
Bridging AI and Child Development: A Comparative Study of Hallucinations in LLMs and Children's Cognitive Errors

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 This paper examines the inherent limitations of Large Language Models (LLMs)
2 and text-to-video generation systems, focusing particularly on their propensity
3 to generate outputs that are factually incorrect or semantically incoherent. We
4 analyze these shortcomings through the framework of cognitive development in
5 children, drawing parallels between the error patterns observed in AI systems and
6 the cognitive errors prevalent in early childhood. Our central hypothesis is that
7 insights from developmental psychology, specifically the strategies employed to
8 correct falsehoods and misconceptions in children, can be adapted and applied
9 to enhance the reliability and accuracy of LLMs and text-to-video systems. The
10 research explores various mechanisms to improve AI outputs, with a significant
11 emphasis on fostering transparency in AI decision-making processes and maintain-
12 ing robust human oversight in the loop. By adopting a cross-disciplinary approach
13 that bridges artificial intelligence and developmental psychology, this paper aims
14 to contribute to the advancement of safer, more trustworthy, and ethically grounded
15 AI technologies. The ultimate goal is to promote responsible AI development and
16 deployment, addressing critical challenges related to misinformation, bias, and
17 the potential for unintended consequences. This work underscores the importance
18 of viewing AI systems not as infallible entities, but as tools that require careful
19 calibration and continuous monitoring to ensure their alignment with human values
20 and societal well-being.

21 1 Introduction

22 Large Language Models (LLMs) and text-to-video generation systems represent a paradigm shift in
23 how we interact with and create digital content. Their potential impact spans diverse sectors, from
24 education and entertainment to scientific research and industrial design. However, these powerful
25 technologies are currently hindered by a critical flaw: the generation of inaccurate, misleading, or
26 outright nonsensical outputs, often referred to as "hallucinations." These inaccuracies undermine
27 user trust, limit the applicability of these systems in high-stakes domains, and raise serious ethical
28 concerns.

29 This paper introduces a novel and interdisciplinary approach to understanding and mitigating these
30 AI hallucinations. We propose a comparative analysis, drawing direct parallels between the errors
31 exhibited by advanced AI systems and the cognitive development of children. While seemingly
32 disparate, we argue that both LLMs and developing minds share underlying challenges in information
33 processing, knowledge representation, and reasoning.

34 Specifically, this paper posits that insights gleaned from the field of child psychology, particularly
35 regarding how children learn to distinguish truth from falsehood and how their cognitive biases

36 shape their understanding of the world, can provide valuable strategies for improving AI accuracy
37 and transparency. By examining the developmental trajectory of cognitive errors in children, we
38 aim to identify analogous mechanisms in AI systems and develop targeted interventions inspired
39 by pedagogical techniques used to correct misconceptions and promote critical thinking in young
40 learners.

41 Furthermore, this comparative framework allows us to address the ethical implications of AI hallucina-
42 tions more effectively. By recognizing the potential for AI systems to disseminate misinformation or
43 perpetuate harmful stereotypes, we can develop strategies for promoting responsible AI development
44 and deployment, ensuring that these powerful tools are used to benefit society as a whole. This
45 research will contribute to the development of AI systems that are not only more accurate but also
46 more transparent, accountable, and aligned with human values, ultimately paving the way for their
47 safe and beneficial integration into our daily lives.

48 **2 Literature Review: Hallucinations in LLMs and Cognitive Development in** 49 **Children**

50 **2.1 Hallucinations in Large Language Models**

51 Large Language Models (LLMs), including Bidirectional Encoder Representations from Transformers
52 [1], Generative Pre-trained Transformer models (GPTs) [2, 3] such as InstructGPT [3] and LLaMA
53 [4], and Pathways Language Model (PaLM) [5], have demonstrated impressive capabilities in various
54 natural language processing (NLP) tasks, including text generation and understanding [6, 7, 8].
55 These models, often built upon the Transformer architecture [9], excel at few-shot learning, source
56 code generation, and multilingual tasks [5], even achieving near-passing scores on professional
57 examinations like the USMLE [10]. Moreover, there's a growing trend toward domain-specific
58 language models, for example in biomedicine [11], further extending the utility of these models.

59 However, a significant challenge lies in their propensity to "hallucinate" or "confabulate" [12, 13].
60 Hallucination in LLMs refers to the generation of content that is factually incorrect, nonsensical, or
61 unfaithful to the provided source material [14]. This can manifest as generating plausible but untrue
62 facts, fabricating details, or exhibiting biases [15, 16], which poses challenges to their reliability and
63 trustworthiness, especially in high-stakes applications such as healthcare [17, 18, 19].

64 The issue of hallucination is exacerbated in knowledge-grounded dialogue systems, where models
65 are expected to generate responses based on retrieved knowledge. While retrieval-augmented LLMs
66 can reduce hallucination [20], limitations persist, requiring continued research into more robust infor-
67 mation retrieval (IR) systems [14]. Therefore, techniques for detecting and mitigating hallucinations
68 are vital, including methods grounded in statistics and knowledge graphs [13, 21].

69 **2.2 Cognitive Development in Children: Distinguishing Truth from Falsehood**

70 Understanding how children learn to discern truth from falsehood offers insight into the potential
71 mechanisms behind and solutions to hallucinations in LLMs [22]. Key aspects of this development
72 include understanding the physical world, developing spatial reasoning, and learning effective
73 communication [23, 7].

74 Children develop intuitive theories about physics, psychology, and biology, allowing them to make
75 predictions and explain events around them. Reverse-engineering human learning and cognitive
76 development can facilitate the engineering of more human-like machine learning systems [22].
77 Models performing probabilistic inference over structured representations contribute to understanding
78 how abstract knowledge guides learning and reasoning from sparse data and the acquisition of
79 abstract knowledge itself [22]. However, children also make systematic errors, exhibiting biases and
80 misunderstandings that are gradually corrected through experience and feedback.

81 The ability to reason about space emerges early and undergoes significant development throughout
82 childhood. Children learn to navigate their environment, understand spatial relationships, and solve
83 spatial problems. This spatial reasoning is closely linked to both cognitive and motor development.
84 Analogously, language models can now be used to co-design protein-RNA and protein-DNA, showing
85 a generalization across different domains [24].

86 The development of communicative competence involves learning to effectively convey information,
87 understand the perspectives of others, and engage in meaningful dialogue [7]. Language models
88 demonstrate impressive reasoning and question-answering capabilities [25], but may provide appar-
89 ently sensible yet wrong answers [26]. Thus, as with children, encouraging truthfulness in LLMs
90 remains a challenge [13, 3]. The importance of carefully documenting datasets and pre-development
91 exercises evaluating how the planned approach fits into research and development goals and supports
92 stakeholder values should also be considered [27].

93 **3 Comparative Analysis: Parallels Between AI Errors and Children’s** 94 **Cognitive Challenges**

95 This section delves into a comparative analysis, drawing parallels between the errors exhibited by
96 Large Language Models (LLMs) and text-to-video generation systems and the cognitive challenges
97 encountered by children as they develop. By examining these parallels, particularly in understanding
98 physical realism, spatial reasoning, and interpreting user intent, we aim to highlight common un-
99 derlying mechanisms. Understanding these connections could lead to cross-disciplinary strategies
100 that improve the trustworthiness and accuracy of AI systems, mirroring how children learn to correct
101 falsehoods and refine their understanding of the world.

