] align 2D and 3D features by projecting them back to reconstructed points for enhanced multimodal understanding. However, both rely on latent feature-level fusion across modalities—point clouds, images, and text—leading to opaque reasoning pipelines that hinder interpretability and analysis. More explicit pipelines, introduced by [28, 42, 49], combine multi-view 2D semantics features with 3D geometric information for improved indoor scene understanding. However, they assume access to high-quality, densely reconstructed point clouds—an assumption that breaks down in real-world settings where input videos span several minutes, as is common in datasets like ScanNet[12], ScanNet++[47] or ARKitScenes[13]. In such cases, sparse sampling or partial reconstructions may miss key objects, making reliance on complete point cloud–based instance detection a major bottleneck.

**Positional encodings for 3D VLMs.** To improve spatial grounding, recent 3D VLMs introduce positional encodings derived from camera poses, depth maps, or back-projected coordinates. These encodings localize visual features within a 3D coordinate system, aiding alignment across views and with language [31]. For instance, LLAVA-3D [56] and VIDEO-3D-LLM [54] embed spatial tokens linking 2D appearances to 3D context. While this enhances object-level spatial sensitivity, they [37, 54, 56] often fail to capture broader geometric structure. Consequently, the resulting representations remain opaque and overfitting, limiting interpretability and generalization in spatial reasoning tasks.

**3D Vision foundational model encodings for 3D VLMs.** More recently, emerging approaches have incorporated pretrained 3D vision foundation models—such as CUT3R [39] and VGGT [38]—into 3D-VLM pipelines to enhance spatial reasoning through learned geometric priors. For example, some recent works [17, 21, 43, 44, 53], process monocular RGB video sequences using geometry encoders that produce implicit 3D tokens, which are fused with visual tokens from existing VLMs. While this fusion improves spatial feature encoding, it still relies on latent representations derived from 2D inputs—lacking explicit met-
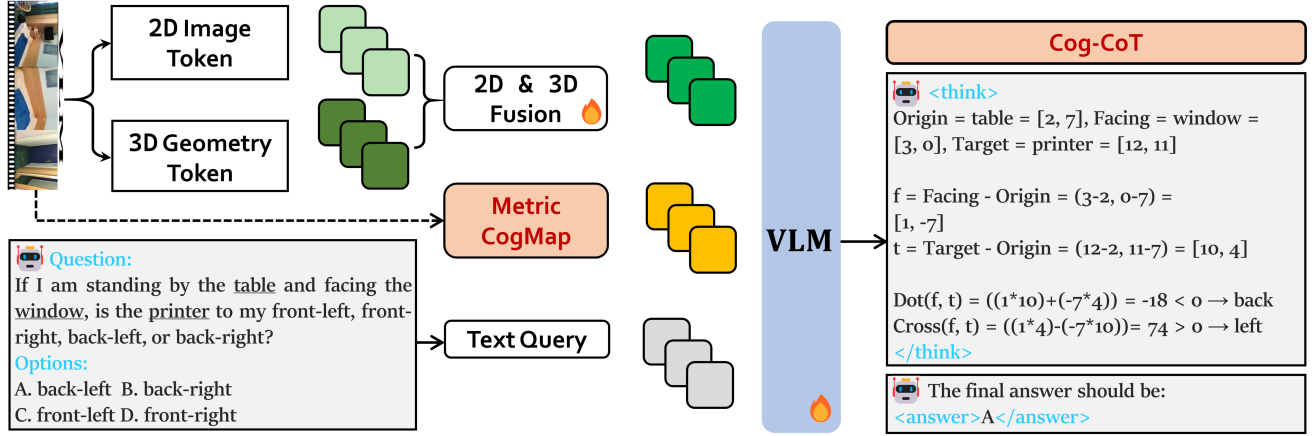
Figure 2. **Overview of our method.** Given an RGB video and a language query, we extract 2D image tokens and 3D geometry-aware tokens, fuse them into 3D-aware visual tokens, and input them to the VLM. The Metric-CogMap (orange blocks) encodes the scene using both a discrete grid and a metric-scale spatial representation. Cog-CoT (right, grey panel) then performs explicit and deterministic geometric reasoning over the map, yielding a transparent and interpretable answer.

ric scale, object-level semantics, and coherent cross-modal alignment. As a result, the encoded spatial structures remain difficult to interpret and limit the model's ability to perform tasks requiring precise geometric reasoning—such as appearance order, object size, or absolute distance in VSI-Bench [45].

**Cognitive maps.** Motivated by principles of human spatial cognition, recent works have explored cognitive map-style representations that discretize 3D environments into symbolic spatial grids encoding relative object positions and inter-object relationships. Approaches like VSI-Bench [45], SpatialMind [50], and MindCube [48] enable spatial reasoning through coarse symbolic cognitive maps based on integer grid coordinates. More recently, this direction has progressed with the integration of scene graphs [4] to model richer relational semantics. While these symbolic representations enhance interpretability and high-level reasoning, they still lack geometric precision and omit metric details essential for fine-grained spatial understanding.

Despite recent advances, 3D-VLMs still lack explicit spatial representations—especially under noisy, subsampled inputs—limiting their ability to support interpretable and transparent reasoning. Prior work tends to prioritize benchmark performance over explainability, offering limited insight into spatial understanding. In this paper, we introduce a framework that constructs metrically grounded cognitive maps and enables chain-of-thought spatial reasoning, enhancing robustness, interpretability, and scalability.

## 3. The Proposed Framework: Map2Thought

**Overview.** We introduce Map2Thought, a unified framework for 3D spatial understanding from monocular RGB videos and language queries, designed to enable explicit and interpretable reasoning grounded in structured geometric representations. As shown in Fig. 2, Map2Thought integrates pretrained vision-language and 3D vision foundation models to jointly process visual, geometric, and textual inputs. The framework is composed of three key components: 1) defining the *Metric-CogMap*, a dual-format representation that encodes both discrete spatial layouts and continuous metric-scale geometry to support relational reasoning and precise spatial grounding; 2) implementing a video-to-map pipeline that constructs the Metric-CogMap from raw video frames, leveraging object detection, tracking, and feed-forward 3D reconstruction; and 3) adopting *Cog-CoT* (Cognitive Chain-of-Thought), a structured reasoning mechanism that performs explicit and verifiable spatial inference over the constructed *Metric-CogMap*.

These components operate in two complementary streams: one for extracting and fusing multimodal tokens for contextual understanding, and another for constructing a structured scene representation (Metric-CogMap) used in explicit reasoning (Cog-CoT). These components work in tandem to produce interpretable and geometry-aware outputs. We detail each part of the architecture below.

### 3.1. Architecture

**Token extraction and fusion.** Map2Thought processes a monocular RGB video and a text query to perform 3D spatial reasoning. Each video frame is resized and passed through a frozen 2D visual encoder (e.g., CLIP-ViT [33]) to extract 2D appearance tokens $Z_{2D}$, capturing local semantic and visual cues. Simultaneously, a pretrained 3D vision foundation model (e.g., CUT3R [39]) produces geometry-aware tokens $Z_{3D}$ that encode the global scene structure
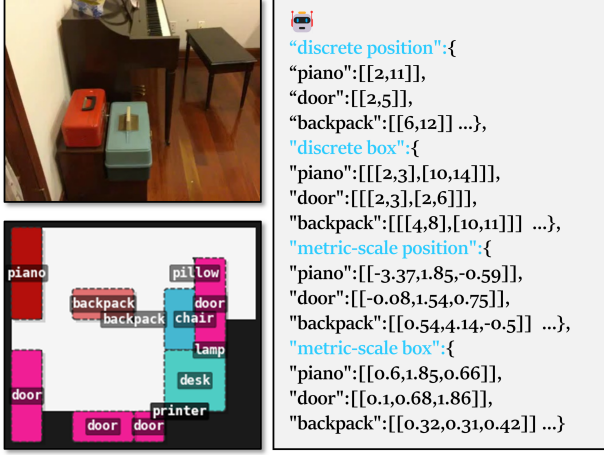
```
"discrete position":{
"piano":[[2,11]],
"door":[[2,5]],
"backpack":[[6,12]] ...},
"discrete box":{
"piano":[[[2,3],[10,14]]],
"door":[[[2,3],[2,6]]],
"backpack":[[[4,8],[10,11]]] ...},
"metric-scale position":{
"piano":[[-3.37,1.85,-0.59]],
"door":[[-0.08,1.54,0.75]],
"backpack":[[0.54,4.14,-0.5]] ...},
"metric-scale box":{
"piano":[[0.6,1.85,0.66]],
"door":[[0.1,0.68,1.86]],
"backpack":[[0.32,0.31,0.42]] ...}
```

Figure 3. Visualization of a *Metric-CogMap* example. Given the reference frame shown at the top left, the discrete grid representation (top two) assigns each object an integer-valued position, while the box occupancy map encodes the grid ranges each object occupies. The metric-scale representation (bottom two) records object centroids and axis-aligned bounding boxes (AABBs) in a globally aligned real-world coordinate system.

across frames. These 2D and 3D tokens are fused using a lightweight cross-attention mechanism [16], where $Z_{2D}$ queries $Z_{3D}$ to generate geometry-enriched visual tokens $Z_f$. The fused tokens are then projected into the VLM's latent space and concatenated with the tokenized language query for downstream spatial reasoning.

**Metric-CogMap and Cog-CoT.** In parallel to token fusion, Map2Thought constructs a *Metric-CogMap* from the input video using off-the-shelf modules for object detection [55], video segmentation [34], and feed-forward 3D reconstruction [41]. This map encodes both discrete spatial grids and continuous metric-scale object geometry, including positions, bounding boxes, and relational layouts. Built atop this representation, *Cog-CoT* (Cognitive Chain-of-Thought) performs explicit reasoning through deterministic operations (e.g., vector relations, occlusion-aware queries) grounded in the map. Together, these modules enable transparent and verifiable 3D reasoning that enhances both accuracy and interpretability under varied supervision regimes.

