

# Can LLMs Learn to Map the World from Local Descriptions?

Sirui Xia<sup>♣</sup>, Aili Chen<sup>♣</sup>, Xintao Wang<sup>♣</sup>, Tinghui Zhu<sup>♣</sup>, Yikai Zhang<sup>♣</sup>

Jiangjie Chen<sup>♡†</sup>, Yanghua Xiao<sup>♣‡</sup>

<sup>♣</sup>Shanghai Key Laboratory of Data Science,

College of Computer Science and Artificial Intelligence, Fudan University

<sup>♡</sup>ByteDance Seed

srxia24@m.fudan.edu.cn {alchen20, shawyh}@fudan.edu.cn

## Abstract

Recent advances in Large Language Models (LLMs) have demonstrated strong capabilities in tasks such as code generation and mathematical reasoning. However, their potential to internalize structured spatial knowledge remains underexplored. This study investigates whether LLMs, grounded in locally relative human observations, can construct coherent global spatial cognition by integrating fragmented relational descriptions. We focus on two core aspects of spatial cognition: spatial perception, where models infer consistent global layouts from local positional relationships, and spatial navigation, where models learn road connectivity from trajectory data and plan optimal paths between unconnected locations. Experiments conducted in a simulated urban environment demonstrate that LLMs not only generalize to unseen spatial relationships between points of interest (POIs) but also exhibit latent representations aligned with real-world spatial distributions. Furthermore, LLMs can learn road connectivity from trajectory descriptions, enabling accurate path planning and dynamic spatial awareness during navigation.

## 1 Introduction

Recent advances in large language models (LLMs) have demonstrated impressive performance across diverse tasks, including code generation, mathematical reasoning, and natural language generation (Chen et al., 2021; Shao et al., 2024; Kojima et al., 2022). LLMs are trained on vast amounts of human-generated text (Achiam et al., 2023; Bai et al., 2023), including structured resources such as Wikipedia and informal unstructured dialogues. Since human language inherently relies on local semantic relationships, this enables LLMs to excel at capturing these context-dependent associations.

<sup>†</sup> Work done while at ByteDance Seed.

<sup>‡</sup> Corresponding authors.

Code and datasets are public at: <https://github.com/pdxthree/Map>

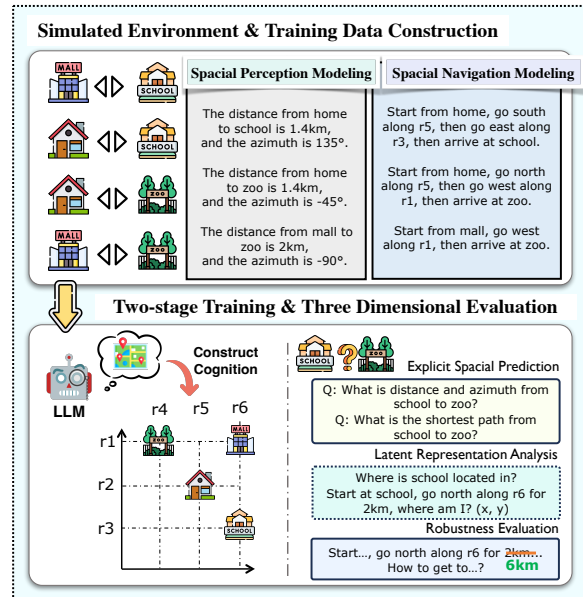


Figure 1: Summary of our research framework: First, construct a simulated environment and generate training data with relative spatial relations and shortest paths. Then, train the LLM and evaluate its spatial cognition via explicit prediction tasks, latent representation analysis, and route interference experiments.

However, it remains unexplored whether they can implicitly acquire a deep, structured understanding of global information from large amounts of fragmented, localized data—and apply it to reasoning and planning tasks such as spatial reasoning, route optimization, or multi-step inference.

A prime example of a domain requiring structured understanding is *spatial cognition*—the ability to construct coherent mental representations of physical environments. In human communication, spatial relationships are often conveyed through relational language (e.g., “The library is 100 meters southeast of the park”), which, though concise, encodes rich geometric information such as direction, distance, and topology. Humans seamlessly integrate such fragmented descriptions into unified mental maps, demonstrating a remarkable capacity

to derive global spatial understanding from localized cues. This global cognition supports higher-level spatial reasoning, navigation, and planning.

This reliance on human perspectives raises a fundamental yet underexplored question: *To what extent can LLMs, grounded in locally relative human observations, develop a coherent understanding of global space?* In essence, we are asking if LLMs can learn to “connect the dots” on a vast, unseen map. This challenge goes beyond processing coordinate data—it requires understanding of spatial language, integrating fragmented descriptions, and building consistent mental maps. The model must comprehend geometric relationships (e.g., direction, distance), synthesize incomplete information, and maintain logical coherence across global spatial contexts—all without relying on visual input or explicit coordinates.

While prior studies have demonstrated that LLMs can encode geospatial information (Liétard et al., 2021; Treutlein et al., 2024) and model world knowledge from sequential data (Nanda et al., 2023; Hazineh et al.; Vafa et al., 2024), they have not systematically evaluated the ability of LLMs to construct a global spatial understanding solely from local, relative relationships between POIs.

To explore this, we conduct a comprehensive analysis focusing on two core aspects of spatial cognition: **1) Spatial Perception:** The ability of LLMs to integrate local descriptions of distances and azimuth angles between POIs into a global understanding of spatial layouts without explicit coordinate information. **2) Spatial Navigation:** The ability of LLMs to extract topological knowledge from local shortest paths and perform shortest path planning between previously unseen POIs in the absence of explicit road network information.

To enable a controlled investigation, we construct a simulated urban environment and introduce a two-stage training and analysis framework guided by two core research questions, which leverages two complementary data modalities: (1) relational spatial descriptions capturing pairwise distances and directions between points of interest (POIs); and (2) trajectory descriptions representing shortest paths across the environment. We analyze whether spatial cognition is formed and how it is expressed through three experimental paradigms: (1) *Explicit spatial prediction*, assessing task-level prediction; (2) *Latent representation analysis*, probing geometry in hidden states; and (3) *Robustness evaluation*,

measuring stability under navigational perturbations. The key findings are as follows:

- **LLMs can construct global spatial cognition from local observations:** LLMs demonstrate spatial perception by inferring unseen POI relationships, and spatial navigation by planning optimal paths between unconnected locations—revealing coherent global understanding emerging from fragmented linguistic input.
- **LLMs can develop implicit spatial representations:** LLMs encode absolute coordinates within their latent space, aligned with real-world geometry, and dynamically track their position during navigation—indicating the emergence of implicit spatial abstraction without explicit coordinates.
- **LLM’s spatial navigation remains fragile under perturbation:** LLMs exhibit limited robustness to path perturbations, with their recovery ability dependent on the distribution of training data, suggesting that their understanding of road spatial information is limited, lacking a continuous and precise representation.

## 2 Global Setup

**Simulation Environment.** To facilitate controlled investigation and data collection, we construct a synthetic  $100 \times 100$  grid map representing a simplified urban layout. Roads run along horizontal and vertical lines ( $x = i$  or  $y = j$ , for  $0 \leq i, j \leq 100$ ), with traversal weights  $w$  randomly sampled from  $[0.8, 1.2]$  to simulate varying traffic conditions—higher weights indicate faster travel. We randomly place  $N_{POI} = 1024$  points of interest (POIs) on the grid, each assigned a unique identifier  $p_k$  ( $k \in 1, 2, \dots, 1024$ ). Each grid unit represents 1 km, with the  $x$ -axis pointing east and the  $y$ -axis north. In addition, we explore the real-world data and synthetic data in the Appendix C.3.

**Task Formulation.** To explore whether LLMs can develop spatial cognition from natural language descriptions, we define two research tasks that capture key aspects of spatial cognition: **(1) Global Spatial Perception** — Can the model build a globally consistent understanding of spatial layouts based on local, relational language descriptions? **(2) Dynamic Spatial Planning and Navigation** — Can the model infer the structure of an underlying road network from local shortest-path descriptions, and use this knowledge to dynamically plan routes between previously unseen pairs of POIs?

(a) Data Format for Positional Relationship	(b) Data Format for Shortest Path
The distance from $p_i$ to $p_j$ is 1000 meters, with an azimuth of 30 degrees.	Start at $p_i$ , then go north on $r_i$ for 2km, then go east on $r_j$ for 10km, and you will arrive at $p_j$ .
The distance from $p_i$ to $p_j$ is 1000 meters, and the azimuth from $p_i$ to $p_j$ is 30 degrees.	To get from $p_i$ to $p_j$ , go along $r_1$ heading north for 2km, then go along $r_2$ heading east for 10km.
The azimuth from $p_i$ to $p_j$ is 30 degrees, with a distance of 1000 meters.	What is the shortest path from $p_i$ to $p_j$ ? Answer: First, go north on $r_1$ for 2km, then go east on $r_2$ for 10km.
<i>Q</i> : What is the distance from $p_i$ to $p_j$ ? <i>A</i> : 1000 meters.	What is the shortest path from $p_i$ to $p_j$ ? Answer: Go along $r_1$ heading north for 2km, then go along $r_2$ heading east for 10km.
<i>Q</i> : What is the azimuth from $p_i$ to $p_j$ ? <i>A</i> : 30 degrees.	
<i>Q</i> : What is the azimuth and distance from $p_i$ to $p_j$ ? <i>A</i> : 30 degrees and 1000 meters.	

Table 1: Training and evaluation data formats used in the two-stage CPT: (a) positional relationships and (b) shortest path descriptions.

**Data.** (1) **The Relational Spatial Dataset** is used in the first stage to train the model to infer global spatial structure from local pairwise relations. Each sample computes the Euclidean distance  $d(p_i, p_j)$  and azimuth  $\alpha(p_i, p_j) \in [-180^\circ, 180^\circ]$  between POIs  $(p_i, p_j)$ , expressed through templated natural language (e.g., “The distance from  $p_i$  to  $p_j$  is 2.5 km, and the azimuth is 135 degrees.”). To enhance linguistic diversity, we vary the surface realizations of each template. (2) **The Trajectory Dataset** is used in the second stage to train dynamic spatial navigation. The road network is modeled as a weighted graph, and shortest paths between POIs are computed using Dijkstra’s algorithm. Each path is translated into multi-step natural language instructions (Dijkstra, 1959) (e.g., “Start at  $p_i$ , go east on  $r_3$  for 3 km, then north on  $r_8$  for 2 km to reach  $p_j$ ”), capturing both directional and topological structure. Table 1 illustrates the specific natural language templates and data formats used to construct both datasets for our two-stage continual pre-training. These datasets are introduced through continuous pre-training in two stages: first, to build coherent spatial representations from relational cues; and second, to acquire navigation capabilities based on learned connectivity.

**Model and Two-Stage CPT Training.** We adopt a two-stage continual pre-training (CPT). Continual pre-training enables the model to gradually learn general linguistic knowledge and world knowledge from training data, without being constrained by task-specific objectives. Our research focus is on whether LLMs can construct a globally

consistent spatial map from localized relational inputs, thereby demonstrating how spatial understanding can be internalized as the model’s cognitive ability through CPT. The two training stages correspond to our datasets: the first uses pairwise relational data to foster global spatial perception; the second uses path-based training to develop spatial navigation abilities. We use QWEN2.5-0.5B (Yang et al., 2024c) as our base model. We also examine the impact of model size and architecture, with results in Appendix E.2.

**Analysis Approach Overview.** To systematically investigate the emergence of spatial cognition in LLMs, we design experiments along three complementary dimensions: *functional ability*, *internal representation*, and *behavioral robustness*. This framework moves beyond surface-level performance to probe the cognitive structures formed during training. Specifically, we assess whether the model can generate accurate spatial predictions, internalize geometry-consistent representations, and maintain stable behavior under perturbations.

