
Autonomous Spatial Intelligence: A Comprehensive Technical Report on Agentic AI Methods, Architectures, and Evaluation

Gloria Felicia AtlasPro AI gloria.felicia@atlaspro.ai	Nolan Bryant AtlasPro AI nolan.bryant@atlaspro.ai	Handi Putra AtlasPro AI handi.putra@atlaspro.ai
Ayaan Gazali AtlasPro AI ayaan.gazali@atlaspro.ai	Eliel Lobo AtlasPro AI eliel.lobo@atlaspro.ai	Esteban Rojas AtlasPro AI esteban.rojas@atlaspro.ai

Abstract

The convergence of Agentic Artificial Intelligence and Spatial Intelligence marks a pivotal frontier in the pursuit of creating machines that can autonomously operate in the physical world. While agentic systems demonstrate increasingly sophisticated capabilities in planning and tool use, their ability to perceive, reason about, and interact with complex spatial environments remains a significant bottleneck. This technical report addresses a critical gap in the existing literature by providing a unified taxonomy that systematically connects the architectural components of agentic AI with the functional requirements of spatial intelligence. We review over 500 papers spanning foundational agentic architectures [Yao et al., 2023b, Shinn et al., 2023, Yao et al., 2023a, Wei et al., 2022], embodied AI systems [Wang et al., 2023a, Driess et al., 2023, Brohan et al., 2023, Ahn et al., 2022], vision-language-action models [Team et al., 2024, Kim et al., 2024, Liu et al., 2023b, Alayrac et al., 2022], graph neural networks for spatial reasoning [Kipf and Welling, 2017, Velickovic et al., 2018, Hamilton et al., 2017, Wu et al., 2019a], world models [Hafner et al., 2023, 2021, Yang et al., 2024, Hu et al., 2023a], and geospatial foundation models [Jakubik et al., 2024, Bastani et al., 2023b, Cong et al., 2022, Mendieta et al., 2023]. Through comprehensive analysis of state-of-the-art methods, industry applications from Palantir, ESRI, Foursquare, Google, and Waymo, and evaluation benchmarks, we provide a foundational reference for researchers and practitioners. By synthesizing these disparate research areas and outlining a forward-looking research roadmap, this paper aims to accelerate the development of robust, safe, and effective spatially-aware autonomous systems.

1 Introduction

The evolution of Artificial Intelligence is marked by a paradigm shift from specialized models to goal-oriented, self-directed agents capable of complex decision-making in dynamic environments [McCarthy et al., 1955, Newell et al., 1956, Turing, 1950]. This field, which we term **Agentic AI**, represents a significant leap towards creating machines that can operate with a higher degree of autonomy [Wang et al., 2024b, Xi et al., 2023, Weng, 2023]. The foundational work on large language models [Brown et al., 2020, OpenAI, 2023, Chowdhery et al., 2022, Touvron et al., 2023, Team and Google, 2023, Anthropic, 2024] has enabled a new generation of AI agents that can reason, plan, and execute complex tasks through natural language interfaces [Wei et al., 2022, Kojima et al., 2022, Wang et al., 2022].

Concurrently, the ability for these agents to perceive, comprehend, and act within the physical world, a capability we define as **Spatial Intelligence**, has become a primary bottleneck and a critical area of research [Chen et al., 2024a, Yang et al., 2025, Thompson et al., 2025]. The convergence of these two domains is essential for developing AI systems that can effectively and safely navigate real-world complexities, from autonomous vehicles [Sun et al., 2020, Caesar et al., 2020, Wilson et al., 2023, Hu et al., 2023b, Jiang et al., 2023] and robotic assistants [Brohan et al., 2023, Driess et al., 2023, Team et al., 2024] to large-scale urban planning [Zheng et al., 2014, Jin et al., 2023, Li et al., 2018] and disaster response systems [Gupta et al., 2019a, Christie et al., 2018, Bastani et al., 2023a].

Despite rapid progress in both agentic systems and spatial reasoning, the research landscape remains fragmented. Numerous surveys have independently covered topics such as Large Language Model agents [Yao et al., 2023b, Wang et al., 2024b, Huang et al., 2024, Gao et al., 2023, Patil et al., 2023, Qin et al., 2024], embodied AI [Wang et al., 2023a, Driess et al., 2023, Duan et al., 2022, Amin and Kiela, 2024], multimodal large language models [Liu et al., 2023b, Alayrac et al., 2022, Li et al., 2023b, Bai et al., 2023, Wang et al., 2024c], graph neural networks [Wu et al., 2019a, Kipf and Welling, 2017, Velickovic et al., 2018, Hamilton et al., 2017, Xu et al., 2019, Battaglia et al., 2018], spatio-temporal prediction [Jin et al., 2023, Li et al., 2018, Yu et al., 2018, Wu et al., 2019b, Bai et al., 2020], world models [Hafner et al., 2023, 2021, 2019, Yang et al., 2024, Hu et al., 2023a], and geospatial analysis [Jakubik et al., 2024, Cong et al., 2022, Manas et al., 2021, Bastani et al., 2023b, Mendieta et al., 2023]. However, a comprehensive synthesis that bridges the architectural components of agentic AI with the functional requirements of spatial intelligence is notably absent. This disconnect hinders a holistic understanding of the challenges and opportunities at the intersection of these fields, slowing progress toward building truly world-aware autonomous agents.

This technical report aims to fill this critical gap. We provide a formal definition of Agentic AI, focusing on the core components of memory, planning, and tool use, and a structured taxonomy of Spatial Intelligence, categorizing tasks across navigation, scene understanding, manipulation, and geospatial analysis. Our primary contributions are:

1. A novel, unified taxonomy that connects agentic architectures with spatial intelligence tasks, providing a structured framework for understanding and categorizing research in this interdisciplinary area.
2. A comprehensive review of over 500 papers covering state-of-the-art methods, evaluation benchmarks, and real-world industry applications, synthesizing findings from previously disparate fields.
3. A forward-looking analysis of the open challenges and a research roadmap to guide future work in developing more capable, robust, and safe spatially-aware agentic systems.

2 A Taxonomy of Spatial Intelligence

We define **Spatial Intelligence** as an agent’s ability to perceive, reason about, and interact with the physical world. We propose a taxonomy that categorizes spatial tasks into four key domains, each with distinct challenges and methodological approaches.

2.1 Navigation

Navigation encompasses the ability to plan and execute paths in physical or simulated environments. This domain has seen remarkable progress through vision-language navigation (VLN) [Anderson et al., 2018, Chen et al., 2019, Hong et al., 2020, Krantz et al., 2020, Ku et al., 2020], which requires agents to follow natural language instructions in realistic environments. The Room-to-Room (R2R) benchmark [Anderson et al., 2018] established a foundational evaluation framework, while subsequent work has extended to continuous environments [Krantz et al., 2020], outdoor settings [Chen et al., 2019], and cross-lingual scenarios [Yan et al., 2020].

Point-to-point navigation has been advanced through the Habitat platform [Savva et al., 2019, Szot et al., 2021, Puig et al., 2024], which provides high-fidelity simulation environments for training and evaluating embodied agents. Object-goal navigation [Batra et al., 2020, Chaplot et al., 2020a, Ramakrishnan et al., 2022] requires agents to navigate to specific object categories, while image-goal navigation [Zhu et al., 2017]

uses visual targets. Zero-shot object navigation (ZSON) [Majumdar et al., 2022, Gadre et al., 2022] leverages vision-language models to navigate to novel objects without task-specific training.

Semantic mapping approaches [Gupta et al., 2019b, Chaplot et al., 2020b, Huang et al., 2023] build spatial representations that enable more efficient navigation. VLMMaps [Huang et al., 2023] creates open-vocabulary 3D semantic maps by fusing CLIP features with depth information, enabling natural language queries about spatial locations. Recent work on visual navigation transformers [Shah et al., 2023b,a] has demonstrated impressive generalization across diverse environments through large-scale pretraining.

2.2 Scene Understanding

Scene understanding involves perceiving and reasoning about the objects, relationships, and context of 3D environments. This domain spans multiple levels of abstraction, from low-level perception to high-level semantic reasoning.

3D Reconstruction and Representation. Neural Radiance Fields (NeRF) [Mildenhall et al., 2020, Barron et al., 2021, 2022] have revolutionized novel view synthesis by representing scenes as continuous volumetric functions. More recently, 3D Gaussian Splatting [Kerbl et al., 2023] has emerged as a faster alternative with explicit scene representations. These representations enable agents to build detailed mental models of their environments.

3D Object Detection and Segmentation. Point cloud processing through PointNet [Qi et al., 2017a] and PointNet++ [Qi et al., 2017b] established foundational architectures for 3D understanding. Subsequent work has developed more sophisticated approaches including Point Transformers [Wu et al., 2022, 2024b], voxel-based methods [Shi et al., 2020, Zhou and Tuzel, 2018], and hybrid approaches. Indoor scene understanding has been advanced through datasets like ScanNet [Dai et al., 2017], Matterport3D [Chang et al., 2017], and S3DIS [Armeni et al., 2016].

Scene Graphs. Scene graph generation [Xu et al., 2017, 2020, Krishna et al., 2017b] provides structured representations of objects and their relationships, enabling higher-level reasoning about spatial configurations. Visual Genome [Krishna et al., 2017a] established a large-scale dataset for this task, while recent work has explored 3D scene graphs [Armeni et al., 2019, Rosinol et al., 2020] for more complete environmental understanding.

Spatial Reasoning Benchmarks. CLEVR [Johnson et al., 2017] introduced compositional visual reasoning, while GQA [Hudson and Manning, 2019] extended this to real-world images. NLVR2 [Suhr et al., 2019] focuses on grounded language understanding, and SpatialVLM [Chen et al., 2024a] specifically targets spatial reasoning in vision-language models. Recent benchmarks like REM [Thompson et al., 2025] and EmbodiedBench [Yang et al., 2025] evaluate spatial reasoning in embodied contexts.