### 3.2. Design of Metric-CogMap

Our *Metric-CogMap* provides a unified spatial representation that integrates two complementary components—(i) a discrete grid-based abstraction and (ii) a continuous metric-scale encoding. Together, these representations support structured, interpretable, and verifiable spatial reasoning by combining symbolic relational structure with precise geometric information.

**Discrete grid representation.** First, the scene is represented on a *discrete position representation*, where each ob-

ject is assigned an integer coordinate (e.g., within a fixed range such as 0–20). Beyond mapping an object to a single cell, we also encode its occupied region using four integer values that specify an axis-aligned bounding box, forming a *discrete box representation* that preserves its spatial extent. This discrete encoding offers an approximate but structured description of position and size, enabling symbolic reasoning and relational parsing in a simplified grid domain while reducing the complexity of spatial inference.

**Continuous metric-scale representation.** In parallel, we preserve real-world spatial fidelity by annotating each object with its *metric-scale position* in the real-world coordinate system. We further associate each object with a *metric-scale axis-aligned bounding box (AABB)*, which provides its physical size and exact spatial boundaries in continuous 3D space. This continuous formulation retains true geometric structure and distance relationships, enabling precise spatial reasoning that discrete layouts cannot capture.

**Unified representation.** Fig. 3 illustrates the integrated *Metric-CogMap*. The discrete grid representation (top two panels) assigns each object an integer-valued position and encodes its occupied grid range through a discrete bounding box, yielding a symbolic and interpretable layout of the scene. Complementarily, the metric-scale representation (bottom two panels) records each object's real-world centroid and its physical axis-aligned bounding box (AABB), capturing precise geometric size and spatial extent in the global coordinate frame. By coherently linking these *discrete* and *continuous* components, the unified map provides a structured relational view aligned with precise geometric measurements, while maintaining an explicit and interpretable inference process.

### 3.3. Construction of Metric Cognitive Map

To support spatially grounded reasoning in our proposed Map2Thought framework, we construct a Metric-CogMap through a multi-stage pipeline that integrates visual semantics with metrically aligned 3D geometry, as shown in Fig. 4. The process involves: 1) Crucial frame selection via key-object anchoring, 2) covisibility map construction, 3) multi-instance-aware object detection (conditional), 4) object-centric video segmentation, 5) geometric transformation and metric alignment, and 6) multi-view segment integration.

**Crucial frame selection via key-object anchoring**: To ensure that all objects referenced in the input question are visually represented in the selected frames, we replace naïve uniform subsampling with a targeted frame selection strategy. We first apply the Detic detector [55] to every 10th frame to identify high-confidence detections of queried objects, retaining the most representative frame for each as a crucial view. To complete a 64-frame input set, we then add uniformly sampled frames from the remaining video
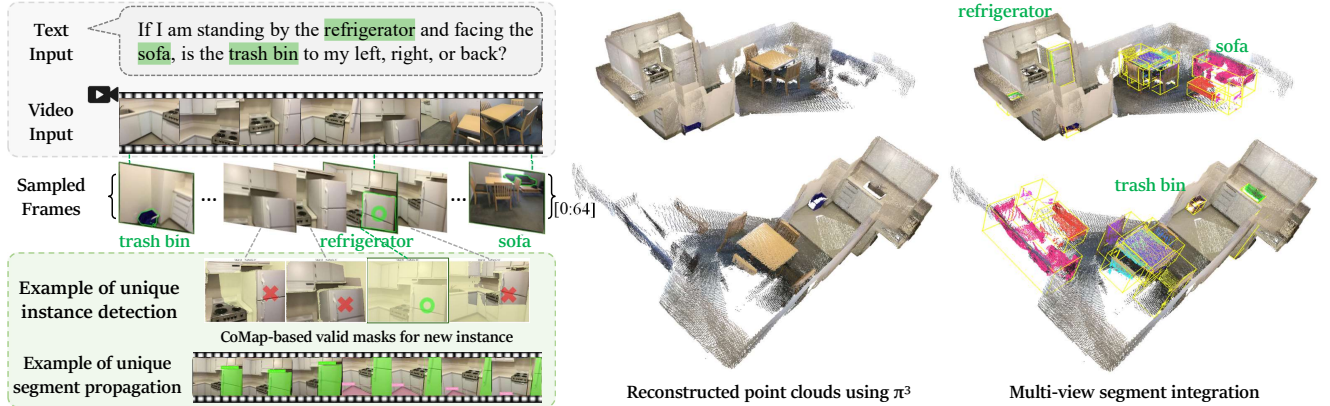
Figure 4. **Overview of the Metric-CogMap construction pipeline.** Given video and text inputs, the pipeline selects 64 uniformly distributed frames while ensuring key frames containing the highest-confidence detections of text-referenced objects (green outlines and annotations). For categories requiring multi-instance detection (e.g., the refrigerator), covisibility maps anchor the highest-confidence frame and suppress redundant detections by checking whether candidate regions appeared in earlier selected views. Only geometrically unique detections within valid mask regions are retained. These detections are propagated across frames, and multi-view segments are merged with consistent IDs during final integration. The example on the right shows the resulting detections for several question-relevant categories (e.g., sofa, trash bin, chair, table) from `scene0651_02` in the ScanNet dataset [12], an evaluation scene in VSI-Bench [45].

to ensure broad spatio-temporal coverage. We then apply Grounding DINO [30], an open-vocabulary object detector, to all selected frames to capture additional object instances that may have been missed by Detic—particularly for underrepresented or category-agnostic entities. This procedure guarantees that key semantic cues are explicitly embedded in the visual stream, overcoming a fundamental limitation of conventional uniform sampling that risks excluding critical objects from the reconstruction pipeline.

**Covisibility map construction**: To establish reliable interframe geometric correspondence, we generate dense pixelwise Covisibility Maps (CoMaps), following the principles of CoMapGS [23]. These maps indicate which regions in one frame are geometrically visible from others, forming the basis for robust multi-view reasoning. Using $\pi^3$ [41], we extract dense 3D point clouds and their associated camera poses. Covisibility is then determined via reprojection consistency: a 3D point observed in a source frame is projected into target views, and covisibility is confirmed if the projection aligns with neighboring 3D structure in the target frame. These pairwise pixel-level covisibility cues are essential for subsequent stages, supporting accurate multiview fusion, disambiguation of duplicate object detections, and spatial validation of instance consistency.

**Multi-instance-aware object detection (conditional)**: For queries requiring multi-instance understanding—such as object counting or appearance ordering—we introduce an additional detection refinement step guided by geometric constraints. Leveraging the previously constructed CoMaps, we enforce spatial distinctiveness by filtering out

redundant detections across views. For each object category, we start with the highest-confidence detection and iteratively suppress other candidates whose masks fall within the covisible regions of already-selected instances. This strategy ensures non-redundant, viewpoint-distinct detections. This deduplication process is applied across both Detic and Grounding DINO [30] outputs, ensuring that only spatially distinct instances are preserved in the map.

**Object-centric video segmentation**: To enhance both the geometric integrity and semantic completeness of object representations, we apply SAM2 [34] to propagate detected instances across frames. This enables the recovery of more complete object regions beyond the initial detection frame, providing a more holistic understanding of each object's shape and spatial extent.

**Geometric transformation and metric alignment**: Leveraging the camera poses and 3D point clouds provided by $\pi^3$ [41], we perform global alignment by anchoring the reconstruction to a dominant horizontal plane (e.g., the floor), or a vertical plane when it provides a more reliable reference. To achieve metric consistency, we rescale the reconstruction based on depth from MoGe-2 [40], producing a metrically accurate, globally aligned 3D scene.

**Multi-view segment integration**: Finally, segmented object instances and their associated 3D points are consolidated across views to form a unified spatial representation. Overlapping regions—defined as 3D segments within a 20% spatial proximity threshold—are merged by retaining the highest-confidence detection, effectively reducing redundancy and addressing potential false duplicates that may

still arise from conventional detection pipelines. Merged objects are then assigned consistent semantic labels and unique instance IDs, producing a compact, semantically coherent, and metrically grounded *Metric-CogMap*.

This integrated map encodes both object-level semantics and spatial topology, providing a robust foundation for structured spatial reasoning via *Cog-CoT* in 3D VLMs.

## 3.4. Cog-CoT: Cognitive Chain-of-Thought

Our *Cog-CoT* (Cognitive Chain-of-Thought) enables explicit and verifiable spatial reasoning—unlike prior approaches that rely on end-to-end predictions over visual inputs with implicit, non-interpretable processes. It operates directly on the structured geometry encoded in the *Metric-CogMap*, executing deterministic operations such as vector analysis, bounding-box distance, and appearance order. This structured reasoning chain allows the VLM to apply its language understanding to interpret geometric evidence and derive the final answer. In the following, we demonstrate how Cog-CoT applies these operations across representative categories of spatial queries, enabling transparent and interpretable 3D reasoning.



Figure 5. **Cog-CoT workflow.** Given a query, the *Metric-CogMap*, and the corresponding *Task-Specific instruction*, the system retrieves the scene-level *Metric-CogMap*, and *Cog-CoT* executes the geometric computations step-by-step, yielding verifiable intermediate numerical evidence used to generate the final answer.

Fig. 5 illustrates a representative *Cog-CoT* trace for the *Absolute Distance Measurement* task:
**Query and instruction:** The language query requests the metric distance between two specified objects, paired with a task-specific instruction that guides the expected reasoning procedure.
**Scene representation:** The relevant *Metric-CogMap* (see Fig. 3) is provided, with the metric-scale representation

used to retrieve object centroids and bounding boxes.
**Geometric reasoning:** The model applies explicit geometric operations (e.g., centroid extraction, distance calculation) to generate verifiable intermediate results.
**Answer generation:** These intermediate results guide the final response in a transparent, verifiable manner.

