

Mind's Eye of LLMs: Visualization-of-Thought Elicits Spatial Reasoning in Large Language Models

Project Page

Wenshan Wu[†], Shaoguang Mao[†], Yadong Zhang^{†,‡}, Yan Xia[†], Li Dong[†], Lei Cui[†], Furu Wei[†]

[†] Microsoft Research [‡] East China Normal University



Introduction

Spatial Reasoning in Human Cognition

- Mental Image: abstract representations from visual perception
- Mind's Eye: mental image manipulation

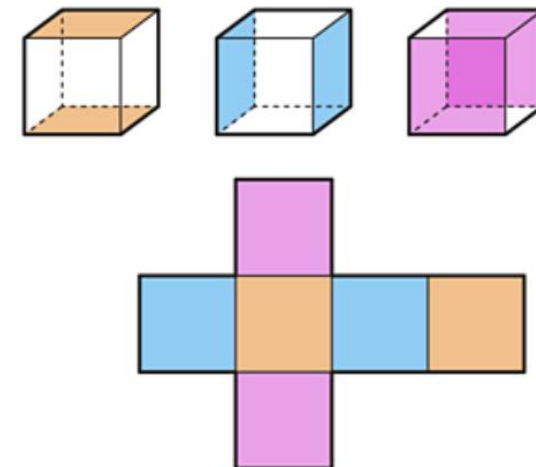
Navigation



Mental rotation



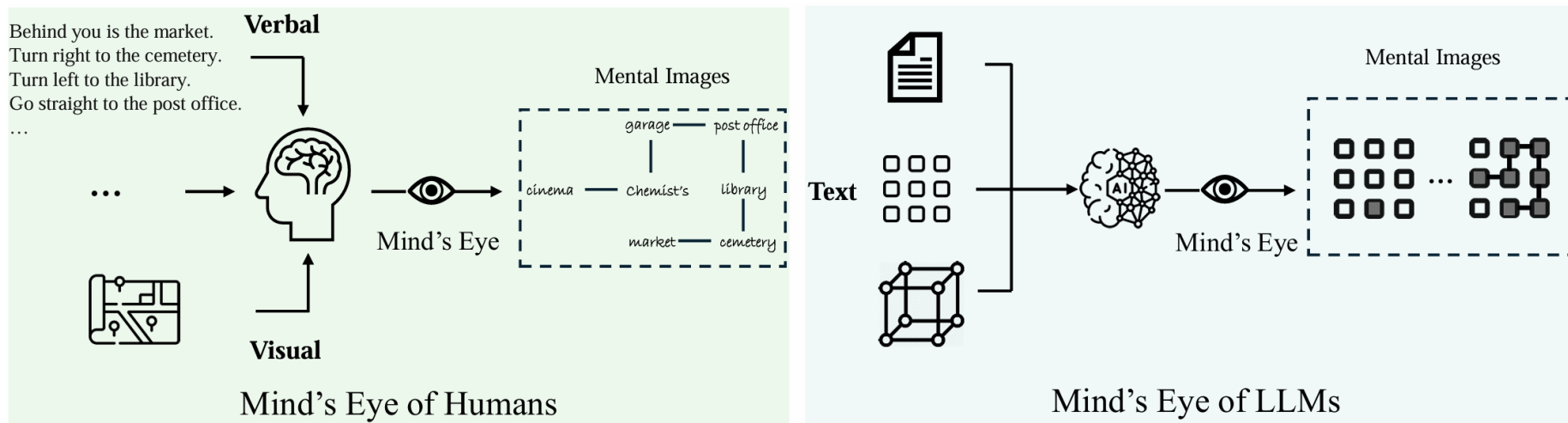
Mental paper folding



Motivation

Similar Mechanism in LLMs: Mind's Eye

- Visualize internal states
- Manipulate mental images to guide subsequent reasoning



Contribution

- We conduct quantitative and qualitative analyses on the mind's eye of LLMs and its limitations. We also explore cues about the **origin of this generalized ability from code pre-training**.
- We develop two tasks of "visual navigation" and "visual tiling", along with corresponding synthetic datasets, emulating various sensory inputs for LLMs. These tasks are structured to support **varying levels of difficulty**.
- We propose **Visualization-of-Thought (VoT) prompting** to elicit the mind's eye of LLMs for spatial reasoning and provide empirical evaluations on three tasks.

