

From Human Cognition to Neural Activations: Probing the Computational Primitives of Spatial Reasoning in LLMs

Anonymous ACL submission

Abstract

As spatial intelligence becomes an increasingly important capability for foundation models, it remains unclear whether large language models’ (LLMs) performance on spatial reasoning benchmarks reflects structured internal spatial representations or reliance on linguistic heuristics. We address this question from a mechanistic perspective by examining how spatial information is internally represented and used. Drawing on computational theories of human spatial cognition, we decompose spatial reasoning into three primitives, relational composition, representational transformation, and stateful spatial updating, and design controlled task families for each. We evaluate multilingual LLMs in English, Chinese, and Arabic under single pass inference, and analyze internal representations using linear probing, sparse autoencoder based feature analysis, and causal interventions. We find that task relevant spatial information is encoded in intermediate layers and can causally influence behavior, but these representations are transient, fragmented across task families, and weakly integrated into final predictions. Cross linguistic analysis further reveals mechanistic degeneracy, where similar behavioral performance arises from distinct internal pathways. Overall, our results suggest that current LLMs exhibit limited and context dependent spatial representations rather than robust, general purpose spatial reasoning, highlighting the need for mechanistic evaluation beyond benchmark accuracy.¹

1 Introduction

Large language models (LLMs) and vision–language models (VLMs) have achieved rapid progress in reasoning, planning, and interactive decision making, and are increasingly deployed in settings that require spatial understanding, such as instruction following, navigation, embodied

¹Our code is available at <It will be published after the paper is accepted.>

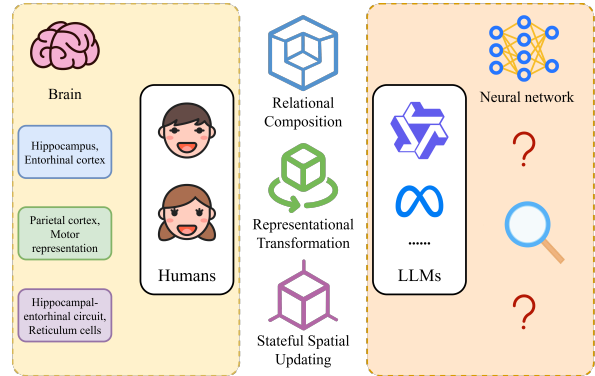


Figure 1: From human spatial cognition to spatial representations in large language models.

interaction, and robotics(Wang et al., 2024a; Brohan et al., 2023; Driess et al., 2023; Guo et al., 2024). At the same time, the study of world models has emphasized the importance of structured internal representations for prediction and control(Yi et al., 2018). These developments raise a fundamental question: do foundation models acquire genuine spatial representations, or do they succeed through surface-level linguistic regularities?(Yamada et al., 2023; Du et al., 2024; Li et al., 2024a; Du et al., 2024)

Recent work has begun to address this question using spatial reasoning benchmarks, reporting steady performance gains with model scale and instruction tuning(Yamada et al., 2023; Li et al., 2024a; Wei et al., 2021). However, benchmark accuracy alone provides limited insight into underlying mechanisms. Correct input–output behavior may arise from linguistic heuristics, memorization, or prompt-induced reasoning strategies, making it unclear whether models rely on structured spatial representations or shallow correlations(Holliday et al., 2014; Geirhos et al., 2020; Xie et al., 2024; Wei et al., 2022).

This stands in contrast to research on human spatial cognition, where spatial ability is characterized

068	not only by behavior but by well-studied internal	2 Related Work	118
069	mechanisms, such as cognitive maps, mental rota-	2.1 Spatial Cognition and Computational	119
070	tion, and path integration. These representations	Mechanisms	120
071	explain generalization, robustness, and systematic	Spatial cognition has long been studied in cogni-	121
072	errors. By comparison, most studies of spatial rea-	tive psychology and neuroscience as a core com-	122
073	soning in large neural models do not directly ex-	ponent of human intelligence. Classical psycho-	123
074	amine internal representations, leaving the basis of	metric frameworks distinguish abilities such as	124
075	their spatial behavior largely opaque(Yamada et al.,	spatial perception, mental rotation, and visualiza-	125
076	2023; Li et al., 2024a; Hewitt and Manning, 2019;	tion(Ekstrom et al., 1976; Shepard and Metzler,	126
077	Elazar et al., 2020).	1971), while later work emphasizes that spatial be-	127
078	The lack of mechanistic understanding has both	havior arises from distinct internal representations	128
079	practical and conceptual implications. As spa-	and computational mechanisms rather than a single	129
080	tial reasoning becomes increasingly embedded in	faculty(Eckardt, 1980; Burgess, 2008).	130
081	downstream systems, limited insight into inter-	Research on cognitive maps shows that humans	131
082	nal representations hinders robustness assessment,	and animals construct structured, allocentric repre-	132
083	interpretability, and principled comparison with	sentations that support relational inference beyond	133
084	biological cognition(Yamins and DiCarlo, 2016).	immediate perception(Tolman, 1948). Studies of	134
085	Without representational evidence, behavioral sim-	mental rotation and perspective taking demonstrate	135
086	ilarities between models and humans remain spec-	that these representations can undergo continuous	136
087	ulative(Firestone, 2020).	geometric transformations(Shepard and Metzler,	137
088	In this work, we argue that progress on spatial	1971), while navigation research highlights mech-	138
089	reasoning in foundation models requires moving	anisms for maintaining and updating spatial state	139
090	beyond benchmark-centric evaluation toward mech-	over time(Etienne and Jeffery, 2004). Together,	140
091	anistic analysis. Rather than asking only whether	these findings motivate decomposing spatial cog-	141
092	a model produces correct answers, we investigate	nition into a small set of core computational prim-	142
093	whether it develops structured, compositional rep-	itives. Our work adopts this computational per-	143
094	resentations that play a functional role in behavior.	spective to guide task design, rather than directly	144
095	Drawing inspiration from cognitive science, we	replicating psychometric categories.	145
096	decompose spatial ability into a small set of core	2.2 Spatial Reasoning in NLP and Large	146
097	computational primitives and design controlled task	Language Models	147
098	families that isolate these primitives under standard	Spatial reasoning in NLP has traditionally focused	148
099	inference settings, without eliciting explicit reason-	on interpreting spatial relations expressed in lan-	149
100	ing traces.	guage, such as prepositions, relative positions, and	150
101	We introduce three families of spatial tasks: rela-	instructions. With the rise of large language mod-	151
102	tional spatial reasoning, perspective transformation,	els (LLMs), recent work has increasingly relied	152
103	and spatial program execution. To disentangle spa-	on benchmarks that evaluate spatial reasoning via	153
104	tial representations from linguistic form, we con-	output accuracy(Suzgun et al., 2022).	154
105	struct parallel tasks in English, Chinese, and Arabic.	These studies show that LLMs can achieve non-	155
106	We analyze both behavior and internal represen-	trivial performance on text-based spatial tasks, es-	156
107	tations using probing, sparse autoencoder-based	pecially with chain-of-thought prompting or ex-	157
108	feature analysis, and causal interventions. This	PLICIT intermediate reasoning. However, they also	158
109	framework enables us to localize spatial informa-	report limitations such as sensitivity to surface	159
110	tion within models, characterize representational	form variation, degraded performance under in-	160
111	differences across task families, and assess their de-	creased compositionality, and poor generaliza-	161
112	pendence on language. By grounding spatial eval-	tion(Wei et al., 2022; Li et al., 2024b). Despite	162
113	uation in cognitive theory and mechanistic inter-	these observations, most work remains behavioral,	163
114	pretability, our study clarifies the nature and limits	offering limited insight into whether correct outputs	164
115	of spatial ability in foundation models, and pro-	reflect structured spatial representations or shal-	165
116	vides a foundation for moving beyond benchmark	low linguistic heuristics(Geirhos et al., 2020). In	166
117	accuracy.	contrast, our work targets internal representations	167