In addition to the distance measurement example, another *Cog-CoT* application—*Relative Direction Reasoning*—is shown in Fig. 2, involving geometric operations such as cross and dot products to determine directional relationships. Further examples and task-specific instructions are provided in the supplementary material.

## 4. Experiments

In this section, we evaluate our framework in terms of spatial reasoning accuracy and data efficiency. Sec. 4.1 outlines the training setup, and Sec. 4.2 introduces the scene-level Metric-CogMap annotations. Benchmark tasks and evaluation metrics are presented in Sec. 4.3, followed by an ablation study in Sec. 4.4.

### 4.1. Training Details

We adopt LLaVA-Next-Video-7B as the backbon, initializing all visual and language components from publicly available checkpoints. Fine-tuning is conducted on the VLM-3R training set, strictly following the same optimization and data protocols from [16]. To ensure efficient adaptation, we employ LoRA [19] to finetune the VLM model. We also employ the same 2D-3D fusion module and projection layer architecture as used in VLM-3R [16].

### 4.2. Training Dataset Construction

Our training process builds upon the VLM-3R [16] instructional dataset, which is curated from ScanNet [12], ScanNet++ [47], and ARKitScenes [13]. For each scene and corresponding question-answer pair, we additionally construct a ground-truth *Metric-CogMap*, enabling explicit spatial computation during Cog-CoT inference. The construction pipeline proceeds as follows:

(1) **Object selection.** We extract the referenced object(s) in each question via text parsing and entity grounding.

(2) **Instance association.** Leveraging scene-level semantic and instance annotations, we locate the corresponding 3D object instances and obtain their axis-aligned bounding boxes (AABBs).

(3) **Metric-CogMap construction.** Each indoor scene is normalized to its room extents and discretized into a $20 \times 20$ spatial grid. For each identified object, we store both a) its discrete position and occupancy on the grid and b) its continuous metric-scale position and AABB parameters in world coordinates, forming the dual-format representation.

Table 1. **Evaluations on VSI-Bench.**

| Methods | Rank | Avg. | Obj. Count | Abs. Dist. | Obj. Size | Room Size | Rel. Dist. | Rel. Dir. | Route Plan | Appr. Order |
|---|---|---|---|---|---|---|---|---|---|---|
| | | | Numerical Answer | | | | Multiple-Choice Answer | | | |
| *Baseline* | | | | | | | | | | |
| Chance Level (Random) | - | - | - | - | - | - | 25.0 | 36.1 | 28.3 | 25.0 |
| Chance Level (Frequency) | - | 34.0 | 62.1 | 32.0 | 29.9 | 33.1 | 25.1 | 47.9 | 28.4 | 25.2 |
| *Proprietary Models (API)* | | | | | | | | | | |
| GPT-4o | 3 | 34.0 | 46.2 | 5.3 | 43.8 | 38.2 | 37.0 | 41.3 | 31.5 | 28.5 |
| Gemini-1.5 Flash | 2 | 42.1 | 49.8 | 30.8 | 53.5 | 54.4 | 37.7 | 41.0 | 31.5 | 37.8 |
| Gemini-1.5 Pro | 1 | 45.4 | 56.2 | 30.9 | 64.1 | 43.6 | 51.3 | 46.3 | 36.0 | 34.6 |
| *Open-sourced VLMs* | | | | | | | | | | |
| LLaVA-OneVision-0.5B | 16 | 28.0 | 46.1 | 28.4 | 15.4 | 28.3 | 28.9 | 36.9 | 34.5 | 5.8 |
| InternVL2-2B | 17 | 27.4 | 21.8 | 24.9 | 22.0 | 35.0 | 33.8 | 44.2 | 30.5 | 7.1 |
| LLaVA-NeXT-Video-7B | 10 | 35.6 | 48.5 | 14.0 | 47.8 | 24.2 | 43.5 | 42.4 | 34.0 | 30.6 |
| InternVL2-8B | 11 | 34.6 | 23.1 | 28.7 | 48.2 | 39.8 | 36.7 | 30.7 | 29.9 | 39.6 |
| LLaVA-OneVision-7B | 12 | 32.4 | 47.7 | 20.2 | 47.4 | 12.3 | 42.5 | 35.2 | 29.4 | 24.4 |
| LongVA-7B | 14 | 29.2 | 38.0 | 16.6 | 38.9 | 22.2 | 33.1 | 43.3 | 25.4 | 15.7 |
| VILA-1.5-8B | 15 | 28.9 | 17.4 | 21.8 | 50.3 | 18.8 | 32.1 | 34.8 | 31.0 | 24.8 |
| LongVILA-8B | 18 | 21.6 | 29.1 | 9.1 | 16.7 | 0.0 | 29.6 | 30.7 | 32.5 | 25.5 |
| InternVL2-40B | 9 | 36.0 | 34.9 | 26.9 | 46.5 | 31.8 | 42.1 | 32.2 | 34.0 | 39.6 |
| VILA-1.5-40B | 13 | 31.2 | 22.4 | 24.8 | 48.7 | 22.7 | 40.5 | 25.7 | 31.5 | 32.9 |
| LLaVA-NeXT-Video-72B | 7 | 40.9 | 48.9 | 22.8 | 57.4 | 35.3 | 42.4 | 36.7 | 35.0 | 48.6 |
| LLaVA-OneVision-72B | 8 | 40.2 | 43.5 | 23.9 | 57.6 | 37.5 | 42.5 | 39.9 | 32.5 | 44.6 |
| Spatial-MLLM-4B | 6 | 48.4 | 65.3 | 34.8 | 63.1 | 45.1 | 41.3 | 46.2 | 33.5 | 46.3 |
| VG-LLM-8B | 4 | 58.2 | 71.7 | 53.8 | 68.8 | 62.1 | 63.8 | 83.0 | 44.3 | 18.4 |
| VLM-3R (100% dataset) | 2 | 60.9 | 70.2 | 49.4 | 69.2 | 67.1 | 65.4 | 80.5 | 45.4 | 40.1 |
| **Ours (100% dataset)** | 1 | 61.0 | 70.8 | 55.0 | 70.1 | 69.4 | 56.9 | 69.8 | 38.1 | 57.4 |
| VLM-3R (25% dataset) | 5 | 54.0 | 69.2 | 41.7 | 67.6 | 62.4 | 58.4 | 47.3 | 43.3 | 41.7 |
| **Ours (25% dataset)** | 3 | 58.8 | 68.2 | 53.0 | 64.4 | 69.1 | 52.1 | 70.9 | 38.1 | 54.5 |

Our model ranks first among open-sourced VLMs, showcasing the effectiveness of *Metric-CogMap* and *Cog-CoT*. For each task within the open-sourced 3D-VLM group, the dark gray cell indicates the best-performing model, while light gray and faint gray highlight the second- and third-best performances, respectively.

## 4.3. Evaluations

**Comparison baselines** Following the evaluation protocol of VSI-Bench [45], we adopt a range of existing video LLM as baselines. These include proprietary models (e.g., Gemini [35], GPT–4o [22]) as well as open-source counterparts such as InternVL2 [11], ViLA [29], LongViLA [51], LongVA [51], LLaVA-OneVision [25], and LLaVA-NeXT-Video [52]. The performance score of these baselines are copied from the VSI-Bench [45] leaderboard.

In addition, we compare against closely related RGB-only spatial reasoning approaches, including VLM-3R [16], VG-LLM [53], and Spatial-MLLM [43]. These methods enhance video-based LMMs by introducing tokens from 3D foundation models to improve spatial understanding.

**Comparison on VSI-Bench** Our model outperforms all existing baselines, including proprietary models such as GPT-4o and Gemini-1.5-Pro, despite using only a 7B-parameter open-source backbone. We achieve state-of-the-art (SOTA) results not only in the overall average performance but also in several individual question categories, including *Obj.Count*, *Abs.Dist*, *Obj.Size*, *Room Size*, and *Appr.Order*. These gains stem from the metric-scale representation in *Metric-CogMap* and geometric reasoning via *Cog-CoT*. Our model excels in tasks like *Abs.Dist* and *Room Size*, where precise real-world understanding is essential and discrete grid maps typically fall short. By combining structured scene representations with interpretable logic, it enables accurate and reliable spatial understanding.

In addition, with the proposed *Metric-CogMap* and *Cog-CoT*, our model enables efficient fine-tuning through a structured, interpretable spatial reasoning process. As shown in the last row of Table 1, using only 25% of the VLM-3R dataset QA pairs, it outperforms VLM-3R trained on the same subset by 4.8% and even surpasses VG-LLM-8B trained on the full dataset. These results underscore not only the data efficiency of our method, but also how explicit metric-scale reasoning contributes to more robust generalization under limited supervision.

Although our model shows slightly lower performance on *Rel.Dist* and *Rel.Dir*, this is primarily due to imperfections in the constructed *Metric-CogMap*, which leads to relatively lower performance in these tasks. To further vali-

date the effectiveness of our *Metric-CogMap* and *Cog-CoT*, we conduct additional evaluations using the ground truth *Metric-CogMap* as input, which are discussed in the subsequent ablation section.

## 4.4. Ablation Study

**Ablation on dataset proportion** We perform ablation studies to investigate the effect of training data proportion on model performance, comparing our approach with VLM-3R. As shown in Tab. 2 and Fig. 1, our model consistently outperforms VLM-3R across all dataset proportions, outperforms state-of-the-art methods by 5.3%, 4.8%, and 4.0% under 10%, 25%, and 50% training subsets, respectively, on the VSI-Bench. When trained on the full 100% dataset, our model achieves performance comparable to VLM-3R (61.0% vs 60.9%). With only 25% of the training data, our model outperforms the VG-LLM-8B model, which is trained on 100% of the VLM-3R training data. These improvements prove that we can achieve higher performance with fewer training data, demonstrating the efficiency and effectiveness of our approach.

We further analyze two representative tasks: *Abs.Dist* and *Rel.Dir*. With only 10% of the training data, our model exceeds VLM-3R by 13.8% and 23.0% on these tasks, respectively. These results demonstrate that the interpretable structure of *Metric-CogMap* and *Cog-CoT* accelerates convergence and enhances performance, especially in geometry-sensitive tasks.

