

---

# 000 AGENTIC AI IN THE WILD: 001 002 FROM HALLUCINATIONS TO RELIABLE AUTONOMY 003 004 005

006 **Website:** <https://hallucination-reliable-agentic-AI.github.io>  
007

## 008 1 SUMMARY AND GOALS 009

010 **When we delegate tasks to AI agents—can we count on them to get it right?** Agentic AI  
011 systems are increasingly stepping beyond static generation tasks into autonomous decision-making:  
012 scheduling meetings, booking travel, managing workflows, and assisting in scientific research. In  
013 these contexts, *reliability is not just important—it is essential*. Yet today’s foundation models remain  
014 prone to a critical failure mode: hallucination, where outputs are factually incorrect, semantically  
015 implausible, or detached from reality (Maynez et al., 2020; Ji et al., 2023; Wachter et al., 2024;  
016 Huang et al., 2025). While hallucinations are concerning in any generative system, these challenges  
017 are *amplified in agentic settings*, where models execute sequences of decisions without continuous  
018 human oversight.

019 The rise of agentic AI marks a critical inflection point: we are no longer simply interacting with  
020 models—we are increasingly entrusting them to act on our behalf. Unlike standard prompt-response  
021 interactions, agentic systems introduce temporal and operational autonomy, where reasoning and  
022 execution unfold independently of the user. Consider a scenario where a user instructs an AI agent to  
023 secure conference registration by a specific deadline. The agent must interpret requirements, research  
024 available options, make selections, and complete transactions. However, if the agent becomes  
025 unreliable during any part of this process—by fabricating confirmation details, misinterpreting the  
026 user’s intent, or falsely reporting task completion—the failure may go unnoticed until much later,  
027 resulting in costly and irreversible consequences.

028 **Topics and key questions.** Despite the growing deployment of such systems, we lack a foundational  
029 understanding of what makes them reliable. This workshop will bring together researchers across  
030 machine learning, NLP, HCI, robotics, and AI safety to explore algorithmic advances, formal  
031 guarantees, empirical insights, and design frameworks for improving agent reliability. We welcome  
032 contributions on topics including hallucination detection, uncertainty estimation, safe planning under  
033 ambiguity, evaluation benchmarks for agents, and interaction protocols for trustworthy human-agent  
034 collaboration. Key questions include:

- 036 • How do hallucinations manifest in multi-step, autonomous agent workflows—and how can  
037 we detect them before harm is done?
- 038 • What forms of uncertainty quantification can forecast model failure in real-time agentic  
039 execution?
- 040 • How can we formalize “reliability” in agentic contexts where success is long-horizon and  
041 ambiguous?
- 042 • What strategies allow agents to defer, escalate, or communicate uncertainty effectively to  
043 human users?
- 044 • What types of benchmarks, platforms and metrics would be sufficient for documenting  
045 unreliability in agentic systems in the wild?

047 **Societal impact.** The rise of autonomous AI agents raises profound societal stakes. In domains such  
048 as healthcare, law, scientific research, and education, hallucinated outputs are not merely technical  
049 failures—they can lead to misinformation, loss of resources, legal liability, and erosion of public trust.  
050 As AI systems shift from tools we supervise to agents we rely on, reliability becomes an urgent public  
051 concern. Equipping agents with the ability to recognize their own uncertainty, avoid false confidence,  
052 and defer when necessary is essential for their safe deployment in high-stakes environments. This  
053 workshop addresses not just a research gap, but a societal responsibility: ensuring that future agentic  
AI systems are not only capable, but cautious, transparent, and trustworthy.

---

## 054 2 INVITED SPEAKERS 055

056 The invited speakers were selected based on their relevance and prior contributions to the field. Each  
057 speaker has authored high-impact publications in machine learning.