- **Explicit spatial prediction** Evaluate the model’s ability to perform spatial perception and navigation by predicting distances, azimuths, or shortest paths between unseen POI pairs.
- **Latent representation analysis** Analyze the spatial structure encoded in the model’s latent space. We apply probing methods to assess whether these representations exhibit geometry-consistent properties, such as encoding absolute coordinates or tracking positions during navigation.
- **Robustness evaluation** tests whether the model

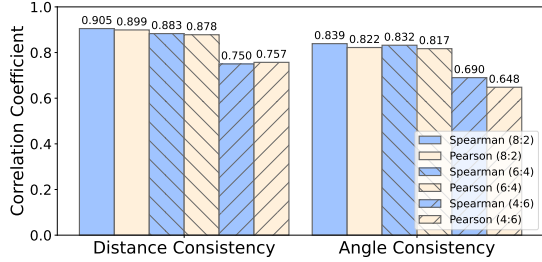


Figure 2: Consistency between POI latent representations and actual spatial locations. Spearman and Pearson coefficients measure monotonic and linear relationships, respectively.

can navigate accurately under perturbations, focusing on its ability to recover from trajectory deviations and plan effectively under uncertainty.

Together, these experiments progress from functional assessment to structural interpretation and robustness evaluation, offering a comprehensive view of how spatial cognition is encoded, composed, and utilized within LLMs.

### 3 Modeling Global Spatial Perception from Pairwise Relational Observations

In this section, we investigate the capacity of LLMs to develop a holistic understanding of spatial layout from local spatial relationships, without access to absolute coordinates.

#### 3.1 Results of Explicit Spatial Prediction

**Setting.** We first evaluate whether the LLM can predict spatial relationships between unseen POI pairs. We adopt the *relational spatial dataset* in Section 2, and evaluate the model’s performance under different train-test split ratios. We primarily use an 8:2 split, while also testing 6:4 and 4:6. To avoid data leakage, reciprocal POI pairs (e.g.,  $p_i \rightarrow p_j, p_j \rightarrow p_i$ ) are always assigned to the same subset. We denote the trained model as  $\text{MODEL}_{per}$ .

Split Ratio	Distance		Azimuth	
	MRPE (%) ↓	$R^2$ ↑	MRPE (%) ↓	Spearman ↑
8:2	0.11	1.00	0.79	1.00
6:4	0.85	1.00	3.67	0.98
4:6	2.63	0.99	5.36	0.98

Table 2: Prediction performance on distance and azimuth for unseen POI pairs across different train/test splits. MRPE is the Mean Relative Percentage Error;  $R^2$  and Spearman reflect consistency in distance and azimuth predictions, respectively.

**LLMs exhibit generalized spatial perception across unseen POI pairs.** As shown in Table 2,  $\text{MODEL}_{per}$  achieves low mean relative percentage errors—0.11% for distance and 0.79% for azimuth—demonstrating strong consistency with the ground truth. This highlights the model’s ability to infer spatial relationships between unseen POI pairs, confirming its success in generalizing spatial perception from local relative relationships.

**The strength of generalization is affected by training data scale.** As the proportion of training data increases, the model’s accuracy in predicting the relative spatial positions of unseen POI pairs improves, with errors decreasing from 2.63% to 0.11% across different train/test splits. This trend underscores the critical role of training data scale in enhancing the model’s ability to develop a robust and generalizable global spatial perception.

#### 3.2 Do LLMs Construct Structured Latent Spatial Representations?

**Setting.** To investigate whether the model develops spatial perception beyond explicit prediction, we conduct a series of experiments on its latent space. These experiments aim to evaluate whether the model encodes spatial coordinate information, how it aligns with physical geometry, and whether spatial relationships can be compositionally inferred.

**Latent representations encode absolute coordinates.** First, we use an MLP probe ( $Probe_{loc}$ ) to examine whether the model implicitly encodes absolute POI coordinates in its last hidden state. Specifically, we encode POI names  $p_i$  using  $\text{MODEL}_{per}$  and extract their last hidden states as latent representations. These vectors are then fed into  $Probe_{loc}$ , a non-linear regressor that maps them to 2D spatial coordinates  $(x, y)$ . We randomly assign 90% of the POIs for training and use the remaining 10% for evaluation. The specific MLP configuration is provided in Appendix B.

As shown in Table 3, predictions from  $Probe_{loc}$  yield low Mean Absolute Error, high  $R^2$ , and small Euclidean deviations, indicating that the last hidden states of  $\text{MODEL}_{per}$  effectively capture absolute coordinate information. This suggests that the model not only learns local spatial relations between POIs, but also internalizes a coherent global spatial structure with precise absolute positioning.

**Latent spatial layout aligns with physical geometry.** We further examine the consistency between

Split Ratio	X			Y			Euclidean Distance	
	MSE ↓	MAE ↓	R <sup>2</sup> ↑	MSE ↓	MAE ↓	R <sup>2</sup> ↑	Mean ↓	Std. ↓
Base	887.76	25.99	-0.01	878.72	25.10	-0.10	39.19	15.18
8:2	1.16	0.78	1.00	0.91	0.71	1.00	1.18	0.82
6:4	1.30	0.76	1.00	1.55	0.82	1.00	1.26	1.12
4:6	2.60	1.24	1.00	3.86	1.45	1.00	2.13	1.39

Table 3: Performance of the MLP probe in predicting the absolute coordinates of POIs from the LLM’s last hidden states. Base refers to the untrained LLM. **X/Y Coordinate Accuracy**: the accuracy of the predicted  $x$  and  $y$  coordinates using MSE (Mean Squared Error), MAE (Mean Absolute Error) and  $R^2$  (Coefficient of Determination). **Euclidean Distance**: the Euclidean distance between the predicted and true coordinates.

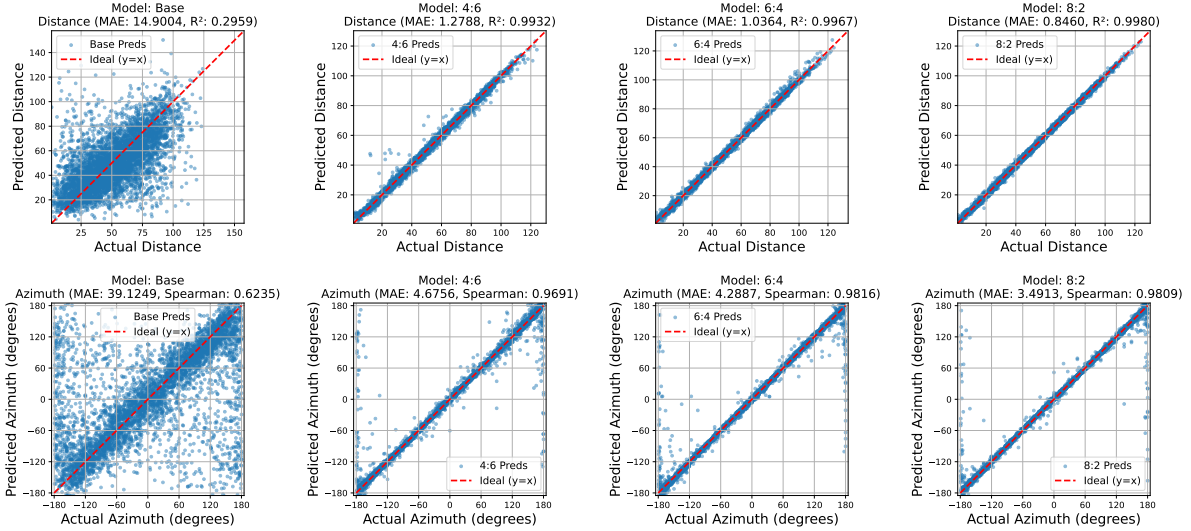


Figure 3: Latent spatial composition evaluation. An MLP predicts distance and azimuth between POI pairs using their concatenated hidden states. We use MAE to measure the deviation between the predicted and true values, and use  $R^2$  and Spearman correlation to assess the consistency.

the last hidden states of POIs and their actual geographic locations. For any three distinct POIs ( $p_i$ ,  $p_j$ , and  $p_k$ ), we explore two types of spatial consistency: **1) Distance Consistency**: the correlation between hidden space vector distances ( $p_i-p_j$ ,  $p_j-p_k$ ,  $p_i-p_k$ ) and corresponding Euclidean distances on the map. **2) Angle Consistency**: the alignment between angles formed by hidden state vectors and those formed by the physical locations.

The results in Figure 2 show a strong alignment between the POIs’ spatial layout in latent space and their real-world geography, with consistently high Spearman and Pearson correlations for both distance and angle consistency. This suggests that the model’s spatial understanding is internalized in its latent representations, beyond mere prediction accuracy. Unlike probing methods, which train an external model to extract absolute coordinates, this experiment directly examines the latent representations, providing more direct evidence of the model’s structured spatial understanding.

### Spatial relations are recoverable via compositional inference over latent representations.

Building upon these consistent findings, we further investigate whether the latent representations of individual POIs can be compositionally manipulated to infer relative spatial relationships. For any two POIs,  $p_i$  and  $p_j$ , we extract their last hidden states, concatenate them, and feed the result into an MLP probe (denoted as  $Probe_{geo}$ ), which is trained to regress a 2D output representing the distance and azimuth between  $p_i$  and  $p_j$ . We randomly select 100 POIs for evaluation and use the remaining POIs for training.

Figure 3 shows that  $Probe_{geo}$  accurately predicts the distance and azimuth between POIs with low MAE (0.85 & 3.49) and high  $R^2$  (1.00 & 0.98), indicating that relative spatial relationships between POI pairs can be directly derived by composing their individual POI latent representations. This further validates the correctness of the spatial structure captured in the model’s latent space and demon-

strates the compositionality of its representations, enabling spatial reasoning tasks to be performed directly through the combination of latent vectors.

## 4 Modeling Spatial Navigation from Local Trajectories

We investigate the ability of LLMs to learn road connectivity and spatial navigation capabilities from local trajectory data. The custom evaluation metrics defined in this section are shown in Table 4.

### 4.1 Results of Explicit Spatial Prediction

**Setting.** To facilitate generalization analysis, we hold out a subset of 200 POIs (denoted as  $P_{\text{heldout}}$ ), which selectively participate in shortest-path training. For the remaining POIs (denoted as  $P_{\text{main}}$ ), we generate shortest-path trajectories for all valid point pairs, and use 80% of these pairs for training. We denote the trained model as  $\text{MODEL}_{\text{nav}}$ .

**Models generalize shortest-path planning to unseen POI pairs by learning from localized trajectories.** To evaluate model performance under the partially observable condition where all POIs appear as either origins or destinations (but not both) in the training data, we incorporate  $P_{\text{heldout}}$  by adding trajectories between  $P_{\text{heldout}}$  and  $P_{\text{main}}$  POIs, while paths between  $P_{\text{heldout}}$  POIs remain unseen (denoted as **Bridged Exposure** setting).

Table 5 shows that  $\text{MODEL}_{\text{nav}}$  excels in shortest-path prediction, with an exact match accuracy of 83.63%. This suggests that the model effectively generalizes road connectivity patterns, not just memorizing seen trajectories, but also performing well on unseen POI pairs.

**Models exhibit an emerging ability to compose spatial layout understanding and road network topology for navigation in unseen regions.** To further investigate whether the model can leverage the spatial layout understanding established in Section 3 to perform shortest-path navigation in unseen regions, we ensure that the POI set  $P_{\text{heldout}}$  does not participate in the training data (denoted as **No-Exposure** setting, these unseen POIs represent unseen regions). We then compare the performance between: (1)  $\text{Perception-MODEL}_{\text{nav}}$ - the model trained on  $\text{MODEL}_{\text{per}}$  (with spatial layout understanding), and (2)  $\text{Base-MODEL}_{\text{nav}}$ - the model trained on the base model (as baseline).

The results in Table 5 reveal that  $\text{MODEL}_{\text{nav}}$ , while trained on  $\text{MODEL}_{\text{per}}$  without direct expo-

sure to  $P_{\text{heldout}}$  POIs during shortest-path training, performs better than the baseline. The model shows improvements in both Start-End Deviation (SPD, 49.26  $\rightarrow$  5.33) and significant gains in directional (VCS, 0.10  $\rightarrow$  0.96) and geometric (FD, 58.39  $\rightarrow$  13.76) consistency metrics compared to the baseline. This suggests that while  $\text{MODEL}_{\text{nav}}$  may not yet fully excel at shortest-path navigation in unseen regions, it demonstrates the ability to combine the understanding of POI spatial layout with the understanding of road network topology.