Spatial Reasoning Tasks

- Existing Benchmarks
 - Spatial semantics are embedded in text, spatial term focused
 - Could be solved by **logic programming** after converting spatial terms to logical forms through LLMs
- Ours
 - Focus on spatial awareness
 - Various aspects: spatial relationships, directions, and **geometric shapes**
 - Essential for action planning in the physical world.
 - Emulating various sensory inputs for LLMs
 - Natural language
 - 2D grid comprising of special text characters

Spatial Reasoning Tasks

- **Natural Language Navigation** [1]

- Square map $W = \{(l_1, o_1), (l_2, o_2), \dots, (l_n, o_n)\}$, each location associated with an object
- Navigation instructions $I = \{i_1, i_2, \dots, i_k\}$
- Task: Find the object o at specific location l determined by navigation instructions

$$o \sim p(o \in W | W = \{(l_1, o_1), (l_2, o_2), \dots, (l_n, o_n)\}, I)$$

- **Visual Navigation**

- Grid map M consisting of k consecutive edges $E = \{e(s_0, s_1), e(s_1, s_2), \dots, e(s_{k-1}, s_k)\}$
- Route planning: generate a sequence of correct directions

$$D \sim p(\{d(s_0, s_1), d(s_1, s_2), \dots, d(s_{k-1}, s_k)\} | M)$$

- Next step prediction: given t navigation instructions, identify the direction of next step

$$d \sim p(d(s_t, s_{t+1}) | M, D_{t, 0 < t < k})$$

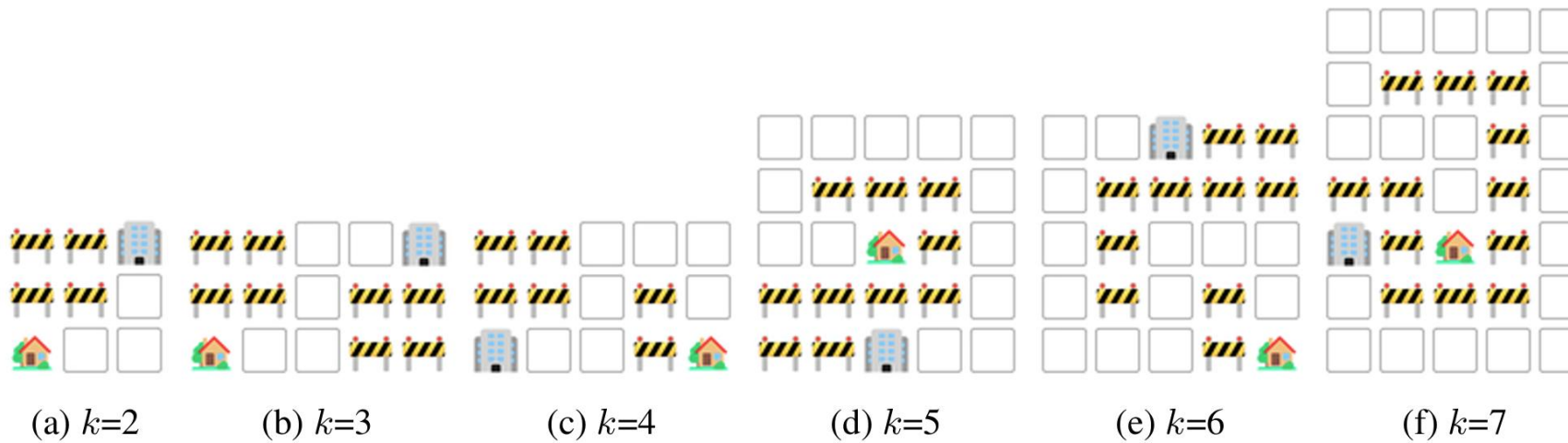
- **Visual Tiling**

- Rectangle R masked with k unique polyominoes $MP = \{mp_1, mp_2, \dots, mp_k\}$
- Two variants of each polyomino $v_{i < k} = \{v_{i1}, v_{i2}\}$, a polyomino query $q \in MP$
- Task: identify the correct variant of q