058 **Dawn Song (UC Berkeley)** [*confirmed*] is a Professor at the University of California, Berkeley. She  
059 is a leading expert in AI, security, and safety. Dawn has received numerous accolades, including the  
060 Alfred P. Sloan Research Fellowship, ACM Fellow, and a MacArthur Fellow. She has recently turned  
061 her focus to the reliability and safety of agentic AI, including foundational work on verifiable and  
062 controllable AI agents, which aligns well with the theme of the workshop.  
063

064 **Mohan Kankanhalli (National University of Singapore)** [*confirmed*] is the Director of NUS AI  
065 Institute, Provost's Chair Professor and is also the Deputy Executive Chairman of AI Singapore  
066 (Singapore's national AI program). Mohan leads major initiatives on human-centric and reliable AI,  
067 including deploying agentic systems in real-world environments.  
068

069 **Hamed Hassani (University of Pennsylvania)** [*confirmed*] is an associate professor at the University  
070 of Pennsylvania, and visiting faculty researcher at Google Research. His works span across LLM  
071 safety and foundations of ML and information theory, and therefore can offer unique perspectives on  
072 uncertainty quantification for LLM agents.  
073

074 **Florian Buettner (Goethe-University Frankfurt, and German Cancer Research Center)** [*con-*  
075 *firmed*] is a professor interested in increasing the precision of ML tools for oncology. Florian has  
076 worked on uncertainty estimation to make models more reliable.  
077

078 **James Zou (Stanford University)** [*confirmed*] works on reliable machine learning with foundation  
079 models and is particularly interested for their applications to the health domain.  
080

081 **Stefano Soatto (UCLA and Vice president of AWS)** [*confirmed*] works on representation and  
082 optimal control, including the control of AI agents and their reliability.  
083

084 Besides the confirmed speakers, there is a list of tentative speakers that we can invite in case of  
085 cancellations (our original schedule posted below is complete):  
086

- 087 • **Csaba Szepesvári** (University of Alberta and team lead at Deepmind, sequential decision  
088 making, uncertainty).
- 089 • **Graham Neubig** (Carnegie Mellon University, question answering, code generation, hallucina-  
090 tion of LLMs).
- 091 • **Sandra Wachter** (University of Oxford, factuality of LLMs).
- 092 • **Yu Su** (Ohio State University, NLP and language agents)

## 093 3 SCHEDULE, CALL FOR PAPERS AND REVIEWING PROCESS 094

095 The workshop will span a **full day** and include a mix of invited talks, spotlight presentations of  
096 selected accepted papers, poster sessions for all accepted submissions, and panel discussion. **Three**  
097 poster sessions will be organized, following the positive feedback from our ICLR workshop (Chrysos  
098 et al., 2025) where we experimented with increased number of poster sessions. This led to an  
099 increased interaction and organic discussions between the interested participants. The tentative  
100 schedule is outlined below.  
101

102 **Call for papers and reviewing:** We invite original research (theoretical, empirical) submissions  
103 on **trustworthy agentic AI and hallucinations**. The process will be managed on *Open-*  
104 *Review* with a **double-blind** policy. The timeline (Table 1) follows the official ICLR sched-  
105 ule. Each paper will receive  $\geq 3$  **expert reviews** plus a meta-review, evaluated on rigor, nov-  
106 elty, impact, reproducibility and relevance to our program on agentic AI and hallucinations. To  
107 ensure a high-quality and fair process, we will enforce a strict conflict-of-interest policy, cap  
108 reviewer workloads. We will select a balanced program across methods, domains, and modal-  
109 ities. We will also conduct targeted outreach (open call, mailing lists, affinity groups), and  
110 reserve slots for replications/negative results. We also have initiatives to foster an inclusive  
111 community, including a **peer mentorship program** for new reviewers and awards for **best**  
112 **paper and best artifact**. To maintain neutrality, the organizers will not give invited talks.  
113

start	duration	event	theme
9:00	0:10	opening remarks	
9:10	0:35	invited talk: Florian Buettner	
9:45	0:35	invited talk: Hamed Hassani	
10:20	1:00	poster session I	Hallucination in agentic systems
11:20	0:35	invited talk: Dawn Song	
11:55	1:00	lunch	
12:55	1:00	poster session II	
13:55	0:30	invited talk: Stefano Soatto	
14:25	0:35	invited talk: James Zou	Towards reliable AI agents
15:00	0:30	invited talk: Mohan Kankanhalli	
15:30	1:00	poster session III	
16:30	0:10	Best paper award	
16:40	1:00	panel	
17:40	0:10	closing remarks	

120

121

122 Similarly, for conflicts of interest, e.g., paper from own group or PhD group or the  
123 same organization, the organizers will recuse themselves from the decision. Workshop  
124 reviewers will not give invited talks to avoid conflict of interest. We will **establish a tiny track of papers** - as workshop  
125 organizers we did this in our ICLR'25 workshop as well.