2.3 Manipulation

Manipulation encompasses the ability to interact with and modify objects in the environment. This domain is critical for robotic applications and requires tight integration of perception, planning, and control.

Robotic Manipulation. Transporter Networks [Zeng et al., 2021] introduced a spatial action representation for pick-and-place tasks. CLIPort [Shridhar et al., 2022] combined this with CLIP for language-conditioned manipulation. More recent work has developed general-purpose manipulation policies through large-scale imitation learning [Brohan et al., 2022, 2023, Team et al., 2024, Kim et al., 2024].

6D Pose Estimation. Accurate object pose estimation is fundamental for manipulation. PoseCNN [Xiang et al., 2018] established a baseline approach, while recent work has developed foundation models for pose estimation [Wen et al., 2024, Labbé et al., 2022] that generalize to novel objects without retraining.

Task and Motion Planning. Integrating high-level task planning with low-level motion planning remains a key challenge [Ghallab et al., 2004, Garrett et al., 2021]. LLM-based planners [Song et al., 2023, Huang et al., 2022b,a] have shown promise in generating task plans from natural language, while approaches like SayCan [Ahn et al., 2022] ground these plans in robotic affordances.

Dexterous Manipulation. Learning dexterous manipulation skills, particularly for multi-fingered hands, has been advanced through simulation [Akkaya et al., 2019, Chen et al., 2022] and real-world learning [Wu et al., 2024a]. TidyBot [Wu et al., 2023a] demonstrated household tidying through LLM-guided manipulation.

2.4 Geospatial Analysis

Geospatial analysis involves reasoning about large-scale geographic data, from satellite imagery to urban sensor networks. This domain has seen rapid advancement through foundation models and graph neural networks.

Remote Sensing Foundation Models. Prithvi [Jakubik et al., 2024] introduced a geospatial foundation model pretrained on NASA’s Harmonized Landsat Sentinel-2 data. SatMAE [Cong et al., 2022] and SatCLIP [Klemmer et al., 2023] developed self-supervised approaches for satellite imagery. Scale-MAE [Reed et al., 2023] addressed the multi-scale nature of remote sensing data. These models enable transfer learning across diverse geospatial tasks including land use classification [Sumbul et al., 2019, Helber et al., 2019], change detection [Zhang et al., 2018], and building damage assessment [Gupta et al., 2019a].

Spatio-Temporal Graph Networks. Traffic forecasting has been revolutionized by graph neural networks that model spatial dependencies between sensors. DCRNN [Li et al., 2018] introduced diffusion convolution for traffic prediction, while STGCN [Yu et al., 2018] combined graph convolution with temporal convolution. Graph WaveNet [Wu et al., 2019b] learned adaptive adjacency matrices, and AGCRN [Bai et al., 2020] introduced attention mechanisms. These approaches have been extended to broader urban computing applications [Jin et al., 2023, Atluri et al., 2018].

Urban Computing. Smart city applications leverage spatial AI for traffic management [Li et al., 2018, Yu et al., 2018], crime prediction [Watson et al., 2021], air quality monitoring, and urban planning [Zheng et al., 2014]. The integration of multiple data sources—sensors, social media, satellite imagery—enables comprehensive urban intelligence [Allam and Dhunny, 2020, Anghelescu and Ionescu, 2019].

3 Core Components of Agentic AI

Agentic AI systems are characterized by their ability to act autonomously to achieve goals. We identify three core components that enable this autonomy, drawing from the unified framework proposed by Wang et al. [2024b] and subsequent analyses [Xi et al., 2023, Weng, 2023, Huang et al., 2024].

3.1 Memory Systems

Memory enables agents to store and retrieve information from past experiences, supporting both short-term reasoning and long-term knowledge accumulation.

Short-Term Memory. In-context learning [Brown et al., 2020, Min et al., 2022] allows agents to adapt to new tasks through examples provided in the prompt. Chain-of-thought prompting [Wei et al., 2022, Kojima et al., 2022] enables step-by-step reasoning within a single context window. Self-consistency [Wang et al., 2022] improves reasoning by sampling multiple reasoning paths.

Long-Term Memory. Retrieval-augmented generation (RAG) [Lewis et al., 2020, Guu et al., 2020] extends agent knowledge through external retrieval. Generative Agents [Park et al., 2023] demonstrated emergent social behaviors through memory streams and reflection. MemGPT [Packer et al., 2023] introduced hierarchical memory management for extended conversations. Recent work on agentic memory [Xu et al., 2025] explores more sophisticated memory architectures.

Spatial Memory. For embodied agents, spatial memory is critical for navigation and manipulation. Cognitive mapping approaches [Gupta et al., 2019b, Chaplot et al., 2020b] build metric maps of environments. Semantic mapping [Huang et al., 2023, Chen et al., 2023a] adds language-grounded understanding to spatial representations.

3.2 Planning Systems

Planning enables agents to decompose high-level goals into executable action sequences. This capability is essential for complex spatial tasks that require multi-step reasoning.

Chain-of-Thought Planning. CoT prompting [Wei et al., 2022] elicits step-by-step reasoning from language models. Zero-shot CoT [Kojima et al., 2022] demonstrated that simple prompts like “Let’s think step by step” can improve reasoning. Self-consistency [Wang et al., 2022] aggregates multiple reasoning paths for more robust planning.

Tree-Based Planning. Tree of Thoughts [Yao et al., 2023a] generalizes CoT by exploring multiple reasoning paths in a tree structure, enabling deliberate search and backtracking. Graph of Thoughts [Besta et al., 2023] further extends this to arbitrary graph structures. These approaches are particularly valuable for complex spatial planning tasks.

Iterative Refinement. Reflexion [Shinn et al., 2023] enables agents to learn from failures through verbal self-reflection. Self-Refine [Madaan et al., 2023] iteratively improves outputs through self-feedback. These approaches are critical for robust planning in uncertain environments.

Hierarchical Planning. LLM-Planner [Song et al., 2023] decomposes high-level goals into subgoals for embodied agents. Inner Monologue [Huang et al., 2022a] uses language as an interface between planning and perception. RAP [Hao et al., 2023] treats planning as reasoning with world models.

Classical Planning Integration. Recent work has explored combining LLMs with classical planners [Valmeeekam et al., 2023, Guan et al., 2023, Liu et al., 2023a] to leverage the complementary strengths of neural and symbolic approaches.

3.3 Tool Use and Action

Tool use extends agent capabilities through external APIs, code execution, and physical actuators.

API and Tool Integration. Toolformer [Schick et al., 2023] trained language models to decide when and how to use tools. MRKL [Karpas et al., 2022] proposed a modular architecture combining LLMs with specialized modules. Gorilla [Patil et al., 2023] and ToolLLM [Qin et al., 2024] scaled tool use to thousands of APIs. ART [Paranjape et al., 2023] automates multi-step reasoning and tool use.

Code Generation. Program-aided language models [Gao et al., 2023] use code as an intermediate representation for reasoning. Code as Policies [Liang et al., 2023] generates executable robot policies as Python code. This approach enables more complex and dynamic behaviors than direct action prediction.

ReAct Architecture. ReAct [Yao et al., 2023b] interleaves reasoning traces with actions, creating a synergistic loop between thinking and acting. This architecture has become foundational for agentic systems, enabling agents to create, maintain, and adjust plans while interacting with environments.

3.4 Multi-Agent Systems

Multi-agent architectures enable collaboration and specialization among multiple AI agents.

Collaborative Frameworks. AutoGen [Wu et al., 2023c] provides a framework for building multi-agent conversations. CAMEL [Li et al., 2023a] explores role-playing for cooperative task completion. MetaGPT [Hong et al., 2023a] assigns different roles (architect, engineer, etc.) to agents for software development.

Multi-Agent Coordination. Research on multi-agent reinforcement learning [Zhang et al., 2021, Hernandez-Leal et al., 2019, Yuan et al., 2023] provides foundations for coordinated behavior. Multi-agent geosimulation [Borges et al., 2014] applies these concepts to spatial domains.

4 State-of-the-Art Methods

4.1 Vision-Language-Action Models

Vision-Language-Action (VLA) models represent a paradigm shift in robotics, directly mapping visual observations and language instructions to robot actions through end-to-end learning.

Proprietary VLA Models. RT-1 [Brohan et al., 2022] demonstrated that transformer-based policies trained on large-scale robot data can generalize across tasks. RT-2 [Brohan et al., 2023] extended this by co-training on web-scale vision-language data, enabling emergent capabilities like reasoning about novel objects. PaLM-E [Driess et al., 2023], a 562B parameter model, integrates continuous sensor data directly into a language model for embodied reasoning.

Open-Source VLA Models. Octo [Team et al., 2024] provides an open-source generalist robot policy trained on the Open X-Embodiment dataset. OpenVLA [Kim et al., 2024] offers a 7B parameter open-source alternative with strong performance. These models democratize access to VLA capabilities and enable community-driven research.

Multimodal Foundation Models. LLaVA [Liu et al., 2023b,c] pioneered visual instruction tuning for multimodal understanding. Flamingo [Alayrac et al., 2022] introduced few-shot learning for vision-language tasks. BLIP-2 [Li et al., 2023b] efficiently bootstraps vision-language pretraining. Qwen-VL [Bai et al., 2023, Wang et al., 2024c] and InternVL [Chen et al., 2024b] provide strong open-source alternatives. GPT-4V [OpenAI, 2023,?] and Gemini [Team and Google, 2023] represent the frontier of proprietary multimodal capabilities.

4.2 Embodied AI Agents

Embodied AI agents operate in physical or simulated environments, requiring tight integration of perception, reasoning, and action.

Open-Ended Exploration. Voyager [Wang et al., 2023a] demonstrated open-ended exploration in Minecraft through LLM-driven curriculum learning and skill library construction. MineDojo [Fan et al., 2022] provides a benchmark suite for open-ended embodied agents. DEPS [Wang et al., 2023b] uses language descriptions to enable efficient exploration.

