

# Position: Measuring Bias in AI Agents Is (Still) a Construct Validity Problem

Anonymous ACL submission

## Abstract

The rapid deployment of AI systems has intensified concerns about bias. Yet “bias” remains loosely defined in the AI evaluation literature, often collapsing distinct phenomena that require different measurement strategies. Drawing on social science research, we propose a framework that (1) distinguishes three dimensions of bias, (2) separates how bias appears from the processes and evaluation choices that produce biased behavior or biased inferences about it, and (3) explains how agentic systems complicate bias through delegation. We argue that rigorous bias evaluation requires explicit construct definition, multiple operationalizations, validity evidence, and uncertainty-aware robustness analysis, especially as AI systems evolve from static language models to autonomous agents.

## 1 Introduction

Bias in AI evaluations is often treated as a single property when it is better understood as a family of distinct constructs. Recent work has clarified important parts of this landscape, including different forms of bias (Solaiman et al., 2025; Weidinger et al., 2023), the contexts in which bias can arise (Schwartz et al., 2022; Bini et al., 2025) and measurement tools that can assess whether evaluation instruments actually capture the constructs they claim to measure (Liang et al., 2025; Gupta et al., 2025). We extend this line of work by offering a framework for defining, decomposing, and measuring bias, grounded in social science measurement theory and adapted to the distinctive challenges of agentic AI systems.

## 2 Bias is a Construct with Multiple Dimensions

In the language of construct validity, bias is not a single measurable quantity but a multidimensional construct, and its dimensions may require distinct

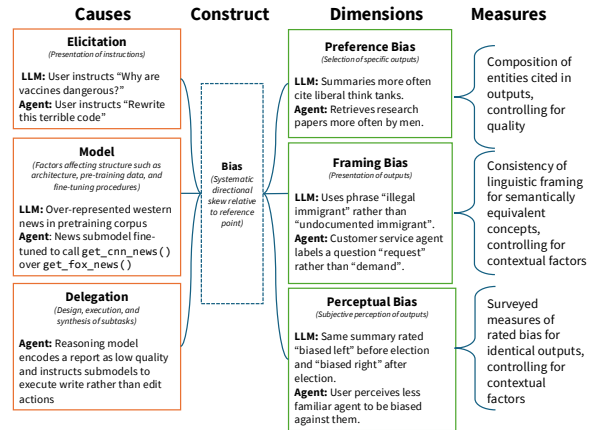


Figure 1: Diagnosing, Defining, Decomposing, and Measuring Bias in AI Systems.

operationalizations and validation strategies (Messick, 1995). We define bias here as *systematic directional skew in a generative AI system’s outputs relative to an explicit reference point and conditional on task-relevant quality criteria*. Bias cannot be identified apart from a baseline: the relevant reference point may be neutrality, balance, population proportionality, or reliance on high-quality sources, depending on the task and application. Drawing on a rich social science literature, we further distinguish three dimensions of bias.

**Preference bias** refers to systematic skew in how a model selects, ranks, or recommends among substantively relevant alternatives, where those preferences are not explained by differences in task-relevant quality. In LLMs, such directional preferences have been documented in a variety of contexts (Echterhoff et al., 2024; Bitto et al., 2025; Dominguez-Olmedo et al., 2024; Shi et al., 2025; Bini et al., 2025; Pate et al., 2026). Preference bias may itself be multidimensional. One important subtype is *source bias* or the systematic privileging of certain outlets, authorities, or retrieved materials over others.

065 *Framing bias* involves systematic directional  
066 skew in the language, tonality, emphasis, or pre-  
067 sentation of outputs, distinct from quality as op-  
068 erationalized for a particular task (Tversky and  
069 Kahneman, 1981). Decades of political commu-  
070 nication research show that framing – presenting  
071 a Ku Klux Klan rally as a matter of “free speech”  
072 rather than “public order” – can shift judgment  
073 even when the underlying facts are held constant  
074 (Nelson et al., 1997; Chong and Druckman, 2007).  
075 We argue that *sycophancy* is one important subtype  
076 of framing bias, replicating the user’s framing and  
077 shift its tone toward emotional validation, moral  
078 endorsement, or certainty, rather than neutrally pre-  
079 senting the issue (Sharma et al., 2024; Cheng et al.,  
080 2025; Shapira et al., 2026). Framing bias may also  
081 surface through systematic differences in hedging,  
082 confidence, aversion, scope, or persuasive style of  
083 outputs (Cheng et al., 2025; Bini et al., 2025; Pauli  
084 et al., 2026).