Table 2. Ablation study across varying training data fractions. We compare our method against VLM-3R. **Bold** denotes better results.

| Data Fraction | Method | Average | Abs. Dist | Rel. Dir |
|---|---|---|---|---|
| 100% Dataset | VLM-3R | 60.9 | 49.4 | **80.5** |
| | Ours | **61.0** | **55.0** | 69.8 |
| 50% Dataset | VLM-3R | 55.5 | 43.7 | 52.5 |
| | Ours | **59.5** | **54.2** | **69.0** |
| 25% Dataset | VLM-3R | 54.0 | 41.7 | 47.3 |
| | Ours | **58.8** | **53.0** | **70.9** |
| 10% Dataset | VLM-3R | 49.1 | 37.3 | 47.0 |
| | Ours | **54.5** | **51.5** | **70.0** |

**Ablation on Metric-CogMap and Cog-CoT** We assess the contributions of the *Metric-CogMap* and *Cog-CoT* components in Tab. 3. All experiments in this table are trained using 25% of the full dataset.

Experiment (1) represents our complete model, which uses a predicted *Metric-CogMap* constructed via our proposed pipeline and applies *Cog-CoT* for stepwise reasoning before answer generation. Experiment (2) disables the *Cog-CoT* module while retaining the predicted *Metric-CogMap*, thus removing explicit cognitive reasoning and relying solely on spatial metric cues. And, experiment (3)

removes both the *Metric-CogMap* and *Cog-CoT*, which is the same as the VLM-3R baseline setting. This confirms that both components play complementary and essential roles in enhancing spatial reasoning, with *Cog-CoT* providing the necessary cognitive reasoning framework and *Metric-CogMap* offering a structured spatial representation. Their combined use allows our model to outperform existing methods, showcasing the importance of integrating explicit spatial and stepwise cognitive inference.

Experiment (4) replaces the learned *Metric-CogMap* with a grid-based cognitive map used in VSI-Bench [45] and MindCube [48], leading to further performance degradation. This degradation indicates that coarse discretized spatial representations are insufficient for capturing continuous metric structure and fine-grained geometric relationships. Finally, experiment (5) uses ground-truth annotations to construct a perfect *Metric-CogMap*, providing an upperbound estimate of our framework's performance. This result demonstrates the potential of our method when supplied with ideal spatial priors, while also underscoring the limitations of current 2D and 3D detection/segmentation algorithms—particularly in tasks like *Rel.Dir* and *Rel.Dist*, as shown in Tab. 1. This points to a promising research direction: bridging the gap between state-of-the-art visual perception and ground-truth-level Metric-CogMap construction to further enhance the Map2Thought framework.

Table 3. Ablation study. We compare different map inputs and the effect of Cog-CoT.

| Method | Components | | Metrics | | |
|---|---|---|---|---|---|
| | Map Type | Cog-CoT | Avg | Abs.Dist | Rel.Dir |
| (1) Ours | Pred | ✓ | 58.8 | 53.0 | 70.9 |
| (2) Variant | Pred | - | 54.0 | 41.9 | 45.9 |
| (3) Baseline | None | - | 54.0 | 41.7 | 47.4 |
| (4) Variant | Grid | - | 49.7 | 19.4 | 46.4 |
| (5) Upper Bound | GT | ✓ | 73.7 | 81.3 | 86.1 |

## 5. Conclusion

In this work, we introduced Map2Thought, a framework that enables explicit and interpretable 3D spatial reasoning for vision-language models. By integrating the *Metric-CogMap* with *Cog-CoT* for structured geometric reasoning, our method supports transparent and verifiable inference. We also proposed a video-to-map extraction pipeline that bridges 2D perception and 3D geometry. Extensive experiments demonstrate that Map2Thought achieves strong data efficiency and superior generalization, advancing the goal of reliable, explainable spatial reasoning in 3D vision-language understanding. While our method shows robust performance, we observed that its effectiveness is bounded by the quality of the reconstructed *Metric-CogMap*. Future work will explore more accurate and robust map construction to further enhance spatial reasoning reliability.

# References

[1] Jean-Baptiste Alayrac, Jeff Donahue, Pauline Luc, Antoine Miech, Iain Barr, Yana Hasson, Karel Lenc, Arthur Mensch, Katherine Millican, Malcolm Reynolds, et al. Flamingo: a visual language model for few-shot learning. 2022. 2

[2] 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. 1

[3] Shuai Bai, Keqin Chen, Xuejing Liu, Jialin Wang, Wenbin Ge, Sibo Song, Kai Dang, Peng Wang, Shijie Wang, Jun Tang, et al. Qwen2.5-vl technical report. *arXiv:2502.13923*, 2025. 2

[4] Yihan Cao, Jiazhao Zhang, Zhinan Yu, Shuzhen Liu, Zheng Qin, Qin Zou, Bo Du, and Kai Xu. Cognav: Cognitive process modeling for object goal navigation with llms. In *ICCV*, 2025. 3

[5] Dave Zhenyu Chen, Angel X Chang, and Matthias Nießner. Scanrefer: 3d object localization in rgb-d scans using natural language. In *European conference on computer vision*, pages 202–221. Springer, 2020. 1

[6] Dave Zhenyu Chen, Qirui Wu, Matthias Nießner, and Angel X Chang. D 3 net: A unified speaker-listener architecture for 3d dense captioning and visual grounding. In *European Conference on Computer Vision*, pages 487–505. Springer, 2022. 1

[7] Shizhe Chen, Pierre-Louis Guhur, Makarand Tapaswi, Cordelia Schmid, and Ivan Laptev. Language conditioned spatial relation reasoning for 3d object grounding. In *NeurIPS*, 2022. 2

[8] 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. 2

[9] Zhenyu Chen, Ali Gholami, Matthias Nießner, and Angel X Chang. Scan2cap: Context-aware dense captioning in rgb-d scans. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 3193–3203, 2021. 1

[10] Zhenyu Chen, Ronghang Hu, Xinlei Chen, Matthias Nießner, and Angel X Chang. Unit3d: A unified transformer for 3d dense captioning and visual grounding. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 18109–18119, 2023. 1

[11] Zhe Chen, Weiyun Wang, Hao Tian, Shenglong Ye, Zhangwei Gao, Erfei Cui, Wenwen Tong, Kongzhi Hu, Jiapeng Luo, Zheng Ma, et al. How far are we to gpt-4v? closing the gap to commercial multimodal models with open-source suites. *arXiv preprint arXiv:2404.16821*, 2024. 7

[12] Angela Dai, Angel X. Chang, Manolis Savva, Maciej Halber, Thomas A. Funkhouser, and Matthias Nießner. Scannet: Richly-annotated 3d reconstructions of indoor scenes. In *CVPR*, 2017. 2, 5, 6, 1

[13] Afshin Dehghan, Gilad Baruch, Zhuoyuan Chen, Yuri Feigin, Peter Fu, Thomas Gebauer, Daniel Kurz, Tal Dimry, Brandon Joffe, Arik Schwartz, and Elad Shulman. Arkitscenes: A diverse real-world dataset for 3d indoor scene understanding using mobile RGB-D data. In *NeurIPS*, 2021. 2, 6, 1

[14] Jiajun Deng, Tianyu He, Li Jiang, Tianyu Wang, Feras Dayoub, and Ian Reid. 3d-llava: Towards generalist 3d lmms with omni superpoint transformer. In *Proceedings of the Computer Vision and Pattern Recognition Conference*, pages 3772–3782, 2025. 2

[15] Mohammed Munzer Dwedari, Matthias Niessner, and Dave Zhenyu Chen. Generating context-aware natural answers for questions in 3d scenes. *arXiv preprint arXiv:2310.19516*, 2023. 1

[16] Zhiwen Fan, Jian Zhang, Renjie Li, Junge Zhang, Runjin Chen, Hezhen Hu, Kevin Wang, Huaizhi Qu, Dilin Wang, Zhicheng Yan, Hongyu Xu, Justin Theiss, Tianlong Chen, Jiachen Li, Zhengzhong Tu, Zhangyang Wang, and Rakesh Ranjan. VLM-3R: Vision-language models augmented with instruction-aligned 3d reconstruction, 2025. 4, 6, 7, 1

[17] Zhiwen Fan, Jian Zhang, Renjie Li, Junge Zhang, Runjin Chen, Hezhen Hu, Kevin Wang, Huaizhi Qu, Dilin Wang, Zhicheng Yan, et al. VLM-3R: Vision-language models augmented with instruction-aligned 3d reconstruction. *arXiv preprint arXiv:2505.20279*, 2025. 2

[18] 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. In *NeurIPS*, 2023. 2

[19] Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, Weizhu Chen, et al. Lora: Low-rank adaptation of large language models. *ICLR*, 1(2):3, 2022. 6

[20] Jiangyong Huang, Silong Yong, Xiaojian Ma, Xiongkun Linghu, Puhao Li, Yan Wang, Qing Li, Song-Chun Zhu, Baoxiong Jia, and Siyuan Huang. An embodied generalist agent in 3d world. In *ICML*, 2024. 1, 2

[21] Xiaohu Huang, Jingjing Wu, Qunyi Xie, and Kai Han. Mllms need 3d-aware representation supervision for scene understanding. *arXiv preprint arXiv:2506.01946*, 2025. 2

[22] Aaron Hurst, Adam Lerer, Adam P Goucher, Adam Perelman, Aditya Ramesh, Aidan Clark, AJ Ostrow, Akila Welihinda, Alan Hayes, Alec Radford, et al. Gpt-4o system card. *arXiv:2410.21276*, 2024. 7

[23] Youngkyoon Jang and Eduardo Pérez-Pellitero. CoMapGS: Covisibility map-based gaussian splatting for sparse novel view synthesis. In *Proceedings of the Computer Vision and Pattern Recognition Conference (CVPR)*, pages 26779–26788, 2025. 5