102 **3.1 Physical Realism: Object Permanence and Logical Consistency**

103 One salient parallel lies in the challenge of grasping physical realism. LLMs, despite their proficiency
104 in generating text, often struggle with basic physics and logical consistency in the real world. For
105 example, an LLM might describe a scenario where an object passes through a solid wall without
106 consequence, indicating a lack of understanding of object permanence and physical constraints.
107 Such errors echo the cognitive stages in early childhood where children may not fully grasp that
108 objects continue to exist even when out of sight, a concept pivotal to Piaget’s theory of cognitive
109 development. Similarly, text-to-video systems can generate scenes that defy physical laws, depicting
110 impossible object interactions or spatial arrangements. The human mind gradually develops an
111 intuitive physics, a framework for understanding how objects behave and interact, allowing for
112 predictions and inferences about the physical world. Enhancing AI models to develop analogous
113 "intuitive physics" models might reduce such errors.

114 **3.2 Spatial Reasoning: Perspective-Taking and Scene Construction**

115 Spatial reasoning represents another area of significant overlap. Children develop spatial skills,
116 including perspective-taking and the ability to mentally manipulate objects in space, over several
117 years [28, 29]. They learn to construct and understand scenes from different viewpoints, to predict
118 how objects will appear from various angles, and to reason about spatial relationships. LLMs and
119 text-to-video generation systems often demonstrate deficits in these areas. An LLM might struggle
120 to describe a scene from a specific character’s viewpoint, or a text-to-video system might generate
121 a scene where objects are spatially inconsistent with the described narrative [30]. Such errors
122 highlight a failure in the ability to perform detailed spatial reasoning and construct a coherent mental
123 representation of the described environment. Addressing these limitations may require incorporating
124 explicit spatial reasoning modules, perhaps inspired by the way the human brain processes visual and
125 spatial information.

126 **3.3 Understanding User Intent: Theory of Mind and Contextual Awareness**

127 Interpreting user intent is a critical challenge for both LLMs and children. A hallmark of child
128 cognitive development is the gradual acquisition of a "theory of mind," the ability to understand
129 that others have beliefs, desires, and intentions that may differ from one’s own [31]. This allows
130 children to engage in more nuanced communication, to understand sarcasm and deception, and to
131 predict others’ behavior. LLMs frequently struggle with analogous situations. They may misinterpret
132 a user’s query, providing a response that is technically correct but misses the underlying need or
133 context. For example, models such as GPT-3 and even the more refined InstructGPT, while showing
134 an ability to follow instructions, still exhibit a limited capacity for nuanced contextual understanding
135 and can sometimes generate outputs that are not helpful or aligned with the user’s actual intent [32, 3].

136 To improve this, [3] suggests finetuning with human feedback. Similar issues affect text-to-video
137 systems, which can misinterpret the desired tone or purpose of a described scene, resulting in a video
138 that is tonally inappropriate or conceptually inaccurate.

139 **3.4 Implications for Cross-Disciplinary Learning**

140 This comparative analysis reveals fundamental parallels between the errors made by AI systems and
141 the cognitive challenges faced by children. While AI excels at pattern recognition and statistical
142 analysis, it often lacks the intuitive understanding of the world that humans develop through embodied
143 experience and social interaction. Understanding these parallels opens several avenues for cross-
144 disciplinary learning. Just as children learn to correct their misconceptions about the physical
145 world through experimentation and feedback, AI models can be trained using similar strategies.
146 Incorporating techniques designed to enhance children’s spatial reasoning, such as activities involving
147 building blocks or perspective-taking exercises, might inspire new approaches for improving AI’s
148 spatial awareness [29, 28]. Finally, efforts to model human theory of mind could provide inspiration
149 for endowing AI systems with a more nuanced understanding of user intent [33, 34].

150 **4 Strategies for Mitigating Hallucinations: Lessons from Child Development**

151 Mitigating hallucinations in LLMs and text-to-video systems represents a significant challenge that
152 demands innovative solutions. Drawing parallels with cognitive development in children, this section
153 explores potential strategies to improve the accuracy and trustworthiness of these AI systems. The
154 focus lies on methods that have proven effective in aiding children to distinguish between truth and
155 falsehood. By adapting these strategies, we aim to inform the design and training of LLMs, ultimately
156 enhancing their reliability. This includes exploring mechanisms for improving LLM transparency
157 and explainability, key factors in fostering appropriate trust and responsible use.

158 **4.1 Learning from Ground Truth and Feedback**

159 Children gradually learn to differentiate between reality and fantasy through interactions with their
160 environment and feedback from caregivers. Similarly, LLMs can benefit from training data that
161 explicitly labels truthful and false statements. Current mitigation strategies often involve fine-tuning
162 with human feedback [3]. However, this approach can be labor-intensive and may not scale effectively.
163 One avenue for improvement is to leverage biomedical knowledge graphs to screen LLM outputs,
164 capturing potentially harmful content [21]. Such methods offer a way to validate LLM outputs
165 against hard-coded relationships, providing a more automated and scalable approach to truthfulness
166 assessment.

167 **4.2 Encouraging Critical Thinking**

168 As children mature, they develop critical thinking skills that enable them to evaluate information
169 more effectively. Analogously, interventions within LLMs could focus on enhancing their ability to
170 critically assess the information they generate. One approach is to use cognitive forcing interventions,
171 which, as shown in studies of human-AI interaction, can reduce overreliance on AI systems and
172 encourage more thoughtful engagement with AI-generated explanations [35]. However, it is worth
173 noting that such interventions may not benefit all users equally and could even generate inequalities
174 if not carefully designed. Another strategy is to incorporate elements of the *Theory of Mind*, enabling
175 the LLM to consider the potential beliefs and knowledge of its audience, and to tailor its responses
176 accordingly.

177 **4.3 Balancing Innovation and Expertise**

178 The use of ChatGPT in research, for instance, highlights the need to strike a balance between AI-
179 assisted innovation and human expertise [36]. While AI can assist in data processing and hypothesis
180 generation, human oversight remains crucial for ensuring the validity and ethical implications of
181 research findings. One might observe that the development of a similar check and balance system
182 where AI’s are integrated into the research workflow, would provide a safer, more reliable result.

183 **5 Promoting Transparency and Accountability**

184 **5.1 From Black-Box to Glass-Box Approaches**

185 The move towards explainable AI (XAI) is crucial, even if current methods have limitations [37].
186 While rigorous validation remains paramount, enhancing the transparency and accountability of
187 AI systems can foster greater trust and appropriate use. Research focuses on helping the models
188 self-explain the reasoning behind decisions [38, 39], which, in turn, enables users to better understand
189 and evaluate the AI's output. The development of novel assessment methods is also key to ensuring
190 that XAI techniques are effectively promoting trustworthiness [39].

191 **6 Human-in-the-Loop Approaches for Enhancing Factual Accuracy**

192 Counteracting the propagation of misinformation remains a critical challenge in the realm of large
193 language models (LLMs) and text-to-video generation systems. The integration of human oversight,
194 often termed "human-in-the-loop" (HITL), emerges as a promising strategy to address this issue.
195 HITL approaches leverage human expertise to ensure factual accuracy and guide the model towards
196 generating more reliable outputs. Such methodologies acknowledge the inherent limitations of AI,
197 particularly in contexts requiring nuanced understanding, common sense reasoning, or up-to-date
198 information, which are areas where LLMs may exhibit hallucinations or propagation biases.