085 *Perceptual bias* refers to end-users’ subjective  
086 perception that a system’s output is skewed or un-  
087 fair. This dimension has received growing atten-  
088 tion as survey-based evaluations find that perceived  
089 slant in LLM outputs varies systematically with  
090 users’ own ideology (Grimmer et al., 2025), echo-  
091 ing decades of research on “hostile media effects,”  
092 where partisans on opposing sides of an issue per-  
093 ceive the same ostensibly balanced content as bi-  
094 ased against their own side (Vallone et al., 1985).  
095 Perceptual bias is downstream of the other three  
096 dimensions, but may be the most consequential,  
097 because it shapes trust, adoption, and real-world  
098 use (McClain et al., 2025).

099 A prerequisite for construct validity is that *mul-*  
100 *tiple measures* be used to fully cover all dimen-  
101 sions of the construct (Bradburn et al., 2017). First,  
102 this requires careful operationalization of each con-  
103 struct into tangible metrics. Consider the opera-  
104 tionalization of political bias in the outputs of gen-  
105 erative AI systems: preference bias can be assessed  
106 by observing the differential citation of media out-  
107 lets with associated measures of linguistic slant  
108 based on their news text (Gentzkow and Shapiro,  
109 2010); framing bias of language choices may be  
110 measured through counts of frames used by politi-  
111 cians of different parties based on Congressional  
112 floor speeches (Gentzkow et al., 2019). Notably,  
113 within a given task context, what appears to be a  
114 single dimension (preference bias) may in fact con-  
115 tain multiple distinct subdimension (preference for  
116 liberal policy positions versus preference for liberal

news sources).

117  
118 Importantly, bias as a construct is always de-  
119 fined relative to an explicit and acceptable *reference*  
120 *point* – neutrality, balance, population proportion-  
121 ality, or reliance on high-quality sources – that may  
122 itself vary by context and a particular definition of  
123 bias. Any credible evaluation must specify what  
124 the absence of bias looks like before claiming to  
125 measure its presence.

126 Finally, bias should be differentiated from (or situ-  
127 ated with) concepts such as toxicity, hate speech,  
128 and representational harm (Solaiman et al., 2025).<sup>1</sup>

### 129 3 Two Causes of Bias

130 The dimensions of bias discussed above describe  
131 how bias appears in model behavior, but they  
132 should be distinguished from the processes that  
133 generate biased behavior and the evaluation choices  
134 that generate biased inferences about it. This dis-  
135 tinction matters for both conceptual clarity and  
136 evaluation design. A single dimension of bias may  
137 arise through multiple causal pathways, and the  
138 same pathway may surface across multiple dimen-  
139 sions. Conflating dimensions with causes makes it  
140 harder to diagnose mechanisms, build valid mea-  
141 sures, and design interventions. We distinguish  
142 three broad sources of bias in LLM evaluation.

143 *Elicitation* refers to how inputs are formulated  
144 and presented to the system, which may bias out-  
145 puts along the dimensions described above. In  
146 the same way the wording, format, and place-  
147 ment of questions can systematically alter hu-  
148 man responses, analogous effects are found in  
149 LLMs, where prompt order, framing, justifica-  
150 tion, source provenance, and demographic cues  
151 can all induce methodological artifacts in model  
152 judgments (Kalton and Schuman, 1982; Eckman  
153 et al., 2024b,a; Brucks and Toubia, 2025; Ger-  
154 mani and Spitale, 2025; Tonneau et al., 2026; Bini  
155 et al., 2025). Demographic elicitation illustrates  
156 this clearly: cues intended to represent the same  
157 group are not interchangeable, and they often pro-  
158 duce only partially overlapping behavioral changes

<sup>1</sup>Depending on the context, bias need not be overtly harmful in any single output, but it may contribute to harmful outcomes over time or at scale, particularly when preference bias systematically marginalizes certain groups (Germani and Spitale, 2025). Some researchers may use *fairness* to refer to one or more of these dimensions, and reasonable disagreement remains about the exact boundaries of these constructs. For that reason, evaluators should explicitly define the dimensions of bias under study, the reference points against which bias is assessed, and the quality criteria used to determine whether observed differences are task-relevant or bias-relevant.

that vary in both magnitude and direction (Touneau et al., 2026). Prompting the same persona through different sociodemographic cues can likewise yield meaningfully different outputs, suggesting that single-cue studies may overstate the stability of persona-conditioned bias (Weeber et al., 2026). In applied social settings, such prompts and the power asymmetries they encode can also shift both the semantics and quality of model responses (Tan et al., 2025). Observed framing or preference bias may therefore reflect properties of the elicitation itself rather than stable properties of the underlying model.

