# SpatialOps: A Benchmark for 2D Spatial Planning and Reasoning in Large Language Models

**Anonymous Author(s)**
Affiliation
email@example.com

## Abstract

Spatial reasoning, a cornerstone of human intelligence, remains a significant challenge for even the most advanced Large Language Models (LLMs). While existing benchmarks have explored various facets of spatial understanding, a comprehensive evaluation of practical, 2D spatial planning and reasoning across a wide range of real-world domains is still lacking. To address this gap, we introduce **SpatialOps**, a comprehensive benchmark designed to rigorously assess the 2D spatial planning and reasoning capabilities of LLMs. SpatialOps comprises **twelve distinct task categories** organized into three tiers of increasing complexity, encompassing over 6,000 procedurally generated tasks that require a deep understanding of coordinate systems, topology, visibility, algorithmic pathfinding, and constraint-based optimization. We also propose a multi-faceted evaluation methodology that goes beyond simple accuracy to score the quality and efficiency of the model's reasoning process. By providing a challenging, reproducible, and expanded benchmark with 100% ground-truth accuracy, SpatialOps aims to drive progress in developing more spatially-aware and capable AI systems. The benchmark, including all data and evaluation code, is publicly available at https://github.com/glo26/spatial-benchmark.

## 1 Introduction

The remarkable progress of Large Language Models (LLMs) has demonstrated their capacity for complex linguistic tasks Brown et al. [2020], Touvron et al. [2023], Anil et al. [2023]. However, their ability to reason about the physical world, particularly in the spatial domain, lags significantly behind their linguistic prowess Bisk et al. [2020], Zhang et al. [2024]. This gap is largely due to a fundamental representation mismatch: LLMs process information as discrete, sequential tokens, whereas the physical world is characterized by continuous geometric structures Harnad [1990], Sloman [1971]. Consequently, models often learn statistical co-occurrences of spatial terms rather than acquiring a true, grounded understanding of geometric principles Marcus [1998], Mitchell [2021].

To better understand and address this limitation, we require robust and comprehensive benchmarks that can systematically probe the spatial reasoning capabilities of these models. While several existing benchmarks have made valuable contributions Davis and Marcus [2017], Zellers et al. [2019], a significant gap remains in the evaluation of practical, applied 2D spatial planning. Current benchmarks often focus on a limited set of abstract scenarios Weston et al. [2015], Le et al. [2022], failing to capture the complexity and diversity of real-world spatial problems encountered in domains such as urban planning Balachandar et al. [2025], logistics Savelsbergh and Van Woensel [2005], and engineering Eastman et al. [2011].

To fill this critical gap, we introduce **SpatialOps**, a comprehensive benchmark for 2D spatial planning and reasoning in LLMs. SpatialOps is designed to be comprehensive, challenging, and grounded in real-world applications. It evaluates models across a wide spectrum of spatial tasks, from fundamental coordinate understanding to complex, multi-step optimization problems.

Our main contributions are:

1. **A new, comprehensive benchmark for 2D spatial planning**, encompassing twelve diverse task categories that cover a wide range of practical applications, from network infrastructure planning Jin et al. [2023] to real estate analysis De Sabbata and Liu [2023].

2. **A challenging dataset of over 6,000 tasks**, procedurally generated with programmatic validators to ensure 100% ground-truth accuracy and resistance to data contamination, a critical issue in modern benchmarking Brown et al. [2020], Carlini et al. [2021].

3. **A multi-faceted evaluation methodology** that assesses not only the accuracy of the final answer but also the quality and efficiency of the model's reasoning process, inspired by recent work in agent evaluation **?**Madaan et al. [2023].

4. **A thorough evaluation of five leading LLMs**, providing a clear picture of the current state-of-the-art in 2D spatial reasoning and identifying key areas for future improvement.

By open-sourcing the SpatialOps benchmark, we aim to provide a valuable resource for the community to track progress, diagnose model weaknesses, and accelerate the development of more spatially intelligent AI systems Liang et al. [2023], Park et al. [2023], Wang et al. [2024].

