

Thinking in Space: How MLLMS See, Remember, Recall Spaces

Ayush Bodade
ayushbodade1@gmail.com

16 January 2025

Abstract

Integrating spatial reasoning in Multimodal Large Language Models (MLLMs) can revolutionize 3D reconstruction and generation by leveraging cognitive maps for perception, alignment, and temporal reasoning. Key challenges in global awareness and egocentric-allocentric shifts must be addressed to unlock applications in robotics, and autonomous systems.

Contents

1	Introduction	2
2	Connections Between Thinking in Space and 3D Systems	2
3	Research Questions	4
4	Potential Applications	5

Introduction

Multimodal Large Language Models (MLLMs) trained on extensive video datasets exhibit emerging abilities to “think in space,” yet they remain far from human-level spatial intelligence. This gap poses significant challenges but also opportunities for enhancing 3D reconstruction and generation. By aligning spatial reasoning with 3D systems, MLLMs can bridge the divide between perception and action, enabling applications in robotics, autonomous navigation, and more. This report explores the interplay between spatial reasoning and 3D technologies, highlighting challenges in egocentric-allocentric transformations, global consistency, and dynamic scene understanding, while identifying future research directions to advance spatial intelligence and its real-world applications.

Connections Between Thinking in Space and 3D Systems

1. Perception and 3D Reconstruction: **Connection:** Spatial reasoning begins with accurate perception, which is the foundation of 3D reconstruction. Both rely on processing raw sensory data (e.g., video, depth, RGB-D) and building representations of objects and their spatial relationships.

Details:

- **Thinking in Space:** MLLMs perceive objects in relation to one another and estimate distances, directions, and sizes from video sequences.
- **3D Reconstruction:** This step corresponds to capturing geometry, texture, and object positions to reconstruct a digital twin of a scene.
- For example:
 - An MLLM’s cognitive map could act as a heuristic to resolve ambiguities in noisy or incomplete scans.
 - Reasoning-based inferences could determine occluded parts of objects or spaces (e.g., estimating a table leg hidden by a chair).

Future Potential: Integrating MLLM-powered reasoning with traditional 3D reconstruction pipelines could reduce the need for perfect visual input, improving performance in complex or cluttered environments.

2. Egocentric-Allocentric Transformation for Global Spatial Consistency:

Connection: Thinking in space involves transforming egocentric (first-person) views into allocentric (global, map-like) perspectives. Similarly, 3D reconstruction involves stitching multiple viewpoints into a unified, global 3D representation.

Details:

- Egocentric-allocentric transformations allow models to reason about spatial layouts irrespective of camera position or movement.
- In 3D reconstruction, aligning frames or point clouds from moving cameras requires global consistency across multiple viewpoints.

Future Potential: Advances in spatial alignment and perspective transformation could improve:

- Reconstruction of large, complex spaces (e.g., multi-room buildings).
- Global navigation in robotics using allocentric maps inferred by MLLMs.

3. Local vs. Global Spatial Awareness: Connection: MLLMs currently excel at *local reasoning* (e.g., understanding nearby object relationships) but struggle with integrating this into *global awareness*, which mirrors challenges in large-scale 3D reconstruction.

Details:

- Cognitive maps could act as modular components that are later stitched into global 3D representations, solving issues of scalability in large environments.
- 3D generation systems could use these local models for finer details while relying on global constraints (e.g., room dimensions or overall layout).

Future Potential: Combining cognitive maps with neural 3D scene reconstruction pipelines could improve the realism and scalability of reconstructed spaces.

4. Dynamic Spaces and Temporal 3D Reasoning: Connection: Videos allow MLLMs to observe spaces dynamically, which aligns with the need for *temporal reasoning* in dynamic 3D scene generation or reconstruction.

Details:

- Thinking in space involves tracking object positions and predicting changes over time (e.g., a door swinging open or a chair being moved).
- Dynamic 3D systems must account for these changes to create interactive environments or simulate realistic scenarios.

Future Potential: Integrating temporal reasoning could enable:

- Time-aware 3D reconstructions (e.g., showing how a room changes throughout the day).
- Interactive AR/VR systems with real-time updates to the virtual environment based on user actions.