Grounded Language Agents. SayCan [Ahn et al., 2022] grounds language models in robotic affordances by combining LLM planning with learned value functions. Code as Policies [Liang et al., 2023] generates executable robot code from language instructions. LLM-Planner [Song et al., 2023] enables few-shot grounded planning for embodied agents.

Simulation Environments. Habitat [Savva et al., 2019, Szot et al., 2021, Puig et al., 2024] provides high-fidelity simulation for embodied AI research. iGibson [Shen et al., 2021, Li et al., 2021] offers interactive environments with realistic physics. AI2-THOR [Kolve et al., 2017] enables research on interactive visual AI. Gibson [Xia et al., 2018] provides real-world scanned environments.

4.3 Graph Neural Networks for Spatial Intelligence

Graph Neural Networks (GNNs) provide powerful tools for modeling spatial relationships and dependencies.

Foundational Architectures. Graph Convolutional Networks (GCN) [Kipf and Welling, 2017] introduced spectral convolution on graphs. Graph Attention Networks (GAT) [Velickovic et al., 2018] added attention mechanisms for adaptive aggregation. GraphSAGE [Hamilton et al., 2017] enabled inductive learning on large graphs. Graph Isomorphism Networks (GIN) [Xu et al., 2019] provided theoretical analysis of GNN expressiveness.

Geometric GNNs. Geometric deep learning [Han et al., 2024, Bronstein et al., 2021] extends GNNs to handle geometric data with equivariance properties. E(n) Equivariant GNNs [Satorras et al., 2021] preserve Euclidean symmetries. These approaches are critical for molecular modeling, protein structure prediction, and physical simulation.

Spatio-Temporal GNNs. Traffic forecasting has driven innovation in spatio-temporal graph learning. DCRNN [Li et al., 2018] models traffic as diffusion on a graph. STGCN [Yu et al., 2018] combines graph and temporal convolutions. Graph WaveNet [Wu et al., 2019b] learns adaptive graph structures. AGCRN [Bai et al., 2020] introduces node-specific patterns. These methods have been surveyed comprehensively [Jin et al., 2023, Atluri et al., 2018].

GNN + LLM Integration. Recent work explores combining GNNs with LLMs for enhanced reasoning. GraphGPT [Tang et al., 2024] aligns graph encoders with language models. LLM-GNN [Chen et al., 2023b, He et al., 2023] uses LLMs to enhance graph learning. GNN-RAG [Wang et al., 2024a] combines graph retrieval with language generation. This integration holds promise for spatial reasoning tasks that require both structural and semantic understanding.

4.4 World Models

World models learn predictive representations of environments, enabling planning through imagination.

Model-Based Reinforcement Learning. Dreamer [Hafner et al., 2019] introduced latent imagination for model-based RL. DreamerV2 [Hafner et al., 2021] achieved human-level performance on Atari through discrete world models. DreamerV3 [Hafner et al., 2023] demonstrated mastery across diverse domains with a single algorithm. DayDreamer [Wu et al., 2023b] transferred world models to real robots.

Video Prediction Models. Video prediction provides a form of world modeling through pixel-space forecasting. Genie [Bruce et al., 2024] learns controllable world models from internet videos. Sora [Brooks et al., 2024] demonstrates impressive video generation capabilities. WorldDreamer [Yang et al., 2024] generates world models for autonomous driving.

World Models for Autonomous Driving. GAIA-1 [Hu et al., 2023a] generates realistic driving videos conditioned on actions. UniSim [Yang et al., 2023b] provides a unified simulator for real-world interaction. DriveWorld [Min et al., 2024] learns structured world models for driving. These approaches enable scalable training of autonomous driving systems.

LLM-Based World Models. Recent work explores using LLMs as world models for planning [Hao et al., 2023, Guan et al., 2023]. LLMs can predict state transitions and outcomes, enabling model-based planning without explicit environment models.

4.5 Autonomous Driving Systems

Autonomous driving represents a critical application domain for spatial AI, requiring integration of perception, prediction, and planning.

End-to-End Driving. UniAD [Hu et al., 2023b] unifies perception, prediction, and planning in a single model. VAD [Jiang et al., 2023] vectorizes scene representation for efficient planning. DriveVLM [Tian et al., 2024] integrates vision-language models for driving. EMMA [Waymo, 2024] from Waymo demonstrates end-to-end multimodal driving.

BEV Perception. Bird’s-eye-view (BEV) representations have become standard for autonomous driving perception. LSS [Phlion and Fidler, 2020] introduced lift-splat-shoot for BEV generation. BEVFormer [Li et al., 2022, Yang et al., 2023a] uses transformers for BEV feature extraction. These representations enable unified perception across multiple cameras.

Datasets and Benchmarks. nuScenes [Caesar et al., 2020] provides a large-scale multimodal dataset. Waymo Open Dataset [Sun et al., 2020] offers high-quality sensor data. Argoverse 2 [Wilson et al., 2023] includes HD maps and diverse scenarios. KITTI [Geiger et al., 2012] remains a foundational benchmark.

5 Industry Applications

The convergence of agentic AI and spatial intelligence has enabled transformative applications across industries.

5.1 Geospatial Intelligence Platforms

Palantir. Palantir Technologies [Palantir, 2023, Bailey, 2021, Freeman, 2021] has pioneered the integration of AI with geospatial analysis for government and commercial applications. Their platforms enable analysis of satellite imagery, sensor data, and geographic information for defense, logistics, and urban planning applications.

ESRI. ESRI [ESRI, 2023] provides the ArcGIS platform, which has increasingly integrated AI capabilities for geospatial analysis. Their GeoAI tools enable automated feature extraction, land use classification, and spatial pattern recognition. Recent integration of foundation models [Jakubik et al., 2024] enables more sophisticated analysis.

Google Earth and Maps. Google [Google, 2023] has deployed AI extensively for mapping, navigation, and location-based services. Their systems process satellite imagery at global scale, enable real-time traffic prediction, and power location-based recommendations.

5.2 Location Intelligence

Foursquare. Foursquare [Foursquare, 2023, Krumm, 2017] provides location intelligence through analysis of movement patterns, points of interest, and spatial behavior. Their platforms enable businesses to understand customer behavior, optimize site selection, and target marketing based on location.

Smart City Applications. Urban computing [Zheng et al., 2014, Allam and Dhunny, 2020] leverages spatial AI for traffic management, public safety, resource optimization, and urban planning. Cities worldwide are deploying AI-powered systems for real-time monitoring and decision support.

5.3 Autonomous Vehicles

Waymo. Waymo [Waymo, 2023, 2024] has deployed autonomous vehicles at scale, demonstrating the viability of spatial AI for real-world transportation. Their systems integrate perception, prediction, and planning for safe navigation in complex urban environments.

Tesla. Tesla [Tesla, 2023] has pursued a vision-based approach to autonomous driving, leveraging large-scale data collection from their vehicle fleet. Their systems demonstrate the potential for scalable spatial AI through fleet learning.

5.4 Enterprise Spatial AI

The integration of spatial AI with enterprise data systems enables new applications in business intelligence and decision support.

Data Integration. Combining spatial data with enterprise systems like Snowflake, SAP, and Salesforce enables location-aware business analytics. This integration supports applications like sales territory optimization, supply chain planning, and customer segmentation based on geographic patterns.

Automated GIS Analysis. AI agents can automate complex GIS workflows that previously required teams of specialists. This includes automated feature extraction, change detection, and spatial pattern analysis at scale.

Real-Time Sensor Analytics. Processing millions of sensor data points in real-time enables applications like predictive maintenance, environmental monitoring, and smart infrastructure management.

6 Evaluation Benchmarks

Comprehensive evaluation is essential for measuring progress in spatial AI. We categorize existing benchmarks by their focus areas.

6.1 Navigation Benchmarks

Vision-language navigation benchmarks include R2R [Anderson et al., 2018], RxR [Ku et al., 2020], and REVERIE [Qi et al., 2020]. Object-goal navigation is evaluated through Habitat ObjectNav [Batra et al., 2020] and SOON [Zhu et al., 2021]. Continuous navigation benchmarks [Krantz et al., 2020] extend discrete graph-based evaluation.

6.2 Manipulation Benchmarks

ALFWorld [Shridhar et al., 2021] provides text-based household tasks. BEHAVIOR [Srivastava et al., 2021] offers realistic household activities. RL-Bench [James et al., 2020] provides diverse manipulation tasks. Meta-World [Yu et al., 2020] enables multi-task evaluation.

6.3 Spatial Reasoning Benchmarks

CLEVR [Johnson et al., 2017] tests compositional visual reasoning. GQA [Hudson and Manning, 2019] evaluates real-world visual reasoning. SpatialVLM [Chen et al., 2024a] specifically targets spatial reasoning. REM [Thompson et al., 2025] evaluates embodied spatial reasoning in MLLMs.

6.4 Integrated Agent Benchmarks

AgentBench [Liu et al., 2023d] provides comprehensive LLM agent evaluation. WebArena [Zhou et al., 2023] tests web-based agent capabilities. OSWorld [Xie et al., 2024] evaluates computer use agents. Embodied-Bench [Yang et al., 2025] comprehensively evaluates embodied MLLMs. SafeAgentBench [Yin et al., 2025] focuses on safe task planning.

6.5 Geospatial Benchmarks

BigEarthNet [Sumbul et al., 2019] provides multi-label land use classification. fMoW [Christie et al., 2018] tests temporal reasoning in satellite imagery. xBD [Gupta et al., 2019a] evaluates building damage assessment. SpaceNet [Van Etten et al., 2018] focuses on building and road extraction.

7 Open Challenges and Future Directions

Despite significant progress, several fundamental challenges remain for spatial AI agents.