168	rather than output accuracy alone.		
169	2.3 Spatial Reasoning Benchmarks for Large		
170	Language Models		
171	A growing body of work proposes benchmarks		
172	to evaluate spatial reasoning in LLMs, focusing		
173	on spatial relations, compositional descriptions,		
174	and navigation-like inference in text(Yamada et al.,		
175	2023; Rizvi et al., 2024). Synthetic and semi-		
176	natural datasets such as StepGame(Shi et al., 2022)		
177	and SpatialEval(Wang et al., 2024b) are widely		
178	used to probe multi-step spatial inference, typically		
179	using accuracy-based evaluation.		
180	Across benchmarks, recent studies report im-		
181	proved performance with prompting strategies, but		
182	also reveal substantial weaknesses, including sen-		
183	sitivity to paraphrasing, compositional complex-		
184	ity, and distribution shift(Sharma, 2023). Sim-		
185	ilar patterns appear in multimodal and vision-		
186	language models, where explicit reference struc-		
187	tures or visualization-style reasoning can boost per-		
188	formance but are often tightly coupled to the evalu-		
189	ation setup(Wu et al., 2024; Liao et al., 2024). Sur-		
190	veys highlight a persistent gap between behavioral		
191	success and evidence for structured, compositional		
192	spatial representations(Liu et al., 2025; Zheng et al.,		
193	2025). Overall, existing benchmarks largely em-		
194	phasize output behavior, providing limited insight		
195	into underlying representations.		
196	2.4 Probing and Mechanistic Analysis of		
197	LLMs		
198	A parallel line of research investigates what in-		
199	formation is encoded inside neural language mod-		
200	els(Maennel et al., 2020; Krause and Reimann,		
201	2024; Yamada et al., 2023; Nanda and Bloom,		
202	2022; Bloom et al., 2024). Linear probing tests		
203	whether variables can be decoded from hidden		
204	states, while more recent work analyzes represen-		
205	tation geometry or identifies interpretable features		
206	using sparse autoencoders(Hewitt and Manning,		
207	2019; Hewitt and Liang, 2019; Gupta et al., 2023;		
208	Raghu et al., 2017; Yan et al., 2024). These meth-		
209	ods also enable causal interventions that modify		
210	internal representations and measure behavioral ef-		
211	fects(Meng et al., 2022).		
212	While mechanistic analysis has been applied to		
213	many linguistic and semantic phenomena, its appli-		
214	cation to spatial reasoning remains limited(Olsson		
215	et al., 2022; Elhage et al., 2022). Our work builds		
216	on these interpretability tools to analyze spatial		
217	representations in LLMs, aiming to move beyond		
	correlational evidence toward functional and causal	218	
	understanding.	219	
	In summary, prior work on spatial reasoning in	220	
	LLMs has largely focused on benchmark-based	221	
	behavioral evaluation, while mechanistic inter-	222	
	pretability studies have rarely targeted spatial cog-	223	
	niton. Our work bridges these lines by grounding	224	
	task design in computational theories of spatial cog-	225	
	niton and applying mechanistic analysis to probe	226	
	whether LLMs implement core spatial primitives.	227	
	3 Task Taxonomy: A Computational	228	
	Decomposition of Spatial Ability	229	
	To distinguish genuine spatial computation from	230	
	surface-level linguistic pattern matching in large	231	
	language models, we introduce a task taxonomy	232	
	grounded in a computational decomposition of spa-	233	
	tial ability. Rather than organizing tasks by domain	234	
	or surface form, we classify them by the minimal	235	
	computational primitives required for correct in-	236	
	ference, drawing on established findings in human	237	
	spatial cognition.	238	
	We identify three irreducible components of spa-	239	
	tial computation: (1) relational composition, (2)	240	
	representational transformation, and (3) stateful	241	
	spatial updating. Each component defines one task	242	
	family, forming a compact and mechanistically ac-	243	
	cessible test suite for probing internal spatial rep-	244	
	resentations beyond linguistic heuristics. Figure 2	245	
	summarizes the three families with representative	246	
	examples.	247	
	3.1 Design Principles	248	
	All task families follow four shared principles. Ab-	249	
	straction: tasks use abstract entities and relations,	250	
	minimizing reliance on world knowledge. Compo-	251	
	sitionality: solving each task requires integrating	252	
	multiple spatial constraints or operations. Parame-	253	
	terizability: difficulty is systematically controlled	254	
	via factors such as entity count, steps, or dimen-	255	
	sionality. Mechanistic Accessibility: intermediate	256	
	spatial variables are explicitly defined, enabling	257	
	probing and causal intervention.	258	
	All tasks are evaluated under standard single-	259	
	pass inference, without eliciting explicit reason-	260	
	ing traces, and are generated using controlled rule-	261	
	-based procedures to ensure computational equiva-	262	
	lence across tasks and languages (Appendix A.1).	263	
	3.2 Relational Spatial Reasoning	264	
	This task family targets relational composition, re-	265	
	quiring models to construct a globally consistent	266	

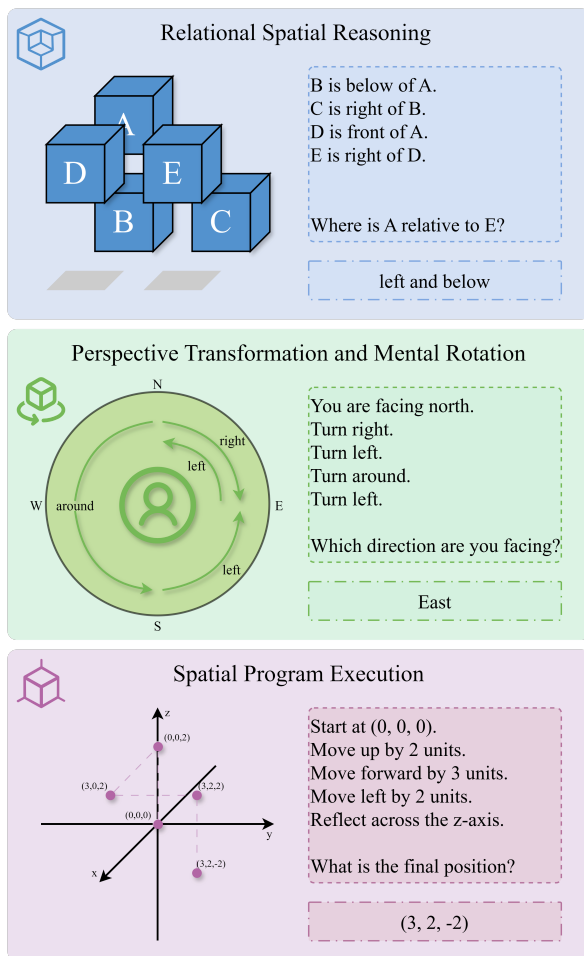


Figure 2: Illustration of the three spatial task families proposed in this work, with examples shown in English for clarity.

spatial structure from multiple pairwise relations (e.g., *A left of B*, *B above C*) and infer relations between indirectly connected entities. No metric information or procedural sequence is provided; the challenge lies entirely in relational structure building. We include both 2D and 3D variants, probing whether models form structured spatial representations rather than relying on memorized linguistic templates.

3.3 Perspective Transformation and Mental Rotation

This family isolates representational transformation. Given an initial spatial configuration or reference frame, models must apply one or more global transformations—such as rotations, reflections, or viewpoint changes—and report the resulting orientation or relation. Unlike relational reasoning, structure is assumed and must be transformed while preserving internal consistency. Both self-centered

and multi-agent perspective-taking variants are included, directly probing equivariance under geometric operations.

3.4 Spatial Program Execution

The third family targets stateful spatial updating. Each instance specifies an initial position and a sequence of movement or transformation commands; the model must compute the final position after executing all steps. Performance depends on maintaining and updating a latent spatial state over time, making errors cumulative and diagnostic of state-tracking failures. This family abstracts the computational core of navigation and path integration and is especially amenable to causal analysis via intervention on intermediate states.

3.5 Alignment with Human Spatial Abilities

Although computationally motivated, the taxonomy aligns naturally with distinctions in human spatial cognition: relational reasoning with cognitive maps and spatial visualization, perspective transformation with mental rotation, and spatial program execution with dynamic updating in navigation. This alignment serves as an interpretive reference rather than a claim of equivalence.

3.6 Cross-Linguistic Task Construction

To disentangle spatial computation from language-specific cues, all task families are independently constructed in English, Chinese, and Arabic. Rather than direct translation, we ensure computational equivalence while allowing natural variation in surface realization, treating language as a controlled variable for analyzing whether spatial representations are shared or language-dependent.

4 Methodology and Experiments

4.1 Experimental Setup

Models. We evaluate spatial reasoning in two multilingual LLM families, Qwen2.5-7B (Yang et al., 2024) and Llama-3-8 (Dubey et al., 2024), including both base and instruction-tuned variants. All experiments use publicly available checkpoints without task-specific fine-tuning, ensuring that observed spatial behaviors reflect general representations rather than adaptation.