[24] Chao Jia, Yinfei Yang, Ye Xia, Yi-Ting Chen, Zarana Parekh, Hieu Pham, Quoc Le, Yun-Hsuan Sung, Zhen Li, and Tom Duerig. Scaling up visual and vision-language representation learning with noisy text supervision. 2021. 2

[25] Bo Li, Yuanhan Zhang, Dong Guo, Renrui Zhang, Feng Li, Hao Zhang, Kaichen Zhang, Yanwei Li, Ziwei Liu, and Chunyuan Li. Llava-onevision: Easy visual task transfer. *arXiv:2408.03326*, 2024. 7

[26] Haoyuan Li, Yanpeng Zhou, Yufei Gao, Tao Tang, Jianhua Han, Yujie Yuan, Dave Zhenyu Chen, Jiawang Bian, Hang Xu, and Xiaodan Liang. Does your 3d encoder really work? when pretrain-sft from 2d vlms meets 3d vlms. *arXiv preprint arXiv:2506.05318*, 2025. 2

[27] Junnan Li, Dongxu Li, Silvio Savarese, and Steven Hoi. Blip-2: Bootstrapping language-image pre-training with frozen image encoders and large language models. In *ICML*, 2023. 2

[28] Zeju Li, Chao Zhang, Xiaoyan Wang, Ruilong Ren, Yifan Xu, Ruifei Ma, Xiangde Liu, and Rong Wei. 3dmit: 3d multi-modal instruction tuning for scene understanding. In *2024 IEEE International Conference on Multimedia and Expo Workshops (ICMEW)*, pages 1–5. IEEE, 2024. 2

[29] Ji Lin, Hongxu Yin, Wei Ping, Pavlo Molchanov, Mohammad Shoeybi, and Song Han. VILA: on pre-training for visual language models. In *CVPR*, 2024. 7

[30] Shilong Liu, Zhaoyang Zeng, Tianhe Ren, Feng Li, Hao Zhang, Jie Yang, Qing Jiang, Chunyuan Li, Jianwei Yang, Hang Su, Jun Zhu, and Lei Zhang. Grounding DINO: Marrying dino with grounded pre-training for open-set object detection. In *ECCV*, page 38–55, 2024. 5

[31] Xianzheng Ma, Brandon Smart, Yash Bhalgat, Shuai Chen, Xinghui Li, Jian Ding, Jindong Gu, Dave Zhenyu Chen, Songyou Peng, Jia-Wang Bian, et al. When llms step into the 3d world: A survey and meta-analysis of 3d tasks via multi-modal large language models. *arXiv preprint arXiv:2405.10255*, 2024. 2

[32] Zhangyang Qi, Zhixiong Zhang, Ye Fang, Jiaqi Wang, and Hengshuang Zhao. Gpt4scene: Understand 3d scenes from videos with vision-language models. *arXiv:2501.01428*, 2025. 2

[33] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. 2021. 2, 3

[34] Nikhila Ravi, Valentin Gabeur, Yuan-Ting Hu, Ronghang Hu, Chaitanya Ryali, Tengyu Ma, Haitham Khedr, Roman Rädle, Chloe Rolland, Laura Gustafson, Eric Mintun, Junting Pan, Kalyan Vasudev Alwala, Nicolas Carion, Chao-Yuan Wu, Ross Girshick, Piotr Dollár, and Christoph Feichtenhofer. SAM 2: Segment anything in images and videos. In *ICLR*, 2025. 4, 5

[35] Gemini Team, Petko Georgiev, Ving Ian Lei, Ryan Burnell, Libin Bai, Anmol Gulati, Garrett Tanzer, Damien Vincent, Zhufeng Pan, Shibo Wang, et al. Gemini 1.5: Unlocking multimodal understanding across millions of tokens of context. *arXiv:2403.05530*, 2024. 7

[36] Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, et al. Llama 2: Open foundation and fine-tuned chat models. *arXiv preprint arXiv:2307.09288*, 2023. 2

[37] Haochen Wang, Yucheng Zhao, Tiancai Wang, Haoqiang Fan, Xiangyu Zhang, and Zhaoxiang Zhang. Ross3d: Reconstructive visual instruction tuning with 3d-awareness, 2025. 2

[38] Jianyuan Wang, Minghao Chen, Nikita Karaev, Andrea Vedaldi, Christian Rupprecht, and David Novotny. VGGT: Visual geometry grounded transformer. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2025. 1, 2

[39] Qianqian Wang, Yifei Zhang, Aleksander Holynski, Alexei A. Efros, and Angjoo Kanazawa. Continuous 3d perception model with persistent state. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 10510–10522, 2025. 1, 2, 3

[40] Ruicheng Wang, Sicheng Xu, Yue Dong, Yu Deng, Jianfeng Xiang, Zelong Lv, Guangzhong Sun, Xin Tong, and Jiaolong Yang. MoGe-2: Accurate monocular geometry with metric scale and sharp details, 2025. 5

[41] Yifan Wang, Jianjun Zhou, Haoyi Zhu, Wenzheng Chang, Yang Zhou, Zizun Li, Junyi Chen, Jiangmiao Pang, Chunhua Shen, and Tong He. $\pi^3$: Scalable permutation-equivariant visual geometry learning, 2025. 4, 5

[42] Zehan Wang, Haifeng Huang, Yang Zhao, Ziang Zhang, and Zhou Zhao. Chat-3d: Data-efficiently tuning large language model for universal dialogue of 3d scenes. *arXiv preprint arXiv:2308.08769*, 2023. 2

[43] Diankun Wu, Fangfu Liu, Yi-Hsin Hung, and Yueqi Duan. Spatial-mllm: Boosting mllm capabilities in visual-based spatial intelligence, 2025. 2, 7

[44] Yueming Xu, Jiahui Zhang, Ze Huang, Yurui Chen, Yanpeng Zhou, Zhenyu Chen, Yu-Jie Yuan, Pengxiang Xia, Guowei Huang, Xinyue Cai, et al. Uniugg: Unified 3d understanding and generation via geometric-semantic encoding. *arXiv preprint arXiv:2508.11952*, 2025. 2

[45] Jihan Yang, Shusheng Yang, Anjali W. Gupta, Rilyn Han, Li Fei-Fei, and Saining Xie. Thinking in space: How multimodal large language models see, remember, and recall spaces. *arXiv:2412.14171*, 2024. 2, 3, 5, 7, 8, 1

[46] Shuquan Ye, Dongdong Chen, Songfang Han, and Jing Liao. 3d question answering. *IEEE Transactions on Visualization and Computer Graphics*, 30(3):1772–1786, 2022. 1

[47] Chandan Yeshwanth, Yueh-Cheng Liu, Matthias Nießner, and Angela Dai. Scannet++: A high-fidelity dataset of 3d indoor scenes. In *ICCV*, 2023. 2, 6, 1

[48] Baiqiao Yin, Qineng Wang, Pingyue Zhang, Jianshu Zhang, Kangrui Wang, Zihan Wang, Jieyu Zhang, Keshigeyan Chandrasegaran, Han Liu, Ranjay Krishna, et al. Spatial mental modeling from limited views. In *Structural Priors for Vision Workshop at ICCV'25*, 2025. 2, 3, 8

[49] Hanxun Yu, Wentong Li, Song Wang, Junbo Chen, and Jianke Zhu. Inst3d-lmm: Instance-aware 3d scene understanding with multi-modal instruction tuning, 2025. 1, 2

[50] Haoyu Zhang, Meng Liu, Zaijing Li, Haokun Wen, Weili Guan, Yaowei Wang, and Liqiang Nie. Spatial understanding from videos: Structured prompts meet simulation data. In *NeurIPS*, 2025. 3

[51] Peiyuan Zhang, Kaichen Zhang, Bo Li, Guangtao Zeng, Jingkang Yang, Yuanhan Zhang, Ziyue Wang, Haoran Tan, Chunyuan Li, and Ziwei Liu. Long context transfer from language to vision. *arXiv preprint arXiv:2406.16852*, 2024. 7

[52] Yuanhan Zhang, Bo Li, haotian Liu, Yong jae Lee, Liangke Gui, Di Fu, Jiashi Feng, Ziwei Liu, and Chunyuan Li. LLaVA-NeXT: A strong zero-shot video understanding model, 2024. 7

[53] Duo Zheng, Shijia Huang, Yanyang Li, and Liwei Wang. Learning from videos for 3d world: Enhancing mllms with 3d vision geometry priors. *arXiv preprint arXiv:2505.24625*, 2025. 2, 7

[54] Duo Zheng, Shijia Huang, and Liwei Wang. Video-3d LLM: learning position-aware video representation for 3d scene understanding. In *CVPR*, 2025. 1, 2

[55] Xingyi Zhou, Rohit Girdhar, Armand Joulin, Philipp Krähenbühl, and Ishan Misra. Detecting twenty-thousand classes using image-level supervision. In *ECCV*, 2022. 4

[56] Chenming Zhu, Tai Wang, Wenwei Zhang, Jiangmiao Pang, and Xihui Liu. Llava-3d: A simple yet effective pathway to empowering lmms with 3d-awareness. *arXiv:2409.18125*, 2024. 1, 2

# Map2Thought: Explicit 3D Spatial Reasoning via Metric Cognitive Maps

## Supplementary Material

## Appendix Overview

This supplementary document provides additional implementation details and results that support and extend the main paper. It is organized as follows:
- **Supp: Additional *Metric-CogMap* examples for evaluation**. Qualitative visualizations of constructed *Metric-CogMap*s across datasets.
- **Supp: *Metric-CogMap* construction for training**. Details of how *Metric-CogMap* annotations are generated for Map2Thought.
- **Supp: *Cog-CoT* construction and Complete QA.** Details of *Cog-CoT* construction. Also, complete examples of questions, maps, and corresponding reasoning traces.
- **Supp: Zero-shot Qualitative Demonstrations**. Qualitative examples showing how *Metric-CogMap* and *Cog-CoT* improve zero-shot spatial reasoning with commercial VLMs.