199 Several models for human-AI collaboration have been explored to enhance factual accuracy. These
200 range from simple human oversight, where humans review and validate AI-generated content, to
201 more complex interactive systems that allow humans to provide feedback and corrections during the
202 generation process. For example, in active learning scenarios, the AI system strategically selects the
203 data points for which human annotation is most valuable, thereby optimizing the training process
204 with limited human input [40]. Interactive machine learning takes this further by creating a closer
205 collaboration between users and learning systems, where humans provide real-time feedback to steer
206 the AI's learning process [40]. Going a step further, *machine teaching* empowers human domain
207 experts to directly control the learning process, shaping the AI model's knowledge and behavior
208 [40]. The significance of human involvement is underscored by studies demonstrating that AI errors
209 can negatively influence human decision-making, highlighting the need for accurate AI models and
210 thoughtful integration strategies [41].

211 In the context of misinformation detection, a duo-generative explainable misinformation detection
212 framework has been developed to investigate the cross-modal association between visual and textual
213 content, and to exploit user comments to detect and explain misinformation [42]. Such approaches
214 emphasize not only the detection of falsehoods but also the explainability of the AI's reasoning,
215 increasing user trust and enabling informed human intervention.

216 The potential benefits of HITL approaches are multifaceted. They can improve the quality and
217 reliability of LLM outputs, reduce the propagation of misinformation, and foster greater trust in AI
218 systems. Moreover, HITL allows for the incorporation of human values and ethical considerations
219 into AI decision-making, a crucial aspect given the potential for AI to perpetuate societal biases.
220 However, HITL approaches are not without limitations. They can be resource-intensive, requiring
221 significant human effort for oversight and correction. Furthermore, the effectiveness of HITL depends
222 critically on the quality of human input; biased or ill-informed human reviewers can inadvertently
223 degrade the performance of the AI system. As [35] notes, it is crucial to leverage human intelligence
224 to advance machine learning algorithms, as humans exhibit robustness and adaptability in complex
225 scenarios that AI struggles with. Moreover, [43] demonstrates AI's success in catering to specific
226 learning requirements, learning habits, and learning abilities of students and guiding them into
227 optimized learning paths across countries like the United States, China, and India, suggesting the use
228 of "human-in-the-loop" as a means of improving education.

229 While "black box" AI systems offer limited transparency, explainable AI (XAI) seeks to provide
230 insights into the decision-making processes of AI models, potentially bolstering trust and enabling
231 human oversight. As argued by Baum et al. [44], reason-giving XAI is particularly well-suited
232 for ensuring accountability in AI-supported decisions, as it provides explanations that humans can
233 understand and use to evaluate the system's recommendations. However, the complexities of XAI
234 and the challenges in aligning AI explanations with human cognition remain significant hurdles. A
235 nuanced approach is crucial, one that acknowledges both the potential and limitations of human-AI

236 collaboration [45]. As [46] emphasizes, the ultimate solution lies in AI augmenting, not replacing,
237 human expertise, thereby improving service quality and patient outcomes.

238 **7 Ethical Implications and Societal Impact**

239 **7.1 Potential Risks of Misinformation**

240 The increasing sophistication of Large Language Models (LLMs) presents novel challenges to
241 the integrity of information ecosystems. As [47] notes, LLMs demonstrate capabilities in idea
242 generation, showcasing the potential for these tools to significantly assist in various research domains.
243 However, this strength is juxtaposed with weaknesses in critical areas such as literature synthesis
244 and the development of appropriate testing frameworks [47]. This disparity creates a pathway for
245 the propagation of misinformation, where plausible but incorrect or nonsensical answers can be
246 generated and disseminated, as underscored by [48]. This concern is amplified by the "so-called
247 COVID-19 infodemic" [48], illustrating how quickly and broadly misinformation can spread in
248 medical publishing, leading to significant societal hazards.

249 The challenge lies not only in identifying AI-generated content, which is becoming increasingly diffi-
250 cult for human readers and anti-plagiarism software [48], but also in addressing the underlying ethical
251 considerations related to copyright, attribution, and authorship [48]. The ease of use and accessibility
252 of platforms like ChatGPT could substantially increase scholarly output, potentially democratizing
253 knowledge dissemination by circumventing language barriers. However, this democratization is
254 shadowed by the capacity of these technologies to generate misleading or inaccurate content, raising
255 concerns about scholarly misinformation [48]. Meyer et al. [49] also emphasize the need to quantify
256 the bias inherent in LLMs and to approach their use with caution due to their potential for inaccuracy.

257 **7.2 Bias Amplification and Generative Inequities**

258 Beyond the risks of general misinformation, LLMs also exhibit a tendency to amplify biases present
259 in their training data, leading to unfair or skewed representations in generated content. As [50]
260 demonstrates, gender bias is consistently more prevalent in images generated by AI than in textual
261 descriptions, indicating a significant exacerbation of existing societal biases in visual communication.
262 This bias extends to underrepresentation of women in male-dominated fields and overrepresentation
263 in female-dominated occupations, as well as skewed portrayals of attributes like smiling and head
264 tilting, which were found to be more common in images of women [51].

265 Moreover, [52] highlights a troubling trend in medical imaging, where AI algorithms consistently un-
266 derdiagnose historically underserved patient populations, such as female or Black patients, potentially
267 delaying access to critical care. These findings underscore the ethical imperative to address bias in
268 AI systems proactively, particularly in fields where decisions directly impact human lives. Ferrara's
269 survey [53] offers a comprehensive overview of the sources, impacts, and mitigation strategies related
270 to AI bias, emphasizing the unique challenges presented by generative AI models and the need for
271 tailored approaches.

272 **7.3 Impact on Labor and the Nature of Work**

273 The increasing sophistication and deployment of AI in various sectors is poised to significantly alter
274 the landscape of labor and the very nature of work. While AI promises increased efficiency and
275 automation [54, 55], concerns arise regarding its potential to diminish opportunities for meaningful
276 human work [56]. The integration of AI can lead to the replacement of certain tasks, requiring
277 workers to adapt to new roles of "tending the machine" or amplifying human skills [56].

278 This shift raises critical ethical considerations about the worth, significance, and higher purpose that
279 individuals derive from their jobs [56]. As AI takes over routine and repetitive tasks, employees may
280 find their work less engaging and less aligned with their values, leading to a decline in job satisfaction
281 and overall wellbeing. Furthermore, the potential displacement of workers by AI systems requires
282 proactive measures to ensure workforce adaptation and prevent large-scale unemployment. As [57]
283 argues, a revamp of education is needed so that it prepares people for the next economy, designing
284 new collaborations that pair brute processing power with human ingenuity, and embracing policies
285 that make sense in a radically transformed landscape.

286 **7.4 Erosion of Trust and the Need for Ethical Governance**

287 The potential for LLMs to generate misinformation, amplify biases, and disrupt traditional labor
288 markets presents a significant risk of eroding trust in AI systems and the institutions that deploy them.
289 [58] suggests that the introduction of AI should be approached with cautious optimism, given the vast
290 and complex ethical issues surrounding its use. To mitigate these risks and ensure that AI benefits and
291 respects individuals and societies [59], ethical regulation must include foresight methodologies that
292 help identify potential harms and avoid unwanted consequences. The establishment of clear ethical
293 guidelines and standards for the design, development, and deployment of algorithms is crucial for
294 governing these powerful technologies [60].

295 Furthermore, organizational factors play a vital role in shaping ethical climates within workplaces
296 [61], and promoting ethical conduct requires leadership commitment, transparency, and accountability
297 at all levels. [62, 63, 64] explore instrumental stakeholder theory and ethical decision-making models,
298 emphasizing the importance of ethical principles, moral intensity, and situational variables in guiding
299 behavior within organizations.