Inference Protocol. All tasks are evaluated using standard single-pass inference without chain-of-thought prompting or explicit reasoning traces. Models select answers in a multiple-choice format

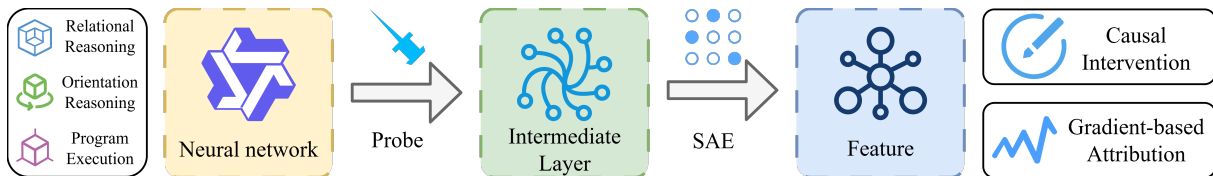


Figure 3: Overview of the proposed framework. We probe intermediate representations of a neural network, extract interpretable features using a sparse autoencoder, and analyze their roles via gradient-based attribution and causal interventions.

Task Family	Train	Test	Difficulty Range
Relational Reasoning	2 000	200	3–10
Orientation Reasoning	2 000	200	2–10
Program Execution	2 000	200	2–10
Total (per language)	6 000	600	–
Total (3 languages)	18 000	1 800	–

Table 1: Dataset statistics across task families. Each task family is constructed in English, Chinese, and Arabic with computational equivalence.

by emitting a single token. This design isolates implicit spatial representations encoded in activations, rather than reasoning strategies expressed in generated text.

Data and Evaluation. For each task family and language, we construct 2,000 training instances and 200 held-out test instances, uniformly distributed across difficulty levels. Test data are never used during probing or SAE training. Unless stated otherwise, results are reported on in-distribution test sets.

Metrics. We report behavioral accuracy and analyze performance as a function of task complexity. For representational analysis, we use linear probe R^2 scores, regression error (MAE/RMSE for continuous variables), and layer-wise trends to characterize where spatial information emerges.

Multilingual Design. All tasks are independently constructed in English, Chinese, and Arabic. Prompts are not translated; instead, computational equivalence is preserved while allowing natural linguistic variation. Models are evaluated across all languages without language-specific tuning.

Table 1 summarizes dataset statistics.

Performance does not degrade monotonically with task complexity. In relational reasoning and perspective transformation, accuracy exhibits non-monotonic fluctuations across step lengths, suggesting regime shifts in internal processing rather than

simple capacity limits. In contrast, spatial program execution shows more stable degradation, indicating that coordinate-based updating is comparatively more learnable from language statistics.

Cross-linguistic evaluation reveals moderate language dependence. English and Chinese show comparable performance overall, while Arabic lags most noticeably on relational reasoning tasks. Notably, spatial program execution exhibits the smallest cross-linguistic gap, suggesting greater language invariance for coordinate-based representations. As shown in Figure 3

4.2 Probing Spatial Representations

We train linear probes to decode task-relevant spatial variables from model hidden states, including relative position vectors (Task 1), orientation vectors (Task 2), and absolute spatial coordinates (Task 3). Probes are trained on the training split and evaluated on held-out test instances.

We train linear probes to decode task-relevant spatial variables from hidden states, including relative position vectors (Task 1), orientation vectors (Task 2), and absolute coordinates (Task 3). Probes are trained on training instances and evaluated on held-out test data. Across models and tasks, spatial information consistently peaks in intermediate layers and declines sharply toward the final layers (Figure 4). For Qwen2.5-7B-Instruct, relational reasoning reaches a maximum R^2 of 0.37 at layer 19, while spatial program execution peaks at $R^2 \approx 0.25$ around layer 16. Orientation variables are weakly decodable throughout ($R^2 < 0.15$). As shown in Table 3.

This inverted-U pattern indicates that spatial representations are constructed during intermediate processing but are not preserved into the final layers responsible for token prediction. Instruction-tuned models consistently show stronger spatial representations than base models, while cross-linguistic comparisons reveal similar layer-wise emergence patterns with varying representational strength.

Model	Lang	Relational	Orientation	Program Execution
Qwen2.5-7B-Base	EN	50.0 (100/200)	23.0 (46/200)	55.5 (111/200)
	ZH	39.5 (79/200)	23.5 (47/200)	58.0 (116/200)
	AR	39.5 (79/200)	24.5 (49/200)	47.0 (94/200)
Qwen2.5-7B-Instruct	EN	49.0 (98/200)	28.0 (56/200)	62.0 (124.0/200)
	ZH	47.5 (95/200)	28.0 (56/200)	52.5 (105/200)
	AR	31.0 (62/200)	26.0 (52/200)	57.5 (115/200)
Llama3-8B-Instruct	EN	4.0 (8/200)	2.0 (4/200)	2.0 (4/200)
	ZH	0.0 (0/200)	2.5 (5/200)	0.0 (0/200)
	AR	3.5 (7/200)	3.0 (6/200)	0.5 (1/200)

Table 2: Accuracy across task families and languages (%). Values are reported as accuracy with raw counts in parentheses. For Program Execution, accuracy may exceed 100% due to partial-credit scoring. Program Execution shows the strongest and most consistent performance across languages, while Orientation Reasoning performs near chance level. Arabic exhibits lower performance on Relational Reasoning, possibly reflecting language-specific encoding challenges.

Model	Lang	Relational Reasoning	Orientation Reasoning	Program Execution
Qwen2.5-7B-Base	EN	$R^2=.382 / L=25 / MAE=.55 / RMSE=.70$	$R^2=-.007 / L=0$	$R^2=.394 / L=26 / MAE=1.83 / RMSE=2.79$
	ZH	$R^2=.254 / L=22 / MAE=.59 / RMSE=.77$	$R^2=-.003 / L=0$	$R^2=.347 / L=24 / MAE=1.93 / RMSE=2.90$
	AR	$R^2=.272 / L=24 / MAE=.59 / RMSE=.76$	$R^2=-.011 / L=0$	$R^2=.347 / L=24 / MAE=1.93 / RMSE=2.90$
Qwen2.5-7B-Instruct	EN	$R^2=.366 / L=19 / MAE=.55 / RMSE=.71$	$R^2=-.008 / L=9$	$R^2=.402 / L=20 / MAE=1.86 / RMSE=2.80$
	ZH	$R^2=.243 / L=19 / MAE=.60 / RMSE=.78$	$R^2=-.006 / L=0$	$R^2=.457 / L=20 / MAE=1.65 / RMSE=2.67$
	AR	$R^2=.264 / L=20 / MAE=.59 / RMSE=.76$	$R^2=-.005 / L=0$	$R^2=.418 / L=20 / MAE=1.82 / RMSE=2.74$
Llama3-8B-Instruct	EN	$R^2=.305 / L=31 / MAE=.57 / RMSE=.75$	$R^2=-.001 / L=17$	$R^2=.168 / L=31 / MAE=2.14 / RMSE=3.27$
	ZH	$R^2=.207 / L=31 / MAE=.60 / RMSE=.80$	$R^2=-.003 / L=0$	$R^2=.156 / L=31 / MAE=2.19 / RMSE=3.27$
	AR	$R^2=.172 / L=31 / MAE=.63 / RMSE=.82$	$R^2=-.002 / L=4$	$R^2=.156 / L=31 / MAE=2.19 / RMSE=3.27$

Table 3: Best-layer probing results by model, language, and task family. We report coefficient of determination (R^2), best-performing layer (L), mean absolute error (MAE), and root mean squared error (RMSE).

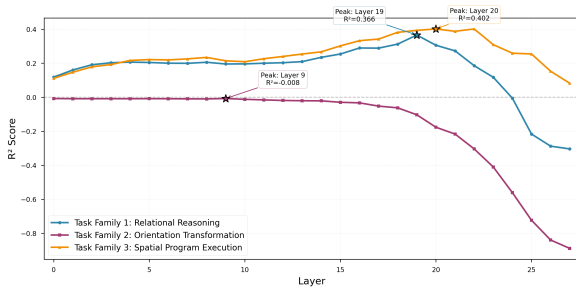


Figure 4: Layer-wise R^2 scores for spatial variable prediction across three task families (Qwen2.5-7B-Instruct, English). All tasks show mid-layer peaks followed by sharp declines in final layers. Task Family 1 and 3 demonstrate strong representational clarity (R^2 up to 0.37 and 0.40 respectively), while Task Family 2 shows minimal spatial encoding.

