

---

# Bridging the Von Neuman Gap: Why LLMs Haven't Made Novel Discoveries

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1 Large language models (LLMs) have been trained on vast data spanning nearly  
2 every scientific discipline, yet they have not produced a single novel discovery.  
3 Human polymaths such as John von Neumann routinely generated breakthroughs  
4 across disparate fields—from game theory to quantum mechanics to the very archi-  
5 tecture of the modern computer—by connecting insights across domains. We  
6 argue this gap reflects a structural limitation of the LLM paradigm rather than  
7 a problem of scale. Drawing on Piaget’s theory of cognitive development and  
8 Gentner’s structure-mapping, we contend novel discovery depends on two core  
9 processes: constructing nuanced internal schemas of the external world and flexibly  
10 redeploying them via analogical mapping. Without embodied data or exploration,  
11 LLMs form shallow world models; and because their architectures optimize for  
12 statistical efficiency, they struggle to extend analogies out of distribution in ways  
13 that capture relational structure across domains. Without rethinking training en-  
14 vironments and architectures, LLMs will remain constrained to weak abstraction  
15 rather than the deep reasoning required for scientific innovation.

## 16 1 Introduction

17 Large language models (LLMs) have reached or exceeded human performance in many specialized  
18 domains, from mathematics and law to protein structure prediction [Abramson et al., 2024, Zhong  
19 et al., 2024]. Yet despite this breadth of competence, no LLM to date has produced a verifiable novel  
20 scientific discovery—a genuine insight not previously known to humans that expands the boundaries  
21 of knowledge [Shojaee et al., 2025]. This absence is often noted with surprise: if LLMs can  
22 solve Olympiad problems, pass graduate-level exams, and synthesize knowledge across disciplines,  
23 why have they not combined these abilities to generate groundbreaking, out-of-distribution (OOD)  
24 findings?

25 We argue that this gap is not surprising at all. Drawing from cognitive science and developmental  
26 psychology, we propose that human novelty generation depends on two essential ingredients current  
27 LLMs lack: (1) the development of robust internal schemas that accurately model the external world,  
28 and (2) the flexible redeployment of these schemas to new contexts through analogical reasoning.  
29 Piaget’s theory of learning shows how embodied experience forms schemas that compress the world  
30 into abstract, reusable models [Beilin, 1992]. Gentner’s Structure-Mapping Theory explains how  
31 these schemas can be redeployed across domains through relational alignment, enabling the deep  
32 analogies behind scientific discovery [Gentner, 1983].

33 Grounding our perspective in Piaget’s and Gentner’s frameworks, we contend that good schemas  
34 lead to good analogies, and good analogies enable novel hypotheses that generalize OOD (Figure 1).  
35 Human history provides vivid illustrations of this process. For example, the hydraulic analogy in  
36 electricity—conceiving current as fluid flow—helped early scientists reason about voltage, resistance,

37 and circuits [Tembrevilla et al., 2019]. James Clerk Maxwell’s vortex analogy in fluid mechanics  
 38 enabled him to formulate the equations of electromagnetism by mapping fluid vortices onto field  
 39 lines [Harman, 1998]. More recently, the discovery of CRISPR–Cas9 gene editing [Jinek et al., 2012]  
 40 emerged from recognizing that bacterial immune systems could be repurposed as programmable  
 41 molecular “scissors”. Si et al. [2024] found that while LLM-generated hypotheses were judged more  
 42 novel than human ones, they were significantly less feasible, reflecting a weak causal model of the  
 43 world. The central challenge lies in moving beyond statistical novelty to schema-based analogical  
 44 reasoning, the cognitive foundation of historical scientific discovery.

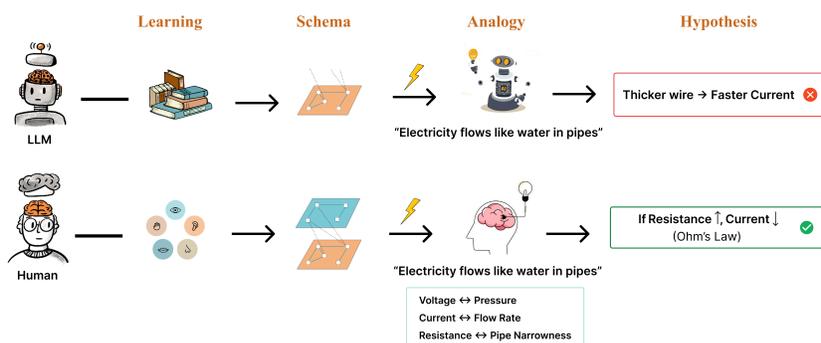


Figure 1: Humans build rich, causal schemas from embodied experience, enabling relational analogies and correct hypotheses (e.g., Ohm’s Law). LLMs, trained primarily on text, form flatter associative schemas, leading to surface analogies and brittle, incorrect hypotheses