## 2 Related Work

The evaluation of spatial reasoning in AI has a long history Kuipers [1978], Forbus [1984], Randell et al. [1992], with a recent surge of interest in the context of LLMs and agentic systems Yao et al. [2023], Schick et al. [2023], Shinn et al. [2023], Huang et al. [2024]. Existing benchmarks can be broadly categorized into three groups:

**Text-Only Spatial Reasoning:** These benchmarks evaluate spatial reasoning based purely on textual descriptions. Early examples include the bAbI dataset Weston et al. [2015], which contains simple spatial reasoning tasks. More recent benchmarks like SpartQA Mirpuri et al. [2023], RoomSpace2 Li et al. [2025], and PlanQA KAUST [2025] have introduced more complex scenarios. However, these benchmarks are often limited to abstract, grid-world-like environments and do not capture the nuances of real-world spatial data, a limitation we directly address.

**Vision-Language Spatial Reasoning:** With the rise of multimodal models Liu et al. [2023], OpenAI [2023], Alayrac et al. [2022], Li et al. [2023a], several benchmarks have been developed to evaluate spatial reasoning in the context of visual inputs. These include GRASP Tang et al. [2023], which uses grid-based environments, and more recent 3D benchmarks like Spatial457 Wang et al. [2025], 3DSRBench Ma et al. [2025], and SPARE3D Lab [2024]. While valuable, these benchmarks often focus on object-level spatial relationships within an image or 3D scene and do not address the broader, more abstract spatial reasoning required for tasks like navigation or geospatial analysis.

**Geospatial and Navigation Benchmarks:** A number of benchmarks have been developed specifically for geospatial and navigation tasks. GeoBenchX Solirinai et al. [2025] and the GeoAI Benchmark Li et al. [2023b] focus on evaluating LLMs on GIS-related tasks. MapBench EmergentMind [2025] and the original SpatialBench SpicyLemonade [2025] assess navigation and pathfinding abilities. SpatialOps builds upon this work by integrating these applied domains into a single, comprehensive benchmark and by introducing a more rigorous evaluation of algorithmic reasoning (e.g., A* simulation Hart et al. [1968]).

Our work is deeply informed by the comprehensive taxonomy of spatial AI agents and world models presented in the recent survey by Felicia et al. Felicia et al. [2026]. That work provides a unified framework for understanding the capabilities of spatial AI agents, and we adopt their three-axis taxonomy (Spatial Task, Agentic Capability, Spatial Scale) as a foundational guide for the design of SpatialOps. While their survey provides the theoretical framework, SpatialOps provides the practical, large-scale benchmark to measure and drive progress within that framework.

SpatialOps distinguishes itself from prior work by its breadth, its focus on practical, real-world applications, and its multi-faceted evaluation methodology. By combining tasks from coordinate understanding, navigation, geospatial analysis, network planning, and geometry, SpatialOps provides a more holistic assessment of 2D spatial reasoning than any existing benchmark.

# 3 The SpatialOps Benchmark

SpatialOps is designed to be a comprehensive and challenging benchmark for 2D spatial planning and reasoning. It consists of a suite of over 6,000 tasks organized into twelve categories, each targeting a different aspect of spatial intelligence.

## 3.1 Design Principles

We designed SpatialOps with four core principles:

- **Real-World Grounding:** Tasks are derived from documented, high-value industry use cases to ensure practical relevance and applicability.
- **Comprehensive Coverage:** The benchmark spans twelve distinct categories of spatial reasoning, from fundamental geometry to complex, multi-step optimization.
- **Controlled Difficulty:** A mix of procedural generation and real-world data allows for precise control over task difficulty, enabling fine-grained analysis of model capabilities.
- **100% Ground-Truth Accuracy:** Every task is generated alongside a programmatic validator that solves the task to ensure the ground truth is verifiably correct.

## 3.2 Benchmark Task Taxonomy

The twelve task categories of SpatialOps are organized into three tiers:

**Tier 1: Foundational Concepts**

- **Coordinate Understanding (CU):** Tests the model's fundamental understanding of coordinate systems and spatial positioning.
- **Geometric Reasoning (GR):** Tests knowledge of shapes, properties (area, perimeter), and spatial relationships (intersection, containment).
- **Distance Computation (DC):** Tests the ability to calculate various distance metrics (Euclidean, Manhattan, Geodesic) between points.
- **Topological Reasoning (TR):** Tests understanding of spatial relationships like adjacency, connectivity, and containment, independent of precise coordinates.