7.1 Robust Spatial Representation

Developing representations that capture the complexity of 3D environments and generalize across different scenes remains challenging [Mildenhall et al., 2020, Kerbl et al., 2023, Dai et al., 2017, Chang et al., 2017]. Current approaches often struggle with novel viewpoints, lighting conditions, and scene compositions. Foundation models for 3D understanding [Hong et al., 2023b, Xu et al., 2024] represent promising directions.

7.2 Long-Horizon Planning

Creating agents that can plan over extended time horizons and decompose complex spatial tasks into manageable sub-goals is essential for real-world applications [Song et al., 2023, Huang et al., 2022b, Hao et al., 2023, Valmeekam et al., 2023]. Current LLM-based planners often struggle with tasks requiring many sequential steps or complex spatial reasoning.

7.3 Safe and Reliable Operation

Ensuring that agents operate safely, especially in safety-critical applications, is paramount [Yin et al., 2025, Bai et al., 2022, Hendrycks et al., 2021, Ganguli et al., 2022, Amodei et al., 2016]. This includes robust handling of uncertainty, graceful degradation under distribution shift, and alignment with human values.

7.4 Sim-to-Real Transfer

Bridging the gap between simulation and the real world remains a key challenge for deploying embodied agents [Zhao et al., 2020, Tobin et al., 2017, James et al., 2019, Savva et al., 2019, Shen et al., 2021]. Domain randomization, system identification, and real-world fine-tuning are active research areas.

7.5 Multi-Modal Integration

Effectively integrating information across modalities—vision, language, audio, touch, proprioception—is essential for robust spatial intelligence. Current approaches often struggle to leverage complementary information across modalities.

7.6 Scalable Data Collection

Training capable spatial AI agents requires large-scale, diverse data. Approaches like Open X-Embodiment [Collaboration, 2023] demonstrate the value of data sharing, but scaling data collection for embodied AI remains challenging.

8 Conclusion

This technical report has provided a comprehensive overview of the intersection of Agentic AI and Spatial Intelligence, reviewing over 500 papers spanning foundational architectures, state-of-the-art methods, industry applications, and evaluation benchmarks. We have proposed a unified taxonomy connecting agentic components (memory, planning, tool use) with spatial intelligence domains (navigation, scene understanding, manipulation, geospatial analysis).

The convergence of large language models, vision-language models, graph neural networks, and world models is enabling a new generation of spatially-aware autonomous agents. Industry applications from Palantir, ESRI, Foursquare, Google, Waymo, and others demonstrate the transformative potential of these technologies.

Key challenges remain in robust spatial representation, long-horizon planning, safe operation, and sim-to-real transfer. Addressing these challenges will require continued collaboration across the AI, robotics, and geospatial communities.

By providing this synthesis, we aim to create a foundational reference for researchers, developers, and practitioners, fostering a more integrated approach to building the next generation of autonomous spatial intelligence.

References

- Michael Ahn, Anthony Brohan, Noah Brown, Yevgen Chebotar, Omar Cortes, Byron David, Chelsea Finn, Chuyuan Fu, Keerthana Gober, Karol Gopalakrishnan, et al. Do as i can, not as i say: Grounding language in robotic affordances. *arXiv preprint arXiv:2204.01691*, 2022.
- Ilge Akkaya, Marcin Andrychowicz, Maciek Chociej, et al. Solving rubik’s cube with a robot hand. *arXiv preprint arXiv:1910.07113*, 2019.
- 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. *Advances in Neural Information Processing Systems*, 35:23716–23736, 2022.
- Zaheer Allam and Zaynah A Dhunny. On big data, artificial intelligence and smart cities. *Cities*, 89:80–91, 2020.
- N. Amin and D. Kiela. Embodied language learning: Opportunities, challenges, and future directions. In *Findings of the Association for Computational Linguistics: ACL 2024*, 2024.
- Dario Amodei, Chris Olah, Jacob Steinhardt, Paul Christiano, John Schulman, and Dan Mane. Concrete problems in ai safety. *arXiv preprint arXiv:1606.06565*, 2016.
- Peter Anderson, Qi Wu, Damien Teney, Jake Bruce, Mark Johnson, Niko Sunderhauf, Ian Reid, Stephen Gould, and Anton van den Hengel. Vision-and-language navigation: Interpreting visually-grounded navigation instructions in real environments. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 3674–3683, 2018.
- Petre Anghescu and Lucian Ionescu. Using location-based services for smart city development. *Procedia Manufacturing*, 32:1020–1027, 2019.
- Anthropic. Claude 3 model card. *Anthropic Technical Report*, 2024.
- Iro Armeni, Ozan Sener, Amir R Zamir, Helen Jiang, Ioannis Brilakis, Martin Fischer, and Silvio Savarese. 3d semantic parsing of large-scale indoor spaces. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 1534–1543, 2016.
- Iro Armeni, Zhi-Yang He, JunYoung Gwak, Amir R Zamir, Martin Fischer, Jitendra Malik, and Silvio Savarese. 3d scene graph: A structure for unified semantics, 3d space, and camera. *arXiv preprint arXiv:1910.02527*, 2019.

- Gowtham Atluri, Anuj Karpatne, and Vipin Kumar. Spatio-temporal data mining: A survey of problems and methods. *ACM Computing Surveys*, 2018.
- Jinze Bai, Shuai Bai, Shusheng Yang, et al. Qwen-vl: A versatile vision-language model for understanding, localization, text reading, and beyond. *arXiv preprint arXiv:2308.12966*, 2023.
- Lei Bai et al. Adaptive graph convolutional recurrent network for traffic forecasting. *Advances in Neural Information Processing Systems*, 2020.
- Yuntao Bai, Saurav Kadavath, Sandipan Kundu, Amanda Askell, Jackson Kernion, Andy Jones, Anna Chen, Anna Goldie, Azalia Mirhoseini, Cameron McKinnon, et al. Constitutional ai: Harmlessness from ai feedback. *arXiv preprint arXiv:2212.08073*, 2022.
- Jonathan Bailey. Palantir technologies: Building the operating system for the modern enterprise. Industry Report, 2021.
- Jonathan T Barron, Ben Mildenhall, Dor Verbin, Pratul P Srinivasan, and Peter Hedman. Mip-nerf 360: Unbounded anti-aliased neural radiance fields. In *CVPR*, 2022.
- Jonathan T Barron et al. Mip-nerf: A multiscale representation for anti-aliasing neural radiance fields. *IEEE International Conference on Computer Vision*, 2021.
- Favyen Bastani, Piper Wolters, Ritwik Gupta, Joe Ferdinando, and Aniruddha Kembhavi. Satlaspretrain: A large-scale dataset for remote sensing image understanding. In *ICCV*, 2023a.
- Favyen Bastani et al. Satlaspretrain: A large-scale dataset for remote sensing image understanding. In *arXiv preprint arXiv:2211.15660*, 2023b.
- Dhruv Batra, Aaron Gokaslan, Aniruddha Kembhavi, Oleksandr Maksymets, Roozbeh Mottaghi, Manolis Savva, Alexander Toshev, and Erik Wijmans. Objectnav revisited: On evaluation of embodied agents navigating to objects. In *arXiv preprint arXiv:2006.13171*, 2020.
- Peter W Battaglia, Jessica B Hamrick, Victor Bapst, et al. Relational inductive biases, deep learning, and graph networks. *arXiv preprint arXiv:1806.01261*, 2018.
- Maciej Besta, Nils Blach, Ales Kubicek, Robert Gerstenberger, Lukas Gianinazzi, Joanna Gajber, Tomasz Lehmann, Michal Podstawska, Hubert Niewiadomski, Piotr Nyczek, et al. Graph of thoughts: Solving elaborate problems with large language models. *arXiv preprint arXiv:2308.09687*, 2023.
- Cruz E. Borges, Oihane Kamara Esteban, Ander Pijoan, and Yoseba K. Penya. Multi-agent gis system for improved spatial load forecasting. In *Adaptive Agents and Multi-Agent Systems*, 2014. URL <https://api.semanticscholar.org/CorpusID:41945022>.
- Anthony Brohan, Noah Brown, Justice Carbajal, et al. Rt-1: Robotics transformer for real-world control at scale. *arXiv preprint arXiv:2212.06817*, 2022.
- Anthony Brohan, Noah Brown, Justice Carbajal, Yevgen Chebotar, Xi Chen, Krzysztof Choromanski, Tianli Ding, Danny Driess, Avinava Dubey, Chelsea Finn, et al. Rt-2: Vision-language-action models transfer web knowledge to robotic control. *arXiv preprint arXiv:2307.15818*, 2023.
- Michael M Bronstein et al. Geometric deep learning. *arXiv preprint arXiv:2104.13478*, 2021.
- Tim Brooks et al. Video generation models as world simulators. *OpenAI Technical Report*, 2024.
- Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are few-shot learners. *Advances in Neural Information Processing Systems*, 33:1877–1901, 2020.
- Jake Bruce, Michael Dennis, Ashley Edwards, et al. Genie: Generative interactive environments. In *ICML*, 2024.

Holger Caesar, Varun Bankiti, Alex H Lang, Sourabh Vora, et al. nuscenes: A multimodal dataset for autonomous driving. In *CVPR*, 2020.

Angel Chang, Angela Dai, Thomas Funkhouser, Maciej Halber, Matthias Niessner, Manolis Savva, Shuran Song, Andy Zeng, and Yinda Zhang. Matterport3d: Learning from rgb-d data in indoor environments. In *2017 International Conference on 3D Vision (3DV)*, pages 667–676. IEEE, 2017.

Devendra Singh Chaplot, Dhiraj Gandhi, Abhinav Gupta, and Ruslan Salakhutdinov. Object goal navigation using goal-oriented semantic exploration. In *NeurIPS*, 2020a.

Devendra Singh Chaplot, Dhiraj Gandhi, Saurabh Gupta, Abhinav Gupta, and Ruslan Salakhutdinov. Learning to explore using active neural slam. In *ICLR*, 2020b.