126

127

## 4 AUDIENCE, ACCESS AND INTERACTIONS

128

129

130 **Audience:** The topic of agentic AI and hallucinations in foundation models has received increasing  
131 attention, with numerous papers accepted at ICLR, ICML, and NeurIPS over the past two years.  
132 Hallucination has also attracted the attention in the popular media (Weisse & Metz, 2023; Weiser,  
133 2023; Kaye, 2023; Verma & Oremus, 2023; Kelsey-Sugg & Carrick, 2024; Legg & McNamara, 2024;  
134 Brittain, 2025; Rahman-Jones, 2025; Weaver, 2025), as explicitly identified by judges (IANS, 2023) or  
135 in the medical domain (Hughes, 2024), with the problem becoming worse with stronger models (Metz  
136 & Weisse, 2025). Consequently, we anticipate substantial interest from the community in a dedicated  
137 workshop on this subject. We expect participation from both early career researchers and students, as  
138 well as established researchers seeking fresh insights into this emerging field. Therefore, we estimate  
139 approximately 300 participants for the workshop.

140

141

142 **Website and Access:** Our website <https://hallucination-reliable-agentic-AI.github.io>.  
143 will release all accepted papers at least one week prior to the workshop, allowing  
144 participants to review the papers in advance. Following ICLR guidelines, our primary focus will be  
145 on in-person attendance. However, to accommodate extenuating circumstances such as visa issues or  
146 other rare exceptions, we will allow online participation for people that have a valid reason.

147

148

149 **Cross-community interaction:** We will intentionally engineer mixing between machine learning,  
150 HCI, planning/agents, LLM and safety communities: (i) paper–discussant pairing that assigns each  
151 accepted paper a discussant from a different field, with a short cross-disciplinary response after  
152 the talk/poster; (ii) guided poster walks where moderators lead mixed-background groups through  
153 themed clusters (tool use, evaluation, human studies, theory), using prompt cards to surface divergent  
154 assumptions; (iii) cross-review: at least one reviewer per paper comes from an adjacent area; and (iv)  
155 ongoing Slack/Discord with channels by theme (ML and AI, theory, human studies, planning, eval)  
156 to coordinate post-workshop projects and shared task proposals.

157

158

159 **5 RELATED WORKSHOPS**  
160 Beyond ICLR, there have been no workshops on reliable agentic AI at ICML or NeurIPS. The  
161 most closely related event was the ICLR workshop on uncertainty quantification in foundation  
162 models (Chrysos et al., 2025). That workshop shared some organizers as the current proposal, but had a  
163 different focus. Many of the organizers of the previous workshop focus on uncertainty quantification,  
164 which is different than the present proposal. The theme of the workshop is also different with the  
165 current proposal focusing on the significant and emerging topic of agentic AI that was not explored or

Table 1: Timeline for contributed work submissions

December 5 <sup>th</sup> (2025)	Portal opens for submissions
December 15 <sup>th</sup> (2025)	Mentoring sessions
January 30 <sup>th</sup> (2026)	Submission deadline
February 20 <sup>th</sup> (2026)	Reviewing deadline
February 26 <sup>th</sup> (2026)	Notification date
April 20 <sup>th</sup> (2026)	Camera-ready deadline