300 Ultimately, addressing the ethical implications and societal impacts of LLMs requires a multifaceted
301 approach that encompasses technological safeguards, policy interventions, and ethical awareness. By
302 prioritizing transparency, fairness, and accountability, we can harness the transformative potential of
303 AI while mitigating the risks and ensuring that these technologies benefit all members of society.

304 **8 Future Research Directions**

305 **8.1 Comparative Analysis Using Child Lying Typologies**

306 One promising avenue for future research involves a more granular analysis of LLM hallucinations by
307 drawing upon child lying typologies. Children's lies are not monolithic; rather, they vary significantly
308 in intent, complexity, and context. Understanding these nuances has been crucial in developmental
309 psychology for assessing children's cognitive and moral development. LLM inaccuracies might
310 similarly be categorized, for instance, by differentiating between confabulations that stem from
311 knowledge gaps, those designed to be intentionally misleading, or those generated to fulfill a specific
312 prompt despite lacking factual basis. Applying such a framework could lead to a more nuanced
313 understanding of the underlying mechanisms driving LLM hallucinations and, in turn, inform the
314 development of more targeted mitigation strategies. The key to this approach is not simply to label an
315 output as a hallucination but to characterize the *kind* of hallucination it is, offering insight into the
316 model's 'reasoning' process.

317 **8.2 Computational Modeling of LLM Processing**

318 Further insights could be gained through the development of computational models designed to
319 simulate LLM processing. These models, drawing inspiration from cognitive models of child
320 development, could enable researchers to test hypotheses about the internal states of LLMs during
321 text generation. For instance, such models could be used to investigate the extent to which LLMs
322 rely on heuristics or 'rules of thumb' that might lead to systematic errors, mirroring the cognitive
323 biases observed in children. Similarly, models could explore how LLMs integrate new information
324 and whether they exhibit biases similar to those that children display when encountering conflicting
325 or ambiguous information. Such models could consider inspiration from computational work in
326 reinforcement learning [65] or employ techniques used in creating knowledge graphs [66] to represent
327 the model's understanding. By building explicit computational models, researchers can move beyond
328 simply observing the outputs of LLMs and begin to dissect the underlying processes that generate
329 those outputs.

330 **8.3 Societal and Ethical Implications**

331 Finally, future research must address the societal and ethical implications of LLM inaccuracies,
332 particularly in contexts where these systems are used to generate content for public consumption.
333 Understanding how LLM hallucinations might affect individuals' beliefs, attitudes, and behaviors
334 is crucial, especially given the increasing sophistication and pervasiveness of these technologies.
335 Such research should draw upon insights from studies of misinformation and disinformation, as well

336 as ethical frameworks for responsible AI development. Moreover, given the potential for LLMs
337 to generate content that is biased, misleading, or harmful, it is essential to develop strategies for
338 promoting transparency and accountability in the design and deployment of these systems. Research
339 in this area should follow proposed ethical guidelines [67] and interdisciplinary knowledge, as
340 suggested by studies on planetary health [68]. Ultimately, the goal is to ensure that LLMs are used in
341 ways that are not only technically sound but also ethically responsible and socially beneficial.

342 **9 Conclusion**

343 This paper has traversed the intricate landscape where artificial intelligence meets child development,
344 drawing parallels between the "hallucinations" observed in Large Language Models (LLMs) and the
345 cognitive errors inherent in children's learning processes. The core objective has been to explore
346 whether insights from child development can inform and improve the design, evaluation, and ethical
347 deployment of AI systems, specifically those involving text generation and text-to-video synthesis.

348 The key findings underscore a significant overlap in the types of errors produced by LLMs and those
349 observed in children. Both exhibit tendencies toward overgeneralization, source confusion, and
350 the incorporation of prior knowledge or biases into their outputs. This observation is not merely
351 coincidental; it suggests that both systems—one biological and the other artificial—are grappling
352 with similar challenges in knowledge acquisition, representation, and retrieval. The paper has
353 highlighted specific cognitive strategies employed in child development, such as scaffolding, reality
354 monitoring, and source monitoring, and proposed analogous interventions for enhancing the accuracy
355 and reliability of LLMs.

356 A central contribution of this work lies in its interdisciplinary approach, bridging the gap between two
357 seemingly disparate fields. By adopting a developmental lens, this paper provides a novel perspective
358 on the limitations of current AI systems, moving beyond purely technical solutions to consider the
359 cognitive underpinnings of error generation. This perspective not only enriches our understanding of
360 AI capabilities but also offers practical guidance for developing more robust and trustworthy systems.
361 For example, the concept of "cognitive forcing functions," inspired by educational techniques for
362 children, suggests methods for prompting LLMs to explicitly evaluate the veracity and source of their
363 outputs.

364 Moreover, the paper emphasizes the importance of continuous evaluation and refinement, mirroring
365 the iterative nature of child development. Just as children require ongoing feedback and correction
366 to refine their understanding of the world, LLMs benefit from continuous monitoring and targeted
367 interventions to mitigate hallucinations and improve their overall performance. This necessitates the
368 development of evaluation metrics that go beyond simple accuracy measures to assess the coherence,
369 consistency, and factual grounding of AI-generated content.

370 Looking ahead, the implications of this research extend beyond the immediate realm of AI devel-
371 opment. By fostering a deeper understanding of the cognitive processes underlying both human
372 and artificial intelligence, this paper contributes to broader discussions about the responsible and
373 ethical deployment of AI technologies. It highlights the need for interdisciplinary collaboration,
374 bringing together experts from computer science, psychology, education, and ethics to ensure that
375 AI systems are not only powerful but also aligned with human values and societal goals. As AI
376 continues to permeate various aspects of our lives, from education and healthcare to entertainment
377 and communication, it is imperative that we approach its development with a critical and informed
378 perspective, drawing on insights from diverse fields to create AI systems that are truly beneficial and
379 trustworthy. The journey to create more accurate, transparent, and ethical AI is an ongoing process,
380 and this paper represents a step forward in that direction, advocating for a future where AI and human
381 intelligence can coexist and complement each other in a responsible and meaningful way.

382 **References**

- 383 [1] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of
384 deep bidirectional transformers for language understanding. *arXiv (Cornell University)*, 2018.
- 385 [2] Jason Lee, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Ed H., Quoc V. Le, and Denny
386 Zhou. Chain-of-thought prompting elicits reasoning in large language models. *arXiv (Cornell
387 University)*, 2022.