$$v \sim p(v_q | R, \{mp_1, \dots, mp_k\}, \{v_{11}, v_{12}, \dots, v_{k1}, v_{k2}\}, q)$$

Spatial Reasoning Tasks

Various Levels of Difficulties



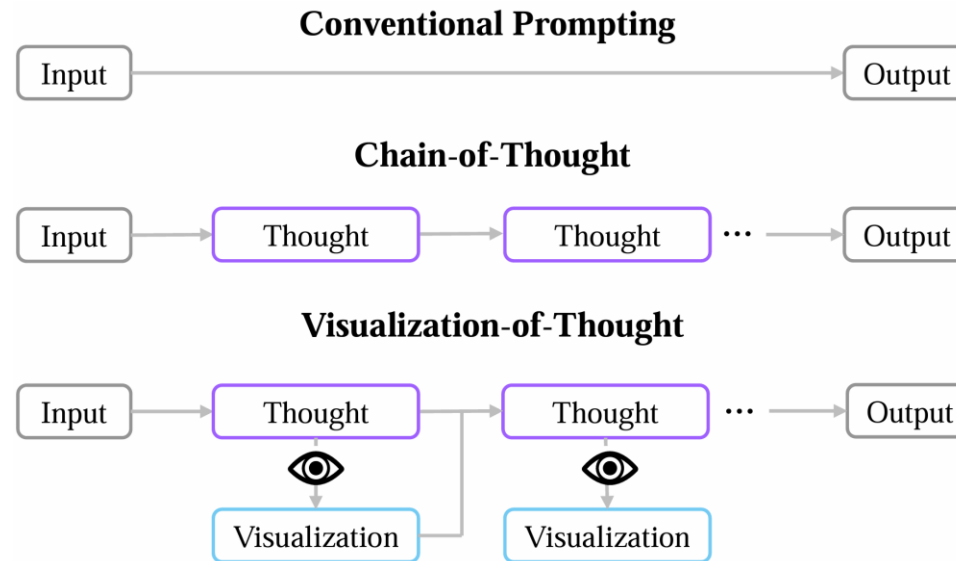
(a) Fit 2 pieces into a masked rectangle



(b) Fit 3 pieces into a masked rectangle

Visualization-of-Thought Prompting

Visualize the state after each reasoning Step



- x : text sequence of input
- v : visualization sequence in text form
- z : language sequence of intermediate steps



$$v_i \sim p_{\theta}(v_i \mid \text{prompt}_{VOT}, x, z_{1\dots i}, v_{1\dots i-1})$$

$$z_{i+1} \sim p_{\theta}(z_{i+1} \mid \text{prompt}_{VOT}, x, z_{1\dots i}, v_{1\dots i})$$

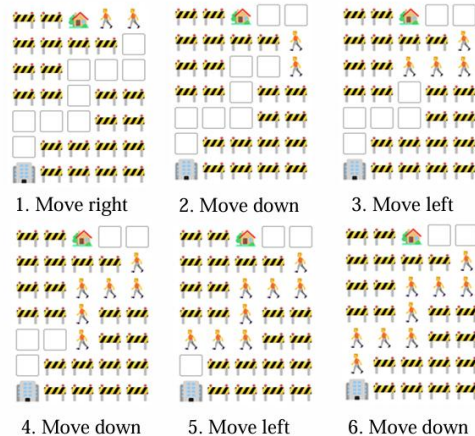
Qualitative Results

Visual Navigation






Starting from , provide the steps to navigate to .