162 mentioned at all in ICLR’25. The ICLR’25 workshop accepted 46 papers and experienced significant  
163 attendance, with invited talks drawing audiences beyond the room’s 200-person capacity. It briefly  
164 mentioned hallucinations in text models (in the panel questions) and did not cover AI agents. By  
165 contrast, our workshop puts higher emphasis on agentic settings, where the agentic AI systems can per-  
166 ceive, plan, and act with tools, memory, and multi-step reasoning. This moves uncertainty beyond next-  
167 token probability to the whole decision pipeline—reasoning traces, tool calls, retrieved evidence, and  
168 downstream actions, which is both fundamentally significant and practically important to the broader  
169 audience. As such, we do not consider that there is a significant overlap with the proposed content.  
170

171 Other related workshops include workshops on LLM agents for science (Chen & et al., 2024;  
172 Koutra & et al., 2025), reinforcement learning agents (Jiang & et al., 2022; de Witt & et al., 2023),  
173 out-of-distribution learning (Deshmukh & et al., 2023), or specific types of generative models (Thana-  
174 palasingam & et al., 2023). The closest workshops are those on trustworthy ML (Papernot & et al.,  
175 2020; Beirami & et al., 2021; Xiao & et al., 2022; Cheng & et al., 2023; Bansal & et al., 2024).  
176 None of them focus on hallucinations and they mostly refer to adversarial robustness, which is a  
177 more narrow topic than the questions we will face with agentic systems. Therefore, we believe the  
178 proposed workshop does not have a significant overlap with the past workshops.  
179

## 6 DIVERSITY OF TEAM AND TOPICS

180 **Organizing Committee Diversity:** The organizing committee reflects a broad range of backgrounds,  
181 with deliberate attention to diversity in gender, career stage, institutional affiliation, and geographic  
182 representation. The core team consists of 1 applied scientist, 2 assistant professors, 1 associate  
183 professor and 1 research professor. Two additional volunteers will assist with workshop logistics.  
184 The committee includes 5 women (representing 70% of the organizing committee), underscoring a  
185 commitment to gender representation. Members are affiliated with institutions in the United States,  
186 and Singapore, spanning academia, research institutes and industry. This composition is intended to  
187 ensure a balanced and inclusive environment.  
188

189 **Invited Speaker Diversity** The invited speakers represent a mix of senior and early-career researchers.  
190 Their contributions cover both theoretical and empirical research with practical applications. A  
191 common focus across all speakers is the development of reliable and trustworthy AI systems for  
192 deployment in real-world contexts, while the invited speakers are affiliated with institutes in both US  
193 and Europe.  
194

195 **Theme Diversity** The workshop is structured to promote interdisciplinary engagement, drawing  
196 on expertise from applied mathematics, statistics, computer science, HCI, and machine learning.  
197 Emphasis will be placed on the practical integration of these disciplines, with the goal of fostering  
198 new collaborations and advancing understanding across traditional boundaries.  
199

## 7 ORGANIZERS

200 Below, we provide further details on the core team and the program committee of the workshop.  
201

### 7.1 CORE TEAM

202 The organizing team has made significant contributions to various areas of out-of-distribution pre-  
203 diction, trustworthy models, uncertainty quantification and generative models. In addition to their  
204 scholarly work, they bring extensive organizational experience, having led more than a **fifteen work-  
205 shops and tutorials** over the past decade. These efforts have successfully engaged a wide range of  
206 sub-communities within the broader machine learning. Several members of the team have also held  
207 prominent leadership roles, including **serving as Program Chairs for ICML**.  
208

209 **Grigoris Chrysos (UW-Madison)** is an Assistant Professor at the University of Wisconsin-Madison.  
210 Before that, Grigoris was a postdoctoral fellow at EPFL following the completion of his PhD at  
211 Imperial College London. Previously, he graduated from the National Technical University of  
212 Athens with a Diploma/MEng in Electrical and Computer Engineering. Grigoris has co-organized  
213 workshops in top conferences (**ICCV’15, CVPR’17, ICCV’17, NeurIPS’24, AAAI’25, ICLR’25**).  
214 The most recent NeurIPS workshop (**‘Fine-Tuning in Modern Machine Learning: Principles and  
215 Scalability’**) is complementary to this workshop and focuses on different perspectives of foundation