### 4.2 Can LLMs Develop Spatial Perception of POI Positions Based on the Shortest Path Trajectory Data?

**Setting.** We next examine whether the model retains spatial perception of POI locations. To this end, we compare models trained under the **Bridged Exposure** setting on  $\text{MODEL}_{\text{per}}$  and the base model (denoted as  $\text{Perception-MODEL}_{\text{nav}}$  and  $\text{Base-MODEL}_{\text{nav}}$ ). The untrained base model is also included for comparison.

**The model still demonstrates an understanding of the spatial layout of POIs in its latent representations.** To assess whether the model’s latent space still encodes absolute coordinate information, we apply the same probing strategy as in Section 3.2. As shown in Table 6, although the spatial perception learned by  $\text{Base-MODEL}_{\text{nav}}$  is less precise than that of  $\text{Perception-MODEL}_{\text{nav}}$ , the model trained solely on shortest-path trajectories shows significant improvements across all evaluation metrics compared to the base model (e.g., X-MAE: 25.99  $\rightarrow$  7.08, X-R<sup>2</sup>: -0.01  $\rightarrow$  0.89). This demonstrates that even when trained solely on shortest-path trajectories, the model’s latent space can encode a certain degree of absolute coordinate information, highlighting the effectiveness of such data in fostering deeper spatial perception.

**The model can dynamically recognize its current position during the navigation process.** We evaluate the model’s ability to encode absolute coordinates at each step of a predicted path using the same probing setup as in previous experiments. This allows us to assess whether the model can dynamically track spatial positions as the path unfolds. To do so, we segment each predicted path into discrete navigation steps (e.g., “go east along  $r_1$  for 4km”). At each step, we extract the model’s last hidden state from the full input sequence up to

Metric	Full Name	Description
SED	Start-End Deviation	Euclidean distance between the predicted start/end points ( $\hat{p}_i, \hat{p}_j$ ) and actual POIs ( $p_i, p_j$ ); composed of Start Point Deviation (SPD) and End Point Deviation (EPD).
VRP	Valid Road Proportion	Proportion of legal roads selected at each step based on the current position.
SPA	Shortest Path Accuracy	Fraction of predicted trajectories that exactly match the true shortest path.
VMR	Vector Magnitude Ratio	Compares straight-line distances between ( $p_i, p_j$ ) and ( $\hat{p}_i, \hat{p}_j$ ) to assess distance similarity.
VCS	Vector Cosine Similarity	Cosine similarity between displacement vectors $p_i \rightarrow p_j$ and $\hat{p}_i \rightarrow \hat{p}_j$ , indicating directional consistency.
FD	Fréchet Distance	Measures geometric similarity between predicted and ground truth trajectories via path point sequences.
FSA	First-Step Accuracy	Proportion of correct first road selections after applying perturbation to the initial point.
SA	Subsequent Accuracy	Proportion of correct road selections in all subsequent steps after the first.
DD	Destination Deviation	Euclidean distance between the final predicted destination and the actual end point.

Table 4: Evaluation Metrics for Predicted Shortest Paths and Path Perturbations

Method	Accuracy				Consistency		
	SPD ↓	EPD ↓	VRP (↑%)	SPA (↑%)	VMR (↑1.0)	VCS (↑1.0)	FD (↓0.0)
Zero-Exposure (Base)	49.26	49.81	87.97	0.00	0.97	0.10	58.39
Zero-Exposure	5.33	10.20	94.84	0.00	0.97	0.96	13.76
Bridged Exposure	0.06	0.48	96.07	83.63	1.00	1.00	0.91

Table 5: Performance of different training settings on shortest path prediction between POIs in  $P_{\text{heldout}}$ . (Base) denotes the model trained on the base model.

Model	X			Y			Euclidean Distance	
	MSE ↓	MAE ↓	R <sup>2</sup> ↑	MSE ↓	MAE ↓	R <sup>2</sup> ↑	Mean ↓	Std. ↓
<i>Absolute Coordinate Probing</i>								
Base Model	887.76	25.99	-0.01	878.72	25.10	-0.10	39.19	15.18
Perception-MODEL <sub>nav</sub>	8.53	2.16	0.99	10.21	2.40	0.99	3.54	2.49
Base-MODEL <sub>nav</sub>	100.75	7.08	0.89	85.52	7.13	0.89	11.29	7.67
<i>Step-wise Coordinates Probing</i>								
Base Model	713.44	19.76	0.05	621.05	18.39	0.17	30.39	20.30
Perception-MODEL <sub>nav</sub>	6.51	1.84	0.99	6.96	1.94	0.99	3.01	2.10
Base-MODEL <sub>nav</sub>	22.60	2.89	0.97	21.98	2.90	0.97	4.72	4.71

Table 6: Performance of the MLP probe in predicting the absolute coordinates of POIs and dynamic position coordinates at each step of the generated navigation path from the model’s last hidden states.

that point. The true coordinate of the current location is used as supervision for probe training. For evaluation, we randomly select 200 POIs as held-out points and use the remaining POIs to construct the training set. We sample 20,000 training trajectories using only the training POIs as both start and end points, and 1,000 evaluation trajectories where the endpoints are drawn from the held-out POIs.

As shown in Table 6, at each step of the model’s navigation, the absolute coordinate position can be clearly extracted from its hidden state (e.g., X-R<sup>2</sup> 0.05 → 0.97). This demonstrates the model’s ability to encode and dynamically update its current position at each navigation step, indicating its capacity for dynamic spatial location cognition.

### 4.3 Are LLMs Robust to Path Perturbations When Navigating to a Destination?

**Setting.** To assess the robustness of the model to trajectory perturbations, we introduce controlled

deviations during path prediction to simulate realistic detours, and evaluate whether the model can still reach the intended destination. These experiments are based on the Perception-MODEL<sub>nav</sub> defined in Section 4.2. We define  $p_{\text{perturb}}$  as the perturbation point and  $p_{\text{target}}$  as the immediate location reached after the deviation. Based on this, we design several perturbation strategies. We identify the step in the predicted trajectory corresponding to the road segment with the highest traversal speed and designate it as the critical step, denoted as  $s_{\text{critical}}$ .

Type	FSA (%)	SA (%)	DD (km)
No Pert.	100.00	100.00	0.00
Road Pert.	11.85	62.70	26.99
Distance Pert.	58.79	77.71	20.24
Direction Pert.	43.61	74.87	56.08

Table 7: Navigation performance under various perturbation strategies applied at critical path steps.

**The model exhibits poor robustness against random disturbances.** We apply the following types of random perturbations to  $s_{\text{critical}}$ : **1) Road Perturbation:** Replace the original road name in  $s_{\text{critical}}$  with a different road (the direction should be modified accordingly); **2) Distance Perturbation:** Randomly adjust the distance at  $s_{\text{critical}}$ , ensuring that it does not exceed the remaining distance to the destination. **3) Direction Perturbation:** Invert the heading direction in  $s_{\text{critical}}$  (e.g., “east”  $\rightarrow$  “west”).

We select 10,000 cases with original correct predictions by the model for evaluation. The experimental results in Table 7 show that the model performs poorly when confronted with random perturbations, and its robustness varies across different types of disturbances. Specifically, in the road perturbation scenario, the model only has an 11.85% chance of selecting the first valid passable road, indicating that it does not have a precise understanding of its current location, or it lacks clarity on the available roads at its current position. This suggests that the model’s understanding of the road network is not coherent.

The experimental results in Table 7 show that the model performs poorly when confronted with random perturbations, and its robustness varies across different types of disturbances. Specifically, in the road perturbation scenario, the model only has an 11.85% chance of selecting the first valid passable road, indicating that it does not have a precise understanding of its current location, or it lacks clarity on the available roads at its current position. This suggests that the model’s understanding of the road network is not coherent.

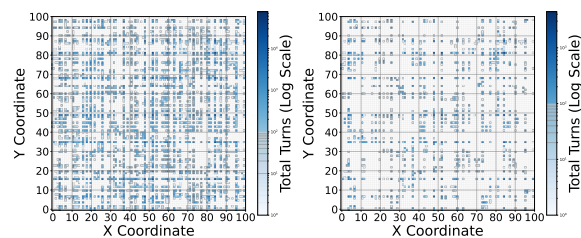


Figure 4: Heatmap of turning point frequencies in shortest paths. The left side shows the training data statistics, while the right side shows the test data statistics.

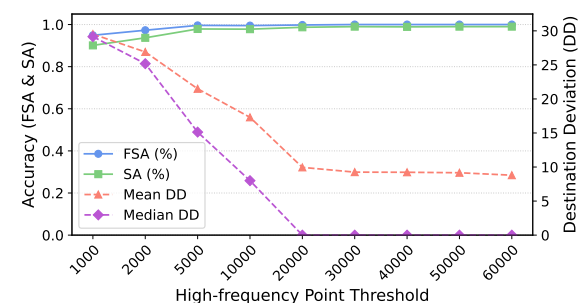


Figure 5: The model’s performance under different frequency thresholds. A higher frequency threshold indicates that the new starting point after the interference appears more frequently in the paths of the training data.

**The model’s robustness to disturbances largely depends on the distribution of the training data.** The results reveal a distinction between distance

(or direction) perturbation and road perturbation. When subjected to distance or direction perturbations, the model remains on the original high-speed road, such as those within  $s_{\text{critical}}$ . In contrast, road perturbations often randomly cause the model to deviate onto slower roads. The roads within  $s_{\text{critical}}$ , which are frequently selected in the shortest-path training data and feature higher-frequency entry and exit points along with corresponding turning patterns, are more robust to disturbances.

To analyze the impact of turning points, we count the frequency of points in the datasets where the direction of movement changes. High-frequency turning points generally correspond to transitions between high-speed and regular roads (Figure 4). We control the selection of  $p_{\text{perturb}}$  and  $p_{\text{target}}$  by ensuring the location frequency exceeds a threshold  $\tau$ . We select 8,464 cases and analyze how the model’s performance varies across different  $\tau$ .

As shown in Figure 5, the model’s performance improves with an increase in the frequency threshold ( $\tau$ ) for selecting  $p_{\text{perturb}}$  and  $p_{\text{target}}$ . This suggests that the model is more robust to perturbations at high-frequency turning points—enabling it to recover more effectively and reorient towards the correct destination. We provide further analysis and examples in Appendix D.2.

These results suggest that although the model exhibits a degree of robustness to perturbations, its recovery ability is highly dependent on the frequency of turning points encountered during training. This reliance implies that the model’s understanding of the road network is likely fragmented and localized, rather than comprehensive and global.

## 5 Related Work

**World Cognition** Previous studies have demonstrated that LLMs can encode real-world geospatial (Liétard et al., 2021; Treutlein et al., 2024) information and temporal (Gurnee and Tegmark) information within their internal representations. However, most of these studies use pretrained LLMs in non-anonymized experiments and have not fully explored the source of these capabilities. Concurrently, many works have focused on the ability of LLMs to learn and internalize rules of the physical world or form a “world cognition” of specific tasks from sequential data, such as in board games (Nanda et al., 2023; Li et al., 2023; Hazineh et al.) or simulated navigation (Jin and Rinard, 2024; Martorell, 2025; Vafa et al., 2024). Unlike

predicting the next token based on sequential data, our work focuses on whether LLMs can create a global understanding from natural language descriptions of local observations.

**Urban Space Reasoning** Some works focus on evaluating and enhancing the geospatial reasoning capabilities of LLMs (Feng et al., 2024a,b; Li et al., 2024). These studies enhance LLMs through knowledge training, external information, or tool use to adapt to spatial reasoning tasks in urban scenarios. In contrast, we focus on evaluating whether LLMs can reconstruct a global spatial understanding from local descriptions.

**Spatial Cognition** Spatial cognition capabilities are essential for LLMs to understand physical environments and perform tasks involving spatial reasoning. Many works focus on evaluating and enhancing the spatial cognition capabilities of LLMs (Mirzaee et al., 2021; Momennejad et al., 2023; Ramakrishnan et al., 2024), particularly in MLLM settings involving spatial memory and path reasoning (Yang et al., 2024b; Wu et al., 2024; Yu et al., 2025). Our work examines text-only LLMs’ ability to construct global spatial cognition from localized natural language observations, without relying on global information or coordinates.