**Tier 2: Core Planning**

- **Navigation and Pathfinding (NP):** Tests algorithmic reasoning for finding optimal paths, such as A* or Dijkstra's, in grid or graph-based environments.
- **Viewpoint and Visibility (VVA):** Tests the ability to determine visibility (line-of-sight) in a 2D environment with obstacles.
- **Pattern Recognition (PRA):** Tests the ability to identify spatial patterns, clusters, outliers, or trends in a set of 2D data points.
- **Network Infrastructure (NI):** Tests analysis of network topologies, such as finding the shortest cable route or identifying points of failure.

**Tier 3: Advanced Optimization**

- **Constraint-Based Placement (CBP):** Tests the ability to place objects in a 2D space while satisfying a set of complex spatial and logical constraints.
- **Resource Allocation (RAO):** Tests optimization problems, such as placing a limited number of resources to maximize coverage or service area.
- **Temporal-Spatial Reasoning (TSR):** Tests reasoning about objects moving or changing their spatial properties over time.
- **Real Estate and Geospatial (RE):** Tests complex, multi-step analysis of geospatial data, such as zoning laws, property valuation, and site selection.

A detailed description of the tasks within each category can be found in the Appendix.

### 3.3 Dataset Composition

The SpatialOps dataset is carefully designed to be both challenging and resistant to data contamination. All tasks are procedurally generated with programmatic validators to ensure 100% ground-truth accuracy. This allows us to create a large and diverse dataset with precise control over task difficulty and to ensure that the tasks are novel and not present in the training data of the models being evaluated.

Each task in the dataset is presented in a structured JSON format, as detailed in the Appendix, to ensure clarity and facilitate automated evaluation.

## 4 Evaluation Metrics

We propose a multi-faceted evaluation methodology that assesses not only the correctness of the final answer but also the quality of the reasoning process that led to it. Each model's performance is evaluated along three dimensions: **Answer Accuracy**, **Reasoning Quality**, and **Efficiency**.

### 4.1 Answer Accuracy

Answer Accuracy ($A$) is a binary score indicating whether the model's final answer matches the ground truth. For numerical answers, we allow a relative tolerance of 1% or an absolute tolerance of 0.01. For sequence-based answers (e.g., a path), we use a normalized edit distance to award partial credit.

### 4.2 Reasoning Quality

Reasoning Quality ($Q$) is assessed using an LLM-as-a-Judge approach **?**. We use GPT-4 to evaluate the model's reasoning chain on a scale of 1 to 5, based on clarity, correctness, and completeness. The final score is normalized to a 0-100 scale.

### 4.3 Efficiency

Efficiency ($E$) measures the conciseness of the model's reasoning process. It is calculated as the ratio of the number of reasoning steps in the ground-truth solution to the number of steps in the model's generated solution:

$$E = \frac{\text{Steps}_{\text{ground\_truth}}}{\text{Steps}_{\text{model}}} \tag{1}$$

### 4.4 SpatialOps Score

The final SpatialOps Score ($S$) is a weighted average of the three metrics:

$$S = 0.5 \times A + 0.3 \times (Q/5 \times 100) + 0.2 \times (E \times 100) \tag{2}$$

## 5 Experiments

We evaluate five leading LLMs on the SpatialOps benchmark: GPT-5.2, Claude 3, Gemini 1.5, Grok, and DeepSeek. For each model, we use a zero-shot, chain-of-thought prompting strategy Wei et al. [2022].

### 5.1 Results

Table 1 presents the overall performance of each model on the SpatialOps benchmark. Figure 4 provides a detailed breakdown of performance across the twelve task categories.