---

216 models. Grigoris also organized tutorials on tensors and architecture design (CVPR'22, AAAI'23,  
217 DSA'24, NeurIPS'25, Dagstuhl'26) and deep learning theory (CVPR'23, ISIT'24). His research  
218 interests lie in trustworthy machine learning, publishing many results in top-tier conferences (CVPR,  
219 NeurIPS, ICLR, ICML). Grigoris serves as an Associate Editor for TMLR and an Area Chair for  
220 top-tier ML conferences (ICLR, ICML, NeurIPS). [\[Google Scholar\]](#)[\[Email\]](#)

221 **Sharon Li (UW-Madison)** is an Associate Professor in the Department of Computer Sciences at the  
222 University of Wisconsin-Madison. She received a Ph.D. from Cornell University in 2017, advised by  
223 John E. Hopcroft. Subsequently, she was a postdoctoral fellow in the Computer Science department  
224 at Stanford University. Her research focuses on the algorithmic and theoretical foundations of reliable  
225 machine learning in the open world, as well as developing responsible foundation models, including  
226 large language models and vision-language models. Sharon has served as the founding organizer and  
227 Program Chair for the *ICML Workshop on Uncertainty and Robustness in Deep Learning* (2019 and  
228 2020), co-organized multiple other workshops including the *ICML Workshop on Distribution-free  
229 Uncertainty Quantification* in 2021 and 2022, *NeurIPS'22 Workshop on Robustness in Sequence  
230 Modeling*, and *ICCV'23 Tutorial on Reliability of Deep Learning for Real-World Deployment*, and  
231 *ICLR'25 Workshop on Quantify Uncertainty and Hallucination in Foundation Models: The Next  
232 Frontier in Reliable AI*. She has served as Area Chair and Senior Program Committee for top-tier  
233 ML conferences including ICLR, ICML and NeurIPS between 2020 and 2024. She is serving as the  
234 Program Chair of ICML 2026. [\[Google Scholar\]](#)[\[Email\]](#)

235 **Etsuko Ishii (Amazon)** is an Applied Scientist at Amazon Web Services, focusing on developing  
236 and advancing Agentic AI applications. Her research expertise lies in contextual understanding  
237 of dialogues and building agentic systems that can effectively interact with users. Before joining  
238 AWS, she completed her Ph.D. in Electronic & Computer Engineering at the Hong Kong University  
239 of Science and Technology in 2024. She serves as Area Chair for ACL Rolling Review. [\[Google  
240 Scholar\]](#)[\[Email\]](#)

241 **Sean Du (Nanyang Technological University - Singapore)** is an Assistant Professor at College of  
242 Computing and Data Science (CCDS), Nanyang Technological University, Singapore. He obtained his  
243 Ph.D. in Computer Sciences at UW-Madison. His research interest is in reliable machine learning and  
244 the applications to foundation models and AI safety. His first-author papers have been recognized with  
245 multiple oral and spotlight presentations at NeurIPS and CVPR. He is a recipient of the Jane Street  
246 Graduate Research Fellowship, and Rising Stars in Data Science award. Sean has co-organized the  
247 *IJCNN'25 Special Session on Responsible Foundation Models in the Wild*. [\[Google Scholar\]](#)[\[Email\]](#)

248 **Katia Sycara (CMU)** is a Research Professor at Carnegie Mellon University. She is a Fellow of  
249 IEEE and AAAI, with over 300 publications in multi-agent systems, semantic web, and human-  
250 agent interaction. Her work has been funded by DARPA, NASA, NSF, and others, including the  
251 development of the RETSINA multi-agent infrastructure. She has held leadership roles in major  
252 conferences and contributed to web standards through W3C and OASIS. Sycara has received multiple  
253 honors, including the ACM/SIGART Agents Research Award and an honorary doctorate. Prof. Sycara  
254 is a founding member and member of the Board of Directors of the International Foundation of  
255 Multiagent Systems (IFMAS). She is a founding Editor-in-Chief of the journal “Autonomous Agents  
256 and Multiagent Systems”; an Editor-in-Chief of the Springer Series on Agents; on the Editorial Board  
257 of the Kluwer book series on “Multiagent Systems, Artificial Societies and Simulated Organizations”.  
258 [\[Google Scholar\]](#)[\[Email\]](#)

## 259 7.2 PROGRAM COMMITTEE

260 **Volunteers:** Two student (one PhD student and one undergraduate student) will help the organizers  
261 hosting the workshop: Andrea Tseng (University of Wisconsin-Madison) and Yiheng Zhang  
262 (University of Wisconsin-Madison).