Visualize the state after each reasoning step.

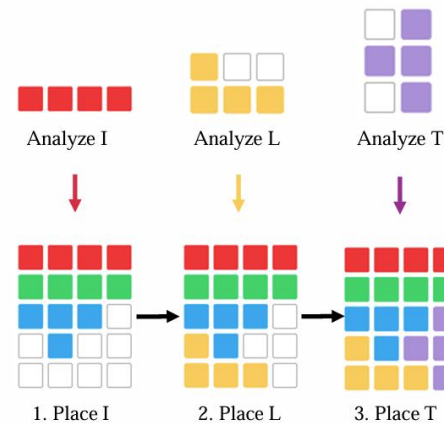


Visual Tiling



Provided: I  T  L 
To fit all the provided polyominoes into the empty squares, what's the correct variation of Tetromino T?

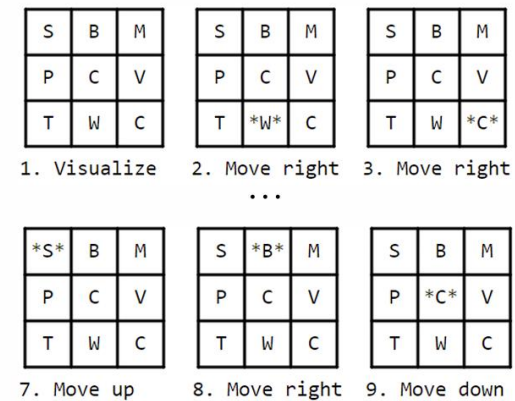
Visualize the state after each reasoning step.



Natural Language Navigation

You have been given a 3 by 3 square grid. Initially, you are at the bottom-left corner...find a cassette player...go right...a wool, go right...a conch, go up...a moving van, go left...a confectionery store, go left...a pot pie, go up...a siamang, go right...a black-and-white colobus, go right...a minivan. Now you have all the information on the map. You start at where the cassette player is located, then you go right by one step, go right...go up...go left...go left...go up...go right...go down by one step. What will you find?

Visualize the state after each reasoning step.



Dataset

Natural Language Navigation

- 200 square maps of size 3x3

Visual Navigation

- 496 navigation maps and 2520 QA instances
- Map size up to 7×9 and 9×7

• Visual Tiling

- 5 x 4 rectangle with 2 or 3 polyomino masked

Task	<i>K</i> Step						Total
	2	3	4	5	6	7	
Route Planning	8	16	32	64	128	248	496
Next Step Prediction	8	32	96	256	640	1488	2520

	Mask count		Total
	2	3	
Configuration	248	124	376
QA Instance	489	307	796

Experiments

- Settings
 - GPT-4 CoT: Let's think step by step.
 - GPT-4 w/o Viz: Don't use visualization. Let's think step by step.
 - GPT-4V CoT: Let's think step by step.
 - GPT-4 VoT: Visualize the state after each reasoning step.

Settings	Visual Navigation			Visual Tiling	Natural-Language Navigation
	Route Planning		Next Step Prediction		
	Completing Rate	Succ Rate			
GPT-4 CoT	37.02	9.48	48.61	54.15	54.00
GPT-4 w/o Viz	37.17	10.28	48.49	46.98	35.50
GPT-4V CoT	33.36	5.65	46.59	49.62	/
GPT-4 VoT	<u>40.77</u>	<u>14.72</u>	<u>55.28</u>	<u>63.94</u>	<u>59.00</u>