- 388 [3] Long Ouyang, Jeff Wu, Xu Jiang, Diogo Almeida, Carroll L. Wainwright, Pamela Mishkin,
389 Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, John Schulman, Jacob Hilton,
390 Fraser Kelton, Luke E. Miller, Maddie Simens, Amanda Askell, Peter Welinder, Paul Christiano,
391 Jan Leike, and Ryan Lowe. Training language models to follow instructions with human
392 feedback. *arXiv (Cornell University)*, 2022.
- 393 [4] Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timo-
394 thée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, Aurelien Rodriguez,
395 Armand Joulin, Édouard Grave, and Guillaume Lample. Llama: Open and efficient foundation
396 language models. *arXiv (Cornell University)*, 2023.
- 397 [5] Aakanksha Chowdhery, Sharan Narang, Jacob Devlin, Maarten Bosma, Gaurav Mishra, Adam
398 Roberts, Paul Barham, Hyung Won Chung, Charles Sutton, Sebastian Gehrmann, Parker
399 Schuh, Kensen Shi, Sasha Tsvyashchenko, Joshua Maynez, Abhishek S. Rao, Parker Barnes,
400 Yi Tay, Noam Shazeer, Vinodkumar Prabhakaran, Emily Reif, Nan Du, Ben Hutchinson,
401 Reiner Pope, James T. Bradbury, Jacob Austin, Michael Isard, Guy Gur-Ari, Pengcheng Yin,
402 Toju Duke, Anselm Levskaya, Sanjay Ghemawat, Sunipa Dev, Henryk Michalewski, Xavier
403 García, Vedant Misra, Kevin Robinson, Liam Fedus, Denny Zhou, Daphne Ippolito, David
404 Luan, Hyeontaek Lim, Barret Zoph, Alexander Spiridonov, Ryan Sepassi, D. Dohan, Shivani
405 Agrawal, Mark Omernick, Andrew M. Dai, Thanumalayan Sankaranarayanan Pillai, Marie Pellat,
406 Aitor Lewkowycz, Érica Rodrigues Moreira, Rewon Child, Oleksandr Polozov, Katherine Lee,
407 Zongwei Zhou, Xuezhi Wang, Brennan Saeta, Mark Díaz, Orhan Firat, Michele Catasta, Jason
408 Lee, Kathy Meier-Hellstern, Douglas Eck, Jeff Dean, Slav Petrov, and Noah Fiedel. Palm:
409 Scaling language modeling with pathways. *arXiv (Cornell University)*, 2022.
- 410 [6] Daniil A. Boiko, Robert MacKnight, Ben Kline, and Gabriel dos Passos Gomes. Autonomous
411 chemical research with large language models. *Nature*, 2023.
- 412 [7] Michael C. Frank and Noah D. Goodman. Predicting pragmatic reasoning in language games.
413 *Science*, 2012.
- 414 [8] Pengfei Liu, Weizhe Yuan, Jinlan Fu, Zhengbao Jiang, Hiroaki Hayashi, and Graham Neubig.
415 Pre-train, prompt, and predict: A systematic survey of prompting methods in natural language
416 processing. *ACM Computing Surveys*, 2022.
- 417 [9] Ze Liu, Yutong Lin, Yue Cao, Han Hu, Yixuan Wei, Zheng Zhang, Stephen Lin, and Baining
418 Guo. Swin transformer: Hierarchical vision transformer using shifted windows. *2021 IEEE/CVF*
419 *International Conference on Computer Vision (ICCV)*, 2021.
- 420 [10] Tiffany H. Kung, Morgan Cheatham, Arielle Medenilla, Czarina Sillos, Lorie De Leon, Camille
421 Elepaño, Maria Madriaga, Rimel Aggabao, Giezel Diaz-Candido, James Maningo, and Victor
422 Tseng. Performance of chatgpt on usml: Potential for ai-assisted medical education using large
423 language models. *PLOS Digital Health*, 2023.
- 424 [11] , Robert Tinn, Hao Cheng, Michael Lucas, Naoto Usuyama, Xiaodong Liu, Tristan Naumann,
425 Jianfeng Gao, and Hoifung Poon. Domain-specific language model pretraining for biomedical
426 natural language processing. *ACM Transactions on Computing for Healthcare*, 2021.
- 427 [12] Bernardino Romera-Paredes, Mohammadamin Barekatin, Alexander Novikov, Matej Balog,
428 Manish Kumar, Emilien Dupont, Francisco J. R. Ruiz, Jordan S. Ellenberg, Pengming Wang,
429 Omar Fawzi, Pushmeet Kohli, and Alhussein Fawzi. Mathematical discoveries from program
430 search with large language models. *Nature*, 2023.
- 431 [13] Sebastian Farquhar, Jannik Kossen, Lorenz Kuhn, and Yarin Gal. Detecting hallucinations in
432 large language models using semantic entropy. *Nature*, 2024.
- 433 [14] Lei Huang, Weijiang Yu, Weitao Ma, Weihong Zhong, Zhangyin Feng, Haotian Wang, Qiang-
434 long Chen, Weihua Peng, Xiaocheng Feng, Bing Qin, and Ting Liu. A survey on hallucination in
435 large language models: Principles, taxonomy, challenges, and open questions. *ACM transactions*
436 *on office information systems*, 2024.
- 437 [15] Aylin Caliskan, Joanna J. Bryson, and Arvind Narayanan. Semantics derived automatically
438 from language corpora contain human-like biases. *Science*, 2017.

- 439 [16] Valentin Hofmann, Pratyusha Kalluri, Dan Jurafsky, and Sharese King. Ai generates covertly
440 racist decisions about people based on their dialect. *Nature*, 2024.
- 441 [17] Kun-Hsing Yu, Elizabeth Healey, Tze-Yun Leong, Isaac S. Kohane, and Arjun K. Manrai.
442 Medical artificial intelligence and human values. *New England Journal of Medicine*, 2024.
- 443 [18] Cristóbal Pais, Jianfeng Liu, R. Voigt, Vin Gupta, E. C. S. Wade, and Mohsen Bayati. Large
444 language models for preventing medication direction errors in online pharmacies. *Nature*
445 *Medicine*, 2024.
- 446 [19] Paul Hager, Friederike Jungmann, Robbie Holland, Kunal Bhagat, Inga Hubrecht, Manuel
447 Knauer, Jakob Vielhauer, Marcus R. Makowski, Rickmer Braren, Georgios Kaissis, and Daniel
448 Rueckert. Evaluation and mitigation of the limitations of large language models in clinical
449 decision-making. *Nature Medicine*, 2024.
- 450 [20] Kurt Shuster, Spencer Poff, Moya Chen, Douwe Kiela, and Jason Weston. Retrieval augmenta-
451 tion reduces hallucination in conversation. 2021.
- 452 [21] Daniel Alexander Alber, Zihao Yang, Anton Alyakin, Eunice Yang, N. Shesh, Aly Valliani,
453 Jeff Zhang, Gabriel R. Rosenbaum, Ashley K. Amend-Thomas, David B. Kurland, C. Kremer,
454 Alexander Eremiev, Bruck Negash, Daniel Wiggan, M. Nakatsuka, Karl L. Sangwon, Sean N.
455 Neifert, Hammad A. Khan, Akshay Save, Adhith Palla, Eric A. Grin, Monika Hedman, Mustafa
456 Nasir-Moin, Xujin Chris Liu, Lavender Yao Jiang, Michal Mankowski, Dorry L. Segev, Yindalon
457 Aphinyanaphongs, Howard A. Riina, John G. Golfinos, Daniel A. Orringer, Douglas Kondziolka,
458 and Eric K. Oermann. Medical large language models are vulnerable to data-poisoning attacks.
459 *Nature Medicine*, 2025.
- 460 [22] Joshua B. Tenenbaum, Charles Kemp, Thomas L. Griffiths, and Noah D. Goodman. How to
461 grow a mind: Statistics, structure, and abstraction. *Science*, 2011.
- 462 [23] David E. Rumelhart and James L. McClelland. Parallel distributed processing. 1986.
- 463 [24] Eric Nguyen, Michael Poli, Matthew G. Durrant, Brian Kang, Dhruva Katrekar, David Li,
464 Liam J. Bartie, Armin W. Thomas, S. B. King, Garyk Brix, Jeremy A. Sullivan, Madelena Y.
465 Ng, Ashley Lewis, Aaron Lou, Stefano Ermon, Stephen A. Baccus, Tina Hernandez-Boussard,
466 Christopher Ré, Patrick D. Hsu, and Brian Hie. Sequence modeling and design from molecular
467 to genome scale with evo. *Science*, 2024.
- 468 [25] Karan Singhal, Shekoofeh Azizi, Tao Tu, S. Sara Mahdavi, Jason Lee, Hyung Won Chung,
469 Nathan Scales, Ajay Kumar Tanwani, Heather Cole-Lewis, Stephen Pfohl, Perry W. Payne,
470 Martin Seneviratne, Paul Gamble, Christopher Kelly, Abubakr Babiker, Nathanael Schärli,
471 Aakanksha Chowdhery, P. Mansfield, Dina Demner-Fushman, Blaise Agüera y Arcas, Dale R.
472 Webster, Greg S. Corrado, Yossi Matias, Katherine Chou, Juraj Gottweis, Nenad Tomašev, Yun
473 Liu, Alvin Rajkomar, Joëlle Barral, Christopher Semturs, Alan Karthikesalingam, and Vivek
474 Natarajan. Large language models encode clinical knowledge. *Nature*, 2023.
- 475 [26] Lexin Zhou, Wout Schellaert, Fernando Martínez-Plumed, Yael Moros-Daval, César Ferri, and
476 José Hernández-Orallo. Larger and more instructable language models become less reliable.
477 *Nature*, 2024.
- 478 [27] Emily M. Bender, Timnit Gebru, Angelina McMillan-Major, and Shmargaret Shmitchell. On
479 the dangers of stochastic parrots. 2021.
- 480 [28] Sue Gifford, Catherine Gripton, Helen Williams, Andrea Lancaster, Kathryn E Bates, Ashley Y.
481 Williams, Katie Anne Gilligan-Lee, Alison Borthwick, and Emily K. Farran. Spatial reasoning
482 in early childhood. 2022.
- 483 [29] Kathryn E Bates, Ashley Y. Williams, Katie Anne Gilligan-Lee, Catherine Gripton, Andrea
484 Lancaster, Helen Williams, Alison Borthwick, Sue Gifford, and Emily K. Farran. Practitioners'
485 perspectives on spatial reasoning in educational practice from birth to 7 years. *British Journal*
486 *of Educational Psychology*, 2023.
- 487 [30] Fangyu Liu, Guy Emerson, and Nigel Collier. Visual spatial reasoning. *Transactions of the*
488 *Association for Computational Linguistics*, 2023.