263 In addition, the following people have agreed to serve in the program committee:

264 Elias Rocamora (EPFL), Blerina Gkotse (University of Wisconsin-Madison), Justin Deschenaux  
265 (EPFL), Andrea Tseng (University of Wisconsin-Madison), Thomas Pethick (EPFL), Stratis Skoulakis  
266 (Aarhus university), Dimitris Halatsis (Imperial College London), Muhammad Ashiq (University  
267 of Wisconsin-Madison), Zhiyuan Wu (University of Oslo), Aggelina Chatziagapi (Stonybrook  
268 University), Jiankang Deng (Imperial College London), Seongheon Park (University of Wisconsin-  
269

---

270 Madison), Changdae Oh (University of Wisconsin-Madison), Froilan Choi (University of Wisconsin-  
271 Madison), Shaokun Zhang (Penn State University), Jiachen (Tianhao) Wang (Princeton University),  
272 Yiran Wu (Penn State University), Wenyue Hua (Microsoft Research), Zhouxing Shi (UC Riverside).  
273

274 Depending on the number of submissions received, it may become necessary to expand the program  
275 committee. Given the track record in the past year, the core team's extensive experience in organizing  
276 scientific events will be instrumental in mobilizing their broad professional network to support  
277 the review process. **Concretely, as organizers we have organized before workshops scaling to**  
278 **hundreds of submissions and we recruited additional reviewers from our networks.**  
279

280

281

282

283

284

285

286

287

288

289

290

291

292

293

294

295

296

297

298

299

300

301

302

303

304

305

306

307

308

309

310

311

312

313

314

315

316

317

318

319

320

321

322

323

---

## 324 REFERENCES

## 325

326 Bansal and et al. Workshop on ‘reliable and responsible foundation models’. In *International Conference*  
327 *on Learning Representations (ICLR) Workshops*, 2024. [URL: [https://iclr.cc/virtual/2024/](https://iclr.cc/virtual/2024/workshop/20588)  
328 [workshop/20588](https://iclr.cc/virtual/2024/workshop/20588)].

329 Beirami and et al. Workshop on ‘responsible ai (rai)’. In *International Conference on Learning Representations*  
330 *(ICLR) Workshops*, 2021. [URL: <https://iclr.cc/virtual/2021/workshop/2132>].

331 Blake Brittain. Anthropic’s lawyers take blame for ai ‘hallucination’ in music publishers’ lawsuit, 2025. [Data: [https://www.reuters.com/legal/legalindustry/](https://www.reuters.com/legal/legalindustry/anthropic-lawsyers-take-blame-ai-hallucination-music-publishers-lawsuit-2025-05-15/)  
332 [anthropic-lawsyers-take-blame-ai-hallucination-music-publishers-lawsuit-2025-05-15/](https://www.reuters.com/legal/legalindustry/anthropic-lawsyers-take-blame-ai-hallucination-music-publishers-lawsuit-2025-05-15/).  
333 Status: Online; accessed 9-October-2025].

335 Chen and et al. Workshop on ‘large language models for agents’. In *International Conference on Learning*  
336 *Representations (ICLR) Workshops*, 2024. [URL: [https://iclr.cc/virtual/2024/workshop/](https://iclr.cc/virtual/2024/workshop/20575)  
337 [20575](https://iclr.cc/virtual/2024/workshop/20575)].

338 Cheng and et al. Workshop on ‘trustworthy and reliable large-scale machine learning models’. In *International*  
339 *Conference on Learning Representations (ICLR) Workshops*, 2023. [URL: <https://iclr.cc/virtual/2023/workshop/12827>].