4.3 Sparse Autoencoder Analysis

To identify interpretable spatial features, we train sparse autoencoders (SAEs) on layers with peak probe performance for each task family. SAEs reveal a small subset of spatially selective features,

typically accounting for 3–5% of all discovered features. As shown in Figure 5.

Across tasks, we observe clear axis- or direction-selective features, including units specialized for cardinal orientations (Task 2), coordinate ranges (Task 3), and relational axes (Task 1). Gradient-based attribution shows that features with stronger spatial selectivity contribute disproportionately to spatial predictions.

Feature overlap across task families is limited (12–18%), suggesting that different forms of spatial computation rely on largely distinct feature sets. Cross-linguistic analysis reveals partial feature sharing based on activation frequency, but substantially lower overlap when features are ranked by causal attribution, indicating language-specific mechanistic pathways supporting functionally similar behavior.

4.4 Causal Interventions

We perform activation patching and SAE feature ablation to assess whether spatial representations

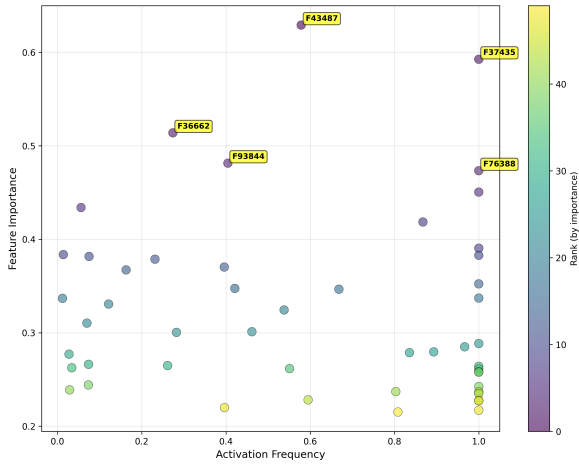


Figure 5: Feature importance versus activation frequency for Task Family 2 (Qwen2.5-7B-Instruct, Chinese). Important features are sparse and not aligned with activation frequency, indicating dissociation between usage and causal contribution.

causally influence model behavior.

In spatial program execution, patching intermediate-layer activations with counterfactual spatial states systematically shifts model outputs toward counterfactual trajectories, with strongest effects in layers 14–18. Interventions at final layers show minimal impact, despite being closest to output generation. As shown in Figure 6.

Complementary SAE feature ablation confirms functional specificity. Removing top spatial features reduces accuracy by 29% in orientation tasks and 14.5% in spatial program execution, while ablating non-spatial control features has negligible effect. These results establish that identified spatial representations are not merely decodable, but causally contribute to behavior.

5 Analyses and Discussion

5.1 What Kind of Spatial Representations Do LLMs Develop?

A central question of this work is whether large language models develop structured spatial representations analogous to cognitive maps, or whether spatial behavior arises from shallow linguistic heuristics. Our results suggest an intermediate regime: LLMs do encode spatial information, but these representations are conditional, fragmented, and weakly integrated into decision-making.

Linear probing reveals that task-relevant spatial variables are decodable from intermediate layers, with peak R^2 values up to 0.37–0.40 for relational reasoning and spatial program execution.

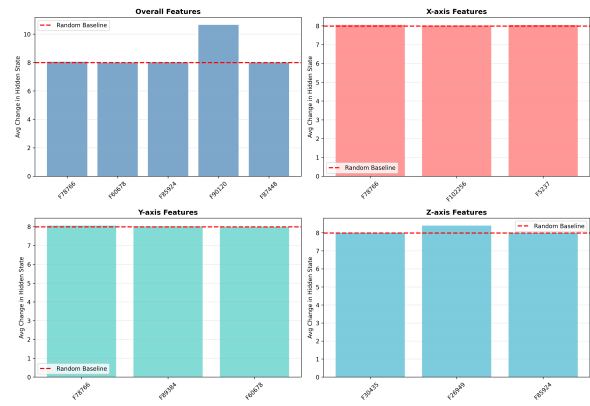


Figure 6: Intervention effects of spatial features across dimensions (Qwen2.5-7B-Instruct, Task Family 3, English). Bars show the average hidden-state change caused by intervening on spatially selective features, with the dashed line indicating a random-feature baseline.

These representations exhibit geometric coherence, supporting multi-axis integration rather than isolated symbolic associations. Crucially, causal interventions confirm functional relevance: activation patching at mid-layers reliably shifts model outputs toward counterfactual spatial trajectories.

At the same time, spatial representations are fragile and transient. Probe performance drops sharply in final layers, indicating that spatial structure constructed during intermediate processing is not preserved through output generation. Moreover, spatial features are largely task-specific, with limited overlap across relational, orientation, and coordinate-based reasoning. This fragmentation contrasts with biological cognitive maps, which support multiple spatial computations through a shared representational substrate.

Taken together, these findings indicate that LLMs exhibit *conditional and partial spatial encoding*: spatial structure is constructed when task format and linguistic cues make it accessible, but is neither persistent nor unified across tasks.

5.2 Failure Modes Across Spatial Tasks

Distinct failure patterns across task families further illuminate the limits of LLM spatial reasoning.

Relational reasoning: compositional collapse.

In relational tasks, accuracy degrades non-monotonically with chain length, with sharp drops at intermediate complexity. Errors frequently preserve one spatial axis while losing another, suggesting partial breakdowns in multi-constraint composition. Performance recovery at higher complexity

492	likely reflects a shift toward coarse heuristics rather	reasoning. Bottlenecks arise in propagating and	541
493	than robust relational inference.	integrating spatial information into final decision	542
494	Orientation transformation: representational	layers.	543
495	absence. Orientation tracking performs near	Implicit vs. explicit reasoning. Evaluating mod-	544
496	chance and exhibits near-zero probe R^2 across all	els without chain-of-thought reveals latent spatial	545
497	layers. Error patterns reveal strong distributional	structure, but this structure is fragile. We hypothe-	546
498	bias toward frequent directions (e.g., north, east),	size that explicit reasoning stabilizes spatial repre-	547
499	indicating reliance on linguistic priors rather than	sentations by externalizing intermediate states, act-	548
500	geometric state tracking. Unlike position or coord-	ing as a compensatory scaffold rather than merely	549
501	inates, heading direction does not appear to be	revealing hidden competence.	550
502	represented as a manipulable internal variable.		
503	Spatial programs: arithmetic dominance. Spa-	Task format matters. Performance differences	551
504	tial program execution shows the strongest behav-	across task families highlight the role of format.	552
505	ioral performance, but error analysis reveals that	Coordinate-based tasks benefit from arithmetic	553
506	a substantial fraction of failures are purely arith-	pathways, while relational and orientation tasks	554
507	metic. Once arithmetic errors are factored out, spa-	demand abstract composition that LLMs struggle	555
508	tial performance aligns more closely with relational	to sustain. High performance on one format does	556
509	reasoning. This suggests that success in coordinate-	not imply general spatial understanding.	557
510	-based tasks partly reflects numerical competence		
511	rather than dedicated spatial computation.		
512	5.3 Language Dependence of Spatial	6 Conclusion	558
513	Representations		
514	Our multilingual evaluation reveals a layered pat-	Large language models achieve strong performance	559
515	tern of language dependence. Behaviorally, per-	on spatial benchmarks, yet it remains unclear	560
516	formance varies by up to 14 points across lan-	whether this reflects genuine spatial reasoning or	561
517	guages, most prominently in relational reasoning.	reliance on linguistic heuristics. Motivated by cog-	562
518	Coordinate-based tasks show the smallest cross-	nitve neuroscience, we adopt a mechanistic per-	563
519	linguistic gap, consistent with numerical formats	spective and introduce a compact task taxonomy	564
520	being more language-invariant.	that isolates three core spatial primitives to analyze	565
521	At the representational level, spatial information	multilingual models across languages. We find	566
522	emerges at similar layers across languages, indicat-	that task relevant spatial information is encoded	567
523	ing a shared computational architecture. However,	in intermediate layers and can causally influence	568
524	SAE analysis shows limited overlap in causally im-	behavior, but these representations are fragile, tran-	569
525	portant features across languages. While different	sient, and highly task specific. Spatial structure	570
526	languages achieve comparable probe performance,	typically emerges in mid layers but is weakly pre-	571
527	they do so using distinct internal features, reflect-	served in final decisions, and different spatial tasks	572
528	ing <i>mechanistic degeneracy</i> : similar outputs supported	rely on largely disjoint internal features, providing	573
529	by different internal pathways.	limited evidence for unified spatial representations.	574
530	These findings challenge the assumption that	Cross linguistic analysis reveals mechanistic degen-	575
531	spatial reasoning in LLMs is inherently language-	eracy, where the emergence of spatial information	576
532	agnostic. Linguistic encoding substantially shapes	is largely language invariant but the internal path-	577
533	both the accessibility and implementation of spatial	ways supporting spatial computation vary across	578
534	computation.	languages, with greater invariance observed in co-	579
535	5.4 Implications for Spatial Reasoning in	ordinate based tasks. Overall, these findings place	580
536	LLMs	current LLMs between shallow pattern matching	581
537	Spatial representations are necessary but insuffi-	and robust spatial cognition as characterized in neu-	582
538	cient. The gap between representational strength	roscience, suggesting that more human like spatial	583
539	and behavioral accuracy indicates that encoding	reasoning will require approaches that better pre-	584
540	spatial variables alone does not guarantee reliable	serve and integrate spatial representations across	585
		processing stages and motivating mechanistic eval-	586
		uation beyond benchmark accuracy.	587