## 45 2 Cognitive Science Frameworks

### 46 2.1 Defining Novel Discovery

47 In a trivial sense, language models constantly produce “new” ideas: every draft completion, bug fix  
 48 suggestion, or paraphrase is technically novel. Genuine breakthroughs, however, provide a testable  
 49 insight about the external world and introduce a principle or relation that extends human knowledge  
 50 beyond existing frameworks [Kuhn, 1962]. By these criteria, discoveries such as Newton’s law of  
 51 gravitation, Maxwell’s equations, or the CRISPR–Cas9 gene-editing system qualify: they reveal  
 52 underlying structures of reality that were previously unknown and enabled entire domains of inquiry.  
 53 Advances such as protein folding exemplify optimization within established frameworks rather than  
 54 out-of-distribution novelty. It’s rare to find examples of LLMs synthesizing their massive training  
 55 dataset into a new conceptual discovery, an imaginative leap like Kekulé’s dream of a snake biting its  
 56 tail [Rocke, 2010] that revealed benzene’s ring structure.

### 57 2.2 Piaget’s Theory of Cognitive Development

58 According to Piaget, the goal of learning is to construct the most accurate internal model of the  
 59 world available at a given time [Beilin, 1992]. In the sensorimotor stage (0–2 years), children build  
 60 embodied schemas through direct interaction with the environment, discovering object permanence  
 61 and forming habits grounded in physical action. In the preoperational stage (2–7 years), schemas  
 62 become symbolic: words, gestures, and images represent objects and events. By the concrete  
 63 operational stage (7–11 years), children can run internal “simulations” of their schemas, applying  
 64 logical operations and reversible reasoning to concrete scenarios. Finally, in the formal operational  
 65 stage (12+ years), these foundations enable abstract thought, hypothetical reasoning, and systematic  
 66 problem-solving. As humans learn, schemas reorganize to become increasingly abstract, hierarchical,  
 67 and nuanced, functioning as cognitive priors for interpreting the world.

### 68 2.3 Gentner’s Structure-Mapping Theory

69 Gentner’s Structure-Mapping Theory provides a cognitive account of how analogy supports ab-  
 70 straction and discovery [Gentner, 1983]. Unlike surface similarity, analogy depends on aligning

71 relational structures across domains. In this process, a familiar “base” domain is mapped onto a less  
72 familiar “target” domain, with correspondences drawn between underlying causal relations. Gentner  
73 formalized this through the systematicity principle, which holds that analogies preserving coherent,  
74 interconnected relations are more powerful than those based on isolated features. High-quality  
75 analogies are therefore indispensable for novel discovery, because they enable relational structures  
76 from well-understood domains to be systematically redeployed in unfamiliar ones. Gentner’s theory  
77 identifies the cognitive mechanism that allows humans to reason out-of-distribution, moving beyond  
78 rote pattern recognition toward flexible, relational inference.

### 79 **3 Why Current LLMs Fail?**

#### 80 **3.1 Internal World Models**

81 Large Language Models do not merely memorize text—they form implicit internal world models that  
82 guide their predictions [Li et al., 2024]. Evidence from mechanistic interpretability shows that even  
83 small transformers learn structured representations of game states rather than just token statistics [Li  
84 et al., 2023, Karvonen et al., 2024]. Recent work further demonstrates that LLMs encode linear spatial  
85 world models [Teheenan et al., 2025] and can apply simple heuristics in physical reasoning tasks such  
86 as pulleys [Robertson and Wolff, 2025]. Yet, as Robertson and Wolff [2025] emphasize, these models  
87 lack the facility to reason over nuanced structural connectivity, failing when problems demand deeper  
88 relational understanding. More broadly, evaluations reveal that world-model coherence often breaks  
89 down under perturbation [Vafa et al., 2024].