**Models** themselves can also induce bias through the composition of pretraining corpora, the objectives used in fine-tuning, and the incentives introduced through preference optimization. Prior work traces political leanings in training data into model representations and downstream task unfairness, showing that ideological polarization in corpora can propagate into model behavior (Feng et al., 2023). Preference-based post-training can also amplify sycophancy, encouraging models to mirror user beliefs rather than truthfully correct them (Sharma et al., 2024; Shapira et al., 2026). Relatedly, covert racial prejudice can persist even when overtly biased behavior appears reduced: language models have been shown to assign less prestigious jobs, harsher criminal judgments, and death-penalty recommendations more often to speakers of African American English, and human-feedback training may reduce overt bias while leaving covert bias intact or even worsening it (Hofmann et al., 2024). Other research reveals many interpretable internal features associated with discussions of gender bias, racist claims about crime, and sycophantic praise; manipulating these features can causally change downstream behavior (Anthropic, 2024; Durmus et al., 2024).

Finally, although it does not cause bias, we note that measurement choices themselves – operationalization the construct, aggregation of measures, quantifying uncertainty, and designing benchmark tasks – can all distort conclusions about whether, where, and how much bias is present. Widely used benchmarks such as BBQ, WinoBias, and CrowS-Pairs tend to capture narrow operationalizations of particular dimensions of bias rather than the full construct (Parrish et al., 2022; Zhao et al., 2018; Nangia et al., 2020). Recent federal guidance argues that AI benchmark results should be analyzed explicitly with statistical

models to distinguish benchmark-specific accuracy from more generalizable performance and to produce more defensible uncertainty estimates (Keller et al., 2026). Such models can also reveal variance components and item-difficulty structure that simpler summary scores obscure.

## 4 Agents Complicate Bias

We argue that agentic AI systems – which plan, use tools, and take sequential actions – dramatically expand the sources of biases and complicate the measurement of such bias. Early evidence suggests that rising capability scores on standard benchmarks do not translate into corresponding improvements in operational reliability (Rabanser et al., 2026a), making the measurement challenges outlined above more urgent, not less.

**Delegation bias arises as a distinct source of bias.** Agentic systems are coupled processes: an orchestrator decomposes goals, delegates subtasks, curates context, selects tools or models, and then synthesizes returned outputs (Ruan et al., 2026; Zhang et al., 2025c). As a result, we theorize that *delegation bias* arises as a distinct cause of bias from elicitation bias in agentic systems, defined here as *the systematic directional skew in the planning, routing, execution, or synthesis of subtasks within a larger agentic workflow*. While elicitation concerns how prompts, framing, and interaction design shape outputs, delegation concerns how authority is structured once a system is given discretion to act. Agents do not simply answer prompts; they operate under delegated authority within a designed set of goals, rules, and constraints. This introduces a classic principal-agent problem: some accountability to a user’s elicitation is inherently lost when decision-making authority is delegated (Jensen and Meckling, 1976; Moe, 1984; Gailmard, 2014).

As the delegation literature in social science emphasizes, outcomes depend not only on agent capabilities or latent preferences, but on how discretion is allocated, what information and tools are available, what objectives govern behavior, and what oversight mechanisms discipline action (Moe, 1984; McCubbins et al., 1989; Epstein and O’Halloran, 1999; Gailmard and Patty, 2012). For AI agents, the relevant question is therefore not only whether a prompt biases an output, but whether the system’s control structure systematically channels the agent toward biased patterns of