588 Limitations

589 This study analyzes spatial representations in two
590 multilingual language model families at the 7–8B
591 scale, and findings may differ for larger models or
592 alternative architectures. The proposed tasks isolate
593 three spatial primitives and do not cover more
594 complex spatial abilities such as planning or embodied
595 interaction. Experiments are limited to English,
596 Chinese, and Arabic, and spatial representations
597 may interact differently with other languages.
598 Causal conclusions rely on activation patching and
599 feature ablation, which provide partial causal
600 characterization, and representational analyses focus on
601 linearly decodable information, potentially missing
602 nonlinear or distributed structure. Finally, tasks are
603 abstract and text-only, limiting conclusions about
604 real-world or multimodal spatial reasoning.

605 References

606 Joseph Bloom, Curt Tigges, Anthony Duong, and David
607 Chanin. 2024. Saelens. [https://github.com/
608 decoderresearch/SAELens](https://github.com/decoderresearch/SAELens).

609 Anthony Brohan, Noah Brown, Justice Carbajal, Yevgen
610 Chebotar, Krzysztof Choromanski, Tianli Ding,
611 Danny Driess, Kumar Avinava Dubey, Chelsea Finn,
612 Peter R. Florence, Chuyuan Fu, Montse Gonzalez
613 Arenas, Keerthana Gopalakrishnan, Kehang Han,
614 Karol Hausman, Alexander Herzog, Jasmine Hsu,
615 Brian Ichter, Alex Irpan, and 27 others. 2023. [Rt-2:
616 Vision-language-action models transfer web knowl-
617 edge to robotic control](#). *ArXiv*, abs/2307.15818.

618 Neil Burgess. 2008. Spatial cognition and the brain.
619 *Annals of the New York Academy of Sciences*,
620 1124(1):77–97.

621 Danny Driess, F. Xia, Mehdi S. M. Sajjadi, Corey
622 Lynch, Aakanksha Chowdhery, Brian Ichter, Ayzaan
623 Wahid, Jonathan Tompson, Quan Ho Vuong, Tianhe
624 Yu, Wenlong Huang, Yevgen Chebotar, Pierre Ser-
625 manet, Daniel Duckworth, Sergey Levine, Vincent
626 Vanhoucke, Karol Hausman, Marc Toussaint, Klaus
627 Greff, and 3 others. 2023. [Palm-e: An embodied
628 multimodal language model](#). In *International Con-
629 ference on Machine Learning*.

630 Mengfei Du, Binhao Wu, Zejun Li, Xuanjing Huang,
631 and Zhongyu Wei. 2024. [Embspatial-bench: Bench-
632 marking spatial understanding for embodied tasks
633 with large vision-language models](#). In *Annual Meet-
634 ing of the Association for Computational Linguistics*.

635 Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey,
636 Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman,
637 Akhil Mathur, Alan Schelten, Amy Yang, Angela
638 Fan, Anirudh Goyal, Anthony S. Hartshorn, Aobo

639 Yang, Archi Mitra, Archie Sravankumar, Artem Ko-
640 renev, Arthur Hinsvark, Arun Rao, Aston Zhang, and
641 510 others. 2024. [The llama 3 herd of models](#).

642 Michael J. Eckardt. 1980. [The hippocampus as a cog-
643 nitive map](#). *Journal of Nervous and Mental Disease*,
644 168:191–192.

645 Ruth B. Ekstrom, John W. French, and Harry H. Harman.
646 1976. [Manual for kit of factor-referenced cognitive
647 tests](#).

648 Yanai Elazar, Shauli Ravfogel, Alon Jacovi, and Yoav
649 Goldberg. 2020. [Amnesic probing: Behavioral expla-
650 nation with amnesic counterfactuals](#). *Transactions of
651 the Association for Computational Linguistics*, 9:160–
652 175.

653 Nelson Elhage, Tristan Hume, Catherine Olsson,
654 Nicholas Schiefer, Tom Henighan, Shauna Kravec,
655 Zac Hatfield-Dodds, Robert Lasenby, Dawn Drain,
656 Carol Chen, and 1 others. 2022. [Toy models of su-
657 perposition](#). *arXiv preprint arXiv:2209.10652*.

658 A. S. Etienne and Kate J. Jeffery. 2004. [Path integration
659 in mammals](#). *Hippocampus*, 14.

660 Chaz Firestone. 2020. [Performance vs. competence in
661 human–machine comparisons](#). *Proceedings of the
662 National Academy of Sciences*, 117:26562 – 26571.

663 Robert Geirhos, Jörn-Henrik Jacobsen, Claudio
664 Michaelis, Richard S. Zemel, Wieland Brendel,
665 Matthias Bethge, and Felix Wichmann. 2020. [Short-
666 cut learning in deep neural networks](#). *Nature Ma-
667 chine Intelligence*, 2:665 – 673.

668 Taicheng Guo, Xiuying Chen, Yaqi Wang, Ruidi Chang,
669 Shichao Pei, N. Chawla, Olaf Wiest, and Xiangliang
670 Zhang. 2024. [Large language model based multi-
671 agents: A survey of progress and challenges](#). In
672 *International Joint Conference on Artificial Intelli-
673 gence*.

674 Sharut Gupta, Joshua Robinson, Derek Lim, Soledad
675 Villar, and Stefanie Jegelka. 2023. [Structuring rep-
676 resentation geometry with rotationally equivariant
677 contrastive learning](#). *ArXiv*, abs/2306.13924.

678 John Hewitt and Percy Liang. 2019. [Designing
679 and interpreting probes with control tasks](#). *ArXiv*,
680 abs/1909.03368.

681 John Hewitt and Christopher D. Manning. 2019. [A
682 structural probe for finding syntax in word representa-
683 tions](#). In *North American Chapter of the Association
684 for Computational Linguistics*.

685 Trenton W Holliday, Joanna R. Gautney, and Lukas
686 Friedl. 2014. [Right for the wrong reasons](#). *Current
687 Anthropology*, 55:696 – 724.

688 Renate Krause and Stefan Reimann. 2024. [Items or
689 relations - what do artificial neural networks learn?](#)
690 *ArXiv*, abs/2404.12401.