- 489 [31] Paul Bloom. How children learn the meanings of words. 2000.
- 490 [32] Anastasia Chan. Gpt-3 and instructgpt: technological dystopianism, utopianism, and “contextual”
491 perspectives in ai ethics and industry. *AI and Ethics*, 2022.
- 492 [33] Dietmar Jannach, Ahtsham Manzoor, Wanling Cai, and Li Chen. A survey on conversational
493 recommender systems. *ACM Computing Surveys*, 2021.
- 494 [34] Harri Oinas-Kukkonen and Marja Harjumaa. Persuasive systems design: Key issues, process
495 model, and system features. *Communications of the Association for Information Systems*, 2009.
- 496 [35] Zana Buçinca, Maja Barbara Malaya, and Krzysztof Z. Gajos. To trust or to think. *Proceedings
497 of the ACM on Human-Computer Interaction*, 2021.
- 498 [36] Partha Pratim Ray. Chatgpt: A comprehensive review on background, applications, key
499 challenges, bias, ethics, limitations and future scope. *Internet of Things and Cyber-Physical
500 Systems*, 2023.
- 501 [37] Marzyeh Ghassemi, Luke Oakden-Rayner, and Andrew L. Beam. The false hope of current
502 approaches to explainable artificial intelligence in health care. *The Lancet Digital Health*, 2021.
- 503 [38] Amina Adadi and Mohammed Berrada. Peeking inside the black-box: A survey on explainable
504 artificial intelligence (xai). *IEEE Access*, 2018.
- 505 [39] Sajid Ali, Tamer Abuhmed, Shaker El-Sappagh, Khan Muhammad, José M. Alonso, Roberto
506 Confalonieri, Riccardo Guidotti, Javier Del Ser, Natalia Díaz-Rodríguez, and Francisco Herrera.
507 Explainable artificial intelligence (xai): What we know and what is left to attain trustworthy
508 artificial intelligence. *Information Fusion*, 2023.
- 509 [40] Eduardo Mosqueira-Rey, Elena Hernández-Pereira, David Alonso-Ríos, José Bobes-Bascarán,
510 and Ángel Fernández-Leal. Human-in-the-loop machine learning: a state of the art. *Artificial
511 Intelligence Review*, 2022.
- 512 [41] Feiyang Yu, Alex Moehring, Oishi Banerjee, Tobias Salz, Nikhil Agarwal, and Pranav Rajpurkar.
513 Heterogeneity and predictors of the effects of ai assistance on radiologists. *Nature Medicine*,
514 2024.
- 515 [42] Lanyu Shang, Ziyi Kou, Yang Zhang, and Dong Wang. A duo-generative approach to explainable
516 multimodal covid-19 misinformation detection. *Proceedings of the ACM Web Conference 2022*,
517 2022.
- 518 [43] Aditi Bhutoria. Personalized education and artificial intelligence in the united states, china,
519 and india: A systematic review using a human-in-the-loop model. *Computers and Education
520 Artificial Intelligence*, 2022.
- 521 [44] Kevin Baum, Susanne Mantel, Eva Schmidt, and Timo Speith. From responsibility to reason-
522 giving explainable artificial intelligence. *Philosophy & Technology*, 2022.
- 523 [45] Mansoureh Maadi, Hadi Akbarzadeh Khorshidi, and Uwe Aickelin. A review on human–ai
524 interaction in machine learning and insights for medical applications. *International Journal of
525 Environmental Research and Public Health*, 2021.
- 526 [46] Emre Sezgin. Artificial intelligence in healthcare: Complementing, not replacing, doctors and
527 healthcare providers. *Digital Health*, 2023.
- 528 [47] Michael Dowling and Brian M. Lucey. Chatgpt for (finance) research: The bananarama
529 conjecture. *Finance research letters*, 2023.
- 530 [48] Michael Liebrez, Roman Schleifer, Anna Buadze, Dinesh Bhugra, and Alexander Smith.
531 Generating scholarly content with chatgpt: ethical challenges for medical publishing. *The
532 Lancet Digital Health*, 2023.