90 Piaget’s Theory of Cognitive Development makes clear why LLMs fail to form robust world models.  
91 The ability to build higher-order abstractions in the formal operational stage depends on foundations  
92 laid in earlier stages: spatial grounding in the sensorimotor and preoperational stages, and exploratory  
93 simulation in the concrete operational stage. Many of our most powerful scientific analogies—such  
94 as electricity flowing like water or the atom resembling a solar system—are rooted in embodied  
95 interaction. Exploration enables the counterfactual reasoning that underpins robust schemas—for  
96 example, asking *what would happen if resistance increased?* LLMs, by contrast, encounter only  
97 linguistic descriptions of these mappings, not the embodied patterns themselves. While text corpora  
98 encode valuable abstractions in mathematics, physics, and scientific reasoning, they lack the embodied  
99 variation necessary to anchor relational concepts [Bisk et al., 2020]. One can read about conservation  
100 of energy, but until they interact with systems of push and pull, the concept remains fragile. LLMs  
101 lack both of these developmental foundations: spatial understanding of the physical world and  
102 self-generated exploration to refine internal schemas.

#### 103 **3.2 Limits of Analogical Mapping**

104 The core problem isn’t that LLMs can’t do analogical mapping—they often succeed at surface-level  
105 reasoning [Musker et al., 2025]. Transformer attention excels at capturing token co-occurrences  
106 and statistical dependencies [Geva et al., 2023, Vig and Belinkov, 2019], but this strength becomes  
107 brittle out of distribution. The deeper issue is that the next-token prediction paradigm is structurally  
108 myopic: trained under teacher forcing, models exploit local token correlations rather than constructing  
109 generalizable rules. As Bachmann and Nagarajan [2024] show, this leads to failures even on simple  
110 lookahead planning tasks, and Nagarajan et al. [2025] demonstrate similar breakdowns on algorithmic  
111 problems requiring novel pattern construction. This brittleness is evident in analogy itself: Lewis  
112 and Mitchell [2024] find that while LLMs handle standard analogy problems, they collapse on  
113 counterfactual variants. Puranam et al. [2025] confirm this gap empirically, showing that GPT-4 often  
114 applies incorrect analogies based on superficial features, while humans generate fewer but causally  
115 grounded mappings.

116 Humans rely on hierarchical schemas grounded in multiple levels of abstraction. Understanding  
117 electricity, for example, involves mathematical formalism, intuitive “flow” metaphors, and mechanical  
118 analogies simultaneously. This layered structure supports analogical reasoning that is deeply tied to  
119 spatial and causal grounding—capacities LLMs lack. Shani et al. [2025] show that LLMs instead  
120 perform “aggressive statistical compression,” prioritizing efficiency over preserving the fine-grained  
121 distinctions essential for human-like reasoning. While this yields broad categorical alignment with  
122 human concepts, it erases the typicality gradients and internal semantic structure that enable flexible  
123 analogical mapping across domains. By relying on statistical similarity from next-token prediction,

124 LLMs fail to produce the relational, out-of-distribution analogies that Gentner identifies as essential  
125 for novel discovery.

## 126 4 Towards Solutions

### 127 4.1 Training Environment

128 **Training with spatial data.** Exposing LLMs to spatially grounded information can significantly  
129 improve their reasoning about the physical world. SpatialLM, trained on large-scale synthetic 3D  
130 indoor scenes, improves layout estimation and 3D object detection [Mao et al., 2025]. SpatialVLM  
131 scales relational spatial reasoning through billions of spatially grounded vision–language examples  
132 [Chen et al., 2024]. Han et al. [2025] highlights the importance of incorporating computational  
133 geometry, 3D point clouds (e.g., LiDAR scans), and CAD/architectural models into LLM training  
134 pipelines to advance spatial reasoning. Moving forward, researchers should prioritize 3D scene data,  
135 multimodal robotics traces, and physics-based simulation datasets, since these expose models to the  
136 structural relationships necessary for robust internal world models.

137 **Training for exploration.** Exploration is critical for building robust world models, enabling LLMs  
138 to internalize causal dynamics. Evidence supports this claim: fine-tuning LLMs on embodied experi-  
139 ences in environments such as VirtualHome yields over 60% improvements in reasoning tasks by  
140 grounding models in object permanence and causal regularities [Xiang et al., 2023]. Similarly, embod-  
141 ied agents like STEVE [Zhao et al., 2024] in Minecraft and Voyager [Wang et al., 2023] demonstrate  
142 how autonomous exploration can accumulate transferable skills, while S2ERS [Zhang et al., 2025]  
143 reduces spatial hallucinations in maze-like planning through reinforcement learning. These results  
144 highlight the broader principle: models must actively probe their environments to uncover invariants  
145 such as conservation laws and stability. We propose two pathways forward—scaling open-ended  
146 virtual environments like Minecraft and MuJoCo to approximate embodied play, and leveraging scien-  
147 tific simulators such as SPICE or fluid dynamics solvers for controlled, intervention-rich exploration.  
148 We see particular promise in physics-based sandboxes, where models must manipulate environments  
149 to rediscover scientific laws from first principles. By repeatedly altering parameters and observing  
150 outcomes, models can learn the causal rules that govern systems, developing robust priors for analogy  
151 and transfer across domains.