691	Fangjun Li, David C. Hogg, and Anthony G. Cohn. 2024a. Reframing spatial reasoning evaluation in language models: A real-world simulation benchmark for qualitative reasoning. <i>ArXiv</i> , abs/2405.15064.	Zhengxiang Shi, Qiang Zhang, and Aldo Lipani. 2022. Stepgame: A new benchmark for robust multi-hop spatial reasoning in texts. In <i>AAAI Conference on Artificial Intelligence</i> .	747 748 749 750
695	Zhiming Li, Yushi Cao, Xiufeng Xu, Junzhe Jiang, Xu Liu, Yon Shin Teo, Shang-Wei Lin, and Yang Liu. 2024b. Llms for relational reasoning: How far are we? In <i>Proceedings of the 1st International Workshop on Large Language Models for Code</i> , pages 119–126.	Mirac Suzgun, Nathan Scales, Nathanael Scharli, Sebastian Gehrmann, Yi Tay, Hyung Won Chung, Aakanksha Chowdhery, Quoc V. Le, Ed H. Chi, Denny Zhou, and Jason Wei. 2022. Challenging big-bench tasks and whether chain-of-thought can solve them. In <i>Annual Meeting of the Association for Computational Linguistics</i> .	751 752 753 754 755 756 757
701	Yuan-Hong Liao, Rafid Mahmood, Sanja Fidler, and David Acuna. 2024. Reasoning paths with reference objects elicit quantitative spatial reasoning in large vision-language models. In <i>Conference on Empirical Methods in Natural Language Processing</i> .	Edward Chace Tolman. 1948. Cognitive maps in rats and men. <i>Psychological review</i> , 55 4:189–208.	758 759
706	Weichen Liu, Qiyao Xue, Haoming Wang, Xiangyu Yin, Boyuan Yang, and Wei Gao. 2025. Spatial reasoning in multimodal large language models: A survey of tasks, benchmarks and methods.	Jiaqi Wang, Zihao Wu, Yiwei Li, Hanqi Jiang, Peng Shu, Enze Shi, Huawen Hu, Chong-Yi Ma, Yi-Hsueh Liu, Xuhui Wang, Yincheng Yao, Xuan Liu, Huaqin Zhao, Zheng Liu, Haixing Dai, Lin Zhao, Bao Ge, Xiang Li, Tianming Liu, and Shu Zhang. 2024a. Large language models for robotics: Opportunities, challenges, and perspectives. <i>ArXiv</i> , abs/2401.04334.	760 761 762 763 764 765 766
710	Hartmut Maennel, Ibrahim M. Alabdulmohsin, Ilya O. Tolstikhin, Robert J. N. Baldock, Olivier Bousquet, Sylvain Gelly, and Daniel Keysers. 2020. What do neural networks learn when trained with random labels? <i>ArXiv</i> , abs/2006.10455.	Jiayu Wang, Yifei Ming, Zhenmei Shi, Vibhav Vineet, Xin Wang, and Neel Joshi. 2024b. Is a picture worth a thousand words? delving into spatial reasoning for vision language models. <i>ArXiv</i> , abs/2406.14852.	767 768 769 770
715	Kevin Meng, David Bau, Alex Andonian, and Yonatan Belinkov. 2022. Locating and editing factual associations in gpt. In <i>Neural Information Processing Systems</i> .	Jason Wei, Maarten Bosma, Vincent Zhao, Kelvin Guu, Adams Wei Yu, Brian Lester, Nan Du, Andrew M. Dai, and Quoc V. Le. 2021. Finetuned language models are zero-shot learners. <i>ArXiv</i> , abs/2109.01652.	771 772 773 774
719	Neel Nanda and Joseph Bloom. 2022. Transformerlens. https://github.com/TransformerLensOrg/TransformerLens .	Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Ed H. Chi, F. Xia, Quoc Le, and Denny Zhou. 2022. Chain of thought prompting elicits reasoning in large language models. <i>ArXiv</i> , abs/2201.11903.	775 776 777 778
722	Catherine Olsson, Nelson Elhage, Neel Nanda, Nicholas Joseph, Nova Dassarma, T. J. Henighan, Benjamin Mann, Amanda Askell, Yuntao Bai, Anna Chen, Tom Conerly, Dawn Drain, Deep Ganguli, Zac Hatfield-Dodds, Danny Hernandez, Scott Johnston, Andy Jones, John Kernion, Liane Lovitt, and 7 others. 2022. In-context learning and induction heads. <i>ArXiv</i> , abs/2209.11895.	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. <i>Advances in Neural Information Processing Systems</i> 37.	779 780 781 782 783
730	Maithra Raghu, Justin Gilmer, Jason Yosinski, and Jascha Narain Sohl-Dickstein. 2017. Sycca: Singular vector canonical correlation analysis for deep learning dynamics and interpretability. In <i>Neural Information Processing Systems</i> .	Chulin Xie, Yangsibo Huang, Chiyuan Zhang, Da Yu, Xinyun Chen, Bill Yuchen Lin, Bo Li, Badih Ghazi, and Ravi Kumar. 2024. On memorization of large language models in logical reasoning. <i>ArXiv</i> , abs/2410.23123.	784 785 786 787 788
735	Md Imbesat Hassan Rizvi, Xiaodan Zhu, and Iryna Gurevych. 2024. Sparc and sparp: Spatial reasoning characterization and path generation for understanding spatial reasoning capability of large language models. In <i>Annual Meeting of the Association for Computational Linguistics</i> .	Yutaro Yamada, Yihan Bao, Andrew Kyle Lampinen, Jungo Kasai, and Ilker Yildirim. 2023. Evaluating spatial understanding of large language models. <i>ArXiv</i> , abs/2310.14540.	789 790 791 792
741	Manasi Sharma. 2023. Exploring and improving the spatial reasoning abilities of large language models. <i>ArXiv</i> , abs/2312.01054.	Daniel Yamins and James J. DiCarlo. 2016. Using goal-driven deep learning models to understand sensory cortex. <i>Nature Neuroscience</i> , 19:356–365.	793 794 795
744	Roger N. Shepard and Jacqueline Metzler. 1971. Mental rotation of three-dimensional objects. <i>Science</i> , 171:701 – 703.	Hanqi Yan, Yanzheng Xiang, Guangyi Chen, Yifei Wang, Lin Gui, and Yulan He. 2024. Encourage or inhibit monosemanticity? revisit monosemanticity from a feature decorrelation perspective. In <i>Conference on Empirical Methods in Natural Language Processing</i> .	796 797 798 799 800 801

802 Qwen An Yang, Baosong Yang, Beichen Zhang,
803 Binyuan Hui, Bo Zheng, Bowen Yu, Chengyuan
804 Li, Dayiheng Liu, Fei Huang, Guanting Dong, Hao-
805 ran Wei, Huan Lin, Jian Yang, Jianhong Tu, Jianwei
806 Zhang, Jianxin Yang, Jiaxin Yang, Jingren Zhou, Jun-
807 yang Lin, and 25 others. 2024. [Qwen2.5 technical
808 report](#). *ArXiv*, abs/2412.15115.

809 Fengji Yi, Wenlong Fu, and Huan Liang. 2018. [Model-
810 based reinforcement learning: A survey](#).

811 Xu Zheng, Zihao Dongfang, Lutao Jiang, Boyuan
812 Zheng, Yulong Guo, Zhenquan Zhang, Giuliano Al-
813 banese, Runyi Yang, Mengjiao Ma, Zixin Zhang,
814 Chenfei Liao, Dingcheng Zhen, Yuanhuiyi Lyu,
815 Yuqian Fu, Bin Ren, Linfeng Zhang, Danda Pani
816 Paudel, Niculae Sebe, Luc van Gool, and Xum-
817 ing Hu. 2025. [Multimodal spatial reasoning in the
818 large model era: A survey and benchmarks](#). *ArXiv*,
819 abs/2510.25760.

A Supplementary Materials 820

This appendix provides additional details on task 821
generation, experimental configurations, and ex- 822
tended results that complement the main text. 823

A.1 Task Generation Details 824

We provide detailed algorithms for generating each 825
task family. All algorithms use controlled random- 826
ization with fixed seeds to ensure reproducibility 827
and computational equivalence across languages. 828

A.1.1 Algorithm 1: Relational Spatial Reasoning 829
830

Algorithm 1 generates multi-hop spatial reasoning 831
instances. The key design principle is to ensure that 832
answering the query question requires composing 833
multiple pairwise relations, with no direct relation 834
between the source and target entities. 835

Spatial Encoding: We use a 3D coordinate sys- 836
tem where each atomic relation corresponds to a 837
unit vector: left/right = $(\pm 1, 0, 0)$, behind/front = 838
 $(0, \pm 1, 0)$, below/above = $(0, 0, \pm 1)$. The vector 839
representation enables systematic computation of 840
transitive relations. 841

Reasoning Requirement: To guarantee multi- 842
hop reasoning, the target entity (last in sequence) 843
cannot directly reference the source entity (first in 844
sequence), ensuring at least two inference steps. 845

A.1.2 Algorithm 2: Orientation Reasoning 846

Algorithm 2 generates orientation reasoning in- 847
stances that require tracking heading direction 848
through a sequence of turn actions. 849