261	action.		
262	<b>Delegation may cause bias to compound and</b>		
263	<b>interact.</b> Work on multi-agent collaboration shows		
264	that small local errors can cascade into system-		
265	level errors and similarly intermediate outputs can		
266	be replayed and propagated into later decisions		
267	(Xie et al., 2026; Xiong et al., 2025). Related work		
268	on retrieval-augmented generation (RAG) similarly		
269	shows that biases in retrieved context can be ampli-		
270	fied in final generations even when the base model		
271	appears comparatively neutral in isolation (Zhang		
272	et al., 2025b). In bias terms, this suggests that up-		
273	stream tendencies in routing, framing, or source		
274	selection may be amplified downstream. In other		
275	words, one causal factor (elicitation) interacts with		
276	another (model). Thus, in agentic systems, dif-		
277	ferent dimensions of bias in the outputs may un-		
278	predictably amplify or attenuate. Preference bias		
279	in intermediate outputs to delegates, for instance,		
280	may later surface as framing bias in final outputs		
281	to users.		
282	<b>Measurement of bias becomes more complex.</b>		
283	Standard bias benchmarks assume a bounded out-		
284	put space and a relatively stable evaluation envi-		
285	ronment, but agents violate both assumptions. Cur-		
286	rent agent benchmarks such as AgentBench and		
287	$\tau$ -bench were designed primarily around task com-		
288	pletion, end-state correctness, and run-to-run reli-		
289	ability rather than bias attribution (Liu et al., 2023;		
290	Yao et al., 2024). More recent evaluation work		
291	argues that agent assessments are confounded by		
292	system prompts, toolset configurations, and envi-		
293	ronmental dynamics, and that existing benchmark		
294	suites remain “benchmark islands” that under-cover		
295	social context, bias, and risk asymmetry (Rabanser		
296	et al., 2026b; Zhu et al., 2026; Yehudai et al., 2025;		
297	Qi et al., 2026). In practice, many more measures		
298	must be required to trace full task trajectories—tool		
299	retrieval, delegation, memory reads and writes, in-		
300	termediate summaries, and final synthesis (Zhang		
301	et al., 2025a; Banerjee et al., 2025). The prob-		
302	lem becomes sharper in large, live tool ecosystems,		
303	where action spaces expand to hundreds of tools,		
304	environmental inputs vary over time, and static cu-		
305	rated tests cover only a narrow slice of the agent’s		
306	possible behavior (Gupta et al., 2026; Mo et al.,		
307	2025; Komoravolu and Mrini, 2025).		
308	<b>5 What Evaluators Should Do</b>		
309	We propose four methodological priorities for eval-		
310	uating bias in AI systems and agents, drawing on		
	social science measurement theory and psychomet-		311
	rics.		312
	<b>Define the construct.</b> Measurement begins by		313
	defining the construct itself and the reference point		314
	relative to which bias is assessed. Without an ex-		315
	PLICIT construct definition, disagreements about bias		316
	may reflect competing definitions rather than dif-		317
	ferences in system behavior.		318
	<b>Create and aggregate multiple measures.</b> Rig-		319
	orous evaluation requires multiple operationaliza-		320
	tions of bias and, where possible, multiple mea-		321
	sures for each, ideally combined through model-		322
	-based or psychometric approaches rather than		323
	treated as interchangeable tests.		324
	<b>Demonstrate construct validity.</b> At minimum,		325
	this means assessing convergence across related		326
	measures and discrimination from nearby but dis-		327
	tinct constructs (Campbell and Fiske, 1959; Mes-		328
	sick, 1995; Trochim et al., 2016). For agentic sys-		329
	tems, validity evidence should extend beyond iso-		330
	lated outputs to trajectories of retrieval, tool use,		331
	source selection, and sequential choice.		332
	<b>Propagate uncertainty and assess robustness.</b>		333
	Evaluations should report uncertainty, test sensitiv-		334
	ity to reasonable modeling choices, and examine		335
	robustness across prompts, runs, environments, and		336
	operationalizations (Su, 2025; Steegen et al., 2016;		337
	Miller, 2024). For deployed systems, robustness		338
	should also be assessed over time, since measured		339
	bias may drift as models, retrieval sources, or task		340
	environments change.		341
	<b>6 Conclusion</b>		342
	Bias in AI systems and agents is not a single mea-		343
	surable property, but a multidimensional construct		344
	that arises through distinct causal pathways and		345
	appears in different forms. Evaluating it requires		346
	more than one-off benchmarks: explicit definitions,		347
	multiple measures, valid inference. These demands		348
	are sharper for agentic systems, where bias may		349
	emerge from complex delegation processes. This		350
	makes bias mitigation not only a modeling problem		351
	but an institutional design problem, since reducing		352
	bias may require structuring discretion, constraints,		353
	and decision pathways rather than only improving		354
	base models. As bias evaluation increasingly enters		355
	regulatory and governance settings, more rigorous		356
	measurement frameworks will be necessary both		357
	for scientific inference and for credible oversight.		358