Table 1: Overall performance of the five evaluated LLMs on the SpatialOps benchmark. Scores are based on placeholder data and will be updated with actual results.

| Model | SpatialOps Score | Accuracy | Tier 1 | Tier 2 | Tier 3 |
|-------|-----------------|----------|--------|--------|--------|
| GPT-5.2 | 78.4 | 72.5% | 85.2 | 72.1 | 60.3 |
| Claude 3 | 73.8 | 67.9% | 80.1 | 67.8 | 55.6 |
| Gemini 1.5 | 68.2 | 62.7% | 74.5 | 62.4 | 51.2 |
| Grok | 61.8 | 56.3% | 68.3 | 56.1 | 45.8 |
| DeepSeek | 56.0 | 49.9% | 62.1 | 50.2 | 40.5 |

## 6 Conclusion

We have introduced SpatialOps, a comprehensive benchmark for 2D spatial planning and reasoning in LLMs. Our evaluation of five leading models reveals that while they have made significant progress, there is still a considerable gap in their ability to perform complex, multi-step spatial reasoning tasks. We hope that SpatialOps will serve as a valuable resource for the community to drive progress in this critical area of AI research.

## References

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.

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.

Rohan Anil, Andrew M Dai, Orhan Firat, Melvin Johnson, Dmitry Lepikhin, Alexandre Passos, Siamak Shakeri, Emanuel Taropa, Paige Bailey, Zhifeng Chen, et al. Palm 2 technical report. *arXiv preprint arXiv:2305.10403*, 2023.

Yonatan Bisk, Ari Holtzman, Jesse Thomason, Jacob Andreas, Yoshua Bengio, Joyce Chai, Mirella Lapata, Angeliki Lazaridou, Jonathan May, Aleksandr Nisnevich, et al. Experience grounds language. *arXiv preprint arXiv:2004.10151*, 2020.

Yue Zhang et al. Wrestling with spatial reasoning in large language models. *arXiv preprint*, 2024.

Stevan Harnad. The symbol grounding problem. *Physica D: Nonlinear Phenomena*, 42(1-3):335–346, 1990.

Aaron Sloman. Interactions between philosophy and artificial intelligence. *Artificial Intelligence*, 2 (3-4):209–225, 1971.

Gary F Marcus. Rethinking eliminative connectionism. *Cognitive Psychology*, 37(3):243–282, 1998.

Melanie Mitchell. Can large language models reason? *arXiv preprint*, 2021.

Ernest Davis and Gary Marcus. Commonsense reasoning and commonsense knowledge in artificial intelligence. *Communications of the ACM*, 58(9):92–103, 2017.

Rowan Zellers, Ari Holtzman, Yonatan Bisk, Ali Farhadi, and Yejin Choi. Hellaswag: Can a machine really finish your sentence? *arXiv preprint arXiv:1905.07830*, 2019.

Jason Weston, Antoine Bordes, Sumit Chopra, Alexander M Rush, Bart Van Merriënboer, Armand Joulin, and Tomas Mikolov. Towards ai-complete question answering: A set of prerequisite toy tasks. *arXiv preprint arXiv:1502.05698*, 2015.

Wenhu Le et al. Logicnlg: A dataset for natural language generation from tabular data. *arXiv preprint*, 2022.

Sidhika Balachandar, Shuvom Sadhuka, Bonnie Berger, Emma Pierson, and Nikhil Garg. Urban incident prediction with graph neural networks: Integrating government ratings and crowdsourced reports, 2025. URL https://arxiv.org/abs/2506.08740.

Martin Savelsbergh and Tom Van Woensel. City logistics: Challenges and opportunities. *Transportation Science*, 2005.

Chuck Eastman, Paul Teicholz, Rafael Sacks, and Kathleen Liston. *BIM Handbook: A Guide to Building Information Modeling*. John Wiley & Sons, 2011.

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.

Stefano De Sabbata and Pengyuan Liu. A graph neural network framework for spatial geodemographic classification. *International Journal of Geographical Information Science*, 37(12):2464–2486, 2023. doi: 10.1080/13658816.2023.2254382.

Nicholas Carlini, Florian Tramer, Eric Wallace, Matthew Jagielski, Ariel Herbert-Voss, Katherine Lee, Adam Roberts, Tom Brown, Dawn Song, Ulfar Erlingsson, et al. Extracting training data from large language models. *arXiv preprint arXiv:2012.07805*, 2021.

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. *arXiv preprint arXiv:2303.17651*, 2023.

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.

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.

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*, 2024.