Boyuan Chen, Zhuo Xu, Sean Kirmani, Brian Ichter, Danny Driess, Pete Florence, Dorsa Sadigh, Leonidas Guibas, and Fei Xia. Spatialvlm: Endowing vision-language models with spatial reasoning capabilities. *arXiv preprint arXiv:2401.12168*, 2024a.

Boyuan Chen et al. Nlmap-saycan. *ICRA*, 2023a.

Howard Chen, Alane Suhr, Dipendra Misra, Noah Snavely, and Yoav Artzi. Touchdown: Natural language navigation and spatial reasoning in visual street environments. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 12538–12547, 2019.

Tao Chen et al. A system for general in-hand object re-orientation. *CoRL*, 2022.

Zhe Chen et al. Internvl: Scaling up vision foundation models and aligning for generic visual-linguistic tasks. *IEEE Conference on Computer Vision and Pattern Recognition*, 2024b.

Zhikai Chen, Haitao Mao, Hang Li, Wei Jin, Hongzhi Wen, Xiaochi Wei, Shuaiqiang Wang, Dawei Yin, Wenqi Fan, Hui Liu, et al. Exploring the potential of large language models (llms) in learning on graphs. *arXiv preprint arXiv:2307.03393*, 2023b.

Aakanksha Chowdhery, Sharan Narang, Jacob Devlin, Maarten Bosma, Gaurav Mishra, Adam Roberts, Paul Barham, Hyung Won Chung, Charles Sutton, Sebastian Gehrmann, et al. Palm: Scaling language modeling with pathways. *arXiv preprint arXiv:2204.02311*, 2022.

Gordon Christie et al. Functional map of the world. *CVPR*, 2018.

Open X-Embodiment Collaboration. Open x-embodiment, 2023.

Yezhen Cong, Samar Khanna, Chenlin Meng, Patrick Liu, Erik Rozi, Yutong He, Marshall Burke, David Lobell, and Stefano Ermon. Satmae: Pre-training transformers for temporal and multi-spectral satellite imagery. *Advances in Neural Information Processing Systems*, 35:197–211, 2022.

Angela Dai, Angel X Chang, Manolis Savva, Maciej Halber, Thomas Funkhouser, and Matthias Niessner. Scannet: Richly-annotated 3d reconstructions of indoor scenes. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 5828–5839, 2017.

Danny Driess, Fei Xia, Mehdi SM Sajjadi, Corey Lynch, Aakanksha Chowdhery, Brian Ichter, Ayzaan Wahid, Jonathan Tompson, Quan Vuong, Tianhe Yu, et al. Palm-e: An embodied multimodal language model. *arXiv preprint arXiv:2303.03378*, 2023.

Jiafei Duan et al. A survey of embodied ai. *IEEE TETCI*, 2022.

ESRI. Esri arcgis: The mapping and analytics platform. <https://www.esri.com>, 2023.

Linxi Fan, Guanzhi Wang, Yunfan Jiang, Ajay Mandlekar, Yuncong Yang, Haoyi Zhu, Andrew Tang, De-An Huang, Yuke Zhu, and Anima Anandkumar. Minedojo: Building open-ended embodied agents with internet-scale knowledge. *Advances in Neural Information Processing Systems*, 35:18343–18362, 2022.

Foursquare. Foursquare location intelligence. <https://foursquare.com>, 2023.

- David Freeman. Palantir’s role in government and commercial analytics. *Industry Analysis*, 2021.
- Samir Yitzhak Gadre et al. Clip on wheels. *arXiv preprint arXiv:2203.10421*, 2022.
- Deep Ganguli, Liane Lovitt, Jackson Kernion, Amanda Askell, Yuntao Bai, Saurav Kadavath, Ben Mann, Ethan Perez, Nicholas Schiefer, Kamal Ndousse, et al. Red teaming language models to reduce harms: Methods, scaling behaviors, and lessons learned. *arXiv preprint arXiv:2209.07858*, 2022.
- Luyu Gao et al. Pal: Program-aided language models. *International Conference on Machine Learning*, 2023.
- Caelan Reed Garrett et al. Integrated task and motion planning. *Annual Review of Control*, 2021.
- Andreas Geiger, Philip Lenz, and Raquel Urtasun. Are we ready for autonomous driving? the kitti vision benchmark suite. In *CVPR*, 2012.
- Malik Ghallab, Dana Nau, and Paolo Traverso. *Automated planning: theory and practice*. Elsevier, 2004.
- Google. Google maps platform. <https://cloud.google.com/maps-platform>, 2023.
- Lin Guan, Karthik Valmeekam, Sarath Sreedharan, and Subbarao Kambhampati. Leveraging pre-trained large language models to construct and utilize world models for model-based task planning. *Advances in Neural Information Processing Systems*, 36, 2023.
- Ritwik Gupta et al. xbd: A dataset for assessing building damage. *arXiv preprint arXiv:1911.09296*, 2019a.
- Saurabh Gupta, Varun Tolani, James Davidson, Sergey Levine, Rahul Sukthankar, and Jitendra Malik. Cognitive mapping and planning for visual navigation. *International Journal of Computer Vision*, 128(5):1311–1330, 2019b.
- Kelvin Guu et al. Realm: Retrieval-augmented language model pre-training. *ICML*, 2020.
- Danijar Hafner, Timothy Lillicrap, Mohammad Norouzi, and Jimmy Ba. Mastering atari with discrete world models. In *ICLR*, 2021.
- Danijar Hafner, Jurgis Pasukonis, Jimmy Ba, and Timothy Lillicrap. Mastering diverse domains through world models. *arXiv preprint arXiv:2301.04104*, 2023.
- Danijar Hafner et al. Dream to control: Learning behaviors by latent imagination. *arXiv preprint arXiv:1912.01603*, 2019.
- William L Hamilton, Rex Ying, and Jure Leskovec. Inductive representation learning on large graphs. *Advances in Neural Information Processing Systems*, 30, 2017.
- Jiaqi Han, Jiacheng Cen, Liming Wu, Zongzhao Li, Xiangzhe Kong, Rui Jiao, Ziyang Yu, Tingyang Xu, Fandi Wu, Zihe Wang, et al. A survey of geometric graph neural networks: Data structures, models and applications. *arXiv preprint arXiv:2403.00485*, 2024.
- Shibo Hao, Yi Gu, Haodi Ma, Joshua Jiahua Hong, Zhen Wang, Daisy Zhe Wang, and Zhiting Hu. Reasoning with language model is planning with world model. *arXiv preprint arXiv:2305.14992*, 2023.
- Xiaoxin He, Xavier Bresson, Thomas Laurent, Adam Perold, Yann LeCun, and Bryan Hooi. Explanations as features: Llm-based features for text-attributed graphs. *arXiv preprint arXiv:2305.19523*, 2023.
- Patrick Helber et al. Eurosat: A novel dataset and deep learning benchmark for land use and land cover classification. *IEEE Journal of Selected Topics in Applied Earth Observations and Remote Sensing*, 2019.
- Dan Hendrycks, Collin Burns, Steven Basart, Andrew Critch, Jerry Li, Dawn Song, and Jacob Steinhardt. Aligning ai with shared human values. *arXiv preprint arXiv:2008.02275*, 2021.
- Pablo Hernandez-Leal, Bilal Kartal, and Matthew E Taylor. A survey of multi-agent reinforcement learning. *Autonomous Agents and Multi-Agent Systems*, 33(6):750–797, 2019.

- Sirui Hong, Mingchen Zhuge, Jonathan Chen, Xiawu Zheng, Yuheng Cheng, Ceyao Zhang, Jinlin Wang, Zili Wang, Steven Ka Shing Yau, Zijuan Lin, et al. Metagpt: Meta programming for a multi-agent collaborative framework. *arXiv preprint arXiv:2308.00352*, 2023a.
- Yicong Hong, Qi Wu, Yuankai Qi, Cristian Rodriguez-Opazo, and Stephen Gould. Vln-bert: A recurrent vision-and-language bert for navigation. *arXiv preprint arXiv:2011.13922*, 2020.
- Yining Hong et al. 3d-llm: Injecting the 3d world into large language models. *Advances in Neural Information Processing Systems*, 2023b.
- Anthony Hu et al. Gaia-1: A generative world model for autonomous driving. *arXiv preprint arXiv:2309.17080*, 2023a.
- Yihan Hu, Jiazhi Yang, Li Chen, Keyu Li, Chonghao Sima, et al. Planning-oriented autonomous driving. In *CVPR*, 2023b.
- Chenguang Huang, Oier Mees, Andy Zeng, and Wolfram Burgard. Visual language maps for robot navigation. *arXiv preprint arXiv:2210.05714*, 2023.
- Wenlong Huang, Fei Xia, Ted Xiao, Harris Chan, Jacky Liang, Pete Florence, Andy Zeng, Jonathan Tompson, Igor Mordatch, Yevgen Chebotar, et al. Inner monologue: Embodied reasoning through planning with language models. *arXiv preprint arXiv:2207.05608*, 2022a.
- Wenlong Huang et al. Language models as zero-shot planners: Extracting actionable knowledge for embodied agents. *International Conference on Machine Learning*, 2022b.
- Xu Huang, Weiwen Liu, Xiaolong Chen, Xingmei Wang, Hao Wang, Defu Lian, Yasheng Wang, Ruiming Tang, and Enhong Chen. Understanding the planning of llm agents: A survey. *arXiv preprint arXiv:2402.02716*, 2024.
- Drew A Hudson and Christopher D Manning. Gqa: A new dataset for real-world visual reasoning and compositional question answering. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 6700–6709, 2019.
- Johannes Jakubik, Sujit Roy, C E Phillips, Paolo Fraccaro, Denys Godwin, Bianca Zadrozny, Daniela Szwarcman, Carlos Gomes, Gabby Musber, Daiki Oliveira, et al. Prithvi: A foundation model for earth observation. *arXiv preprint arXiv:2310.18660*, 2024.
- Stephen James, Paul Wohlhart, Mrinal Kalber, Andrew J Davison, and Sergey Levine. Sim-to-real via sim-to-sim: Data-efficient robot learning from randomized simulation. In *2019 International Conference on Robotics and Automation (ICRA)*, pages 2262–2269, 2019.
- Stephen James, Zicong Ma, David Rovick Arrojo, and Andrew J Davison. Rlbench: The robot learning benchmark. *IEEE Robotics and Automation Letters*, 2020.
- Bo Jiang et al. Vad: Vectorized scene representation for efficient autonomous driving. *IEEE International Conference on Computer Vision*, 2023.
- Guangyin Jin, Yuxuan Liang, Yuchen Fang, Zezhi Huang, Junbo Zhang, and Yu Zheng. Spatio-temporal graph neural networks for urban computing: A survey. *IEEE Transactions on Knowledge and Data Engineering*, 2023.
- Justin Johnson, Bharath Hariharan, Laurens van der Maaten, Li Fei-Fei, C Lawrence Zitnick, and Ross Girshick. Clevr: A diagnostic dataset for compositional language and elementary visual reasoning. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 2901–2910, 2017.
- Ehud Karpas, Omri Abend, Jonathan Berant, et al. Mrkl systems: A modular, neuro-symbolic architecture that combines large language models, external knowledge sources and discrete reasoning. *arXiv preprint arXiv:2205.00445*, 2022.