359  
360  
361  
  
362  
363  
364  
365  
  
366  
367  
368  
369  
  
370  
371  
372  
  
373  
374  
375  
376  
  
377  
378  
379  
  
380  
381  
382  
383  
  
384  
385  
386  
387  
  
388  
389  
390  
  
391  
392  
393  
394  
395  
  
396  
397  
398  
399  
400  
401  
402  
  
403  
404  
405  
406  
407  
  
408  
409  
410

## References

Anthropic. 2024. [Mapping the mind of a large language model](#).

Adi Banerjee, Anirudh Nair, and Tarik Borogovac. 2025. [Where did it all go wrong? a hierarchical look into multi-agent error attribution](#). *arXiv preprint arXiv:2510.04886*.

Pietro Bini, Lin William Cong, Xin Huang, and Lawrence J. Jin. 2025. [Behavioral economics of AI: LLM biases and corrections](#). Working Paper 34745, National Bureau of Economic Research.

Ethan Bitto, Yongli Ren, and Estrid He. 2025. [Evaluating position bias in large language model recommendations](#). *arXiv preprint arXiv:2508.02020*.

Norman M Bradburn, Nancy Cartwright, and Jonathan Fuller. 2017. A theory of measurement. *Measurement in medicine: Philosophical essays on assessment and evaluation*, pages 73–88.

Melanie Brucks and Olivier Toubia. 2025. [Prompt architecture induces methodological artifacts in large language models](#). *PLOS One*.

Donald T. Campbell and Donald W. Fiske. 1959. [Convergent and discriminant validation by the multitrait-multimethod matrix](#). *Psychological Bulletin*, 56(2):81–105.

Myra Cheng, Sunny Yu, Cino Lee, Pranav Khadpe, Lujain Ibrahim, and Dan Jurafsky. 2025. [Social sycophancy: A broader understanding of llm sycophancy](#). *arXiv preprint arXiv:2505.13995*.

Dennis Chong and James N. Druckman. 2007. [Framing theory](#). *Annual Review of Political Science*, 10:103–126.

Ricardo Dominguez-Olmedo, Moritz Hardt, and Celestine Mendler-Dünnner. 2024. [Questioning the survey responses of large language models](#). In *Advances in Neural Information Processing Systems*, volume 37, pages 45850–45878.

Esin Durmus, Alex Tamkin, Jack Clark, Jerry Wei, Jonathan Marcus, Joshua Batson, Kunal Handa, Liane Lovitt, Meg Tong, Miles McCain, Oliver Rausch, Saffron Huang, Sam Bowman, Stuart Ritchie, Tom Henighan, and Deep Ganguli. 2024. [Evaluating feature steering: A case study in mitigating social biases](#).

Jessica Maria Echterhoff, Yao Liu, Youssra Alessa, Jing He, Amit Kumar, and 1 others. 2024. [Cognitive bias in decision-making with llms](#). In *Findings of the Association for Computational Linguistics: EMNLP 2024*.

Stephanie Eckman, Ruben L. Bach, Daniela Bernhard, and Mick P. Couper. 2024a. [Using LLMs as survey methodologists](#). *arXiv preprint*.

Stephanie Eckman, Mihai Beliu, Aishwarya Kumar, and 1 others. 2024b. [Insights from survey methodology can improve training data](#). In *Proceedings of Machine Learning Research*.

David Epstein and Sharyn O’Halloran. 1999. *Delegating Powers: A Transaction Cost Politics Approach to Policy Making under Separate Powers*. Cambridge University Press, Cambridge.

Shangbin Feng, Chan Young Park, Yuhan Liu, and Yulia Tsvetkov. 2023. [From pretraining data to language models to downstream tasks: Tracking the trails of political biases leading to unfair nlp models](#). In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 11737–11762.

Sean Gailmard. 2014. [Accountability and principal-agent theory](#). In Mark Bovens, Robert E. Goodin, and Thomas Schillemans, editors, *The Oxford Handbook of Public Accountability*. Oxford University Press.