## A. Additional *Metric-CogMap* examples for evaluation

While the main paper provides an overview of our *Metric-CogMap* construction pipeline, space constraints limited the inclusion of more visual examples. To address this, we include additional visualizations as shown in Fig. S1, showcasing scene-level *Metric-CogMaps* produced by our evaluation-time pipeline across the three datasets in VSI-Bench [45]. These examples demonstrate how our pipeline identifies and integrates query-relevant objects into a unified, metrically accurate representation, which supports downstream spatial reasoning tasks across multiple queries within each scene. Note that the construction pipeline described in the main manuscript (Sec.3.3) is specific to the evaluation phase; a separate, ground-truth annotation-based procedure is used during training, independent of the evaluation-time construction process, as detailed in Sec. B.

## B. *Metric-CogMap* construction for training

We build our training dataset upon VLM-3R [16], which provides instructional question-answer pairs from Scan-Net [12], ScanNet++ [47], and ARKitScenes [13]. As noted earlier, while Sec.3.3 in the main manuscript outlines the *Metric-CogMap* construction process for evaluation, this section details the complementary pipeline used for the training dataset construction. To facilitate explicit spatial computation in Cog-CoT inference, we augment each training scene with a ground-truth *Metric-CogMap*. The construction process involves the following steps:

**Metadata extraction.** To construct our training dataset, we extract structured metadata from three RGB-D scene datasets: ScanNet [12], ScanNet++ [47], and ARKitScenes [13]. Each scene's annotated 3D point cloud includes spatial coordinates $(x, y, z)$, RGB values, semantic labels, and instance IDs. From this, we extract geometric and semantic information essential for training. We extract the following metadata for each scene:
- *Scene-level attributes:* (1) Room area, computed via the convex hull of floor points; (2) Room center, defined as the centroid of all 3D points; (3) Video path, linking each scene to its corresponding RGB-D sequence.
- *Object-level attributes:* (1) Object counts per semantic category; (2) Axis-aligned 3D bounding boxes for each instance, parameterized by centroid, spatial extents, normalized rotation axes, and min/max coordinates.

For label mapping, we handle two annotation schemes: ScanNetV2 (20 categories) and ScanNet200 (200 categories). Raw semantic labels from the point clouds are mapped to these standardized labels using the official label mapping files. We construct unified label mapping functions $f : \mathcal{L}_{\text{raw}} \to \mathcal{L}_{\text{target}}$ that convert raw label IDs to target category indices.

Bounding boxes are computed per instance by determining the min/max coordinates of each instance's point cluster along the three spatial axes. These boxes efficiently capture object extents and are used throughout our spatial reasoning pipeline. The extraction process is parallelized across scenes using multiprocessing to handle the large-scale datasets efficiently. Invalid or incomplete scenes (e.g., missing PLY files, corrupted annotations) are automatically filtered out during processing. Extracted metadata is stored in JSON format for efficient loading during training.

**Load scene XY range.** To facilitate spatial reasoning and metric map construction, we extract the spatial boundaries of each scene from the point cloud. For each scene, we load the mesh-aligned point cloud and compute the axis-aligned bounding box that encompasses all points in the scene. Specifically, we determine the minimum and maximum coordinates along the $x$ and $y$ axes: $[x_{\min}, x_{\max}]$ and $[y_{\min}, y_{\max}]$. To ensure uniform spatial representation and simplify downstream map discretization, we normalize these bounds to form a square region. We compute the extent along each axis: $\Delta x = x_{\max} - x_{\min}$ and $\Delta y = y_{\max} - y_{\min}$. If $\Delta x > \Delta y$, we expand the $y$ bounds symmetrically by $(\Delta x - \Delta y)/2$ on each side; conversely, if $\Delta y > \Delta x$, we expand the $x$ bounds by $(\Delta y - \Delta x)/2$.

This normalization ensures that the scene occupies

Figure S1. **Examples of evaluation *Metric-CogMap* outputs on the VSI-Bench [45].** Each row showcases a representative scene from one of the three datasets used in VSI-Bench: (top) ScanNet, (middle) ScanNet++, and (bottom) ARKitScenes. Above each scene visualization, we display a sample query with highlighted object mentions (green background). Within the corresponding 3D visualization, the mentioned objects are labeled in green text and enclosed in yellow axis-aligned bounding boxes, as produced by the *Metric-CogMap* construction pipeline. Additionally, objects mentioned in other queries from the same scene are shown in gray text. This reflects the unified nature of the *Metric-CogMap*, which integrates all detected and reconstructed objects relevant to multiple queries within a single scene-level map.

a square spatial region, which is essential for constructing grid-based metric cognitive maps with consistent resolution. The extracted boundaries are stored as `room_x_bound` and `room_y_bound` for each scene, representing the final square-normalized spatial extent. This preprocessing is applied uniformly across all three datasets (ScanNet, ScanNet++, and ARKitScenes) for both training and validation splits. The boundary information is saved in JSON format and later used during metric map construction to discretize continuous 3D space into a structured grid representation.

**Make *Metric-CogMap* for each scene** To support both symbolic spatial reasoning and precise metric computation, we construct a *Metric-CogMap* for each scene that encodes a dual-format representation: discrete grid coordinates for efficient relational queries and continuous metric-scale coordinates for accurate geometric analysis. Given the normalized scene boundaries $[x_{\min}, x_{\max}]$ and $[y_{\min}, y_{\max}]$, we represent each object in two parallel forms. First, for the discrete representation, we map each object's 3D centroid $(x, y, z)$ to a 2D grid cell $(g_x, g_y)$ using:

$$g_x = \left\lfloor \frac{x - x_{\min}}{x_{\max} - x_{\min}} \times N \right\rfloor, \quad g_y = \left\lfloor \frac{y - y_{\min}}{y_{\max} - y_{\min}} \times N \right\rfloor \tag{1}$$

, where $N$ is the grid size (default: 20), and grid indices are clipped to $[0, N-1]$. Similarly, we also compute discrete bounding boxes $([g_{x_{\min}}, g_{x_{\max}}], [g_{y_{\min}}, g_{y_{\max}}])$ by projecting the object's axis-aligned bounding box corners to the grid.

Second, for the metric-scale representation, we store the original 3D centroid $(x, y, z)$ in meters and extract the object's physical dimensions $(w, h, d)$ from its axis-aligned bounding box. This dual-format design allows the model to leverage discrete grids for efficient (symbolic) spatial relationship reasoning (e.g., relative direction) while enabling fine-grained metric reasoning (e.g., absolute distance) when precision is required.

For each scene, the *Metric-CogMap* encodes the following structured information:

(1) **Scene-level metadata:** The scene identifier (`scene_id`), grid configuration (`grid_size`, default: $20 \times 20$), and room dimensions (`room_size`) specifying the physical extent of the scene in meters along each axis.

(2) **Object-level information:** Objects are organized by semantic category (e.g., *door*, *towel*, *mirror*), with each category containing a list of object instances. Each object instance encodes both discrete and continuous spatial attributes:
- `instance_id`: A unique identifier for each object instance within the scene.

- `grid_position`: Discrete 2D coordinates $[g_x, g_y]$ representing the grid cell occupied by the object's plane centroid, enabling efficient spatial indexing and relationship queries.
- `grid_box_position`: Grid-aligned bounding box $[[g_{x_{\min}}, g_{x_{\max}}], [g_{y_{\min}}, g_{y_{\max}}]]$ representing the discrete spatial extent of the object in the discrete grid coordinates.
- `centroid`: Continuous 3D coordinates $[x, y, z]$ in meters, preserving precise metric-scale location information from the original point cloud.
- `bounding_box_size`: Physical dimensions $[w, h, d]$ in meters, computed from the axis-aligned bounding box, enabling accurate size-based reasoning.

This hybrid structure enables the model to bridge symbolic spatial reasoning and metric-scale computation, supporting both qualitative spatial queries and quantitative distance estimations during Cog-CoT inference. The resulting *Metric-CogMap* is serialized in JSON format for seamless integration into the training pipeline.

**Object selection for each QA pair** To manage reasoning complexity, we identify the object categories referenced in each question–answer pair. This step filters the relevant entities from the scene's *Metric-CogMap*, enabling the model to focus computation on task-relevant objects rather than processing all instances in the scene. The extraction procedure relies on category-aware text parsing. For each question, we apply word–boundary–based pattern matching using the category vocabulary predefined for each dataset. Specifically, we use regex word-boundary matching ($\backslash$b*category*$\backslash$b) to ensure exact category matching and avoid spurious detections (e.g., preventing "table" from matching "vegetable").

To address cases where a shorter category name is embedded within a longer one (e.g., "box" within "tissue box"), we apply a positional filtering strategy. A category is retained only if at least one of its matched spans in the text is not entirely contained within the span of another category. This prevents us from incorrectly counting a shorter category (e.g., "box") when it appears only as part of a longer one (e.g., "tissue box"), while still allowing both categories to be detected when they genuinely appear separately in the question. Additionally, we normalize common linguistic variations through plural-to-singular mapping (e.g., "doors" → "door") and synonym resolution (e.g., "staircase" → "stairs")improving robustness to phrasing differences.

The final output for each QA pair includes: (1) question identifier and scene name; (2) the original question text; and (3) the list of identified object categories (`found_categories`). This structured representation enables the model to retrieve only the relevant object in-

stances from the scene's *Metric-CogMap* during inference, significantly reducing computational overhead and improving precision in spatially grounded reasoning.

**Create *Metric-CogMap* for each QA pair.** Once the relevant object categories are identified for a given question, we construct a question-specific *Metric-CogMap* by extracting only the corresponding object instances from the full scene-level map. This targeted subset serves as the spatial context for the question, reducing noise and focusing the model's attention on task-relevant elements. The resulting map is embedded into the input prompt alongside task-specific reasoning instructions, enabling efficient and interpretable spatial reasoning.