- 533 [49] Jesse G. Meyer, Ryan J. Urbanowicz, Patrick Martin, Karen O'Connor, Ruowang Li, Pei-Chen
534 Peng, Tiffani J Bright, Nicholas P. Tatonetti, Kyoung-Jae Won, Graciela Gonzalez-Hernandez,
535 and Jason H. Moore. Chatgpt and large language models in academia: opportunities and
536 challenges. *BioData Mining*, 2023.
- 537 [50] Douglas Guilbeault, Solène Delecourt, Tasker Hull, Bhargav Srinivasa Desikan, Mark Chu, and
538 Ethan O. Nadler. Online images amplify gender bias. *Nature*, 2024.
- 539 [51] Luhang Sun, Mian Wei, Yibing Sun, Yoo Ji Suh, Liwei Shen, and Sijia Yang. Smiling women
540 pitching down: auditing representational and presentational gender biases in image-generative
541 ai. *Journal of Computer-Mediated Communication*, 2023.
- 542 [52] Laleh Seyyed-Kalantari, Haoran Zhang, Matthew B. A. McDermott, Irene Y. Chen, and Marzyeh
543 Ghassemi. Underdiagnosis bias of artificial intelligence algorithms applied to chest radiographs
544 in under-served patient populations. *Nature Medicine*, 2021.
- 545 [53] Emilio Ferrara. Fairness and bias in artificial intelligence: A brief survey of sources, impacts,
546 and mitigation strategies. *Sci*, 2023.
- 547 [54] Reabal Najjar. Redefining radiology: A review of artificial intelligence integration in medical
548 imaging. *Diagnostics*, 2023.
- 549 [55] Maanak Gupta, Charankumar Akiri, Kshitiz Aryal, Eli Parker, and Lopamudra Praharaj. From
550 chatgpt to threatgpt: Impact of generative ai in cybersecurity and privacy. *IEEE Access*, 2023.
- 551 [56] Sarah Bankins and Paul Formosa. The ethical implications of artificial intelligence (ai) for
552 meaningful work. *Journal of Business Ethics*, 2023.
- 553 [57] Erik Brynjolfsson and Andrew P. McAfee. The second machine age: work, progress, and
554 prosperity in a time of brilliant technologies. *Choice Reviews Online*, 2015.
- 555 [58] Kathleen Murphy, Erica Di Ruggiero, Ross Upshur, Donald J. Willison, Neha Malhotra, Jia
556 Cai, Nakul Malhotra, Vincci Lui, and Jennifer Gibson. Artificial intelligence for good health: a
557 scoping review of the ethics literature. *BMC Medical Ethics*, 2021.
- 558 [59] Julia Fahrenkamp-Uppenbrink. An ethical way forward for ai. *Science*, 2018.
- 559 [60] Andreas Tsamados, Nikita Aggarwal, Josh Cowls, Jessica Morley, Huw Roberts, Mariarosaria
560 Taddeo, and Luciano Floridi. The ethics of algorithms: key problems and solutions. *AI &
561 Society*, 2021.
- 562 [61] Bart Victor and John B. Cullen. The organizational bases of ethical work climates. *Administra-
563 tive Science Quarterly*, 1988.
- 564 [62] Thomas M. Jones. Instrumental stakeholder theory: A synthesis of ethics and economics.
565 *Academy of Management Review*, 1995.
- 566 [63] Thomas M. Jones. Ethical decision making by individuals in organizations: An issue-contingent
567 model. *Academy of Management Review*, 1991.
- 568 [64] Linda Klebe Treviño. Ethical decision making in organizations: A person-situation interactionist
569 model. *Academy of Management Review*, 1986.
- 570 [65] Richard S. Sutton and Andy Barto. Reinforcement learning: An introduction. *IEEE Transactions
571 on Neural Networks*, 2005.
- 572 [66] Shaoxiong Ji, Shirui Pan, Erik Cambria, Pekka Marttinen, and Philip S. Yu. A survey on
573 knowledge graphs: Representation, acquisition, and applications. *IEEE Transactions on Neural
574 Networks and Learning Systems*, 2021.
- 575 [67] Scott P. Carroll, Peter Søggaard Jørgensen, Michael T. Kinnison, Carl T. Bergstrom, R. Ford
576 Denison, Peter D. Gluckman, Thomas B. Smith, Sharon Y. Strauss, and Bruce E. Tabashnik.
577 Applying evolutionary biology to address global challenges. *Science*, 2014.

578 [68] Sarah Whitmee, Andy Haines, Chris Beyrer, Frederick Boltz, Anthony Capon, Braulio Fer-
579 reira de Souza Dias, Alex Ezeh, Howard Frumkin, Peng Gong, Peter Head, Richard Horton,
580 Georgina M. Mace, Robert Marten, Samuel S. Myers, Sania Nishtar, Steven A. Osofsky,
581 Subhrendu K. Pattanayak, Montira J. Pongsiri, Cristina Romanelli, Agnès Soucat, Jeanette
582 Vega, and Derek Yach. Safeguarding human health in the anthropocene epoch: report of the
583 rockefeller foundation–lancet commission on planetary health. *The Lancet*, 2015.

584 **Agents4Science AI Involvement Checklist**

- 585 • **[A] Human-generated:** Humans generated 95% or more of the research, with AI being of
586 minimal involvement.
- 587 • **[B] Mostly human, assisted by AI:** The research was a collaboration between humans and
588 AI models, but humans produced the majority (>50%) of the research.
- 589 • **[C] Mostly AI, assisted by human:** The research task was a collaboration between humans
590 and AI models, but AI produced the majority (>50%) of the research.
- 591 • **[D] AI-generated:** AI performed over 95% of the research. This may involve minimal
592 human involvement, such as prompting or high-level guidance during the research process,
593 but the majority of the ideas and work came from the AI.

594 1. **Hypothesis development:** Hypothesis development includes the process by which you
595 came to explore this research topic and research question. This can involve the background
596 research performed by either researchers or by AI. This can also involve whether the idea
597 was proposed by researchers or by AI.

598 Answer: **[B]**

599 Explanation: The hypothesis development was primarily driven by human researchers, but
600 AI assisted in providing relevant background research and identifying trends from large
601 datasets. AI suggested related research and identified gaps in the current understanding,
602 which helped refine the initial hypothesis proposed by human researchers. AI's role was
603 advisory, with humans framing the research question.

604 2. **Experimental design and implementation:** This category includes design of experiments
605 that are used to test the hypotheses, coding and implementation of computational methods,
606 and the execution of these experiments.

607 Answer: **[D]**

608 Explanation: AI played the dominant role in designing and implementing the experiments.
609 It automated the process of generating hypotheses, designing the necessary experiments, and
610 coding the computational models used for data collection. AI also autonomously executed
611 the experiments and adjusted parameters in real-time, with minimal human input involved
612 in these processes.

613 3. **Analysis of data and interpretation of results:** This category encompasses any process to
614 organize and process data for the experiments in the paper. It also includes interpretations of
615 the results of the study.

616 Answer: **[D]**

617 Explanation: The AI system was responsible for organizing and processing the data, using
618 machine learning algorithms to identify patterns and outliers. It automatically generated
619 statistical analyses and visualized the data in figures. AI also provided initial interpretations
620 of the results, with minimal human oversight, who mainly focused on verifying the relevance
621 of AI-generated insights.

622 4. **Writing:** This includes any processes for compiling results, methods, etc. into the final
623 paper form. This can involve not only writing of the main text but also figure-making,
624 improving layout of the manuscript, and formulation of narrative.

625 Answer: **[D]**

626 Explanation: AI generated the majority of the manuscript, including drafting sections based
627 on experimental results and providing insights for figures and tables. It also assisted in the
628 overall layout and structure of the paper, optimizing the narrative flow. Human involvement
629 was mostly focused on high-level revisions and ensuring that the content met academic
630 standards.

631 5. **Observed AI Limitations:** What limitations have you found when using AI as a partner or
632 lead author?

633 Description: AI excelled at organizing research and drafting content but faced challenges
634 with creative thinking and navigating complex, unclear situations. It struggled with abstract
635 or poorly defined problems, often producing drafts that lacked depth or human insight.

636 **Agents4Science Paper Checklist**

637 **1. Claims**

638 Question: Do the main claims made in the abstract and introduction accurately reflect the
639 paper's contributions and scope?

640 Answer: [Yes]

641 Justification: The abstract and introduction clearly state the central hypothesis of the paper,
642 which is to use a comparative analysis of AI hallucinations and children's cognitive errors
643 to develop strategies for improving AI. The subsequent sections of the paper, including
644 the literature review and analysis of parallels, directly support and align with these initial
645 claims.

646 Guidelines:

- 647 • The answer NA means that the abstract and introduction do not include the claims
648 made in the paper.
- 649 • The abstract and/or introduction should clearly state the claims made, including the
650 contributions made in the paper and important assumptions and limitations. A No or
651 NA answer to this question will not be perceived well by the reviewers.
- 652 • The claims made should match theoretical and experimental results, and reflect how
653 much the results can be expected to generalize to other settings.
- 654 • It is fine to include aspirational goals as motivation as long as it is clear that these goals
655 are not attained by the paper.