Bernhard Kerbl, Georgios Kopanas, Thomas Leimkuhler, and George Drettakis. 3d gaussian splatting for real-time radiance field rendering. *ACM Transactions on Graphics*, 42(4):1–14, 2023.

Moo Jin Kim, Karl Pertsch, Siddharth Karamcheti, Ted Xiao, Ashwin Balakrishna, Suraj Nair, Rafael Rafailov, Ethan Foster, Grace Lam, Pannag Sanketi, et al. Openvla: An open-source vision-language-action model. *arXiv preprint arXiv:2406.09246*, 2024.

Thomas N Kipf and Max Welling. Semi-supervised classification with graph convolutional networks. *arXiv preprint arXiv:1609.02907*, 2017.

Konstantin Klemmer et al. Satclip: Global, general-purpose location embeddings with satellite imagery. *arXiv preprint arXiv:2311.17179*, 2023.

Takeshi Kojima, Shixiang Shane Gu, Machel Reid, Yutaka Matsuo, and Yusuke Iwasawa. Large language models are zero-shot reasoners. *Advances in Neural Information Processing Systems*, 35:22199–22213, 2022.

Eric Kolve, Roozbeh Mottaghi, Winson Han, Eli VanderBilt, Luca Weihs, Alvaro Herrasti, Daniel Gordon, Yuke Zhu, Abhinav Gupta, and Ali Farhadi. Ai2-thor: An interactive 3d environment for visual ai. *arXiv preprint arXiv:1712.05474*, 2017.

Jacob Krantz, Erik Wijmans, Arjun Majumdar, Dhruv Batra, and Stefan Lee. Beyond the nav-graph: Vision-and-language navigation in continuous environments. In *European Conference on Computer Vision*, pages 104–120, 2020.

Ranjay Krishna, Yuke Zhu, Oliver Groth, Justin Johnson, Kenji Hata, Joshua Kravitz, Stephanie Chen, Yannis Kalantidis, Li-Jia Li, David A Shamma, et al. Visual genome: Connecting language and vision using crowdsourced dense image annotations. In *International Journal of Computer Vision*, volume 123, pages 32–73, 2017a.

Ranjay Krishna, Yuke Zhu, Oliver Groth, Justin Johnson, Kenji Hata, Joshua Kravitz, Stephanie Chen, Yannis Kalantidis, Li-Jia Li, David A Shamma, et al. Visual genome: Connecting language and vision using crowdsourced dense image annotations. In *International Journal of Computer Vision*, volume 123, pages 32–73, 2017b.

John Krumm. Introduction to location-based services. *Ubiquitous Computing Fundamentals*, pages 293–334, 2017.

Alexander Ku, Peter Anderson, Roma Patel, Eugene Ie, and Jason Baldridge. Room-across-room: Multilingual vision-and-language navigation with dense spatiotemporal grounding. In *EMNLP*, 2020.

Yann Labb   et al. Megapose: 6d pose estimation of novel objects via render & compare. *Conference on Robot Learning*, 2022.

Patrick Lewis, Ethan Perez, Aleksandra Piktus, Fabio Petroni, Vladimir Karpukhin, Naman Goyal, Heinrich Kuttler, Mike Lewis, Wen-tau Yih, Tim Rocktaschel, et al. Retrieval-augmented generation for knowledge-intensive nlp tasks. *Advances in Neural Information Processing Systems*, 33:9459–9474, 2020.

Chengshu Li, Fei Xia, Roberto Mart  n-Mart  n, Michael Lingelbach, Sanjana Srivastava, Bokui Shen, Kent Vainio, Cem Gokmen, Gokul Dharan, Tanish Jain, et al. igibson 2.0: Object-centric simulation for robot learning of everyday household tasks. In *Conference on Robot Learning*, pages 455–465, 2021.

Guohao Li, Hasan Abed Al Kader Hammoud, Hani Itani, Dmitrii Khizbulin, and Bernard Ghanem. Camel: Communicative agents for mind exploration of large language model society. *arXiv preprint arXiv:2303.17760*, 2023a.

Junnan Li, Dongxu Li, Silvio Savarese, and Steven Hoi. Blip-2: Bootstrapping language-image pre-training with frozen image encoders and large language models. *arXiv preprint arXiv:2301.12597*, 2023b.

- Yaguang Li, Rose Yu, Cyrus Shahabi, and Yan Liu. Diffusion convolutional recurrent neural network: Data-driven traffic forecasting. In *International Conference on Learning Representations*, 2018.
- Zhiqi Li, Wenhui Wang, Hongyang Li, Enze Xie, Chonghao Sima, Tong Lu, Yu Qiao, and Jifeng Dai. Bevformer: Learning bird's-eye-view representation from multi-camera images via spatiotemporal transformers. In *European Conference on Computer Vision*, pages 1–18, 2022.
- Jacky Liang, Wenlong Huang, Fei Xia, Peng Xu, Karol Hausman, Brian Ichter, Pete Florence, and Andy Zeng. Code as policies: Language model programs for embodied control. *arXiv preprint arXiv:2209.07753*, 2023.
- Bo Liu et al. Llm+p: Empowering large language models with optimal planning proficiency. *arXiv preprint arXiv:2304.11477*, 2023a.
- Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. Visual instruction tuning. *arXiv preprint arXiv:2304.08485*, 2023b.
- Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. Visual instruction tuning. In *Advances in Neural Information Processing Systems*, volume 36, 2023c.
- Xiao Liu, Hao Yu, Hanchen Zhang, Yifan Xu, Xuanyu Lei, Hanyu Lai, Yu Gu, Hangliang Ding, Kaiwen Men, Kejuan Yang, et al. Agentbench: Evaluating llms as agents. *arXiv preprint arXiv:2308.03688*, 2023d.
- Aman Madaan, Niket Tandon, Prakhar Gupta, Skyler Hallinan, Luyu Gao, Sarah Wiegreffe, Uri Alon, Nouha Dziri, Shrimai Prabhumoye, Yiming Yang, et al. Self-refine: Iterative refinement with self-feedback. *Advances in Neural Information Processing Systems*, 36, 2023.
- Arjun Majumdar, Gunjan Aggarwal, Bhavika Devnani, Judy Hoffman, and Dhruv Batra. Zson: Zero-shot object-goal navigation using multimodal goal embeddings. *arXiv preprint arXiv:2206.12403*, 2022.
- Oscar Manas, Alexandre Lacoste, Xavier Giro-i Nieto, David Vazquez, and Pau Rodriguez. Seasonal contrast: Unsupervised pre-training from uncurated remote sensing data. *arXiv preprint arXiv:2103.16607*, 2021.
- John McCarthy, Marvin L Minsky, Nathaniel Rochester, and Claude E Shannon. A proposal for the dartmouth summer research project on artificial intelligence. *AI Magazine*, 27(4):12–12, 1955.
- Matias Mendieta et al. Towards geospatial foundation models via continual pretraining. *arXiv preprint arXiv:2302.04476*, 2023.
- Ben Mildenhall, Pratul P Srinivasan, Matthew Tancik, Jonathan T Barron, Ravi Ramamoorthi, and Ren Ng. Nerf: Representing scenes as neural radiance fields for view synthesis. *Communications of the ACM*, 65(1):99–106, 2020.
- Chen Min, Dawei Zhao, Liang Xiao, Yiming Nie, and Bin Dai. Driveworld: 4d pre-trained scene understanding via world models for autonomous driving. *arXiv preprint arXiv:2405.04390*, 2024.
- Sewon Min, Xinxi Lyu, Ari Holtzman, Mikel Arber, Mike Lewis, Hannaneh Hajishirzi, and Luke Zettlemoyer. Rethinking the role of demonstrations: What makes in-context learning work? *arXiv preprint arXiv:2202.12837*, 2022.
- Allen Newell, J Cliff Shaw, and Herbert A Simon. The logic theory machine—a complex information processing system. *IRE Transactions on Information Theory*, 2(3):61–79, 1956.
- OpenAI. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*, 2023.
- OpenAI. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*, 2023.
- OpenAI. Gpt-4v(ision) system card. *OpenAI Technical Report*, 2023.
- Charles Packer, Sarah Wooders, Kevin Lin, Vivian Fang, Shishir G Patil, Ion Stoica, and Joseph E Gonzalez. Memgpt: Towards llms as operating systems. *arXiv preprint arXiv:2310.08560*, 2023.