341 Grigoris Chrysos, Sharon Li, Anastasios Angelopoulos, Stephen Bates, Barbara Plank, and Emtiyaz Khan.  
342 Workshop on ‘quantify uncertainty and hallucination in foundation models: The next frontier in reliable  
343 ai’. In *International Conference on Learning Representations (ICLR) Workshops*, 2025. [URL: <https://uncertainty-foundation-models.github.io>].

345 Schroeder de Witt and et al. Workshop on ‘ai for agent-based modelling (ai4abm)’. In *International Conference*  
346 *on Learning Representations (ICLR) Workshops*, 2023. [URL: [https://iclr.cc/virtual/2023/](https://iclr.cc/virtual/2023/workshop/12840)  
347 [workshop/12840](https://iclr.cc/virtual/2023/workshop/12840)].

348 Deshmukh and et al. Workshop on ‘what do we need for successful domain generalization?’. In *International*  
349 *Conference on Learning Representations (ICLR) Workshops*, 2023. [URL: [https://iclr.cc/virtual/2023/](https://iclr.cc/virtual/2023/workshop/12824)  
350 [workshop/12824](https://iclr.cc/virtual/2023/workshop/12824)].

351 Lei Huang, Weijiang Yu, Weitao Ma, Weihong Zhong, Zhangyin Feng, Haotian Wang, Qianglong Chen, Weihua  
352 Peng, Xiaocheng Feng, Bing Qin, et al. A survey on hallucination in large language models: Principles,  
353 taxonomy, challenges, and open questions. *ACM Transactions on Information Systems*, 43(2):1–55, 2025.

355 Alex Hughes. Your gp could be using chatgpt to diagnose you, a study finds, 2024. [Data: <https://www.sciencefocus.com/news/gps-use-chatgpt-study-finds>. Status: Online; accessed  
356 9-October-2025].

358 IANS. Us judge orders lawyers not to use chatgpt-drafted content in court, 2023. [Data: [https://www.business-standard.com/world-news/us-judge-orders-lawyers-not-to-use-chatgpt-drafted-content-in-court-123053100238\\_1.html](https://www.business-standard.com/world-news/us-judge-orders-lawyers-not-to-use-chatgpt-drafted-content-in-court-123053100238_1.html). Status: Online; accessed 9-October-2025].

362 Ziwei Ji, Nayeon Lee, Rita Frieske, Tiezheng Yu, Dan Su, Yan Xu, Etsuko Ishii, Ye Jin Bang, Andrea Madotto,  
363 and Pascale Fung. Survey of hallucination in natural language generation. *ACM computing surveys*, 55(12):  
364 1–38, 2023.

365 Jiang and et al. Workshop on ‘agent learning in open-endedness’. In *International Conference on Learning*  
366 *Representations (ICLR) Workshops*, 2022. [URL: [https://iclr.cc/virtual/2022/](https://iclr.cc/virtual/2022/workshop/4547)  
367 [workshop/4547](https://iclr.cc/virtual/2022/workshop/4547)].

368 Byron Kaye. Australian mayor readies world’s first defamation lawsuit over  
369 chatgpt content, 2023. [Data: [https://www.reuters.com/technology/](https://www.reuters.com/technology/australian-mayor-readies-worlds-first-defamation-lawsuit-over-chatgpt-content-2023-04-05/)  
370 [australian-mayor-readies-worlds-first-defamation-lawsuit-over-chatgpt-content-2023-04-05/](https://www.reuters.com/technology/australian-mayor-readies-worlds-first-defamation-lawsuit-over-chatgpt-content-2023-04-05/).  
371 Status: Online; accessed 9-October-2025].