### 152 4.2 Cognitively Aligned Architectures

153 We contend that progress toward analogical reasoning requires neurosymbolic architectures that  
154 combine the explicit relational structures of symbolic systems, with the statistical generalization ca-  
155 pabilities of neural networks [Bougzime et al., 2025]. Early results support this approach: on Raven’s  
156 Progressive Matrices, ARLC achieves near-perfect performance by explicitly modeling relational  
157 rules [Hersche et al., 2024], while [Shah et al., 2022] show that integrating symbolic background  
158 knowledge with neural embeddings enables analogical inferences beyond surface correlations. At the  
159 same time, scalability remains an open engineering challenge. Current approaches either collapse into  
160 fuzzy embeddings or brittle symbolic rules that fail to generalize [Naik et al., 2024]. The next step is  
161 to design architectures that learn structured symbolic schemas at scale while retaining generalization,  
162 enabling cross-domain analogical reasoning.

## 163 5 Conclusion

164 Bridging the von Neumann gap will require more than scale: it demands training regimes and  
165 architectures that embed the principles of human reasoning. Incorporating spatial data and enabling  
166 LLMs to explore counterfactuals can anchor robust internal world models, while neurosymbolic  
167 architectures may provide the scaffolding for deep analogical inference across domains. Imagine  
168 thousands of polymathic systems like Von Neumann, each capable of connecting insights across  
169 disparate fields and generating unexpected analogies that drive scientific breakthroughs. We urge  
170 the AI community to ground future research in cognitive theory, so that AI moves beyond statistical  
171 efficiency and emerges as an engine of unprecedented discovery.