## 6 Conclusion

Our study shows that LLMs can develop a global spatial understanding by training on local relative positions and shortest-path data. This is evident in their ability to generalize to unseen POI-pair relationships and in the strong alignment between latent representations and real-world geographic structures. These findings suggest that the model can autonomously build structured spatial cognition from unstructured language to support spatial reasoning. However, its limited robustness to navigation disturbances reveals the constraints of its understanding of road network structures.

## Limitations

Our study reveals that during the training process, the model develops an understanding of the global spatial distribution of Points of Interest (POIs) through the description of local relative relationships. However, how the model utilizes this spatial understanding when explicitly predicting positional relationships and shortest-path trajectories between unseen point pairs has not been fully analyzed. Our

experiments lack in-depth investigation into the internal mechanisms underlying the model’s explicit predictions of relative positions and shortest-path trajectories, which we will explore in future work.

In this study, to ensure the controllability of experiments, we used simplified synthetic data and constrained the experiments to a grid scenario. We present partial experimental results of real-world POI points in the appendix, and the synthetic data is reasonable and sufficient for our evaluation and analysis scenarios. However, how to verify such capabilities of LLMs in more complex scenarios and apply these capabilities of LLMs to real-world scenarios remains a problem we need to address in the future. Furthermore, our training process has impacted the model’s original general language capabilities. Given that our work is primarily analytical rather than enhancement-oriented, although the issue of catastrophic forgetting does exist, it does not affect our evaluations or conclusions within the context of our assessment scenarios. Nevertheless, in future downstream application scenarios, how to balance the model’s general capabilities with its internal spatial cognitive abilities remains an open research question.

## Ethics Statement

We hereby acknowledge that all authors of this work are aware of the provided ACL Code of Ethics and honor the code of conduct.

**Datasets Source** All studies in this work are based on a simulated, synthetically constructed dataset. The generated data is solely for model analysis research and contains no other usable information. To ensure privacy and ethical compliance, the dataset has been anonymized with placeholder names and contains no real-world information. As a result, the risk of sensitive information leakage is effectively eliminated.

**AI assistants** AI assistants (ChatGPT) were solely used to improve the grammatical structure of the text.

## References

- Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altschmidt, Sam Altman, Shyamal Anadkat, et al. 2023. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*.
- AI@Meta. 2024. [Llama 3 model card](#).

- Jinze Bai, Shuai Bai, Yunfei Chu, Zeyu Cui, Kai Dang, Xiaodong Deng, Yang Fan, Wenbin Ge, Yu Han, Fei Huang, et al. 2023. Qwen technical report. *arXiv preprint arXiv:2309.16609*.
- Mark Chen, Jerry Tworek, Heewoo Jun, Qiming Yuan, Henrique Ponde De Oliveira Pinto, Jared Kaplan, Harri Edwards, Yuri Burda, Nicholas Joseph, Greg Brockman, et al. 2021. Evaluating large language models trained on code. *arXiv preprint arXiv:2107.03374*.
- E. W. Dijkstra. 1959. A note on two problems in connexion with graphs. *Numerische Mathematik*, 1(1).
- Jie Feng, Yuwei Du, Tianhui Liu, Siqi Guo, Yuming Lin, and Yong Li. 2024a. Citygpt: Empowering urban spatial cognition of large language models. *arXiv preprint arXiv:2406.13948*.
- Jie Feng, Jun Zhang, Junbo Yan, Xin Zhang, Tianjian Ouyang, Tianhui Liu, Yuwei Du, Siqi Guo, and Yong Li. 2024b. Citybench: Evaluating the capabilities of large language model as world model. *CoRR*.
- Wes Gurnee and Max Tegmark. Language models represent space and time. In *The Twelfth International Conference on Learning Representations*.
- Dean Hazineh, Zechen Zhang, and Jeffrey Chiu. Linear latent world models in simple transformers: A case study on othello-gpt. In *Socially Responsible Language Modelling Research*.
- Charles Jin and Martin Rinard. 2024. Emergent representations of program semantics in language models trained on programs. In *Proceedings of the 41st International Conference on Machine Learning*, volume 235 of *Proceedings of Machine Learning Research*, pages 22160–22184. PMLR.
- Takeshi Kojima, Shixiang Shane Gu, Machel Reid, Yutaka Matsuo, and Yusuke Iwasawa. 2022. Large language models are zero-shot reasoners. *Advances in neural information processing systems*, 35:22199–22213.
- Kenneth Li, Aspen K. Hopkins, David Bau, Fernanda B. Viégas, Hanspeter Pfister, and Martin Wattenberg. 2023. Emergent world representations: Exploring a sequence model trained on a synthetic task. In *Proceedings of the 11th International Conference on Learning Representations (ICLR)*.
- Wenbin Li, Di Yao, Ruibo Zhao, Wenjie Chen, Zijie Xu, Chengxue Luo, Chang Gong, Quanliang Jing, Haining Tan, and Jingping Bi. 2024. Stbench: Assessing the ability of large language models in spatio-temporal analysis. *arXiv preprint arXiv:2406.19065*.
- Bastien Liétard, Mostafa Abdou, and Anders Søgaard. 2021. Do language models know the way to rome? In *Proceedings of the Fourth BlackboxNLP Workshop on Analyzing and Interpreting Neural Networks for NLP*, pages 510–517.
- Nicolas Martorell. 2025. From text to space: Mapping abstract spatial models in llms during a grid-world navigation task.
- Roshanak Mirzaee, Hossein Rajaby Faghihi, Qiang Ning, and Parisa Kordjamshidi. 2021. SPARTQA: A textual question answering benchmark for spatial reasoning. In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 4582–4598, Online. Association for Computational Linguistics.
- Ida Momennejad, Hosein Hasanbeig, Felipe Vieira Fruteri, Hiteshi Sharma, Nebojsa Jojic, Hamid Palangi, Robert Ness, and Jonathan Larson. 2023. Evaluating cognitive maps and planning in large language models with cogeval. *Advances in Neural Information Processing Systems*, 36:69736–69751.
- Neel Nanda, Andrew Lee, and Martin Wattenberg. 2023. Emergent linear representations in world models of self-supervised sequence models. In *Proceedings of the 6th BlackboxNLP Workshop: Analyzing and Interpreting Neural Networks for NLP*, pages 16–30, Singapore. Association for Computational Linguistics.
- Fabian Pedregosa, Gaël Varoquaux, Alexandre Gramfort, Vincent Michel, Bertrand Thirion, Olivier Grisel, Mathieu Blondel, Peter Prettenhofer, Ron Weiss, Vincent Dubourg, et al. 2011. Scikit-learn: Machine learning in python. *the Journal of machine Learning research*, 12:2825–2830.
- Santhosh Kumar Ramakrishnan, Erik Wijmans, Philipp Kraehenbuehl, and Vladlen Koltun. 2024. Does spatial cognition emerge in frontier models? *arXiv preprint arXiv:2410.06468*.
- Zhihong Shao, Peiyi Wang, Qihao Zhu, Runxin Xu, Junxiao Song, Xiao Bi, Haowei Zhang, Mingchuan Zhang, YK Li, Y Wu, et al. 2024. Deepseekmath: Pushing the limits of mathematical reasoning in open language models. *arXiv preprint arXiv:2402.03300*.
- Gemma Team, Aishwarya Kamath, Johan Ferret, Shreya Pathak, Nino Vieillard, Ramona Merhej, Sarah Perrin, Tatiana Matejovicova, Alexandre Ramé, Morgane Rivière, et al. 2025. Gemma 3 technical report. *arXiv preprint arXiv:2503.19786*.
- Johannes Treutlein, Dami Choi, Jan Betley, Samuel Marks, Cem Anil, Roger B Grosse, and Owain Evans. 2024. Connecting the dots: LLMs can infer and verbalize latent structure from disparate training data. *Advances in Neural Information Processing Systems*, 37:140667–140730.
- Keyon Vafa, Justin Chen, Ashesh Rambachan, Jon Kleinberg, and Sendhil Mullainathan. 2024. Evaluating the world model implicit in a generative model. *Advances in Neural Information Processing Systems*, 37:26941–26975.

- Wenshan Wu, Shaoguang Mao, Yadong Zhang, Yan Xia, Li Dong, Lei Cui, and Furu Wei. 2024. [Mind’s eye of llms: Visualization-of-thought elicits spatial reasoning in large language models](#). In *Neural Information Processing Systems*.
- An Yang, Beichen Zhang, Binyuan Hui, Bofei Gao, Bowen Yu, Chengpeng Li, Dayiheng Liu, Jianhong Tu, Jingren Zhou, Junyang Lin, et al. 2024a. Qwen2.5-math technical report: Toward mathematical expert model via self-improvement. *arXiv preprint arXiv:2409.12122*.
- Jihan Yang, Shusheng Yang, Anjali W Gupta, Rilyn Han, Li Fei-Fei, and Saining Xie. 2024b. Thinking in space: How multimodal large language models see, remember, and recall spaces. *arXiv preprint arXiv:2412.14171*.
- Qwen An Yang, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chengyuan Li, Dayiheng Liu, Fei Huang, Guanting Dong, Haoran Wei, Huan Lin, Jian Yang, Jianhong Tu, Jianwei Zhang, Jianxin Yang, Jiaxin Yang, Jingren Zhou, Junyang Lin, Kai Dang, Keming Lu, Keqin Bao, Kexin Yang, Le Yu, Mei Li, Mingfeng Xue, Pei Zhang, Qin Zhu, Rui Men, Runji Lin, Tianhao Li, Tingyu Xia, Xingzhang Ren, Xuancheng Ren, Yang Fan, Yang Su, Yi-Chao Zhang, Yunyang Wan, Yuqi Liu, Zeyu Cui, Zhenru Zhang, Zihan Qiu, Shanghaoran Quan, and Zekun Wang. 2024c. [Qwen2.5 technical report](#). *ArXiv*, abs/2412.15115.
- Dazhou Yu, Riyang Bao, Gengchen Mai, and Liang Zhao. 2025. Spatial-rag: Spatial retrieval augmented generation for real-world spatial reasoning questions. *arXiv preprint arXiv:2502.18470*.
- Yaowei Zheng, Richong Zhang, Junhao Zhang, Yanhan Ye, and Zheyang Luo. 2024. [LlamaFactory: Unified efficient fine-tuning of 100+ language models](#). In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 3: System Demonstrations)*, pages 400–410, Bangkok, Thailand. Association for Computational Linguistics.

## A Notation Table

Definition	
<i>Task Formulation</i>	
$\mathcal{G}$	A graph where intersections are nodes, roads remain unchanged, and the average travel speed is used as the edge weight.
$p_i$	Names of Points of Interest.
$r_i$	Names of roads.
MODEL <sub>per</sub>	LLM trained on data describing the relative positional relationships between POIs.
MODEL <sub>nav</sub>	LLM trained on data describing the shortest path trajectories.
<i>Metrics</i>	
<b>MSE</b>	Mean Squared Error, a metric quantifying the average squared difference between predictions and actual values.
<b>MRPE</b>	Mean Relative Percentage Error, a metric quantifying the average relative percentage difference between predictions and actual values.
<b>MAE</b>	Mean Absolute Error, a metric quantifying the average absolute difference between predictions and actual values.
<b>RMSE</b>	Root Mean Squared Error, a metric quantifying the square root of the average squared difference between predictions and actual values.
<b>R<sup>2</sup></b>	R-squared, a metric quantifying the proportion of the variance in the dependent variable that is predictable from the independent variable(s).
<b>Spearman</b>	Spearman correlation coefficient, a metric measuring the strength and direction of the monotonic relationship between two variables.
<b>Pearson</b>	Pearson correlation coefficient, a metric measuring the strength and direction of the linear relationship between two variables.
<b>FD</b>	Fréchet Distance, a metric that measures the similarity between two curves by considering the location and ordering of points.

Table 8: The notation table.

In Table 8, we list the notations and abbreviations in this paper, together with their definitions.