372 Anna Kelsey-Sugg and Damien Carrick. Ai hallucinations caused artificial intelligence to falsely de-  
373 scribe these people as criminals, 2024. [Data: [https://www.abc.net.au/news/2024-11-04/](https://www.abc.net.au/news/2024-11-04/ai-artificial-intelligence-hallucinations-defamation-chatgpt/104518612)  
374 [ai-artificial-intelligence-hallucinations-defamation-chatgpt/104518612](https://www.abc.net.au/news/2024-11-04/ai-artificial-intelligence-hallucinations-defamation-chatgpt/104518612).  
375 Status: Online; accessed 9-October-2025].

376 Koutra and et al. Workshop on ‘towards agentic ai for science: Hypothesis generation, comprehension,  
377 quantification, and validation’. In *International Conference on Learning Representations (ICLR) Workshops*,  
378 2025. [URL: [https://iclr.cc/virtual/2025/](https://iclr.cc/virtual/2025/workshop/23991)  
379 [workshop/23991](https://iclr.cc/virtual/2025/workshop/23991)].

---

378 Michael Legg and Vicki McNamara. Ai is creating fake legal cases – and making its way into  
379 real courtrooms with disastrous results, 2024. [Data: <https://theconversation.com/ai-is-creating-fake-legal-cases-and-making-its-way-into-real-courtrooms-with-disastrous-results>]. Status: Online; accessed 9-October-2025].

382 Joshua Maynez, Shashi Narayan, Bernd Bohnet, and Ryan McDonald. On faithfulness and factuality in  
383 abstractive summarization. *arXiv preprint arXiv:2005.00661*, 2020.

384 Cade Metz and Karen Weisse. A.i. hallucinations are getting worse as models get  
385 smarter, 2025. [Data: <https://www.nytimes.com/2025/05/05/technology/ai-hallucinations-chatgpt-google.html>]. Status: Online; accessed 9-October-2025].

388 Papernot and et al. Workshop on ‘towards trustworthy ml: Rethinking security and privacy for ml’. In  
389 *International Conference on Learning Representations (ICLR) Workshops*, 2020. [URL: [https://iclr.cc/virtual\\_2020/workshops\\_6.html](https://iclr.cc/virtual_2020/workshops_6.html)].

390 Imran Rahman-Jones. Man files complaint after chatgpt said he killed his children, 2025. [Data: <https://www.bbc.com/news/articles/c0kgydkr516o>]. Status: Online; accessed 9-October-2025].

393 Thanapalasingam and et al. Workshop on ‘neurosymbolic generative models (nesy-gems)’. In *International Conference on Learning Representations (ICLR) Workshops*, 2023. [URL: <https://iclr.cc/virtual/2023/workshop/12830>].

396 Pranshu Verma and Will Oremus. Chatgpt sometimes makes up facts. for one law prof, it went  
397 too far, 2023. [Data: <https://www.washingtonpost.com/technology/2023/04/05/chatgpt-lies/>]. Status: Online; accessed 9-October-2025].

400 Sandra Wachter, Brent Mittelstadt, and Chris Russell. Do large language models have a legal duty to tell the  
401 truth? *Royal Society Open Science*, 11(8):240197, 2024.

402 Matthew Weaver. Ai chatbots distort and mislead when asked about current affairs, bbc  
403 finds, 2025. [Data: <https://www.theguardian.com/technology/2025/feb/11/ai-chatbots-distort-and-mislead-when-asked-about-current-affairs-bbc-finds>]. Status: Online; accessed 9-October-2025].

406 Benjamin Weiser. Here’s what happens when your lawyer uses chatgpt, 2023. [Data: <https://www.nytimes.com/2023/05/27/nyregion/avianca-airline-lawsuit-chatgpt.html>]. Status: Online; accessed 9-October-2025].

409 Karen Weisse and Cade Metz. When a.i. chatbots hallucinate, 2023. [Data: <https://www.nytimes.com/2023/05/01/business/ai-chatbots-hallucination.html>]. Status: Online; accessed 9-October-2025].

412 Xiao and et al. Workshop on ‘socially responsible machine learning’. In *International Conference on Learning Representations (ICLR) Workshops*, 2022. [URL: <https://iclr.cc/virtual/2022/workshop/4558>].

415

416

417

418

419

420

421

422

423

424

425

426

427

428

429

430

431