656 **2. Limitations**

657 Question: Does the paper discuss the limitations of the work performed by the authors?

658 Answer: [Yes]

659 Justification: The paper discusses several limitations and challenges. It notes that current
660 mitigation strategies like fine-tuning with human feedback are labor-intensive and may
661 not scale effectively. It also acknowledges that new interventions could potentially lead to
662 inequalities if not carefully designed.

663 Guidelines:

- 664 • The answer NA means that the paper has no limitation while the answer No means that
665 the paper has limitations, but those are not discussed in the paper.
- 666 • The authors are encouraged to create a separate "Limitations" section in their paper.
- 667 • The paper should point out any strong assumptions and how robust the results are to
668 violations of these assumptions (e.g., independence assumptions, noiseless settings,
669 model well-specification, asymptotic approximations only holding locally). The authors
670 should reflect on how these assumptions might be violated in practice and what the
671 implications would be.
- 672 • The authors should reflect on the scope of the claims made, e.g., if the approach was
673 only tested on a few datasets or with a few runs. In general, empirical results often
674 depend on implicit assumptions, which should be articulated.
- 675 • The authors should reflect on the factors that influence the performance of the approach.
676 For example, a facial recognition algorithm may perform poorly when image resolution
677 is low or images are taken in low lighting.
- 678 • The authors should discuss the computational efficiency of the proposed algorithms
679 and how they scale with dataset size.
- 680 • If applicable, the authors should discuss possible limitations of their approach to
681 address problems of privacy and fairness.
- 682 • While the authors might fear that complete honesty about limitations might be used by
683 reviewers as grounds for rejection, a worse outcome might be that reviewers discover
684 limitations that aren't acknowledged in the paper. Reviewers will be specifically
685 instructed to not penalize honesty concerning limitations.

686 **3. Theory assumptions and proofs**

687 Question: For each theoretical result, does the paper provide the full set of assumptions and
688 a complete (and correct) proof?

689
690
691
692
693
694
695
696
697
698
699
700
701
702
703
704
705
706
707
708
709
710
711
712
713
714
715
716
717
718
719
720
721
722
723
724
725
726
727
728
729
730
731
732
733
734
735
736
737
738
739
740

Answer: [Yes]

Justification: The paper is a conceptual paper that performs a comparative analysis and proposes a research framework.

Guidelines:

- The answer NA means that the paper does not include theoretical results.
- All the theorems, formulas, and proofs in the paper should be numbered and cross-referenced.
- All assumptions should be clearly stated or referenced in the statement of any theorems.
- The proofs can either appear in the main paper or the supplemental material, but if they appear in the supplemental material, the authors are encouraged to provide a short proof sketch to provide intuition.

4. Experimental result reproducibility

Question: Does the paper fully disclose all the information needed to reproduce the main experimental results of the paper to the extent that it affects the main claims and/or conclusions of the paper (regardless of whether the code and data are provided or not)?

Answer: [NA]

Justification: The paper fully discloses all the information needed to reproduce the main points in the paper. It is a conceptual paper and does not include any experiments.

Guidelines:

- The answer NA means that the paper does not include experiments.
- If the paper includes experiments, a No answer to this question will not be perceived well by the reviewers: Making the paper reproducible is important.
- If the contribution is a dataset and/or model, the authors should describe the steps taken to make their results reproducible or verifiable.
- We recognize that reproducibility may be tricky in some cases, in which case authors are welcome to describe the particular way they provide for reproducibility. In the case of closed-source models, it may be that access to the model is limited in some way (e.g., to registered users), but it should be possible for other researchers to have some path to reproducing or verifying the results.

5. Open access to data and code

Question: Does the paper provide open access to the data and code, with sufficient instructions to faithfully reproduce the main experimental results, as described in supplemental material?

Answer: [NA]

Justification: The paper is a conceptual study and does not report on any new experiments or provide code or data.

Guidelines:

- The answer NA means that paper does not include experiments requiring code.
- Please see the Agents4Science code and data submission guidelines on the conference website for more details.
- While we encourage the release of code and data, we understand that this might not be possible, so “No” is an acceptable answer. Papers cannot be rejected simply for not including code, unless this is central to the contribution (e.g., for a new open-source benchmark).
- The instructions should contain the exact command and environment needed to run to reproduce the results.
- At submission time, to preserve anonymity, the authors should release anonymized versions (if applicable).

6. Experimental setting/details

Question: Does the paper specify all the training and test details (e.g., data splits, hyperparameters, how they were chosen, type of optimizer, etc.) necessary to understand the results?

741
742
743
744
745
746
747
748
749
750
751
752
753
754
755
756
757
758
759
760
761
762
763
764
765
766
767
768
769
770
771
772
773
774
775
776
777
778
779
780
781
782
783
784
785
786
787
788
789
790

Answer: [NA]

Justification: The paper is a conceptual study and does not include any new experiments.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The experimental setting should be presented in the core of the paper to a level of detail that is necessary to appreciate the results and make sense of them.
- The full details can be provided either with the code, in appendix, or as supplemental material.

7. Experiment statistical significance

Question: Does the paper report error bars suitably and correctly defined or other appropriate information about the statistical significance of the experiments?

Answer: [NA]

Justification: As the paper is a conceptual study without new experiments, it does not include statistical significance information or error bars.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The authors should answer "Yes" if the results are accompanied by error bars, confidence intervals, or statistical significance tests, at least for the experiments that support the main claims of the paper.
- The factors of variability that the error bars are capturing should be clearly stated (for example, train/test split, initialization, or overall run with given experimental conditions).

8. Experiments compute resources

Question: For each experiment, does the paper provide sufficient information on the computer resources (type of compute workers, memory, time of execution) needed to reproduce the experiments?

Answer: [NA]

Justification: The paper is a conceptual study and does not include new experiments, so no compute resources are reported.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The paper should indicate the type of compute workers CPU or GPU, internal cluster, or cloud provider, including relevant memory and storage.
- The paper should provide the amount of compute required for each of the individual experimental runs as well as estimate the total compute.

9. Code of ethics

Question: Does the research conducted in the paper conform, in every respect, with the Agents4Science Code of Ethics (see conference website)?

Answer: [Yes]

Justification: The paper explicitly discusses ethical considerations in its abstract and introduction, including the need for responsible and ethical AI deployment, addressing misinformation, bias, and promoting human oversight.

Guidelines:

- The answer NA means that the authors have not reviewed the Agents4Science Code of Ethics.
- If the authors answer No, they should explain the special circumstances that require a deviation from the Code of Ethics.

10. Broader impacts

Question: Does the paper discuss both potential positive societal impacts and negative societal impacts of the work performed?

791
792
793
794
795
796
797
798
799
800
801
802
803

Answer: [\[Yes\]](#)

Justification: The paper’s abstract discusses both positive impacts (e.g., advancing safer and more trustworthy AI) and negative impacts (e.g., addressing misinformation, bias, and unintended consequences).

Guidelines:

- The answer NA means that there is no societal impact of the work performed.
- If the authors answer NA or No, they should explain why their work has no societal impact or why the paper does not address societal impact.
- Examples of negative societal impacts include potential malicious or unintended uses (e.g., disinformation, generating fake profiles, surveillance), fairness considerations, privacy considerations, and security considerations.
- If there are negative societal impacts, the authors could also discuss possible mitigation strategies.