## B Training Parameters

**LLM Training** For the continual pre-training of the LLM, we use  $4 \times$  A800 80G GPUs with a batch size of 128, a learning rate of  $1.0e^{-4}$ , and a warmup ratio of 0.1, training for 10 epochs. Additionally, we designate the POI names  $P = \{p_i\}_{i=1}^{1024}$  and road names  $R = \{r_i\}_{i=1}^{200}$  as special tokens.

For the SFT of the LLM, we train on a single A800 80G GPU with a batch size of 256, a learning rate of  $3.0e^{-5}$ , a warmup ratio of 0.1, and train for 10 epochs. We use Llama-Factory as our training framework (Zheng et al., 2024).

**Probe** We use the MLPRegression model from scikit-learn (Pedregosa et al., 2011). The MLP probe we use consists of two hidden layers, with 128 and 64 neurons, and ReLU activation functions. The model is trained using the Adam optimizer with an initial learning rate of 0.001, and L2 regularization (alpha = 0.0001) with adaptive learning rate adjustment. The maximum number of training epochs is set to 500, and early stopping is enabled based on validation set performance (patience = 100 epochs), with a validation set proportion of 10%. The batch size is adjusted automatically during training, and data is shuffled before each epoch to improve generalization. All models and tools are publicly available for research purposes.

## C Experimental Details in Modeling Spatial Cognition

### C.1 POI Distribution

The spatial distribution of POI points in Section 2 is shown in the Figure 6.

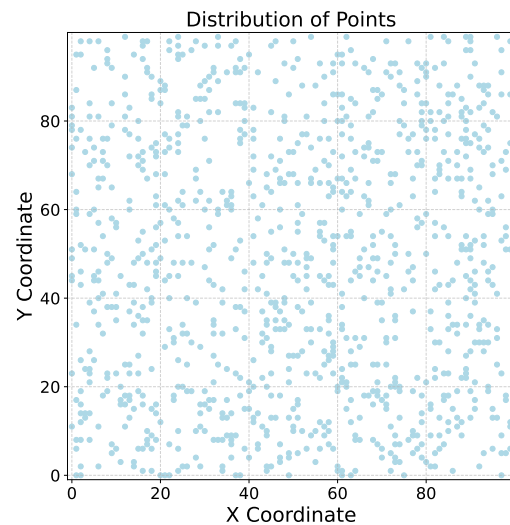


Figure 6: POIs Distribution.

### C.2 Data Size

In our main experiment, a total of 1,047,552 pairwise relational pairs can be formed for all POIs. These pairs are split into a training set and an evaluation set at an 8:2 ratio. After data augmentation using templates, the training set contains 5,028,250 data entries, with a total size of approximately 554 MB.

In terms of the number of trajectories, the training set covers 806,202 POI pair relationships. Com-

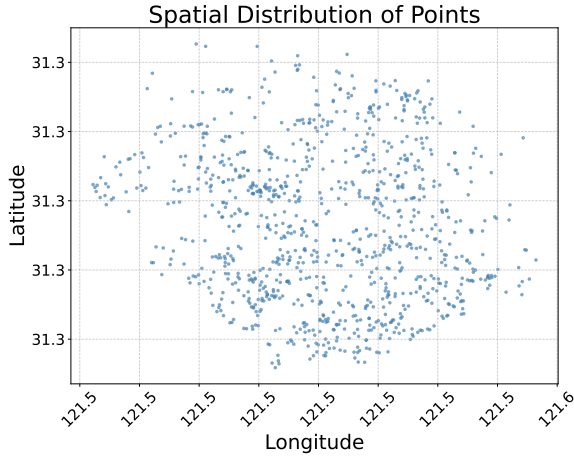


Figure 7: Real-World POIs Distribution.

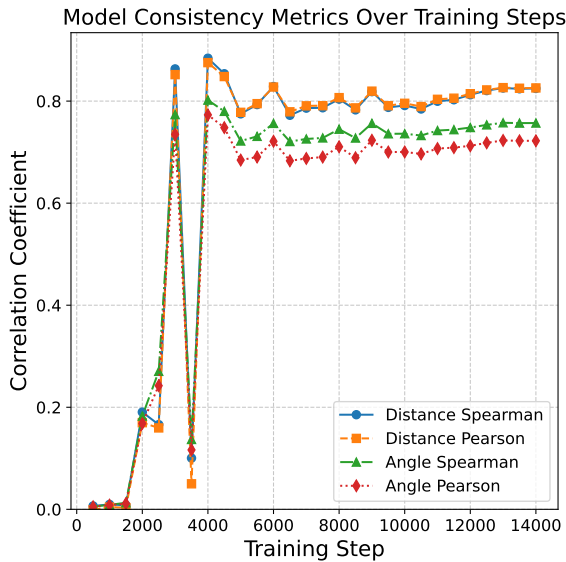


Figure 8: Consistency of POI hidden state vectors with actual spatial locations across training steps.

binning with different templates, 3,224,808 training data entries are generated; these training data have an average word count of 46.37, an average total distance of 67.13 kilometers, and an average number of passing roads of 5.20. The evaluation set includes 241,350 POI pair relationships, and its data has an average word count of 46.40, an average total distance of 67.29 kilometers, and an average number of passing roads of 5.20.

### C.3 Additional Experiment

**POI-in-Area Prediction** We design a simple spatial reasoning task that is inconsistent with the form of the training data in Section 2. We train the  $\text{MODEL}_{per}$  through supervised fine-tuning to determine whether a specific POI lies within a given

region. We consider two types of region descriptions: 1) a circular region defined by a central POI and a given radius; 2) a triangular region formed by three POIs. The LLM is required to provide a “yes” or “no” answer.

Additionally, we reserve a quarter of the POI points in the Map region, which are not included in the region descriptions of the SFT training data and are only used for evaluation. The remaining POIs are randomly sampled and divided into training and testing sets. We directly use prediction accuracy for evaluation.

POI Type	Circle (%)	Triangle (%)
Included (8:2)	98.8	96.2
Excluded (8:2)	97.1	97.9
Included (6:4)	98.9	97.8
Excluded (6:4)	96.1	95.8

Table 9: Prediction accuracy for POI-in-Area experiment.

**Real World POIs** The representation used in our synthetic data is universal and transferable. Real-world geographic data can be represented using our method and used for training, with no substantial differences between synthetic and real-world data.

**The focus of our work is to evaluate whether LLMs can construct global cognition from discrete local descriptions; thus, using synthetic data is appropriate here.** Building spatial cognition in practical scenarios and accomplishing related downstream tasks will be the focus of future work.

Also, collecting real-world data is challenging, especially the shortest path between two POIs, as the shortest path depends not only on distance but also on road conditions of each segment. In our synthetic dataset, we design a road weighting mechanism to simulate real-world road conditions. This weight represents the average driving speed of each segment. For routes with the same straight line distance, driving speeds may vary due to factors such as road roughness or curvature.

To enhance the realism and generalizability of the experiments, we sample 1,000 real-world POIs, represented by their geographic coordinates (latitude and longitude). As with the synthetic data, we anonymize these POI points, compute the pairwise Euclidean distances and azimuths, and split the dataset into training and testing sets (80/20).

The LLM training parameters remain consistent with those used for the synthetic data. The prediction accuracy for distances and azimuths on unseen POI pairs is shown in Table 10.

Distance		Azimuth	
MRPE (%) ↓	R <sup>2</sup> ↑	MRPE (%) ↓	Spearman ↑
0.30	1.00	0.53	1.00

Table 10: Prediction performance for distance and azimuth on unseen POI pairs in real-world scenarios.

The experimental results indicate that in more complex real-world scenarios, the model can also accurately model global positional cognition based on local relative positional relationships. We do not conduct experiments on the shortest path in real-world scenarios. This is because shortest path data is often difficult to collect in real-world settings, and our synthetic data simulates traffic conditions on real roads through weights, which is sufficient for our evaluation scenarios.

## D Experimental Details in Modeling Spatial Navigation

### D.1 Metric Calculations

**Start-End Deviation (SED)** : evaluates the spatial accuracy of the predicted path description by computing the Euclidean distance between predicted and ground truth coordinates at both the start and end points. The predicted trajectory is reconstructed by simulating the movement along a parsed sequence of road-based navigation steps using map information. The final metric is reported as a tuple: Start Deviation (SD) and End Deviation (ED). Detailed computation logic is provided in Algorithm 1.

**Valid Road Proportion (VRP)** : measures the proportion of valid road choices at each step of the predicted path description. The path is parsed into a sequence of steps, and for each step, the algorithm checks if the road and direction are valid according to the map’s connectivity and direction rules. The final metric, VRP, is the ratio of valid steps to the total steps in the path description. If no steps are described, the VRP is defined as 0. Detailed computation logic is provided in Algorithm 2.

**Shortest Path Accuracy (SPA)** : measures the proportion of cases where the model-generated trajectory exactly matches the ground truth shortest

path.

### D.2 Case Study

**Failure Analysis** In terms of distance and azimuth prediction, the model demonstrates high accuracy, with most errors occurring in the shortest path prediction task, especially in the presence of perturbations.

To better understand the failure modes of the model in shortest path prediction, we conducted an error analysis. We categorized prediction errors into three types: 1) start point errors, 2) intermediate path errors, and 3) end point errors. Since end point errors are always a consequence of one of the first two types, we do not report them separately.

The breakdown of errors on the test set (in terms of error count / total number of test cases) is as follows:

- Start point errors: 917 / 39800
- Intermediate path errors: 5656 / 39800

In intermediate path errors, we record the step at which the first error occurs. The distribution is as follows:

Step	1	2	3	4	5	6	7	8	9
Error Count	1819	2254	1049	367	113	39	13	1	1

Table 11: Distribution of the step where the first error occurred in intermediate path errors

We further categorize the causes of intermediate path errors into the following types:

- Direction errors: 894
- Road name errors: 37 (cases where the direction is correct but the road name is incorrect)
- Distance errors: 4725

These results indicate that most errors stem from a single incorrect step in the intermediate path (errors mainly occur in the early steps, primarily because the average number of steps across all cases is 5.2).

It is worth noting that most intermediate path errors are caused by incorrect distance predictions, accounting for 4725 out of 5656 cases.

**Disturbance Case** Figure 12 demonstrates the performance of LLM in handling intermediate disturbances under different turning point frequencies. As the frequency increases, LLM exhibits stronger robustness against disturbances and can reach the

final destination after being disturbed. When the frequency is low, the model is more prone to output interruptions (*e.g.*, not knowing where to go).

### D.3 Additional Experiment

Model	Distance % ↓	Azimuth % ↓
Perception-MODEL <sub>nav</sub>	3.08	5.52
Base-MODEL <sub>nav</sub>	12.03	13.84

Table 12: Evaluation results for distance and azimuth prediction, evaluated using MRPE.

**The model remains capable of performing explicit spatial relationship prediction.** To assess whether the model directly trained on path data can still understand the relative positional relationships between POIs, we fine-tune it with supervised training to predict the distance and azimuth between POI pairs. We use 200 POIs to construct the test set, while the remaining POIs are used to generate the training data (randomly sample 100,000 cases).

The results in Table 12 show that training the base model on shortest-path trajectories (Base-MODEL<sub>nav</sub>) allows it to capture the relative spatial relationships between POI pairs, achieving reasonable performance in both distance and azimuth prediction, with MRPE values of 12.03% and 13.84%, respectively. This suggests that, even without directly relying on local distance and azimuth information between POI pairs, the model is still able to leverage shortest-path trajectories to build a certain level of global spatial perception. This also indicates that shortest-path trajectories, as a topologically structured data format, are effective in constructing an understanding of spatial layout.

## E Additional Experiments and Results

### E.1 Training Strategy

Training Strategy	Distance		Azimuth	
	MRPE ↓	R <sup>2</sup> ↑	MRPE ↓	Spearman ↑
CPT	0.11	1.00	0.79	1.00
SFT	0.003	1.00	0.025	1.00

Table 13: The performance of the model’s prediction of distance and azimuth for unseen POI pairs under different training strategies.