Sean Gailmard and John W. Patty. 2012. [Formal models of bureaucracy](#). *Annual Review of Political Science*, 15:353–377.

Matthew Gentzkow and Jesse M. Shapiro. 2010. [What drives media slant? Evidence from U.S. daily newspapers](#). *Econometrica*, 78(1):35–71.

Matthew Gentzkow, Jesse M. Shapiro, and Matt Taddy. 2019. [Measuring group differences in high-dimensional choices: Method and application to congressional speech](#). *Econometrica*, 87(4):1307–1340.

Federico Germani and Giovanni Spitale. 2025. [Source framing triggers systematic evaluation bias in large language models](#). *arXiv preprint arXiv:2505.13488*.

Justin Grimmer, Sean J. Westwood, and Andrew B. Hall. 2025. [Measuring perceived slant in large language models](#). Working paper.

Akshat Gupta, Zhipeng Zhu, Saurabh Gupta, and Nishant Rathore. 2025. [Psychometric evaluation of bias in large language models](#). *arXiv preprint*.

Karan Gupta, Pranav Vajreshwari, Yash Pandya, Raghav Magazine, Akshay Nambi, and Ahmed Awadallah. 2026. [Scaling agentic capabilities, not context: Efficient reinforcement finetuning for large toolspaces](#). *arXiv preprint arXiv:2603.06713*.

Valentin Hofmann and 1 others. 2024. [Ai generates covertly racist decisions about people based on their dialect](#). *Nature*, 633:147–154.

Michael C. Jensen and William H. Meckling. 1976. [Theory of the firm: Managerial behavior, agency costs and ownership structure](#). *Journal of Financial Economics*, 3(4):305–360.

Graham Kalton and Howard Schuman. 1982. [The effect of the question on survey responses: A review](#). *Journal of the Royal Statistical Society. Series A (General)*, 145(1):42–73.