5. Multi-modal Data for Enhanced Understanding: **Connection:** MLLMs combine language and vision, while 3D systems often rely on multi-modal inputs (e.g., RGB-D, LiDAR, and natural language for semantic labeling).

Details:

- MLLMs could bridge gaps in 3D reconstruction by adding semantic layers (e.g., “This is a chair; it should be near the table”).
- Multi-modal MLLMs trained on 3D datasets could directly infer spatial layouts and semantics from diverse inputs.

Future Potential: Applications like semantic 3D maps or object-aware 3D reconstructions for robotics, where understanding the purpose and placement of objects enhances functionality.

Research Questions

Fundamentals of Spatial Thinking:

1. Egocentric-allocentric alignment: How can MLLMs reliably transform first-person video data into globally consistent 3D reconstructions?
2. Unified cognitive map frameworks: Developing scalable methods for merging local cognitive maps into global spatial models.

3D Reconstruction and Generation:

3. Dynamic 3D reconstruction: Investigating how MLLMs can infer and represent changes in spaces over time for real-time applications.
4. Semantic-aware 3D generation: Combining spatial reasoning with semantic labels to create meaningful, functional 3D spaces (e.g., a “living room” with correctly placed furniture).

Cross-Domain Integration:

5. Vision-language-3D fusion: Exploring multi-modal training paradigms where MLLMs integrate visual and textual data to enhance 3D reconstruction/generation.
6. Temporal 3D reasoning: Building models that can reason about time to reconstruct evolving or interactive environments.

Applications and Deployment:

7. AR/VR for real-world spaces: Creating tools that adapt virtual overlays to real-world environments using MLLM-based spatial reasoning.

8. Spatial AI for autonomous systems: Developing MLLM-based reasoning for autonomous robots to operate in complex and dynamic real-world spaces.

Benchmarking and Evaluation:

9. VSI-enhanced benchmarks: Designing new benchmarks that test MLLMs' ability to reason about global consistency, temporal dynamics, and semantic understanding in 3D environments.
10. Human-in-the-loop systems: Researching collaborative frameworks where humans and MLLMs jointly refine cognitive maps or 3D reconstructions.

Potential Applications

1. Robotics and Autonomous Navigation:

- Robots could combine cognitive maps with real-time 3D reconstructions to navigate cluttered or unfamiliar environments.
- Enhanced reasoning about spaces would enable robots to:
 - Avoid obstacles intelligently.
 - Plan routes based on incomplete or dynamic spatial information.
 - Perform context-aware object manipulation (e.g., fetching a tool from a specific location).

2. AR/VR and Metaverse Applications:

- **Immersive experiences:** Spatially accurate and dynamically adaptive 3D environments can improve virtual training simulations, gaming, or virtual meetings.
- **Real-time updates:** AR glasses or VR headsets could use MLLMs to enhance the realism of overlays by reasoning about spatial relationships in the environment.

3. Architecture and Interior Design:

- Cognitive maps and 3D reconstruction could assist designers by:
 - Simulating how objects fit into spaces.
 - Automatically generating 3D layouts from blueprints or videos of empty rooms.
 - Proposing optimized layouts based on room dimensions and object relationships.

4. Autonomous Driving:

- MLLMs' spatial reasoning could integrate with 3D scene reconstruction to predict object trajectories, improve scene understanding, and navigate complex driving environments.
- Enhancing safety in scenarios requiring quick spatial reasoning, such as avoiding collisions or predicting pedestrian movements.

5. Healthcare and Assistive Technology:

- **Rehabilitation tools:** Virtual spaces for physical therapy that adapt dynamically to the user's progress.
- **Assistive robots:** Guiding visually impaired individuals using 3D reconstructions paired with spatial reasoning for real-time obstacle detection.

References

- [1] Jihan Yang, Shusheng Yang, Anjali W. Gupta, Rilyn Han, Li Fei-Fei, Saining Xie, "Thinking in Space: How Multimodal Large Language Models See, Remember, and Recall Spaces", arXiv:2412.14171.