Direction Encoding: Cardinal directions are 850
encoded as angles in the standard mathematical 851
convention: east = 0° , north = 90° , west = 180° , 852
south = 270° . For probing analysis, directions are 853
further encoded as unit vectors $(\cos \theta, \sin \theta)$ to en- 854
able continuous regression. 855

Turn Actions: We define three turn actions 856
with deterministic effects: “Turn right” (clockwise 857
 90°), “Turn left” (counterclockwise 90°), and “Turn 858
around” (180°). These correspond to rotation ma- 859
trices applied to the current heading vector. 860

Stateful Reasoning: Unlike Task Family 1, 861
which involves relational composition over a static 862
configuration, Task Family 2 requires maintain- 863
ing and updating a latent state (current orientation) 864
across sequential operations. 865

Algorithm 1 Generation of Multi-hop Spatial Reasoning Instances

Require: Number of samples N , hop range $[S_{\min}, S_{\max}]$

Ensure: Dataset \mathcal{D}

- 1: Define atomic spatial relations \mathcal{R} and their vector encodings
 - 2: Initialize empty dataset \mathcal{D}
 - 3: **for** $i = 1$ **to** N **do**
 - 4: Sample hop count $S \sim \text{UNIFORM}(S_{\min}, S_{\max})$
 - 5: Create entity sequence (e_0, e_1, \dots, e_S)
 - 6: Set position $\mathbf{p}(e_0) \leftarrow \mathbf{0}$
 - 7: Initialize fact set $\mathcal{K} \leftarrow \emptyset$
 - 8: **for** $j = 1$ **to** S **do**
 - 9: Sample reference entity $r_j \in \{e_0, \dots, e_{j-1}\}$
 - 10: Sample relation $r \in \mathcal{R}$
 - 11: $\mathbf{p}(e_j) \leftarrow \mathbf{p}(r_j) + \mathbf{v}(r)$
 - 12: Add fact (e_j is r of r_j) to \mathcal{K}
 - 13: **end for**
 - 14: Let source $s \leftarrow e_0$, target $t \leftarrow e_S$
 - 15: Compute relative vector $\Delta \leftarrow \mathbf{p}(t) - \mathbf{p}(s)$
 - 16: Derive gold relation $a \leftarrow \text{REL}(\Delta)$
 - 17: Construct question q from \mathcal{K} and (s, t)
 - 18: Sample distractor options D from $\mathcal{R} \setminus \{a\}$
 - 19: Shuffle options $O \leftarrow \{a\} \cup D$
 - 20: Add instance (q, a, O, Δ, S) to \mathcal{D}
 - 21: **end for**
 - 22: **return** \mathcal{D}
-

A.1.3 Algorithm 3: Spatial Program Execution

Algorithm 3 generates spatial program execution instances that require computing the final position after applying a sequence of geometric transformations.

Action Space: We include five types of operations: (1) *Move* – translate along cardinal directions; (2) *Reflect* – mirror across coordinate axes; (3) *Rotate* – rotate around axes by $90^\circ/180^\circ/270^\circ$; (4) *Scale* – multiply all coordinates by a factor; (5) *Translate* – add offset vector.

Coordinate System: All operations are defined in a 3D Cartesian coordinate system starting at origin $(0, 0, 0)$. The coordinate system uses right-handed convention with x (left/right), y (backward/forward), z (down/up).

Cumulative State: Each operation modifies the current position, and the final answer depends on the cumulative effect of all operations. Errors in in-

Algorithm 2 Generation of Orientation Reasoning Instances

Require: Number of samples N , step range $[S_{\min}, S_{\max}]$

Ensure: Dataset \mathcal{D}

- 1: Define cardinal directions $\mathcal{D} = \{\text{north, east, south, west}\}$
 - 2: Define turn actions \mathcal{A} with rotation offsets
 - 3: Initialize empty dataset \mathcal{D}
 - 4: **for** $i = 1$ **to** N **do**
 - 5: Sample number of steps $S \sim \text{UNIFORM}(S_{\min}, S_{\max})$
 - 6: Sample initial orientation $o_0 \sim \mathcal{D}$
 - 7: Set current orientation $o \leftarrow o_0$
 - 8: Initialize action list $\mathcal{K} \leftarrow \emptyset$
 - 9: **for** $j = 1$ **to** S **do**
 - 10: Sample turn action $a_j \sim \mathcal{A}$
 - 11: Update orientation $o \leftarrow \text{ROTATE}(o, a_j)$
 - 12: Append a_j to \mathcal{K}
 - 13: **end for**
 - 14: Let final orientation $o_S \leftarrow o$
 - 15: Construct question q from (o_0, \mathcal{K})
 - 16: Set answer $a \leftarrow o_S$
 - 17: Encode target vector $\mathbf{t} \leftarrow (\cos \theta, \sin \theta)$ for o_S
 - 18: Sample distractor options O from $\mathcal{D} \setminus \{a\}$
 - 19: Shuffle answer and distractors into multiple-choice options
 - 20: Add instance (q, a, O, \mathbf{t}, S) to dataset
 - 21: **end for**
 - 22: **return** \mathcal{D}
-

termediate steps propagate to the final result, making this task diagnostic of stateful computation and error accumulation.

Quality Control: To prevent extreme coordinate values that could arise from repeated scaling, we implement rejection sampling and constrain maximum coordinate magnitudes to ± 50 units.

A.2 Full Probing Results

A.2.1 Layer-wise Probing Patterns

Across all experiments, we observe consistent patterns in how spatial information is distributed across model layers:

- **Inverted-U profile:** Spatial information peaks in intermediate layers and declines sharply in final layers, indicating that spatial representations are constructed during intermediate processing but not preserved through

Algorithm 3 Generation of Spatial Procedure Execution Instances

Require: Number of samples N , step range $[S_{\min}, S_{\max}]$

Ensure: Dataset \mathcal{D}

```
1: Define action space  $\mathcal{A}$  (move, reflect, rotate,
   scale, translate)
2: Initialize empty dataset  $\mathcal{D}$ 
3: for  $i = 1$  to  $N$  do
4:   Sample number of steps  $S \sim$ 
   UNIFORM( $S_{\min}, S_{\max}$ )
5:   Initialize position  $\mathbf{p}_0 \leftarrow (0, 0, 0)$ 
6:   Set current position  $\mathbf{p} \leftarrow \mathbf{p}_0$ 
7:   Initialize action sequence  $\mathcal{K} \leftarrow \emptyset$ 
8:   for  $j = 1$  to  $S$  do
9:     Sample spatial action  $a_j \sim \mathcal{A}$ 
10:    Update position  $\mathbf{p} \leftarrow \text{APPLY}(\mathbf{p}, a_j)$ 
11:    Append  $a_j$  to  $\mathcal{K}$ 
12:   end for
13:   Let final position  $\mathbf{p}_S \leftarrow \mathbf{p}$ 
14:   Construct question  $q$  from  $(\mathbf{p}_0, \mathcal{K})$ 
15:   Set answer  $a \leftarrow \mathbf{p}_S$ 
16:   Sample distractor positions  $O$  by perturbing
    $\mathbf{p}_S$ 
17:   Shuffle correct answer and distractors into
   multiple-choice options
18:   Add instance  $(q, a, O, \mathbf{p}_S, S)$  to dataset
19: end for
20: return  $\mathcal{D}$ 
```

output generation.

- **Task-specific variations:** Different task families show peak spatial encoding at different layer positions, suggesting distinct computational stages for different types of spatial reasoning.
- **Cross-linguistic consistency:** The layer-wise emergence pattern is similar across languages, though the magnitude of decodable information varies.

Detailed layer-by-layer probing results are presented in the main paper (Table 3).