Analysis

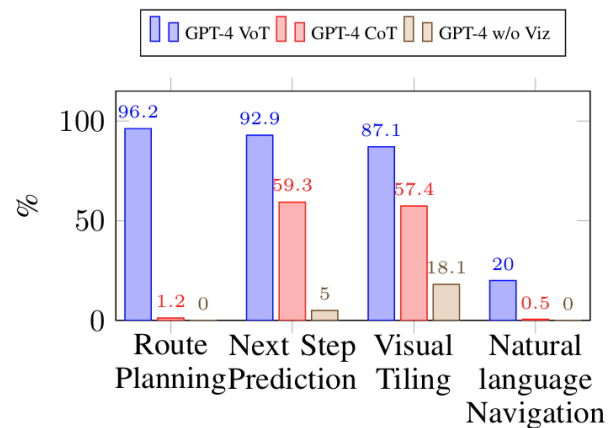
Do visual state tracking behaviors differ among prompting methods?

l_v : length of visualization sequence, l_s : number of reasoning steps

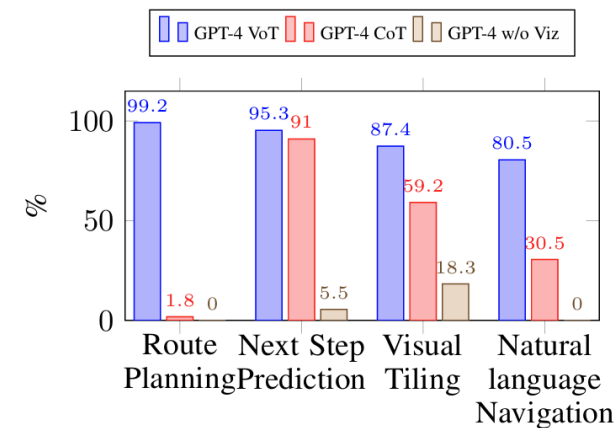
$$\text{Complete tracking rate} = \sum_i^n (l_v == l_s) / n$$

$$\text{Patial tracking rate} = \sum_i^n (l_v < 0) / n$$

- VoT markedly improves the visual tracking rate
- LLMs inherently exhibit the capability of visual state tracking in some tasks.



(a) Complete tracking rate



(b) Partial tracking rate

Analysis

How visualizations enhance final answers?

- Visualization Quality
 - Compliance: visualization satisfies requirements in 51-52% cases
 - Accuracy: visualization aligns with the corresponding state in 24%-26% cases
- Performance enhancement
 - LLMs are able to make correct decisions in 65%-77% of the cases **when accurate internal state visualizations are generated**

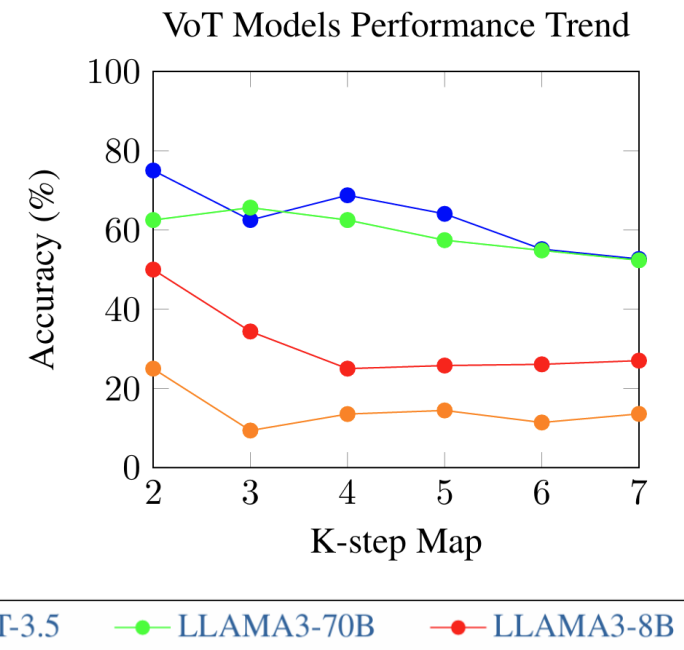
Task	Spatial Visualization		Spatial Understanding
	Compliance	Accuracy	Accuracy
Visual Navigation	51.14	26.48	65.16
Visual Tiling	52.01	24.25	77.20

Analysis

Can VoT benefit less powerful language models?