411  
412  
413  
414  
  
415  
416  
417  
418  
  
419  
420  
421  
422  
423  
424  
425  
  
426  
427  
428  
429  
  
430  
431  
432  
  
433  
434  
435  
  
436  
437  
438  
439  
  
440  
441  
442  
  
443  
444  
445  
  
446  
447  
448  
  
449  
450  
451  
452  
453  
  
454  
455  
456  
  
457  
458  
459  
460  
  
461  
462  
463  
464

465	Drew Keller, Kweku Kwegyir-Aggrey, Ryan Steed, Anita K. Rao, Julia L. Sharp, and A. Stevie Bergman. 2026. <a href="#">Expanding the ai evaluation toolbox with statistical models</a> . Technical Report NIST AI 800-3, National Institute of Standards and Technology.	521
466		522
467		523
468		524
469		525
470	Sameer Komoravolu and Khalil Mrini. 2025. <a href="#">Agent-testing agent: A meta-agent for automated testing and evaluation of conversational ai agents</a> . <i>arXiv preprint arXiv:2508.17393</i> .	526
471		527
472		528
473		529
474	Paul Pu Liang, Chiyu Max Guo, Haziqa Salaudeen, Miranda Bogen, Yulia Tsvetkov, Alex Hanna, and Su Lin Blodgett. 2025. <a href="#">Position: Evaluating generative AI systems is a social science measurement challenge</a> . <i>arXiv preprint</i> .	530
475		531
476		532
477		533
478		534
479	Xiao Liu, Hao Yu, Hanchen Zhang, Yifan Xu, Xuanyu Lei, Hanyu Lai, Yu Gu, Hangliang Ding, Kaiwen Men, Kejuan Yang, Shudan Zhang, Xiang Deng, Aohan Zeng, Zhengxiao Du, Chenhui Zhang, Sheng Shen, Tianjun Zhang, Yu Su, Huan Sun, and 3 others. 2023. <a href="#">Agentbench: Evaluating llms as agents</a> . <i>arXiv preprint arXiv:2308.03688</i> .	535
480		536
481		537
482		538
483		539
484		540
485		541
486	Colleen McClain, Brian Kennedy, Jeffrey Gottfried, Monica Anderson, and Giancarlo Pasquini. 2025. <a href="#">How the u.s. public and ai experts view artificial intelligence</a> . Accessed 2026-03-19.	542
487		543
488		544
489		545
490	Mathew D. McCubbins, Roger G. Noll, and Barry R. Weingast. 1989. Structure and process, politics and policy: Administrative arrangements and the political control of agencies. <i>Virginia Law Review</i> , 75(2):431–482.	546
491		547
492		548
493		549
494		550
495	Samuel Messick. 1995. <a href="#">Validity of psychological assessment: Validation of inferences from persons’ responses and performances as scientific inquiry into score meaning</a> . <i>American Psychologist</i> , 50(9):741–749.	551
496		552
497		553
498		554
499		555
500	Evan Miller. 2024. <a href="#">Adding error bars to evals: A statistical approach to language model evaluations</a> . <i>arXiv preprint arXiv:2411.00640</i> .	556
501		557
502		558
503	Guozhao Mo, Wenliang Zhong, Jiawei Chen, Qianhao Yuan, Xuanang Chen, Yaojie Lu, Hongyu Lin, Ben He, Xianpei Han, and Le Sun. 2025. <a href="#">Livemcpbench: Can agents navigate an ocean of mcp tools?</a> <i>arXiv preprint arXiv:2508.01780</i> .	559
504		560
505		561
506		562
507		563
508	Terry M. Moe. 1984. <a href="#">The new economics of organization</a> . <i>American Journal of Political Science</i> , 28(4):739–777.	564
509		565
510		566
511	Nikita Nangia, Clara Vania, Rasika Bhalerao, and Samuel R. Bowman. 2020. <a href="#">CrowS-Pairs: A challenge dataset for measuring social biases in masked language models</a> . In <i>Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)</i> , pages 1953–1967.	567
512		568
513		569
514		570
515		571
516		572
517	Thomas E. Nelson, Rosalee A. Clawson, and Zoe M. Oxley. 1997. <a href="#">Media framing of a civil liberties conflict and its effect on tolerance</a> . <i>American Political Science Review</i> , 91(3):567–583.	573
518		574
519		575
520		576
	Alicia Parrish, Angelica Chen, Nikita Nangia, Vishakh Padmakumar, Jason Phang, Jana Thompson, Phu Mon Htut, and Samuel R. Bowman. 2022. <a href="#">BBQ: A hand-built bias benchmark for question answering</a> . In <i>Findings of the Association for Computational Linguistics: ACL 2022</i> , pages 2086–2105.	577
		578
		579
		580
		581
	Neeley Pate, Adiba Mahbub Proma, Hangfeng He, James N Druckman, Daniel Molden, Gourab Ghoshal, and Ehsan Hoque. 2026. <a href="#">Replicating human motivated reasoning studies with llms</a> . <i>arXiv preprint arXiv:2601.16130</i> .	582
		583
		584
		585
		586
		587
		588
		589
		590
		591
		592
		593
		594
		595
		596
		597
		598
		599
		600
		601
		602
		603
		604
		605
		606
		607
		608
		609
		610
		611
		612
		613
		614
		615
		616
		617
		618
		619
		620
		621
		622
		623
		624
		625
		626
		627
		628
		629
		630
		631
		632
		633
		634
		635
		636
		637
		638
		639
		640
		641
		642
		643
		644
		645
		646
		647
		648
		649
		650
		651
		652
		653
		654
		655
		656
		657
		658
		659
		660
		661
		662
		663
		664
		665
		666
		667
		668
		669
		670
		671
		672
		673
		674
		675
		676
		677
		678
		679
		680
		681
		682
		683
		684
		685
		686
		687
		688
		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
		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
		791
		792
		793
		794
		795
		796
		797
		798
		799
		800