## 172 References

- 173 J. Abramson, J. Adler, J. Dunger, et al. Accurate structure prediction of biomolecular interactions  
174 with alphafold 3. *Nature*, 630:493–500, 2024. doi: 10.1038/s41586-024-07487-w.
- 175 Gregor Bachmann and Vaishnavh Nagarajan. The pitfalls of next-token prediction. *arXiv preprint*  
176 *arXiv:2403.06963*, 2024.
- 177 Harry Beilin. Piaget’s enduring contribution to developmental psychology. *Developmental Psychol-*  
178 *ogy*, 28(2):191–204, 1992. doi: 10.1037/0012-1649.28.2.191.
- 179 Yonatan Bisk, Ari Holtzman, Jesse Thomason, Jacob Andreas, Yoshua Bengio, Joyce Chai, Mirella  
180 Lapata, Angeliki Lazaridou, Jonathan May, Aleksandr Nisnevich, Nicolas Pinto, and Joseph  
181 Turian. Experience grounds language. In Bonnie Webber, Trevor Cohn, Yulan He, and Yang  
182 Liu, editors, *Proceedings of the 2020 Conference on Empirical Methods in Natural Language*  
183 *Processing (EMNLP)*, pages 8718–8735, Online, November 2020. Association for Computational  
184 Linguistics. doi: 10.18653/v1/2020.emnlp-main.703. URL [https://aclanthology.org/2020.](https://aclanthology.org/2020.emnlp-main.703/)  
185 [emnlp-main.703/](https://aclanthology.org/2020.emnlp-main.703/).
- 186 Oualid Bougzime, Samir Jabbar, Christophe Cruz, and Frédéric Demoly. Unlocking the potential  
187 of generative ai through neuro-symbolic architectures: Benefits and limitations. *arXiv preprint*  
188 *arXiv:2502.11269*, 2025.
- 189 Boyuan Chen, Zhuo Xu, Sean Kirmani, Brain Ichter, Dorsa Sadigh, Leonidas Guibas, and Fei Xia.  
190 Spatialvlm: Endowing vision-language models with spatial reasoning capabilities. In *Proceedings*  
191 *of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 14455–14465,  
192 2024.
- 193 Dedre Gentner. Structure-mapping: A theoretical framework for analogy. *Cognitive Science*, 7  
194 (2):155–170, 1983. ISSN 0364-0213. doi: 10.1016/S0364-0213(83)80009-3. URL <https://www.sciencedirect.com/science/article/pii/S0364021383800093>.
- 196 Mor Geva, Jasmijn Bastings, Katja Filippova, and Amir Globerson. Dissecting recall of factual  
197 associations in auto-regressive language models. *arXiv preprint arXiv:2304.14767*, 2023.
- 198 Bin Han, Robert Wolfe, Anat Caspi, and Bill Howe. Can large language models integrate spatial  
199 data? empirical insights into reasoning strengths and computational weaknesses. *arXiv preprint*  
200 *arXiv:2508.05009*, 2025.
- 201 Peter M. Harman. *The Natural Philosophy of James Clerk Maxwell*. Cambridge University Press,  
202 Cambridge, UK, 1998.
- 203 Michael Hersche, Giacomo Camposampiero, Roger Wattenhofer, Abu Sebastian, and Abbas Rahimi.  
204 Towards learning to reason: Comparing llms with neuro-symbolic on arithmetic relations in abstract  
205 reasoning. *arXiv preprint arXiv:2412.05586*, 2024.
- 206 Martin Jinek, Krzysztof Chylinski, Ines Fonfara, Michael Hauer, Jennifer A. Doudna, and Em-  
207 manuelle Charpentier. A programmable dual-rna-guided dna endonuclease in adaptive bacterial  
208 immunity. *Science*, 337(6096):816–821, August 2012. doi: 10.1126/science.1225829. Epub 2012  
209 Jun 28.
- 210 Adam Karvonen, Benjamin Wright, Can Rager, Rico Angell, Jannik Brinkmann, Logan Smith,  
211 Claudio Mayrink Verdun, David Bau, and Samuel Marks. Measuring progress in dictionary learning  
212 for language model interpretability with board game models. *Advances in Neural Information*  
213 *Processing Systems*, 37:83091–83118, 2024.
- 214 Thomas S. Kuhn. *The Structure of Scientific Revolutions*. University of Chicago Press, Chicago, IL,  
215 1962. ISBN 9780226458086.
- 216 Martha Lewis and Melanie Mitchell. Using counterfactual tasks to evaluate the generality of  
217 analogical reasoning in large language models. *arXiv preprint arXiv:2402.08955*, 2024.
- 218 Kenneth Li, Aspen K Hopkins, David Bau, Fernanda Viégas, Hanspeter Pfister, and Martin Watten-  
219 berg. Emergent world representations: Exploring a sequence model trained on a synthetic task.  
220 *ICLR*, 2023.