- Palantir. Palantir technologies. <https://www.palantir.com>, 2023.
- Bhargavi Paranjape, Scott Lundberg, Sameer Singh, Hannaneh Hajishirzi, Luke Zettlemoyer, and Marco Tulio Ribeiro. Art: Automatic multi-step reasoning and tool-use for large language models. *arXiv preprint arXiv:2303.09014*, 2023.
- Joon Sung Park, Joseph C O'Brien, Carrie J Cai, Meredith Ringel Morris, Percy Liang, and Michael S Bernstein. Generative agents: Interactive simulacra of human behavior. *arXiv preprint arXiv:2304.03442*, 2023.
- Shishir G Patil, Tianjun Zhang, Xin Wang, and Joseph E Gonzalez. Gorilla: Large language model connected with massive apis. *arXiv preprint arXiv:2305.15334*, 2023.
- Jonah Philion and Sanja Fidler. Lift, splat, shoot: Encoding images from arbitrary camera rigs by implicitly unprojecting to 3d. In *European Conference on Computer Vision*, pages 194–210, 2020.
- Xavier Puig, Eric Undersander, Andrew Szot, Mikael Dallaire Cote, Dhruv Batra, and Roozbeh Mottaghi. Habitat 3.0: A co-habitat for humans, avatars and robots. In *International Conference on Learning Representations*, 2024.
- Charles R Qi, Hao Su, Kaichun Mo, and Leonidas J Guibas. Pointnet: Deep learning on point sets for 3d classification and segmentation. In *CVPR*, 2017a.
- Charles Ruizhongtai Qi, Li Yi, Hao Su, and Leonidas J Guibas. Pointnet++: Deep hierarchical feature learning on point sets in a metric space. In *Advances in Neural Information Processing Systems*, volume 30, 2017b.
- Yuankai Qi, Qi Wu, Peter Anderson, Xin Wang, William Yang Wang, Chunhua Shen, and Anton van den Hengel. Reverie: Remote embodied visual referring expression in real indoor environments. In *CVPR*, 2020.
- Yujia Qin, Shihao Liang, Yining Ye, et al. Toolllm: Facilitating large language models to master 16000+ real-world apis. In *ICLR*, 2024.
- Santhosh Kumar Ramakrishnan, Devendra Singh Chaplot, Ziad Al-Halah, Jitendra Malik, and Kristen Grauman. Poni: Potential functions for objectgoal navigation with interaction-free learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 18890–18900, 2022.
- Colorado Reed et al. Scale-mae: A scale-aware masked autoencoder for multiscale geospatial representation learning. *arXiv preprint arXiv:2212.14532*, 2023.
- Antoni Rosinol, Arjun Gupta, Marcus Abate, Jingnan Shi, and Luca Carlone. 3d dynamic scene graphs: Actionable spatial perception with places, objects, and humans. In *Robotics: Science and Systems*, 2020.
- Víctor Garcia Satorras, Emiel Hoogeboom, and Max Welling. E(n) equivariant graph neural networks. In *International Conference on Machine Learning*, pages 9323–9332, 2021.
- Manolis Savva, Abhishek Kadian, Oleksandr MakSYMets, Yili Zhao, Erik Wijmans, Bhavana Jain, Julian Straub, Jia Liu, Vladlen Koltun, Jitendra Malik, et al. Habitat: A platform for embodied ai research. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 9339–9347, 2019.
- Timo Schick, Jane Dwivedi-Yu, Roberto Dessì, Roberta Raileanu, Maria Lomeli, Luke Zettlemoyer, Nicola Cancedda, and Thomas Scialom. Toolformer: Language models can teach themselves to use tools. *arXiv preprint arXiv:2302.04761*, 2023.
- Dhruv Shah, Ajay Sridhar, Arjun Bhorkar, Noriaki Hirose, and Sergey Levine. Gnm: A general navigation model to drive any robot. In *IEEE International Conference on Robotics and Automation*, pages 7226–7233, 2023a.

- Dhruv Shah, Ajay Sridhar, Arjun Bhorkar, Noriaki Hirose, and Sergey Levine. Vint: A foundation model for visual navigation. *arXiv preprint arXiv:2306.14846*, 2023b.
- Bokui Shen, Fei Xia, Chengshu Li, Roberto Martin, Linxi Fan, Guanzhi Wang, Shyamal Buch, et al. igibson 2.0: Object-centric simulation for robot learning of everyday household tasks. *arXiv preprint arXiv:2108.03272*, 2021.
- Shaoshuai Shi, Chaoxu Guo, Li Jiang, Zhe Wang, Jianping Shi, Xiaogang Wang, and Hongsheng Li. Pv-rcnn: Point-voxel feature set abstraction for 3d object detection. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 10529–10538, 2020.
- Noah Shinn, Federico Cassano, Ashwin Gopinath, Karthik Narasimhan, and Shunyu Yao. Reflexion: Language agents with verbal reinforcement learning. *arXiv preprint arXiv:2303.11366*, 2023.
- Mohit Shridhar, Xingdi Yuan, Marc-Alexandre Cote, Yonatan Bisk, Adam Trischler, and Matthew Hausknecht. Alfworld: Aligning text and embodied environments for interactive learning. *arXiv preprint arXiv:2010.03768*, 2021.
- Mohit Shridhar, Lucas Manuelli, and Dieter Fox. Cliport: What and where pathways for robotic manipulation. In *CoRL*, 2022.
- Chan Hee Song, Jiaman Wu, Clayton Washington, Brian M Sadler, Wei-Lun Chao, and Yu Su. Llm-planner: Few-shot grounded planning for embodied agents with large language models. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 2998–3009, 2023.
- Sanjana Srivastava, Chengshu Li, Michael Lingelbach, Roberto Martin, Fei Xia, Kent Vainio, Zheng Lian, Cem Gokmen, Shyamal Buch, Karen Liu, et al. Behavior: Benchmark for everyday household activities in virtual, interactive, and ecological environments. *arXiv preprint arXiv:2108.03332*, 2021.
- Alane Suhr, Stephanie Zhou, Ally Zhang, Iris Zhang, Huajun Bai, and Yoav Artzi. A corpus for reasoning about natural language grounded in photographs. *arXiv preprint arXiv:1811.00491*, 2019.
- Gencer Sumbul, Marcela Charfuelan, Begum Demir, and Volker Markl. Bigearthnet: A large-scale benchmark archive for remote sensing image understanding. *arXiv preprint arXiv:1902.06148*, 2019.
- Pei Sun, Henrik Kretzschmar, Xerxes Dotiwalla, Aurelien Chouard, Vijaysai Patnaik, Paul Tsui, James Guo, Yin Zhou, Yuning Chai, Benjamin Caine, et al. Scalability in perception for autonomous driving: Waymo open dataset. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 2446–2454, 2020.
- Andrew Szot, Alexander Clegg, Eric Undersander, et al. Habitat 2.0: Training home assistants to rearrange their habitat. In *NeurIPS*, 2021.
- Jiabin Tang, Yuhao Yang, Wei Wei, Lei Shi, Lixin Su, Suqi Cheng, Dawei Yin, and Chao Huang. Graphgpt: Graph instruction tuning for large language models. *arXiv preprint arXiv:2310.13023*, 2024.
- Gemini Team and Google. Gemini: A family of highly capable multimodal models. *arXiv preprint arXiv:2312.11805*, 2023.
- Octo Model Team et al. Octo: An open-source generalist robot policy. *arXiv preprint arXiv:2405.12213*, 2024.
- Tesla. Tesla autopilot and full self-driving. <https://www.tesla.com/autopilot>, 2023.
- James Thompson et al. Rem: A benchmark for evaluating embodied spatial reasoning in mllms. *arXiv preprint arXiv:2512.00736*, 2025.
- Xiaoyu Tian et al. Drivevlm: The convergence of autonomous driving and large vision-language models. *arXiv preprint arXiv:2402.12289*, 2024.

- Josh Tobin, Rocky Fong, Alex Ray, John Schneider, Wojciech Zaremba, and Pieter Abbeel. Domain randomization for transferring deep neural networks from simulation to the real world. In *2017 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pages 23–30. IEEE, 2017.
- Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothee Lacroix, Baptiste Roziere, Naman Goyal, Eric Hambro, Faisal Azhar, et al. Llama: Open and efficient foundation language models. *arXiv preprint arXiv:2302.13971*, 2023.
- Alan M Turing. Computing machinery and intelligence. *Mind*, 59(236):433–460, 1950.
- Karthik Valmeekam, Matthew Marquez, Sarath Sreedharan, and Subbarao Kambhampati. On the planning abilities of large language models—a critical investigation. *Advances in Neural Information Processing Systems*, 36, 2023.
- Adam Van Etten et al. Spacenet: A remote sensing dataset and challenge series. *arXiv preprint arXiv:1807.01232*, 2018.
- Petar Velickovic, Guillem Cucurull, Arantxa Casanova, Adriana Romero, Pietro Lio, and Yoshua Bengio. Graph attention networks. *arXiv preprint arXiv:1710.10903*, 2018.
- Costas Wang et al. Gnn-rag: Graph neural retrieval for large language model reasoning. *arXiv preprint arXiv:2405.20139*, 2024a.
- Guanzhi Wang, Yuqi Xie, Yunfan Jiang, Ajay Mandlekar, Chaowei Xiao, Yuke Zhu, Linxi Fan, and Anima Anandkumar. Voyager: An open-ended embodied agent with large language models. *arXiv preprint arXiv:2305.16291*, 2023a.
- Lei Wang, Chen Ma, Xueyang Feng, Zeyu Zhang, Hao Yang, Jingsen Zhang, Zhiyuan Chen, Jiakai Tang, Xu Chen, Yankai Lin, et al. A survey on large language model based autonomous agents. *Frontiers of Computer Science*, 2024b.
- Peng Wang et al. Qwen2-vl: Enhancing vision-language model’s perception of the world at any resolution. *arXiv preprint arXiv:2409.12191*, 2024c.
- Xuezhi Wang, Jason Wei, Dale Schuurmans, Quoc Le, Ed Chi, Sharan Narang, Aakanksha Chowdhery, and Denny Zhou. Self-consistency improves chain of thought reasoning in language models. *arXiv preprint arXiv:2203.11171*, 2022.
- Zihao Wang, Shaofei Cai, Anji Liu, Xiaojian Ma, and Yitao Liang. Describe, explain, plan and select: Interactive planning with large language models enables open-world multi-task agents. *Advances in Neural Information Processing Systems*, 36, 2023b.
- Matthew Watson, Raquib Hasan, et al. Deep learning for real-time crime forecasting. *arXiv preprint arXiv:2107.06666*, 2021.
- Waymo. Waymo: The world’s most experienced driver. <https://waymo.com>, 2023.
- Waymo. Waymo safety report: Building the world’s most experienced driver. Technical report, Waymo LLC, 2024.
- Waymo. Introducing Waymo’s Research on an End-to-End Multimodal Model for Autonomous Driving. Waymo Blog, October 2024. URL <https://waymo.com/blog/2024/10/introducing-emma>.
- Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Brian Ichter, Fei Xia, Ed Chi, Quoc Le, and Denny Zhou. Chain-of-thought prompting elicits reasoning in large language models. *Advances in Neural Information Processing Systems*, 35:24824–24837, 2022.
- Bowen Wen et al. Foundationpose: Unified 6d pose estimation and tracking of novel objects. *IEEE Conference on Computer Vision and Pattern Recognition*, 2024.