**(1) Question-specific map construction.** For each QA pair, we retrieve all instances of the identified categories from the scene's complete *Metric-CogMap* to create a compact, question-specific JSON representation. This tailored map includes four components:

- **Cognitive map:** Discrete grid positions $[g_x, g_y]$ for each object instance, e.g., `{"trash bin": [[5, 13]], "nightstand": [[14, 11]]}`.
- **Cognitive box map:** Grid-aligned bounding boxes $[[g_{x_{\min}}, g_{x_{\max}}], [g_{y_{\min}}, g_{y_{\max}}]]$, e.g., `{"trash bin": [[[4, 5], [12, 13]]]}`.
- **Box centroid:** Continuous metric 3D coordinates $[x, y, z]$ in meters, e.g., `{"trash bin": [[-1.5, 0.58, 0.12]]}`.
- **Box size:** Physical dimensions $[w, h, d]$ in meters, e.g., `{"trash bin": [[0.36, 0.30, 0.37]]}`.

**(2) Task-specific reasoning Instructions.** The question-specific *Metric-CogMap* is serialized in JSON format and appended to the original question text. To enable explicit spatial reasoning, we additionally inject task-specific computational instructions that explicitly instruct the model on the required mathematical operations. Below are representative examples:

- *Relative direction: "Task: Determine relative direction (front-left/front-right/back-left/back-right) using vector analysis based on cognitive map. Format: Origin = [x,y], Facing = [x,y], Target = [x,y] → f = Facing-Origin → t = Target-Origin → Dot(f,t) for front/back → Cross(f,t) for left/right. Output: Mathematical calculations with explicit coordinate subtraction, dot/cross product computation, and directional conclusion."*
- *Relative distance: "Task: Compare distances between multiple candidates and a target to find the closest one based on cognitive map and cognitive box map. Format: Target = [x,y], Candidates = [obj1, obj2, ...] → Distance(obj1,target) = calculation → Distance(obj2,target) = calculation → Compare all distances → Closest object = result. Output: Show distance calculations for each*

candidate with explicit mathematical steps and final comparison."

- *Absolute distance: "Task: Calculate absolute distance between two specific objects based on box centroid and box size. Format: Object1: Centroid c1 = [x,y,z], Size = [w,h,d], Half size s1 → Object2: Centroid c2, Size, Half size s2 → Centroid difference → Half size sum → Rough distance calculation. Output: Step-by-step box data extraction, centroid calculations, and rough distance approximation with explicit numerical values."*
- *Room size: "Task: Estimate room area using the number of valid cells calculation. Format: Number valid cells → Multiply by each cell area (0.36 m²). Answer with format like* `<answer>final_room_size</answer>`*. Output: Explicit calculations with multiplication, rough room size approximation."*

This structured representation provides the model with both the spatial data and explicit computational guidance needed for grounded spatial reasoning.

## C. *Cog-CoT* construction and Complete QA.

**Generating *Cog-CoT* for each QA pair.** To enable the model to perform explicit spatial reasoning over the *Metric-CogMap*, we augment each training QA pair with programmatically generated chain-of-thought (Cog-CoT) supervision. This supervision encodes step-by-step spatial computations tailored to the reasoning requirements of each question type.

**Programmatic CoT generation.** For each question category, we design dedicated algorithms that execute the relevant geometric operations on the *Metric-CogMap* to generate verifiable intermediate reasoning steps.

(a) **Relative direction CoT:** The algorithm extracts object coordinates from the cognitive map and computes the facing vector $\mathbf{f} = \text{Facing} - \text{Origin}$ and target vector $\mathbf{t} = \text{Target} - \text{Origin}$ with explicit coordinate subtraction. It then calculates the dot product $\mathbf{f} \cdot \mathbf{t} = f_x \cdot t_x + f_y \cdot t_y$ to determine whether the target is in front ($> 0$) or behind ($< 0$), and the 2D cross product $f_x t_y - f_y t_x$ to determine whether the target lies to the left ($> 0$) or right ($< 0$), outputting all intermediate numerical values.

The complete question, *Metric-Cogmap*, *Cog-CoT*, and answer for *relative direction* is illustrated in Fig. S2.

(b) **Room size CoT:** The algorithm operates on the $20 \times 20$ occupancy grid, where each cell corresponds to an area of $0.36m^2$. It counts the number of valid (occupied) cells, multiplies by the per-cell area to compute the total room size, and explicitly shows the calculation steps (e.g., "Total cells = 76, total area = $76 \times 0.36 =$

## Relative Direction QA:

**Question:** These are frames of a video.\n If I am standing by the window and facing the trash bin, is the lamp to my left, right, or back?\nAn object is to my back if I would have to turn at least 135 degrees in order to face it.\nOptions:\nA. back\nB. left\nC. right\n

This is a scene information, including object cognitive map, centroid, size:{\

"cognitive map\": {\"trash bin\":[[9,17]],\"lamp\":[[15,10]],\"window\":[[9,0]]},\

"cognitive box map\": {\"trash bin\":[[[8,9],[16,18]]],\"lamp\":[[[15,15],[10,11]]],\"window\":[[[5,14],[0,1]]]},\

"box centroid\": {\"trash bin\":[[-0.34,1.72,0.32]],\"lamp\":[[1.29,0.02,1.23]],\"window\":[[-0.19,-2.62,1.17]]},\

"box size\": {\"trash bin\":[[0.39,0.36,0.49]],\"lamp\":[[0.16,0.15,0.1]],\"window\":[[2.21,0.32,1.09]]}}

**Task:** Determine relative direction (front-left/front-right/back-left/back-right) using vector analysis based on cognitive map.\n

**Format:** Origin = [x,y], Facing = [x,y], Target = [x,y] → f = Facing-Origin → t = Target-Origin → Dot(f,t) for front/back → Cross(f,t) for left/right.\n

**Output:** Mathematical calculations with explicit coordinate subtraction, dot/cross product computation, and directional conclusion.

**Answer:** <think>Origin = window = [9, 0], Facing = trash bin = [9, 17], Target = lamp = [15, 10]\n

f = Facing - Origin = (9-9, 17-0) = [0, 17]\n

t = Target - Origin = (15-9, 10-0) = [6, 10]\n

Dot(f, t) = ((0*6)+(17*10)) = 170 > 0 → front\n

Cross(f, t) = ((0*10)-(17*6)) = -102 < 0 → right</think>

The final answer should be: <answer>C</answer>

Figure S2. **Examples of *Metric-CogMap* and *CogCoT* on relative_direction**

## Room Size QA:

**Question:** These are frames of a video.\nWhat is the size of this room (in square meters)? \nIf multiple rooms are shown, estimate the size of the combined space.\n

This is a scene layout represented by a 20*20 grid. Each valid cell corresponds to 0.36 square meters. The room size can be estimated from the number of valid cells.The scene has 76 cells.

**Task:** Estimate room area using the number of valid cells calculation.\n

**Format:** Number valid cells → Multiply by each cell area (0.36 $m^2$). Answer with format like <answer>final_room_size</answer>\n

**Output:** Explicit calculations with multiplication, rough room size approximation.

**Answer:** <think>The scene has 76 cells. The room size is about 76 * 0.36 = 27.36 square meters. Based on the video and rough room size, I can estimate the final room size.</think>

The final answer should be: <answer>24.9</answer>

Figure S3. **Examples of *Metric-CogMap* and *CogCoT* on room_size**

$27.36m^2$").

The complete question, *Metric-Cogmap*, *Cog-CoT*, and answer for *room size* is illustrated in Fig. S3.

(c) **Absolute distance CoT:** To estimate the distance between two specified objects, the algorithm retrieves their centroids $c_1 = [x_1, y_1, z_1]$ and $c_2 = [x_2, y_2, z_2]$ from the box centroid data, and their sizes from the box size data. It computes half-sizes $s_1 = \text{Size}_1/2$ and $s_2 = \text{Size}_2/2$ component-wise, calculates the absolute centroid difference $\Delta = |c_1 - c_2|$, sums the half-sizes $s = s_1 + s_2$, and estimates the approximate object-to-object distance as dist $= \max(\Delta - s, 0)$ component-wise. All arithmetic operations are explicitly shown with numerical values (e.g., "Centroid difference $= \|[-1.50, 0.58, 0.12] - [1.76, -0.10, 0.34]\| = [3.26, 0.68, 0.22]$").

The complete question, *Metric-Cogmap*, *Cog-CoT*, and answer for *absolute distance* is illustrated in Fig. S4.

(d) **Appearance order CoT:** To infer appearance order, we leverage the sparse 3D point clouds stored in the *Metric-CogMap*, which are quantized per object instance using the predicted camera poses provided by $\pi^3$. Since precise geometry is not strictly required for this task, we use the precomputed quantized point clouds and project them onto each of the 64 subsampled frames used as input to the VLM. For each object instance, we project its 3D points into the image plane of each frame and record the first frame index (ranging from 0 to 63) where the object becomes visibly present. To avoid false positives caused by occlusions (e.g., when a wall or another object blocks the target instance), we apply a visibility filter: projected points are discarded if there exists a reconstructed (unlabeled) point in front of them—closer to the camera by more than a threshold of 0.2 meters (20 centimeters in metric scale). This filtering ensures that only unoccluded instances are considered when identifying appearance frames.

For temporal ordering questions, the algorithm receives frame detection indices for each object (e.g., "door: frame 30, printer: frame 249, mouse: frame 253"), sorts the objects in ascending order by frame index, and determines their relative order (e.g., which appeared first or last). In cases of missing detections (e.g., "kettle not detected"), the reasoning chain explicitly includes a note about the absence, allowing the model to learn inference under incomplete observations. The reasoning concludes by matching the derived order against multiple-choice options.

The complete question, *Metric-Cogmap*, *Cog-CoT*, and answer for *appearance order* is illustrated in Fig. S5.