- 221 Zichao Li, Yanshuai Cao, and Jackie CK Cheung. Do LLMs build world representations? probing  
222 through the lens of state abstraction. In *The Thirty-eighth Annual Conference on Neural Information*  
223 *Processing Systems*, 2024. URL <https://openreview.net/forum?id=lzfzjYuWgY>.
- 224 Yongsen Mao, Junhao Zhong, Chuan Fang, Jia Zheng, Rui Tang, Hao Zhu, Ping Tan, and Zihan  
225 Zhou. Spatiallm: Training large language models for structured indoor modeling. *arXiv preprint*  
226 *arXiv:2506.07491*, 2025.
- 227 Sam Musker, Alex Duchnowski, Raphaël Millière, and Ellie Pavlick. Llms as models for analogical  
228 reasoning. *Journal of Memory and Language*, 145:104676, 2025.
- 229 Vaishnavh Nagarajan, Chen Henry Wu, Charles Ding, and Aditi Raghunathan. Roll the dice &  
230 look before you leap: Going beyond the creative limits of next-token prediction. In *Forty-second*  
231 *International Conference on Machine Learning*, 2025. URL [https://openreview.net/forum?](https://openreview.net/forum?id=Hi0SyHMmkd)  
232 [id=Hi0SyHMmkd](https://openreview.net/forum?id=Hi0SyHMmkd).
- 233 Aaditya Naik, Jason Liu, Claire Wang, Amish Sethi, Saikat Dutta, Mayur Naik, and Eric Wong.  
234 Dolphin: A programmable framework for scalable neurosymbolic learning. *arXiv preprint*  
235 *arXiv:2410.03348*, 2024.
- 236 Phanish Puranam, Prothit Sen, and Maciej Workiewicz. Can llms help improve analogical rea-  
237 soning for strategic decisions? experimental evidence from humans and gpt-4. *arXiv preprint*  
238 *arXiv:2505.00603*, 2025.
- 239 Cole Robertson and Philip Wolff. Llm world models are mental: Output layer evidence of brittle  
240 world model use in llm mechanical reasoning. *arXiv preprint arXiv:2507.15521*, 2025.
- 241 Alan J. Rocke. *Image and Reality: Kekulé, Kopp, and the Scientific Imagination*. University of  
242 Chicago Press, Chicago, IL, 2010. ISBN 9780226723320.
- 243 Vishwa Shah, Aditya Sharma, Gautam Shroff, Lovekesh Vig, Tirtharaj Dash, and Ashwin Srimi-  
244 vasan. Knowledge-based analogical reasoning in neuro-symbolic latent spaces. *arXiv preprint*  
245 *arXiv:2209.08750*, 2022.
- 246 Chen Shani, Dan Jurafsky, Yann LeCun, and Ravid Shwartz-Ziv. From tokens to thoughts: How llms  
247 and humans trade compression for meaning. *arXiv preprint arXiv:2505.17117*, 2025.
- 248 Parshin Shojaee, Ngoc-Hieu Nguyen, Kazem Meidani, Amir Barati Farimani, Khoa D Doan, and  
249 Chandan K Reddy. Llm-srbench: A new benchmark for scientific equation discovery with large  
250 language models. *arXiv preprint arXiv:2504.10415*, 2025.
- 251 Chenglei Si, Diyi Yang, and Tatsunori Hashimoto. Can llms generate novel research ideas? a  
252 large-scale human study with 100+ nlp researchers. *arXiv preprint arXiv:2409.04109*, 2024.
- 253 Matthieu Tehenan, Christian Bolivar Moya, Tenghai Long, and Guang Lin. Linear spatial world  
254 models emerge in large language models. *arXiv preprint arXiv:2506.02996*, 2025.
- 255 Gerald Tembrevilla, Marina Milner-Bolotin, and Stephen Petrina. Electric fluid to electric current:  
256 The problematic attempts of abstraction to concretization, 2019. URL [https://arxiv.org/](https://arxiv.org/abs/1910.02762)  
257 [abs/1910.02762](https://arxiv.org/abs/1910.02762).
- 258 Keyon Vafa, Justin Y Chen, Ashesh Rambachan, Jon Kleinberg, and Sendhil Mullainathan. Evaluating  
259 the world model implicit in a generative model. *Advances in Neural Information Processing*  
260 *Systems*, 37:26941–26975, 2024.
- 261 Jesse Vig and Yonatan Belinkov. Analyzing the structure of attention in a transformer language  
262 model. *arXiv preprint arXiv:1906.04284*, 2019.
- 263 Guanzhi Wang, Yuqi Xie, Yunfan Jiang, Ajay Mandlikar, Chaowei Xiao, Yuke Zhu, Linxi Fan, and  
264 Anima Anandkumar. Voyager: An open-ended embodied agent with large language models. *arXiv*  
265 *preprint arXiv:2305.16291*, 2023.
- 266 Jiannan Xiang, Tianhua Tao, Yi Gu, Tianmin Shu, Zirui Wang, Zichao Yang, and Zhiting Hu.  
267 Language models meet world models: Embodied experiences enhance language models. *Advances*  
268 *in neural information processing systems*, 36:75392–75412, 2023.

- 269 H. Zhang, H. Deng, J. Ou, et al. Mitigating spatial hallucination in large language models for path plan-  
270 ning via prompt engineering. *Scientific Reports*, 15:8881, 2025. doi: 10.1038/s41598-025-93601-5.  
271 URL <https://doi.org/10.1038/s41598-025-93601-5>.
- 272 Zhonghan Zhao, Wenhao Chai, Xuan Wang, Boyi Li, Shengyu Hao, Shidong Cao, Tian Ye, and  
273 Gaoang Wang. See and think: Embodied agent in virtual environment. In *European Conference*  
274 *on Computer Vision*, pages 187–204. Springer, 2024.
- 275 Wanjun Zhong, Ruixiang Cui, Yiduo Guo, Yaobo Liang, Shuai Lu, Yanlin Wang, Amin Saied,  
276 Weizhu Chen, and Nan Duan. AGIEval: A human-centric benchmark for evaluating foundation  
277 models. In Kevin Duh, Helena Gomez, and Steven Bethard, editors, *Findings of the Association*  
278 *for Computational Linguistics: NAACL 2024*, pages 2299–2314, Mexico City, Mexico, June 2024.  
279 Association for Computational Linguistics. doi: 10.18653/v1/2024.findings-naacl.149. URL  
280 <https://aclanthology.org/2024.findings-naacl.149/>.