Our primary experiments adopt a continual pre-training approach for LLM training. In addition to this, we explore the use of SFT for

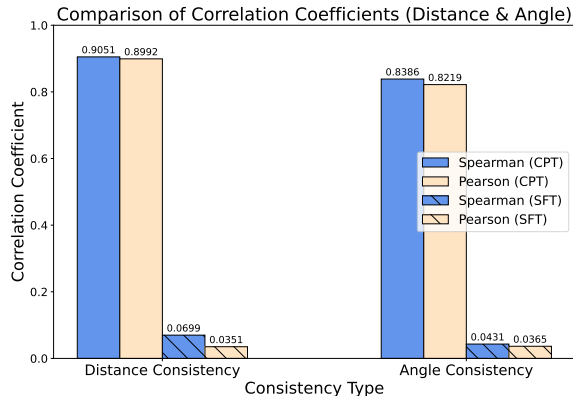


Figure 9: Consistency of POI point last hidden state vector with actual spatial location in terms of distance and angle under different training strategies.

training MODEL<sub>per</sub> and MODEL<sub>nav</sub>. For training MODEL<sub>per</sub>, we retain the question-answer format data from the original complete dataset and adopt an 80/20 split for training and test sets. For MODEL<sub>nav</sub>, we follow the **Bridged Exposure** strategy. We evaluate whether the LLM trained with SFT can perform explicit predictions and construct cognitive representations in the latent space.

We conduct training using 4×A800 80G GPUs, with a batch size of 512 and a learning rate set to 3.0e-5. The LLM is trained for 10 epochs.

**Spatial Perception** The experimental results for evaluating Spatial Cognition are shown in Table 13, Table 14, Table 15 and Figure 9.

Experimental results show that while SFT-trained LLM outperform CPT-trained LLM in distance and azimuth prediction accuracy, they exhibit weaker latent spatial cognition, as evidenced by blurred awareness of absolute coordinates in hidden states and poor alignment between latent vector distributions and actual spatial layouts.

This result is expected, as the POI name tokens in the SFT training process do not directly contribute to the loss calculation. Consequently, their embeddings are not explicitly optimized, leading to a lack of structured distribution in the latent space. This highlights the importance of continual pre-training for fostering deeper internal representations. At the same time, it suggests that a well-structured latent distribution of individual POIs is not strictly necessary for predicting relative relationships between unseen POI pairs.

**Spatial Navigation** The experimental results for evaluating Spatial Navigation are shown in Ta-

Training Strategy	X			Y			Euclidean Distance	
	MSE ↓	MAE ↓	R <sup>2</sup> ↑	MSE ↓	MAE ↓	R <sup>2</sup> ↑	Mean ↓	Std. ↓
Base	887.76	25.99	-0.01	878.72	25.10	-0.10	39.19	15.18
CPT	1.16	0.78	1.00	0.91	0.71	1.00	1.18	0.82
SFT	406.66	15.41	0.46	373.35	14.23	0.53	23.10	15.69

Table 14: Performance of the MLP probe in predicting the absolute coordinates of POIs from the LLM’s last hidden states under different training strategies.

Training Strategy	Distance		Azimuth	
	MAE (km)	R <sup>2</sup>	MAE (°)	Spearman
Base	14.90	0.03	39.12	0.62
CPT	0.85	1.00	3.49	0.98
SFT	31.62	-2.92	66.48	0.38

Table 15: Latent spatial composition evaluation. An MLP predicts distance and azimuth between POI pairs using their concatenated hidden states.

ble 16 and Table 17.

In addition, we further train the continual pre-trained model  $MODEL_{per}$  using the sft approach for the shortest path task, and evaluate its robustness against disturbances. The experimental results are shown in Table 18 and Figure 10.

Experimental results show that *Cognition-ModelTwo* trained via SFT exhibits robustness comparable to that of the CPT-trained counterpart, with both being influenced by the training data distribution—performing better at critical points with larger thresholds. Meanwhile, when facing random disturbances, the SFT-trained model reaches destinations closer to the target on average, but demonstrates a significantly lower proportion of selecting valid roads at each step.

## E.2 Model Architecture and Scale

To investigate the impact of architecture and parameter scale on models’ spatial cognition, in addition to Qwen2.5-0.5B used in the main experiments, we further examine the performance of Qwen2.5-1.5B, Qwen2.5-3B, LLaMA-3.2-1B (AI@Meta, 2024), Gemma-3-1B (Team et al., 2025) and Qwen2.5-Math-1.5B (Yang et al., 2024a).

**Spatial Perception** The results are shown in Table 24, Table 23 and Table 25. Experiments show that for the Qwen2.5 series models, as the model parameter scale increases, no significant improvement is observed in the explicit prediction of distance and azimuth, nor in the probing accuracy of absolute coordinates. Even when the model pa-

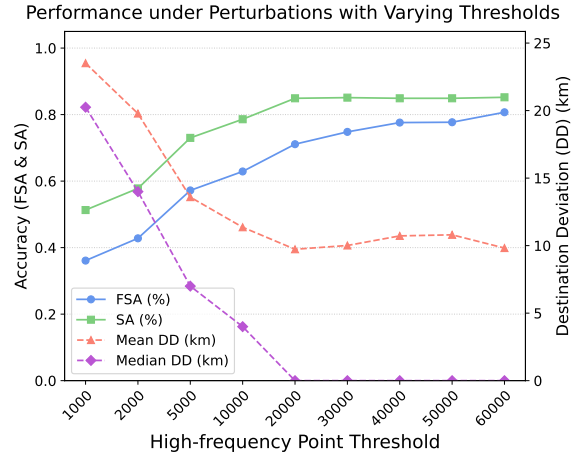


Figure 10: Performance metrics (FSA, SA, Mean/Median DD) versus high-frequency point thresholds. Left y-axis: FSA/SA; Right y-axis: DD (km).

rameter scale is small (0.5B), it already achieves high accuracy. In addition, models with different architectures (LLaMA, Gemma3) also demonstrate highly accurate modeling cognition of relative positions and absolute coordinates, presenting consistent experimental conclusions.

**Spatial Navigation** The results are shown in Table 26 and Figure 11. The experimental results show that for the Qwen2.5 series models, as the scale of the model parameter increases, the prediction accuracy of the shortest path improves (89.0% → 89.9% → 91.9%), but the robustness against path interference does not improve. Moreover, Gemma-3 models also achieve consistently strong performance on the shortest path navigation task, whereas LLaMA models perform poorly in learning local path information and completing shortest path navigation, exhibiting notable bias in identifying the starting point.

It is worth noting that although a larger model size does not lead to better prediction accuracy in some indicators, the prediction accuracy and error of all models fall within a favorable range, and there is no

Training Strategy	Accuracy				Consistency		
	SPD ↓	EPD ↓	VRP (↑%)	SPA (↑%)	VMR (↑1.0)	VCS (↑1.0)	FD (↓0.0)
CPT	0.06	0.48	96.07	83.63	1.00	1.00	0.91
SFT	0.02	0.02	99.65	97.34	1.00	1.00	0.11

Table 16: Performance of different training settings on shortest path prediction between POIs in  $P_{\text{heldout}}$ .

Model	X			Y			Euclidean Distance	
	MSE ↓	MAE ↓	R <sup>2</sup> ↑	MSE ↓	MAE ↓	R <sup>2</sup> ↑	Mean ↓	Std. ↓
<i>Absolute Coordinate Probing</i>								
Base Model	887.76	25.99	-0.01	878.72	25.10	-0.10	39.19	15.18
Cognition-CPT	8.53	2.16	0.99	10.21	2.40	0.99	3.54	2.49
Base-CPT	100.75	7.08	0.89	85.52	7.13	0.89	11.29	7.67
Cognition-SFT	13.05	2.89	0.98	12.88	2.76	0.99	4.39	3.84
Base-SFT	630.21	20.83	0.25	659.85	21.14	0.25	32.55	15.19
<i>Step-wise Coordinates Probing</i>								
Base Model	713.44	19.76	0.05	621.05	18.39	0.17	30.39	20.30
Cognition-CPT	6.51	1.84	0.99	6.96	1.94	0.99	3.01	2.10
Base-CPT	22.60	2.89	0.97	21.98	2.90	0.97	4.72	4.71
Cognition-SFT	11.97	2.53	0.98	12.78	2.56	0.98	4.07	3.56
Base-SFT	39.01	3.91	0.95	80.64	5.13	0.89	7.50	5.21

Table 17: Performance of the MLP probe in predicting the absolute coordinates of POIs and dynamic position coordinates at each step of the generated navigation path from the LLM’s last hidden states.

Method	FSA (%)	SA (%)	DD (km)
No Pert.	100.00	100.00	0.00
Road Pert.	8.14	52.31	12.42
Distance Pert.	14.95	60.29	9.46
Direction Pert.	6.01	59.53	40.89

Table 18: Evaluation Results for Different Types of Perturbations Trained via SFT.

significant prediction deviation among the models. The reason for this phenomenon may be that we use the same training parameters for models of different sizes without conducting separate parameter adjustment and optimization.

### E.3 Linear vs. Non-linear Probe

**Setup** We use the LinearRegression model from scikit-learn. It relies on a direct mathematical solution to find the best-fit line, and we used its default configuration. For the non-linear probe, we use the same MLP configuration as in the main experiment.

**Results** We use  $\text{MODEL}_{\text{per}}$  and  $\text{Base-MODEL}_{\text{nav}}$  to compare linear and non-linear probes in several experiments involving probing. The experimental results are shown in the Table 22, Table 19.

**Conclusion** The results in Table 22 demonstrate that a linear probe can map hidden states to actual coordinate values, indicating the presence of linearly accessible coordinate information within

Probe Type	Distance		Azimuth	
	MAE (km)	R <sup>2</sup>	MAE (°)	Spearman
Non-linear	0.85	1.00	3.49	0.98
Linear	17.89	0.20	51.94	0.78

Table 19: Latent spatial composition evaluation. An MLP predicts distance and azimuth between POI pairs using their concatenated hidden states.

the hidden representations of the LLM. However, non-linear regression achieves higher prediction accuracy. Furthermore, in the LLM trained on shortest-path trajectory data, the performance of the linear probe deteriorates significantly, with the average Euclidean distance increasing from 3.01 to 18.68. This suggests that non-linear probes are better suited for capturing position information in more complex tasks.

The experimental results in Table 19 show that when performing regression to predict distance and azimuth by combining the hidden states of two POIs, the linear probe performs poorly (R<sup>2</sup> of only 0.20 for distance prediction). This suggests that we cannot achieve combined prediction through simple linear regression, which may also be related to how we process the two POI vectors (*e.g.*, concatenation).

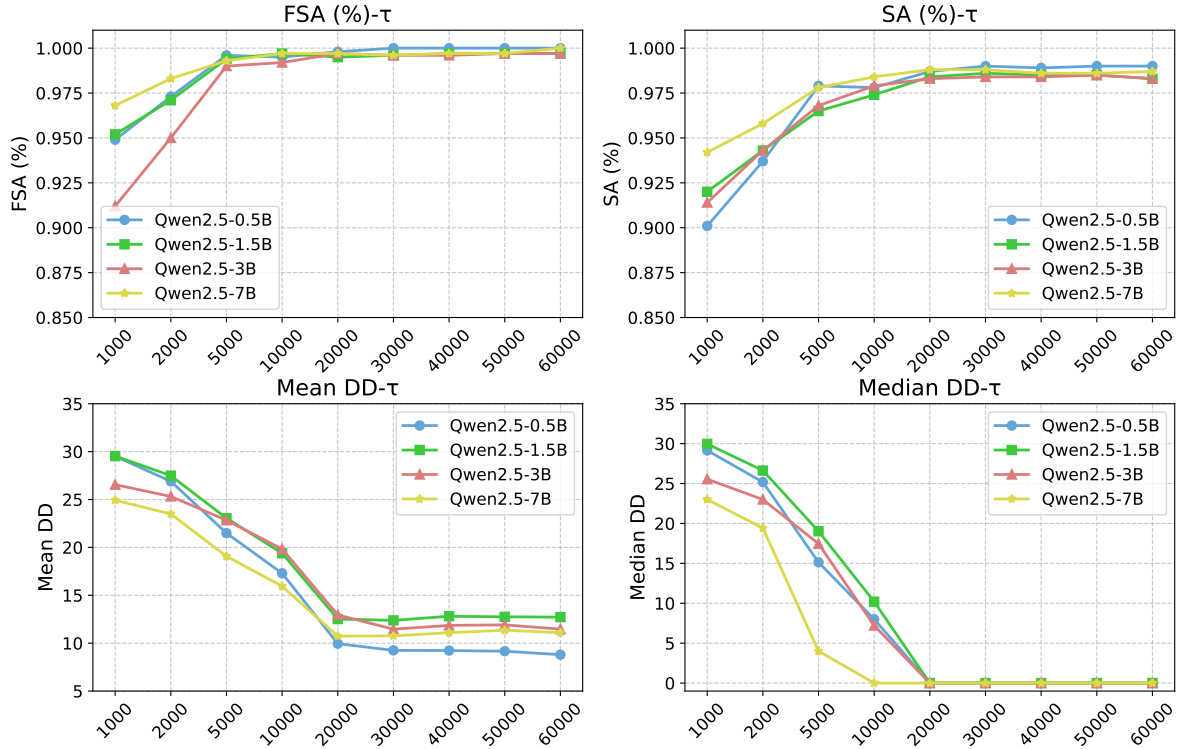


Figure 11: The robustness performance of models with different parameter scales when facing path interference.