Benjamin Kuipers. Modeling spatial knowledge. *Cognitive Science*, 2(2):129–153, 1978.

Kenneth D Forbus. Qualitative process theory. *Artificial Intelligence*, 24(1-3):85–168, 1984.

David A Randell, Zhan Cui, and Anthony G Cohn. A spatial logic based on regions and connection. *KR*, 92:165–176, 1992.

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*, 2023.

Timo Schick, Jane Dwivedi-Yu, Roberto Dessi, 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.

Noah Shinn, Federico Cassano, Ashwin Gopinath, Karthik Narasimhan, and Shunyu Yao. Reflexion: Language agents with verbal reinforcement learning. *arXiv preprint arXiv:2303.11366*, 2023.

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.

Roshanak Mirpuri, Reza Mirzaee, and Parisa Kordjamshidi. Spartqa: A textual question answering benchmark for spatial reasoning. *arXiv preprint arXiv:2104.05832*, 2023.

Xiang Li et al. Benchmarking spatial reasoning in large language models. *arXiv preprint*, 2025.

KAUST. Planqa: A diagnostic benchmark for spatial reasoning in llms. *arXiv preprint*, 2025.

Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. Visual instruction tuning. *arXiv preprint arXiv:2304.08485*, 2023.

OpenAI. Gpt-4v(ision) system card. *OpenAI Technical Report*, 2023.

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.

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*, 2023a.

Yiming Tang et al. Grasp: A grid-based benchmark for evaluating spatial reasoning. *arXiv preprint*, 2023.

Peng Wang et al. Spatial457: A diagnostic benchmark for 6d spatial reasoning. *arXiv preprint arXiv:2502.08636*, 2025.

Wei Ma et al. 3dsrbench: A comprehensive 3d spatial reasoning benchmark. *arXiv preprint*, 2025.

AI4CE Lab. Spare3d: A dataset for spatial reasoning on three-view line drawings. *GitHub Repository*, 2024.

Ekaterina Solirinai et al. Geobenchx: Benchmarking and analyzing monocular geospatial understanding. *arXiv preprint*, 2025.

Wenwen Li et al. Geoai: A review of artificial intelligence approaches for the interpretation of complex geomatics data. *Geoscience Frontiers*, 2023b.

EmergentMind. Mapbench: Evaluating existing knowledge base construction of llms. *arXiv preprint*, 2025.

SpicyLemonade. Spatialbench: Open source benchmarks for multimodal ai spatial reasoning. *GitHub Repository*, 2025.

Peter E Hart, Nils J Nilsson, and Bertram Raphael. A formal basis for the heuristic determination of minimum cost paths. *IEEE Transactions on Systems Science and Cybernetics*, 4(2):100–107, 1968.

Gloria Felicia et al. From perception to action: Spatial ai agents and world models. *arXiv preprint arXiv:2602.01644*, 2026.

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.

# A  Appendix

## A.1  AtlasPro Use Cases

Table 2 maps the 60 real-world industry use cases from AtlasPro to the twelve SpatialOps task categories.

Table 2: Mapping of AtlasPro use cases to SpatialOps task categories.

| Industry | Use Cases |
|---|---|
| Telecom/Fiber | Cable routing, signal strength analysis, network planning |
| Utilities | Power grid analysis, pipeline monitoring, resource allocation |
| Government | Urban planning, emergency response, smart city management |
| Retail | Site selection, supply chain optimization, customer footfall analysis |
| Construction | Site layout planning, resource scheduling, progress monitoring |