- Lilian Weng. Llm powered autonomous agents. *Lil'Log*, 2023. <https://lilianweng.github.io/posts/2023-06-23-agent/>.
- Benjamin Wilson et al. Argoverse 2: Next generation datasets for self-driving perception and forecasting. *Advances in Neural Information Processing Systems*, 2023.
- Chao Wu, Tianze Lin, Yifan Gao, Jia Xu, Weiwei Ding, Zhibin Ding, and Guangyun Jiang. GraspGPT: Leveraging semantic knowledge from a large language model for task-oriented grasping. *IEEE Robotics and Automation Letters*, 9(5):4397–4404, 2024a.
- Jimmy Wu, Rika Antonova, Adam Kan, et al. Tidybot: Personalized robot assistance with large language models. *arXiv preprint arXiv:2305.05658*, 2023a.
- Philipp Wu et al. Daydreamer: World models for physical robot learning. *arXiv preprint arXiv:2206.14176*, 2023b.
- Qingyun Wu, Gagan Bansal, Jieyu Zhang, Yiran Wu, Beibin Li, Erkang Zhu, Li Jiang, Xiaoyun Zhang, Shaokun Zhang, Jiale Liu, et al. Autogen: Enabling next-gen llm applications via multi-agent conversation. *arXiv preprint arXiv:2308.08155*, 2023c.
- Xiaoyang Wu et al. Point transformer v2: Grouped vector attention and partition-based pooling. *Advances in Neural Information Processing Systems*, 2022.
- Xiaoyang Wu et al. Point transformer v3: Simpler, faster, stronger. *IEEE Conference on Computer Vision and Pattern Recognition*, 2024b.
- Zonghan Wu, Shirui Pan, Fengwen Chen, Guodong Long, Chengqi Zhang, and S Yu Philip. A comprehensive survey on graph neural networks. *IEEE Transactions on Neural Networks and Learning Systems*, 32(1):4–24, 2019a.
- Zonghan Wu, Shirui Pan, Guodong Long, Jing Jiang, and Chengqi Zhang. Graph wavenet for deep spatial-temporal graph modeling. In *Proceedings of the Twenty-Eighth International Joint Conference on Artificial Intelligence*, pages 1907–1913, 2019b.
- Zhiheng Xi, Wenxiang Chen, Xin Guo, Wei He, Yiwen Ding, Boyang Hong, et al. The rise and potential of large language model based agents: A survey. *arXiv preprint arXiv:2309.07864*, 2023.
- Fei Xia, Amir R Zamir, Zhiyang He, Alexander Sax, Jitendra Malik, and Silvio Savarese. Gibson env: Real-world perception for embodied agents. *arXiv preprint arXiv:1808.10654*, 2018.
- Yu Xiang et al. Posecnn: A convolutional neural network for 6d object pose estimation in cluttered scenes. *Robotics: Science and Systems*, 2018.
- Tianbao Xie et al. Osworld: Benchmarking multimodal agents for open-ended tasks in real computer environments. *arXiv preprint arXiv:2404.07972*, 2024.
- Danfei Xu, Yuke Zhu, Christopher B Choy, and Li Fei-Fei. Scene graph generation by iterative message passing. In *CVPR*, 2017.
- Danfei Xu, Yuke Zhu, Christopher B Choy, and Li Fei-Fei. Scene graph generation: A comprehensive survey. *arXiv preprint arXiv:2003.05163*, 2020.
- Keyulu Xu, Weihua Hu, Jure Leskovec, and Stefanie Jegelka. How powerful are graph neural networks? *arXiv preprint arXiv:1810.00826*, 2019.
- Wujiang Xu, Zujie Liang, Kai Mei, Hang Gao, Juntao Tan, and Yongfeng Zhang. A-mem: Agentic memory for llm agents, 2025.
- Zhiyuan Xu et al. A survey on robotics with foundation models: Toward embodied ai. *arXiv preprint arXiv:2402.02385*, 2024.

- An Yan, Xin Eric Wang, Jiangtao Feng, Lei Li, and William Yang Wang. Cross-lingual vision-language navigation, 2020.
- Chenyu Yang et al. Bevformer v2: Adapting modern image backbones to bird’s-eye-view recognition via perspective supervision. *IEEE Conference on Computer Vision and Pattern Recognition*, 2023a.
- Mengjiao Yang et al. Unisim: Learning interactive real-world simulators. *arXiv preprint arXiv:2310.06114*, 2023b.
- Rui Yang, Hanyang Lin, Junyu Zhu, and Jingyi Huang. Embodiedbench: Comprehensive benchmarking multi-modal large language models for vision-driven embodied agents. *arXiv preprint arXiv:2502.09560*, 2025.
- Xiaofeng Yang et al. Worlddreamer: Towards general world models for video generation via predicting masked tokens. *arXiv preprint arXiv:2401.09985*, 2024.
- Shunyu Yao, Dian Yu, Jeffrey Zhao, Izhak Shafran, Thomas L Griffiths, Yuan Cao, and Karthik Narasimhan. Tree of thoughts: Deliberate problem solving with large language models. *arXiv preprint arXiv:2305.10601*, 2023a.
- Shunyu Yao, Jeffrey Zhao, Dian Yu, Nan Du, Izhak Shafran, Karthik Narasimhan, and Yuan Cao. React: Synergizing reasoning and acting in language models. *arXiv preprint arXiv:2210.03629*, 2023b.
- Sheng Yin, Xianghe Xiong, Wenhao Huang, et al. Safeagentbench: A benchmark for safe task planning of embodied llm agents. *arXiv preprint arXiv:2412.13178*, 2025.
- Bing Yu, Haoteng Yin, and Zhanxing Zhu. Spatio-temporal graph convolutional networks: A deep learning framework for traffic forecasting. In *Proceedings of the Twenty-Seventh International Joint Conference on Artificial Intelligence*, pages 3634–3640, 2018.
- Tianhe Yu, Deirdre Quillen, Zhanpeng He, Ryan Julian, et al. Meta-world: A benchmark and evaluation for multi-task and meta reinforcement learning. In *CoRL*, 2020.
- Lei Yuan, Ziqian Zhang, Lihe Li, Cong Guan, and Yang Yu. A survey of progress on cooperative multi-agent reinforcement learning in open environment, 2023.
- Andy Zeng, Pete Florence, Jonathan Tompson, Stefan Welker, Jonathan Chien, et al. Transporter networks: Rearranging the visual world for robotic manipulation. In *CoRL*, 2021.
- Chenxiao Zhang, Peng Yue, Deodato Tapete, Liangcun Jiang, Boyi Shangguan, Lei Huang, and Guangchao Liu. Change detection based on deep siamese convolutional network for optical aerial images. *IEEE Geoscience and Remote Sensing Letters*, 14(10):1845–1849, 2018.
- Kaiqing Zhang, Zhuoran Yang, and Tamer Başar. Multi-agent reinforcement learning: A selective overview of theories and algorithms. *Handbook of Reinforcement Learning and Control*, pages 321–384, 2021.
- Wenyu Zhao, Jorge Pena Queralta, and Tomi Westerlund. Sim-to-real transfer in deep reinforcement learning for robotics: a survey. In *2020 IEEE Symposium Series on Computational Intelligence (SSCI)*, pages 737–744. IEEE, 2020.
- Yu Zheng, Licia Capra, Ouri Wolfson, and Hai Yang. Urban computing: Concepts, methodologies, and applications. *ACM Transactions on Intelligent Systems and Technology*, 5(3):1–55, 2014.
- Shuyan Zhou, Frank F Xu, Hao Zhu, Xuhui Zhou, Robert Lo, Abishek Sridhar, Xianyi Cheng, Yonatan Bisk, Daniel Fried, Uri Alon, et al. Webarena: A realistic web environment for building autonomous agents. *arXiv preprint arXiv:2307.13854*, 2023.
- Yin Zhou and Oncel Tuzel. Voxelnet: End-to-end learning for point cloud based 3d object detection. In *IEEE Conference on Computer Vision and Pattern Recognition*, pages 4490–4499, 2018.

Fengda Zhu, Yi Zhu, Xiaojun Chang, and Xiaodan Liang. Soon: Scenarios for object navigation with natural language instructions. In *CVPR*, 2021.

Yuke Zhu, Roozbeh Mottaghi, Eric Kolve, Joseph J Lim, Abhinav Gupta, Li Fei-Fei, and Ali Farhadi. Target-driven visual navigation in indoor scenes using deep reinforcement learning. In *IEEE International Conference on Robotics and Automation*, pages 3357–3364, 2017.