#### Data Format

At an azimuth of 30 degrees from  $p_i$ ,  $p_j$  is located 1000 meters away.

$p_j$  lies 1000 meters from  $p_i$  at an azimuth of 30 degrees.

The azimuth from  $p_i$  to  $p_j$  is 30 degrees, and the separation is 1000 meters.

Q: How far is  $p_j$  from  $p_i$ ?  
A: 1000 meters.

Q: In what direction does  $p_j$  lie relative to  $p_i$ ?  
A: An azimuth of 30 degrees.

Q: What is the direction and separation between  $p_i$  and  $p_j$ ?  
A: An azimuth of 30 degrees and a distance of 1000 meters.

Table 20: An Alternative Template for Training and Evaluation Data of Positional Relationship Description.

Template	Distance		Azimuth	
	MRPE (%) ↓	R <sup>2</sup> ↑	MRPE (%) ↓	Spearman ↑
type 1	0.11	1.00	0.79	1.00
type 2	1.10	1.00	0.92	1.00

Table 21: The model’s prediction performance under different data construction templates: Type 1 represents the original template used in the main experiment, and Type 2 represents the replaced template.

#### E.4 Data Construction Template

**Setup** In addition to the data construction template adopted in the main experiment, we also experiment with other forms of templates to explore the impact of data construction templates on model performance. For the training, prediction, and evaluation of distance and azimuth, the templates we use are shown in Table 20.

**Results** We all adopt an 8:2 split ratio between the training set and the evaluation set. The experimental results of the two different data construction templates are shown in Table 21. In addition, after training the model using the replaced template, we still attempt to use the original template as the input for model evaluation, specifically the question “What is the distance from  $p_i$  to  $p_j$ ?”. The model’s mean relative prediction error (MRPE) for distance prediction remains only 1.30%.

**Conclusion** The experimental results show that after replacing with more diverse templates, the prediction errors of the model are still controlled within a very small range (1.1%), which indicates that the templates have little impact on the model’s construction of such spatial cognitive ability. Moreover, when using an evaluation method different from the templates, the performance of the model

Probe Type	X			Y			Euclidean Distance	
	MSE ↓	MAE ↓	R <sup>2</sup> ↑	MSE ↓	MAE ↓	R <sup>2</sup> ↑	Mean ↓	Std. ↓
<i>Absolute Coordinate Probing</i>								
Non-linear	1.16	0.78	1.00	0.91	0.71	1.00	1.18	0.82
Linear	21.18	3.61	0.97	12.70	2.75	0.99	4.99	2.99
<i>Step-wise Coordinates Probing</i>								
Non-linear	6.51	1.84	0.99	6.96	1.94	0.99	3.01	2.60
Linear	238.65	11.97	0.68	228.98	11.73	0.69	18.68	10.86

Table 22: Performance of the MLP probe in predicting the absolute coordinates of POIs and dynamic position coordinates at each step of the generated navigation path from the LLM’s last hidden states.

Model	X			Y			Euclidean Distance	
	MSE ↓	MAE ↓	R <sup>2</sup> ↑	MSE ↓	MAE ↓	R <sup>2</sup> ↑	Mean ↓	Std. ↓
Qwen2.5-0.5B	1.16	0.78	1.00	0.91	0.71	1.00	1.18	0.82
Qwen2.5-1.5B	6.83	1.96	0.99	3.40	1.47	1.00	2.73	1.66
Qwen2.5-3B	5.84	1.79	0.99	4.72	1.71	0.99	2.90	1.75
Qwen2.5-7B	7.39	1.98	0.99	7.61	2.00	0.99	3.25	2.14
LlaMA-3.2-1B-Base	984.84	26.57	-0.11	1066.15	28.21	-0.28	47.13	16.87
LlaMA-3.2-1B	5.71	1.94	0.99	6.97	1.99	0.99	3.07	1.81
LlaMA-3.2-3B-Base	948.04	25.99	-0.21	1183.55	29.06	-0.43	45.93	17.34
LlaMA-3.2-3B	2.94	1.33	1.00	6.43	1.73	0.99	2.54	1.42

Table 23: Performance of the MLP probe in predicting the absolute coordinates of POIs from the LLM’s last hidden states under different models.

Model	Distance		Azimuth	
	MRPE ↓	R <sup>2</sup> ↑	MRPE ↓	Spearman ↑
Qwen2.5-0.5B	0.11	1.00	0.79	1.00
Qwen2.5-1.5B	0.28	1.00	1.30	0.99
Qwen2.5-3B	0.11	1.00	0.89	1.00
Qwen2.5-7B	0.21	1.00	0.58	1.00
LlaMA-3.2-1B	1.71	1.00	3.99	0.98
LlaMA-3.2-3B	0.75	1.00	1.54	1.00
Gemma3-1B	0.07	1.00	0.96	1.00
Qwen2.5-Math-1.5B	0.25	1.00	1.60	1.00

Model	Distance		Azimuth	
	MAE (km)	R <sup>2</sup>	MAE (°)	Spearman
Qwen2.5-0.5B	0.85	1.00	3.49	0.98
Qwen2.5-1.5B	1.61	0.99	5.81	0.97
Qwen2.5-3B	0.84	1.00	3.81	0.98
Qwen2.5-7B	1.26	0.99	4.13	0.98
LlaMA-3.2-1B	1.18	1.00	4.32	0.96

Table 24: The performance of the model’s prediction of distance and azimuth for unseen POI pairs under different models.

Table 25: Latent spatial composition evaluation. An MLP predicts distance and azimuth between POI pairs using their concatenated hidden states.

Model	Accuracy				Consistency		
	SPD ↓	EPD ↓	VRP (↑%)	SPA (↑%)	VMR (↑1.0)	VCS (↑1.0)	FD (↓0.0)
Qwen2.5-0.5B	0.07	0.47	97.5	89.0	1.00	1.00	0.81
Qwen2.5-1.5B	0.04	0.27	97.6	89.9	1.00	1.00	0.59
Qwen2.5-3B	0.03	0.24	98.0	91.9	1.00	1.00	0.49
Qwen2.5-7B	0.05	0.29	98.6	91.6	1.00	1.00	0.63
LlaMA-3.2-1B	32.18	1.16	96.5	27.4	1.05	0.74	23.30
LlaMA-3.2-3B	59.33	55.91	3.8	0.0	0.18	0.03	53.64
Gemma3-1B	0.04	0.24	98.5	91.9	1.00	1.00	0.68
Qwen2.5-Math-1.5B	0.05	0.35	97.8	88.3	1.00	1.00	0.72

Table 26: Performance of different training settings on shortest path prediction between POIs in  $P_{\text{heldout}}$ .

Paths after Interference under Different Thresholds

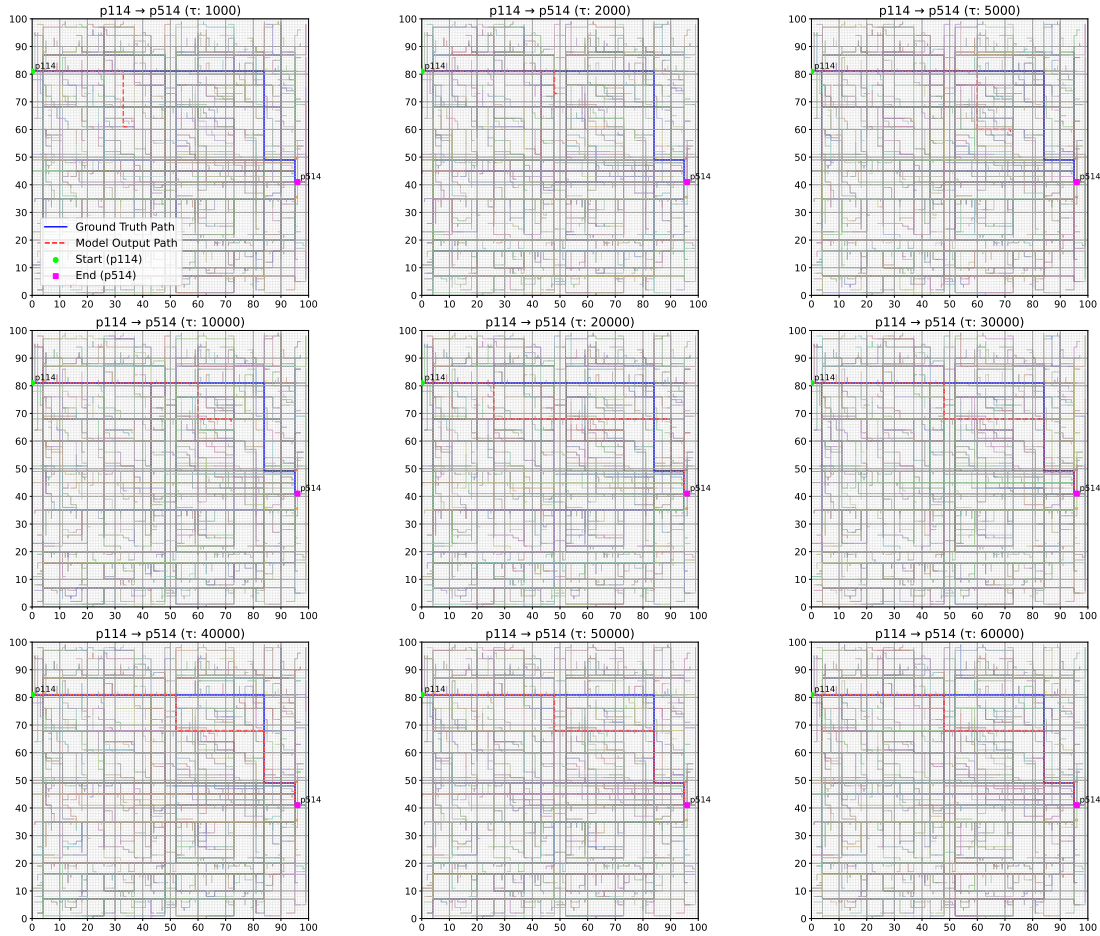


Figure 12: Case study on the model’s behavior under interference during navigation at different statistical frequencies.

is still not significantly affected (1.1%  $\rightarrow$  1.3%), which suggests that LLM has strong generalization ability and can understand texts with the same meaning but different forms.

### E.5 Ablation on Special Tokens and Random Seeds

**Setup** To validate the rationality of special token settings and the robustness of experimental results, we conduct ablation studies from two key perspectives. In the main experiment, POI and road names are treated as special tokens to facilitate probe training and hidden state extraction, preventing numeric identifiers (e.g.,  $p_{1000}$ ) from being split into multiple sub-tokens by the tokenizer. To assess the necessity of this setting, we compare model performance between configurations with and without special tokens. Additionally, to confirm the reliability of the observed high prediction accuracy, we replicate experiments using different random seeds

Settings	Distance		Azimuth	
	MRPE (%) $\downarrow$	R <sup>2</sup> $\uparrow$	MRPE (%) $\downarrow$	Spearman $\uparrow$
Base	0.11	1.00	0.79	1.00
w/o special token	0.93	1.00	1.97	1.00

Table 27: Impact of special token settings on distance and azimuth prediction based on Qwen2.5-0.5B.

to analyze result stability.