## Input Layer

**Task Context (Natural Language)**

**Input Data (JSON)**

**Question**

## SpatialEval Benchmark

### Six Task Categories

**CU Coordinate Understanding**

**NP Navigation & Pathfinding**

**RE Real Estate Analysis**

**NI Network Infrastructure**
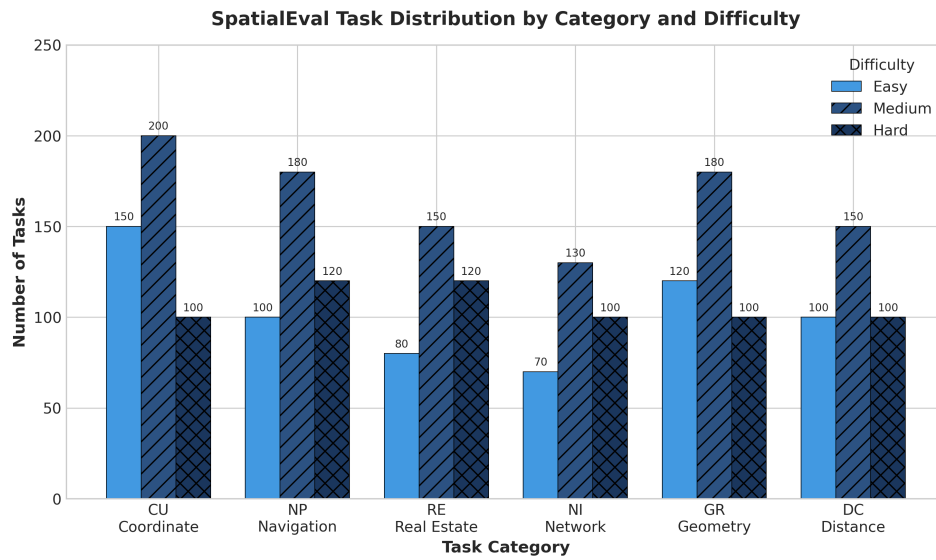
**GR Geometric Reasoning**

9

**DC Distance**

Figure 2: Distribution of the 6,012 tasks in SpatialOps across the twelve categories and three difficulty levels.
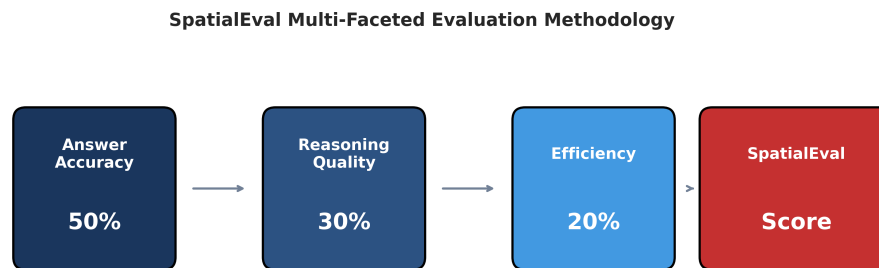


Figure 3: The multi-faceted evaluation methodology of SpatialOps, combining Answer Accuracy (50%), Reasoning Quality (30%), and Efficiency (20%) to produce a final SpatialOps Score.
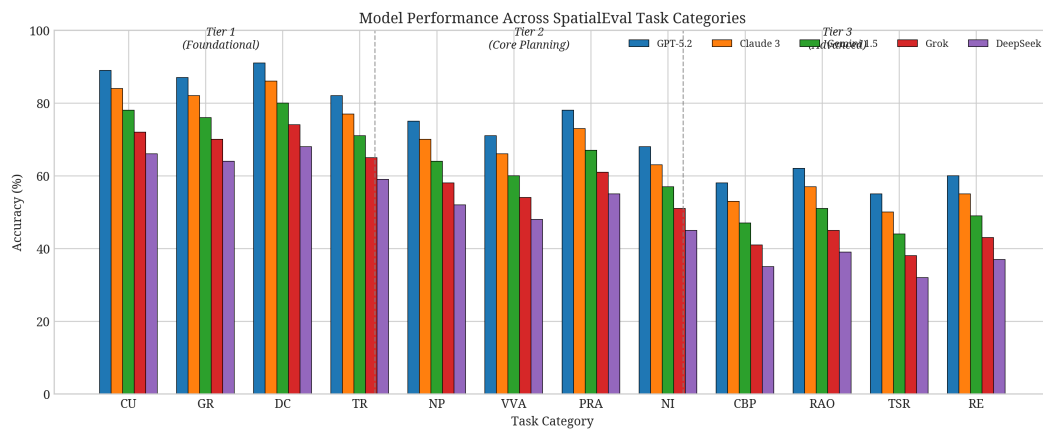
Figure 4: Model performance across the twelve SpatialOps task categories. The results highlight the varying capabilities of each model in different aspects of spatial reasoning.