A.3 SAE Hyperparameters and Visualizations

A.3.1 SAE Training Configuration

All sparse autoencoders were trained using the SAE-Lens library with the following configuration:

- **Architecture:** Standard one-layer SAE with ReLU activation

- **Expansion factor:** $32\times$ (e.g., $3584 \rightarrow 114,688$ features for Qwen2.5-7B) 921-922
- **Training samples:** 2,000 task instances per language 923-924
- **Training tokens:** Approximately 300,000–420,000 tokens depending on language 925-926
- **Batch size:** 4,096 tokens 927
- **Learning rate:** 3×10^{-4} with linear warmup (1,000 steps) and cosine decay 928-929
- **L1 coefficient:** $\lambda = 0.001$ (selected via validation sweep over $\{0.0001, 0.0005, 0.001, 0.005\}$) 930-931-932
- **Optimizer:** AdamW with $\beta_1 = 0.9$, $\beta_2 = 0.999$ 933-934
- **Training steps:** 300–400 batches until convergence 935-936
- **Target layer:** Peak probe layer for each task family 937-938

A.3.2 SAE Quality Metrics 939

Task	Lang	MSE	R^2	Sparsity	L0
Task 1	EN	0.136	0.9996	0.0025	288.3
	ZH	0.798	0.9965	0.0023	268.2
	AR	0.252	0.9989	0.0019	219.5
Task 2	EN	0.147	0.9994	0.0026	301.4
	ZH	0.823	0.9962	0.0024	276.8
	AR	0.268	0.9987	0.0020	227.3
Task 3	EN	0.142	0.9995	0.0027	294.7
	ZH	0.765	0.9968	0.0025	271.5
	AR	0.241	0.9990	0.0021	223.8

Table 4: SAE reconstruction quality and sparsity metrics for Qwen2.5-7B-Instruct. MSE: mean squared error; R^2 : explained variance; Sparsity: fraction of active features; L0: average number of active features per sample. Best values per task family in bold.

Reconstruction MSE ranges from 0.136 (EN, Task 1) to 0.798 (CN, Task 1), with explained variance consistently above 0.996 except for Chinese Task 1 (0.9965). Average L0 (number of active features) ranges from 219 (AR, Task 1) to 301 (EN, Task 2), indicating successful sparsity enforcement. 940-941-942-943-944-945

A.3.3 Feature Activation Distributions 946

Most features activate rarely, with a heavy-tailed distribution where a small proportion of features account for the majority of activation mass. This 947-948-949

950 confirms that SAEs successfully learn sparse repre-
951 sentations.

952 A.3.4 Feature Selectivity Analysis

953 We compute feature selectivity using gradient-
954 based attribution (gradient \times activation). Analy-
955 sis of SAE features reveals that spatially selective
956 features are not necessarily the most frequently
957 activated, indicating dissociation between usage
958 frequency and causal importance for spatial reason-
959 ing.

960 A.4 Prompt Templates (English)

961 All tasks follow a multiple-choice format with four
962 options. Below are representative prompt templates
963 for each task family in English.

964 A.4.1 Task Family 1: Relational Spatial 965 Reasoning

[English]

System: You are a helpful assistant.

User: Given the following spatial facts, an-
answer the question by selecting one option.

D is behind C.
B is left of A.
C is behind B.
E is right of D.

Where is E relative to A?

- A. right and behind
- B. left and front
- C. behind
- D. left and behind

Assistant: The answer is

A.4.2 Task Family 2: Orientation Reasoning

967

[English]

System: You are a helpful assistant.

User: Follow the instructions step by step and
answer the question by selecting one option.

You are facing north.

Turn right.

Turn left.

Turn around.

Turn right.

Which direction are you facing now?

- A. north
- B. east
- C. south
- D. west

Assistant: The answer is

968

A.4.3 Task Family 3: Spatial Program Execution

969

970

[English]

System: You are a helpful assistant.

User: Execute the spatial operations step by
step and answer the question by selecting one
option.

Start at (0, 0, 0).

Move forward by 3 units.

Move right by 2 units.

Reflect the position across the z-axis.

Move up by 1 unit.

What is the final position?

- A. (2, 3, -1)
- B. (2, 3, 1)
- C. (-2, 3, 1)
- D. (2, -3, 1)

Assistant: The answer is

971

A.4.4 Inference Protocol

972

For all experiments:

973

- Models generate a single token continuation
after the prompt 974
975
- The most likely single-character token
(A/B/C/D) is selected as the answer 976
977
- No chain-of-thought or explicit reasoning is
elicited 978
979

966

980	• Temperature is set to 0 (greedy decoding)	• Model inference (2,000 samples): 10–15 minutes	1021
981	• Maximum generation length is 1 token	• Linear probe training (all layers): 30–45 minutes	1022
982	This protocol ensures that evaluation targets implicit spatial representations encoded in activations rather than verbalized reasoning strategies.	• SAE training (one layer, one task): 2–3 hours	1023
983		• Intervention experiments: 1–2 hours per configuration	1024
984		Total compute: approximately 500 GPU-hours across all models, languages, and tasks.	1025
985	A.5 Dataset Statistics and Quality Control	A.6.2 Reproducibility	1026
986	A.5.1 Data Generation Process	We provide:	1027
987	All datasets were generated using rule-based procedures with controlled randomization. For each task family:	• Complete data generation code with fixed random seeds	1028
988		• Exact model checkpoints and inference configurations	1029
989		• Probe training scripts with hyperparameter specifications	1030
990	1. Entity/Action Sampling: Entities (Task 1), directions (Task 2), and operations (Task 3) are sampled uniformly at random.	• SAE training configurations and learned feature dictionaries	1031
991		• Intervention protocols and analysis notebooks	1032
992		All code will be released upon publication to ensure full reproducibility.	1033
993	2. Constraint Verification: Generated instances are verified to ensure they require multi-step reasoning and do not allow shortcut solutions.	A.7 Intervention Experiment Details	1034
994		A.7.1 Activation Patching Protocol	1043
995		For each intervention experiment:	1044
996	3. Answer Distribution: Correct answers are balanced across option positions (A/B/C/D) with uniform distribution (25% each).	1. Baseline forward pass: Run model on original input, collect activations at target layer	1045
997		2. Counterfactual forward pass: Run model on counterfactual input (e.g., different spatial trajectory), collect activations at same layer	1046
998		3. Patching: Replace activations at target layer with counterfactual activations	1047
999	4. Difficulty Stratification: Instances are stratified by complexity level (number of steps) with uniform distribution across difficulty range.	4. Measurement: Continue forward pass and measure output shift toward counterfactual prediction	1048
1000		We report:	1049
1001		• Shift magnitude: KL divergence between original and patched output distributions	1050
1002		• Prediction change: Whether patching changes the top-1 prediction	1051
1003	A.5.2 Cross-Linguistic Consistency	• Layer specificity: Effect size as a function of intervention layer	1052
1004	To ensure computational equivalence across languages, we verify:		1053
1005			1054
1006	• Structural isomorphism: Same reasoning steps and answer patterns		1055
1007			1056
1008	• Token count variance: Chinese prompts are typically 15–20% shorter due to character-level encoding; Arabic prompts are 10–15% longer		1057
1009			1058
1010			1059
1011			1060
1012	• Template diversity: All languages use the same number of prompt templates		1061
1013			1062
1014	• Random seed alignment: Parallel instances across languages use the same random seed, ensuring matched difficulty levels		
1015			
1016			
1017	A.6 Additional Experimental Details		
1018	A.6.1 Computational Resources		
1019	All experiments were conducted on NVIDIA A100 GPUs (40GB). Typical resource requirements:		
1020			

1063 **A.7.2 SAE Feature Ablation Protocol**

1064 For each feature ablation experiment:

- 1065 1. **Feature identification:** Rank SAE features
1066 by gradient-based attribution
- 1067 2. **Selective ablation:** Zero out top- k features
1068 for varying k values
- 1069 3. **Reconstruction:** Decode modified SAE acti-
1070 vations back to hidden states
- 1071 4. **Evaluation:** Measure accuracy drop on test
1072 set

1073 Feature ablation experiments confirm that iden-
1074 tified spatial features causally contribute to model
1075 behavior, with accuracy systematically decreasing
1076 as more top-ranked features are removed.

1077 **A.7.3 Counterfactual Construction**

1078 For spatial program execution (Task 3), we con-
1079 struct counterfactuals by:

- 1080 • Flipping a single operation in the sequence
1081 (e.g., “move left” \rightarrow “move right”)
- 1082 • Computing the resulting position divergence
1083 $\|\mathbf{P}_{\text{actual}} - \mathbf{P}_{\text{counterfactual}}\|$
- 1084 • Selecting pairs with divergence in range [3,
1085 10] units (neither too similar nor too different)

1086 For relational reasoning (Task 1), counterfactu-
1087 als are constructed by:

- 1088 • Flipping one spatial relation (e.g., “A left of
1089 B” \rightarrow “A right of B”)
- 1090 • Verifying that the modified configuration is
1091 logically consistent
- 1092 • Computing answer change (same vs. different
1093 final relation)