**Results** All ablation experiments adopt the same 8:2 train-evaluation split and consistent training configurations. The impact of special tokens on distance and azimuth prediction is presented in Table 27, while their influence on shortest path prediction is shown in Table 28. Model performance under different random seeds is reported in Table 29.

**Conclusion** Experimental results demonstrate that special tokens are not a prerequisite for the model to achieve high-performance spatial predic-

Settings	Accuracy			Consistency			
	SPD ↓	EPD ↓	VRP (↑%)	SPA (↑%)	VMR (↑1.0)	VCS (↑1.0)	FD (↓0.0)
Base	0.07	0.47	97.5	89.0	1.00	1.00	0.81
w/o special token	0.10	0.55	98.8	87.2	1.00	1.00	0.84

Table 28: Impact of special token settings on shortest path prediction based on Qwen2.5-0.5B.

Settings	Distance		Azimuth	
	MRPE (%) ↓	R <sup>2</sup> ↑	MRPE (%) ↓	Spearman ↑
seed=42	0.11	1.00	0.79	1.00
seed=123	0.89	0.98	1.76	1.00
seed=456	0.46	0.98	1.53	1.00

Table 29: Prediction stability of Qwen2.5-0.5B under different random seeds.

tion. Even without special tokens, the model maintains reliable prediction accuracy with only a slight increase in errors. Under different random initialization seeds, minor fluctuations in prediction errors are observed, but all results remain within a highly accurate range and consistent with the main experimental findings.

## F More Discussions

**Research Focus** The core focus of our research is not to propose a method for solving spatial reasoning tasks. Although external tools (*e.g.*, map algorithms) can solve spatial tasks with high accuracy, this study centers on a unique and complementary objective: evaluating whether LLMs can independently construct coherent global mental representations solely through fragmented local observations, just like humans—*i.e.*, their intrinsic spatial cognitive ability that does not rely on external tools. This motivation stems from the following core insight: human-like autonomy. Humans do rely on tools such as Google Maps, but in familiar environments, they can obtain location, route information, and other spatial data without external tools like maps (*e.g.*, finding a friend’s home in the neighborhood). Instead, they integrate local spatial memories (*e.g.*, "The café is 200 meters east of the park") to form a global mental map, thereby acquiring the ability for global location awareness and autonomous navigation. We aim to explore whether LLMs can also exhibit such autonomy, which would enable them to potentially move toward higher-level intelligence.

**Related Research** While numerous existing studies focus on large-scale spatial reasoning and the

formation of global spatial understanding, our research is entirely based on text-only input, unaffected by visual signals or predefined spatial structures. It drives spatial cognition and constructs global spatial perception solely through natural language descriptions of local relative relationships (*e.g.*, distance and direction), with a focus on the inherent spatial modeling capability of language itself. Its core lies in the implicit modeling ability to connect points into lines and lines into planes based on local cues. We delve into the model’s intrinsic implicit encoding mechanisms—such as verifying whether absolute coordinates can be decoded from hidden states, whether explicit relationships can be inferred by combining individual representations, and whether spatial states are dynamically encoded during navigation—and focus on the core question of “whether the model implicitly constructs spatial cognition.”

## G Other Statements

Our use of existing artifacts are consistent with their intended use, and we follow their license and terms.

---

**Algorithm 1** *SED*: Start-End Deviation Calculation

---

```
1: Input: Ground truth start coordinates  $P_{start\_gt}$ , Ground truth end coordinates  $P_{end\_gt}$ , LLM-
   generated textual path description  $\mathcal{A}$ , Map information  $\mathcal{M}_{map}$ 
2: Output: Start-End Deviation  $SED$   $\triangleright$  Euclidean distance between predicted and ground truth points
3:  $\mathcal{S} \leftarrow \text{ParsePathDescription}(\mathcal{A})$   $\triangleright$  Parse  $\mathcal{A}$  into sequence of steps  $\mathcal{S} = [(r_1, d_1, l_1), \dots, (r_n, d_n, l_n)]$ 
4: if  $|\mathcal{S}| < 2$  then
5:    $P_{start\_pred} \leftarrow P_{start\_gt}$   $\triangleright$  Use ground truth start if path description has fewer than 2 steps
6: else
7:   Let  $(r_1, d_1, l_1) = \mathcal{S}[1]$   $\triangleright$  First step details
8:   Let  $(r_2, d_2, l_2) = \mathcal{S}[2]$   $\triangleright$  Second step details
9:    $P_{intersect} \leftarrow \text{FindIntersection}(r_1, r_2, \mathcal{M}_{map})$   $\triangleright$  Find intersection of the first two roads (position
   after the first step)
10:  if  $P_{intersect}$  is valid then  $\triangleright$  Check if a valid intersection was found
11:     $P_{start\_pred} \leftarrow \text{MoveAlongRoad}(P_{intersect}, r_1, \text{Opposite}(d_1), l_1, \mathcal{M}_{map})$   $\triangleright$  Backtrack from
   intersection to estimate start
12:  else
13:     $P_{start\_pred} \leftarrow P_{start\_gt}$   $\triangleright$  Fallback to ground truth start if intersection is indeterminate
14:  end if
15: end if
16:  $P_{current} \leftarrow P_{start\_pred}$   $\triangleright$  Initialize current position
17: if  $|\mathcal{S}| > 0$  then  $\triangleright$  Simulate the path if steps exist
18:   for each step  $(r_i, d_i, l_i)$  in  $\mathcal{S}$  do
19:      $P_{current} \leftarrow \text{MoveAlongRoad}(P_{current}, r_i, d_i, l_i, \mathcal{M}_{map})$   $\triangleright$  Update position
20:   end for
21: end if
22:  $P_{end\_pred} \leftarrow P_{current}$   $\triangleright$  The final position is the predicted end position
23:  $SD \leftarrow \text{EuclideanDistance}(P_{start\_pred}, P_{start\_gt})$   $\triangleright$  Calculate Start Deviation
24:  $ED \leftarrow \text{EuclideanDistance}(P_{end\_pred}, P_{end\_gt})$   $\triangleright$  Calculate End Deviation
25: return  $(SD, ED)$   $\triangleright$  Return deviations at both start and end points
```

$\triangleright$  Helper Functions:

$\triangleright$  -  $\text{ParsePathDescription}(\mathcal{A})$ : Parses the textual path description  $\mathcal{A}$  into a structured list  $\mathcal{S}$  of tuples, where each tuple is  $(road\_id, direction, length)$ .

$\triangleright$  -  $\text{FindIntersection}(r_a, r_b, \mathcal{M}_{map})$ : Returns the geographic coordinates of the intersection between road segment  $r_a$  and road segment  $r_b$  based on  $\mathcal{M}_{map}$ . Returns an invalid/null state if no relevant intersection exists.

$\triangleright$  -  $\text{MoveAlongRoad}(P_{origin}, r, d, l, \mathcal{M}_{map})$ : Calculates the coordinates resulting from starting at  $P_{origin}$ , moving along road  $r$  in direction  $d$  for distance  $l$ , according to  $\mathcal{M}_{map}$ .

$\triangleright$  -  $\text{Opposite}(d)$ : Returns the direction directly opposite to  $d$  (e.g.,  $\text{Opposite}(\text{North}) = \text{South}$ ).

$\triangleright$  -  $\text{EuclideanDistance}(P_1, P_2)$ : Computes the L2 norm (straight-line distance)  $\|P_1 - P_2\|_2$ .

---

---

**Algorithm 2** *VRP*: Valid Road Proportion Calculation

---

```
1: Input: Ground truth start coordinates  $P_{start\_gt}$ , LLM-generated textual path description  $\mathcal{A}$ , Map
   information  $\mathcal{M}_{map}$ 
2: Output: Valid Road Proportion VRP  $\triangleright$  Proportion of steps choosing a valid next road
3:  $\mathcal{S} \leftarrow \text{ParsePathDescription}(\mathcal{A})$   $\triangleright$  Parse  $\mathcal{A}$  into sequence of steps  $\mathcal{S} = [(r_1, d_1, l_1), \dots, (r_n, d_n, l_n)]$ 
4: if  $|\mathcal{S}| < 2$  then
5:    $P_{start\_pred} \leftarrow P_{start\_gt}$   $\triangleright$  Use ground truth start if path description has fewer than 2 steps
6: else
7:   Let  $(r_1, d_1, l_1) = \mathcal{S}[1]$   $\triangleright$  First step details
8:   Let  $(r_2, d_2, l_2) = \mathcal{S}[2]$   $\triangleright$  Second step details
9:    $P_{intersect} \leftarrow \text{FindIntersection}(r_1, r_2, \mathcal{M}_{map})$   $\triangleright$  Find intersection of the first two roads (position
   after the first step)
10:  if  $P_{intersect}$  is valid then  $\triangleright$  Check if a valid intersection was found
11:     $P_{start\_pred} \leftarrow \text{MoveAlongRoad}(P_{intersect}, r_1, \text{Opposite}(d_1), l_1, \mathcal{M}_{map})$   $\triangleright$  Backtrack from
   intersection to estimate start
12:  else
13:     $P_{start\_pred} \leftarrow P_{start\_gt}$   $\triangleright$  Fallback to ground truth start if intersection is indeterminate
14:  end if
15: end if
16:  $P_{current} \leftarrow P_{start\_pred}$   $\triangleright$  Initialize current position
17:  $valid\_steps \leftarrow 0$   $\triangleright$  Initialize counter for valid road choices
18:  $total\_steps \leftarrow |\mathcal{S}|$   $\triangleright$  Total number of steps in the described path
19: if  $total\_steps > 0$  then  $\triangleright$  Simulate the path if steps exist
20:   for each step  $(r_i, d_i, l_i)$  in  $\mathcal{S}$  do
21:      $\mathcal{R}_{valid} \leftarrow \text{GetValidNextRoads}(P_{current}, \mathcal{M}_{map})$   $\triangleright$  Get set of valid (road_name, road_direct)
22:     if  $(r_i, d_i) \in \mathcal{R}_{valid}$  then  $\triangleright$  Check if the chosen road and direction are valid options
23:        $valid\_steps \leftarrow valid\_steps + 1$   $\triangleright$  Increment valid step count
24:     end if
25:      $P_{current} \leftarrow \text{MoveAlongRoad}(P_{current}, r_i, d_i, l_i, \mathcal{M}_{map})$   $\triangleright$  Update position
26:   end for
27: end if
28: if  $total\_steps == 0$  then
29:    $VRP \leftarrow 0$   $\triangleright$  Define VRP as 0 for empty paths
30: else
31:    $VRP \leftarrow valid\_steps / total\_steps$   $\triangleright$  Calculate the proportion of valid steps
32: end if
33: return VRP
```

$\triangleright$  Helper Functions:

$\triangleright$  -  $\text{ParsePathDescription}(\mathcal{A})$ : Parses the textual path description  $\mathcal{A}$  into a structured list  $\mathcal{S}$  of tuples  $(road\_id, direction, length)$ .

$\triangleright$  -  $\text{MoveAlongRoad}(P_{origin}, r, d, l, \mathcal{M}_{map})$ : Calculates coordinates after moving from  $P_{origin}$  along road  $r$  in direction  $d$  for distance  $l$ .

$\triangleright$  -  $\text{Opposite}(d)$ : Returns the direction opposite to  $d$ .

$\triangleright$  -  $\text{GetValidNextRoads}(P_{pos}, \mathcal{M}_{map})$ : Returns a set of valid next moves as  $(road\_id, direction)$  tuples accessible from position  $P_{pos}$ . This considers connectivity and travel direction rules based on map data  $\mathcal{M}_{map}$ .

$\triangleright$  -  $\text{EuclideanDistance}(P_1, P_2)$ : Computes the L2 norm  $\|P_1 - P_2\|_2$ . (Included for consistency, though not used in VRP calculation itself).

---