- VoT offers a scaling advantage when applied to more advanced models
- Less capable models tend to rely on random guessing

Settings	Visual Navigation			Visual Tiling	Natural-Language Navigation
	Route Planning		Next Step Prediction		
	Completing Rate	Succ Rate			
GPT-3.5 CoT GPT-3.5 VoT	16.10 19.02	2.62 1.61	17.42 13.10	44.10 47.99	8.50 9.00
LLAMA3-8B CoT LLAMA3-8B VoT	4.65 4.97	0 0.2	28.73 26.75	47.24 46.73	16.50 15.50
LLAMA3-70B CoT LLAMA3-70B VoT	19.90 30.24	2.62 5.85	49.01 54.09	56.41 56.03	26.00 32.50

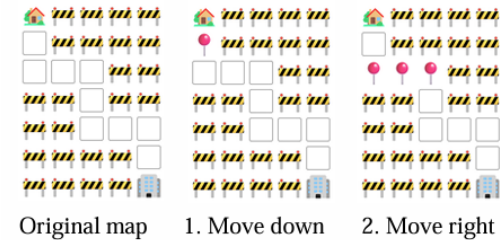


Appendix

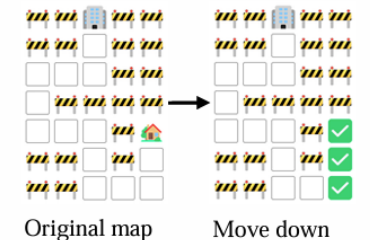
Mental Images for State Tracking

- Mark the path with unique symbols
- Mark path and direction with arrows
- Mark path with temporal steps
- Remove road: turning roads into obstacles

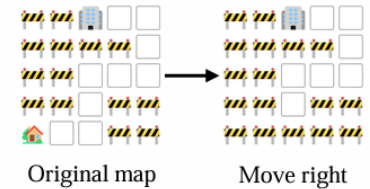
Use round pin



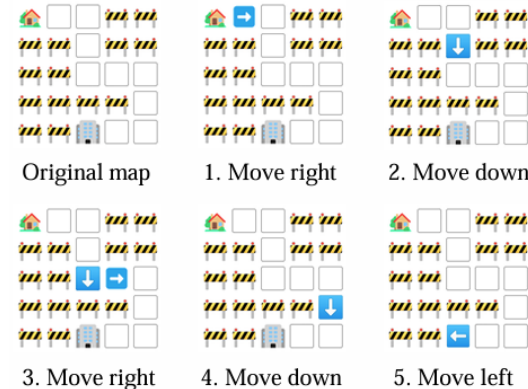
Use checklist



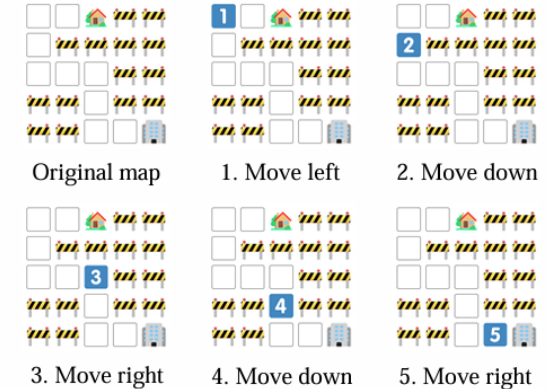
Remove road to avoid turning back



Use arrows to reflect direction



Use numbers for temporal steps



Appendix

Ascii-art in Code Comments

- Represents data structure, diagram, geometry
- Illustrates how an algorithm works or simulates an operation
 - Spatial Causality: [Double-ended queue in Rust](#), [Scrolling web pages](#), [tree rotation](#) present triplets of previous visual state, instruction, and updated state of instruction following.
 - Temporal Causality: [Undo systems from emacs](#) provides various temporal states of the undo system when undo operation happens in different timelines and corresponding visualizations in an interleaved manner. Each visualization reflects the temporal causality of the system state.