(e) **Relative distance CoT:** For queries that ask which object is closest among several candidates, the algorithm computes axis-aligned bounding box (AABB) distances. For each candidate category, it retrieves all object instances from the *Metric-CogMap*. Then, for each instance–target pair, it calculates per-axis separation: if the boxes do not overlap on axis $x$, then $d_x = \max(x_{\text{cand}}^{\min} - x_{\text{tgt}}^{\max}, x_{\text{tgt}}^{\min} - x_{\text{cand}}^{\max})$; otherwise, $d_x = 0$. The squared AABB distance $d^2 = d_x^2 + d_y^2$ is computed for each instance, showing detailed calculations (e.g., "dx = 15 - 14 = 1, dy = 0 (y-axis overlap), AABB distance $= 1^2 + 0^2 = 1.0$"). The minimum distance is selected per category, and all candidates are compared to identify the closest object.

The complete question, *Metric-Cogmap*, *Cog-CoT*, and answer for *relative distance* is illustrated in Fig. S6.

**Structured response format.** Each training sample's ground-truth response is structured as:

```
<think>step-by-step reasoning
with numerical computations</think>
The final answer should be:
<answer>answer</answer>.
```

This format enables the model to learn both the spatial reasoning process and the final answer prediction jointly. By providing explicit computational demonstrations through Cog-CoT, the model learns to decompose complex spatial questions into verifiable symbolic operations over the *Metric-CogMap* representation, effectively bridging natural language understanding with geometric reasoning.

## D. Zero-shot qualitative demonstrations

We conducted a small-scale zero-shot evaluation to assess whether augmenting a generic vision-language model with our proposed *Metric-CogMap* and *Cog-CoT* prompts improves spatial reasoning, even without task-specific fine-tuning. To this end, we tested a few randomly selected examples using publicly available commercial models, following the same input protocol as VSI-Bench. Specifically, we used a zero-shot setting in which only 10 subsampled video frames (the maximum allowed for upload) and the test query were provided.

Figs. S7 and S8 illustrate comparisons between two setups for Gemini-3-Pro and GPT-5-Thinking, respectively: (a) using only images and the question as input, and (b) the same models augmented with *Metric-CogMap* and *Cog-CoT* prompts. In Fig. S7, the Gemini-3-Pro baseline fails to answer a relative direction question, while the ours successfully infers the correct spatial relation. Similarly, in Fig. S8, the GPT-5-Thinking baseline struggles with distance esti-

**Absolute Distance QA:**

Question: These are frames of a video.\nMeasuring from the closest point of each object, what is the direct distance between the trash bin and the nightstand (in meters)?\n

This is a scene information, including object cognitive map, centroid, size: {\

"cognitive map\":{\"trash bin\":[[5,13]],\"nightstand\":[[14,11]]},\

"cognitive box map\":{\"trash bin\":[[[4,5],[12,13]]],\"nightstand\":[[[13,15],[10,12]]]},\

"box centroid\":{\"trash bin\":[[-1.5,0.58,0.12]],\"nightstand\":[[1.76,-0.1,0.34]]},\

"box size\":{\"trash bin\":[[0.36,0.3,0.37]],\"nightstand\":[[0.61,0.88,0.72]]}}

Task: Calculate absolute distance between two specific objects based on box centroid and box size.\n

Format: Object1: Centroid $c_1$ = [x,y,z], Size = [w,h,d], Half size $s_1$ → Object2: Centroid $c_2$, Size, Half size $s_2$ → Centroid difference → Half size sum → Rough distance calculation.\n

Output: Step-by-step box data extraction, centroid calculations, and rough distance approximation with explicit numerical values.

Answer: <think>Objects: trash bin and nightstand\n

trash bin:\n  Centroid $c_1$ = [-1.50, 0.58, 0.12]\n  Size = [0.36, 0.30, 0.37]\n  Half size $s_1$ = [0.18, 0.15, 0.18]\n

nightstand:\n  Centroid $c_2$ = [1.76, -0.10, 0.34]\n Size = [0.61, 0.88, 0.72]\n  Half size $s_2$ = [0.30, 0.44, 0.36]\n

Centroid difference = $|c_1 - c_2|$ = [3.26, 0.68, 0.22]\n

Half size sum = $s_1 + s_2$ = [0.48, 0.59, 0.54]\n

Rough distance = |Relative difference - Half size sum| = [2.77, 0.09, 0.00]\n

Based on the Rough distance and the input rgb video, I can estimate the final answer.</think>

The final answer should be: <answer>2.9</answer>

Figure S4. **Examples of *Metric-CogMap* and *CogCoT* on absolute_distance**

**Appearance Order QA:**

Question: These are frames of a video.\nWhat will be the first-time appearance order of the following categories in the video: cup, suitcase, cabinet, jar?\n[Each object was first detected at frame: cup: frame 49, suitcase: frame 3, cabinet: frame 9] (Note: jar not detected)\nOptions:\nA. cabinet, cup, jar, suitcase\nB. cabinet, jar, cup, suitcase\nC. cup, jar, cabinet, suitcase\nD. suitcase, cabinet, cup, jar

Answer: <think>\nLet me analyze the frame indices to determine the appearance order.\n\nSorting by frame index (earliest to latest):\n1. suitcase (frame 3) - appears first\n2. cabinet (frame 9)\n3. cup (frame 49) - appears last among detected objects\nThe correct answer is D.\n</think>\n

<answer>D</answer>

Figure S5. **Examples of *Metric-CogMap* and *CogCoT* on appearance_order**

Figure S6. **Examples of *Metric-CogMap* and *CogCoT* on relative_distance**

mation, whereas our zero-shot-enhanced prompt enables it to reason accurately.

Due to the limitations of commercial APIs, we could not perform full VSI-Bench evaluation. However, these preliminary examples demonstrate that our Map2Thought framework improves spatial understanding in zero-shot scenarios—even when applied to state-of-the-art proprietary models—highlighting the value of explicit spatial reasoning guidance for models that otherwise lack this capability.

(a) Baseline zero-shot result from Gemini-3-Pro using only image and text query inputs. The model fails to interpret the spatial relationship.



(b) Our zero-shot result using the same model augmented with *Metric-CogMap* and *Cog-CoT*. The model correctly reasons and answers the spatial query.

Figure S7. **Qualitative comparison of zero-shot spatial direction reasoning.** This example demonstrates the effect of embedding *Metric-CogMap* and *Cog-CoT* into a frozen Gemini-3-Pro model. While the baseline fails to resolve the `relative_direction` query using only image and text input, our prompt-enhanced approach enables successful geometric reasoning and accurate answer prediction.

GPT 5 accept frames and question as input.
Ground Truth is 1.1 m.
GPT 5 prediction is far from gt with
Metric-CogMap and CogCoT

Measuring from the closest point of each object, what is the distance between the table and the radiator (in meters)?

Thought for 11s ›

~0.45 m (approx. 45 cm).

(Estimated visually from the gap between the table edge and the radiator; without camera calibration or a reference object, this can't be measured exactly.)

(a) Baseline zero-shot result from ChatGPT-5-Thinking using only image and text query inputs. The model fails to measure metric distance.

Thought for 9s ›

GPT 5 accept frames and question as input.
Ground Truth is 1.1 m.
GPT 5 prediction is close to gt with
Metric-CogMap and CogCoT

Measuring from the closest point of each object, what is the distance between the table and the radiator (in meters)?\nThis is a scene information, including object cognitive map, centroid, size: {\"cognitive map\":{\"radiator\":[[10,15]],\"table\":[[10,11]]},\"cognitive box map\": {\"radiator\":[[[9,12],[14,15]],\"table\":[[[9,11],[9,13]]]},\"box centroid\": {\"radiator\":[[0.37,3.05,0.92]],\"table\":[[0.41,0.72,0.71]]},\"box size\": {\"radiator\":[[1.71,0.3,0.4]],\"table\":[[1.47,2.24,0.24]]}}Task: Calculate absolute distance between two specific objects based on box centroid and box size.\nFormat: Object1: Centroid c1 = [x,y,z], Size = [w,h,d], Half size s1 → Object2: Centroid c2, Size, Half size s2 → Centroid difference → Half size sum → Rough distance calculation.\nOutput: Step-by-step box data extraction, centroid calculations, and rough distance approximation with explicit numerical values.

Object1 (table):

Centroid $c_1$ = [0.41, 0.72, 0.71], Size = [1.47, 2.24, 0.24] → Half-size $s_1$ = [0.735, 1.12, 0.12]

Object2 (radiator):

Centroid $c_2$ = [0.37, 3.05, 0.92], Size = [1.71, 0.30, 0.40] → Half-size $s_2$ = [0.855, 0.15, 0.20]

Centroid difference $|c_2 - c_1|$ = [|0.37−0.41|, |3.05−0.72|, |0.92−0.71|] = [0.04, 2.33, 0.21]

Half-size sums (axis-aligned): $s_1 + s_2$ = [0.735+0.855, 1.12+0.15, 0.12+0.20] = [1.59, 1.27, 0.3

Axis separations (closest points):

$\max(|\Delta x| - 1.59, 0) = \max(0.04 - 1.59, 0) = 0$

$\max(|\Delta y| - 1.27, 0) = \max(2.33 - 1.27, 0) = 1.06$

$\max(|\Delta z| - 0.32, 0) = \max(0.21 - 0.32, 0) = 0$

Rough distance between boxes (Euclidean of separations):

$\sqrt{(0^2 + 1.06^2 + 0^2)}$ = 1.06 m

(b) Our zero-shot result using the same model augmented with Metric-CogMap and Cog-CoT. The model correctly calculates the distance and answers the spatial query.

Figure S8. **Qualitative comparison of zero-shot absolute distance estimation.** This example shows the impact of embedding Metric-CogMap and Cog-CoT into a frozen GPT-5-Thinking model. The baseline model fails to compute correct spatial distance in the `absolute_distance` task, while our zero-shot-enhanced prompt enables metric scale reasoning and correct inference.