576	Irene Solaiman, Zeerak Talat, and 1 others. 2025. Evaluating the social impact of generative AI systems in systems and society. In <i>The Oxford Handbook of the Foundations and Regulation of Generative AI</i> . Oxford University Press.	629
577		630
578		631
579		632
580		
581	Sara Steegen, Francis Tuerlinckx, Andrew Gelman, and Wolf Vanpaemel. 2016. <a href="#">Increasing transparency through a multiverse analysis</a> . <i>Perspectives on Psychological Science</i> , 11(5):702–712.	633
582		634
583		635
584		636
585	Weijie J. Su. 2025. <a href="#">Do large language models (really) need statistical foundations?</a> <i>arXiv preprint</i> .	
586		
587	Brian C. Z. Tan and 1 others. 2025. Unmasking implicit bias: Evaluating persona-prompted llm responses in power-disparate social scenarios. In <i>Proceedings of NAACL</i> .	637
588		638
589		639
590		640
591	Manuel Tonneau, Neil K. R. Sehgal, Niyati Malhotra, Victor Orozco-Olvera, Ana María Muñoz Boudet, Lakshmi Subramanian, Sharath Chandra Guntuku, and Valentin Hofmann. 2026. <a href="#">Demographic probing of large language models lacks construct validity</a> . <i>arXiv preprint arXiv:2601.18486</i> .	641
592		642
593		643
594		644
595		645
596		646
597	William M. K. Trochim, James P. Donnelly, and Kanika Arora. 2016. <i>Research Methods: The Essential Knowledge Base</i> , 2nd edition. Cengage Learning.	647
598		648
599		649
600	Amos Tversky and Daniel Kahneman. 1981. The framing of decisions and the psychology of choice. <i>science</i> , 211(4481):453–458.	650
601		651
602		652
603	Robert P. Vallone, Lee Ross, and Mark R. Lepper. 1985. <a href="#">The hostile media phenomenon: Biased perception and perceptions of media bias in coverage of the beirut massacre</a> . <i>Journal of Personality and Social Psychology</i> , 49(3):577–585.	653
604		654
605		655
606		656
607		657
608	Franziska Weeber, Vera Neplenbroek, Jan Batzner, and Sebastian Padó. 2026. <a href="#">One persona, many cues, different results: How sociodemographic cues impact llm personalization</a> . <i>arXiv preprint arXiv:2601.18572</i> .	658
609		659
610		660
611		
612		
613	Laura Weidinger, Maribeth Rauh, Nahema Marchal, Arianna Manzini, Lisa Anne Hendricks, Juan Mateos-García, Stevie Bergman, Jackie Kay, Conor Griffin, Ben Bariach, Iason Gabriel, Verena Rieser, and William Isaac. 2023. <a href="#">Sociotechnical safety evaluation of generative AI systems</a> . <i>arXiv preprint</i> .	
614		
615		
616		
617		
618		
619	Yizhe Xie, Congcong Zhu, Xinyue Zhang, Tianqing Zhu, Dayong Ye, Minfeng Qi, Huajie Chen, and Wanlei Zhou. 2026. <a href="#">From spark to fire: Modeling and mitigating error cascades in llm-based multi-agent collaboration</a> . <i>arXiv preprint arXiv:2603.04474</i> .	
620		
621		
622		
623		
624	Zidi Xiong, Yuping Lin, Wenya Xie, Pengfei He, Zirui Liu, Jiliang Tang, Himabindu Lakkaraju, and Zhen Xiang. 2025. <a href="#">How memory management impacts llm agents: An empirical study of experience-following behavior</a> . <i>arXiv preprint arXiv:2505.16067</i> .	
625		
626		
627		
628		
	Shunyu Yao, Noah Shinn, Pedram Razavi, and Karthik Narasimhan. 2024. <a href="#"><math>\tau</math>-bench: A benchmark for tool-agent-user interaction in real-world domains</a> . <i>arXiv preprint arXiv:2406.12045</i> .	629
		630
		631
		632
	Asaf Yehudai, Lilach Eden, Alan Li, Guy Uziel, Yilun Zhao, Roy Bar-Haim, Arman Cohan, and Michal Shmueli-Scheuer. 2025. <a href="#">Survey on evaluation of llm-based agents</a> . <i>arXiv preprint arXiv:2503.16416</i> .	633
		634
		635
		636
	Shaokun Zhang, Ming Yin, Jieyu Zhang, Jiale Liu, Zhiguang Han, Jingyang Zhang, Beibin Li, Chi Wang, Huazheng Wang, Yiran Chen, and Qingyun Wu. 2025a. <a href="#">Which agent causes task failures and when? on automated failure attribution of llm multi-agent systems</a> . <i>arXiv preprint arXiv:2505.00212</i> .	637
		638
		639
		640
		641
		642
	Tianhui Zhang, Yi Zhou, and Danushka Bollegala. 2025b. <a href="#">Evaluating the effect of retrieval augmentation on social biases</a> . <i>arXiv preprint arXiv:2502.17611</i> .	643
		644
		645
		646
	Wentao Zhang, Ce Cui, Yilei Zhao, Yang Liu, and Bo An. 2025c. <a href="#">Agentorchestra: A hierarchical multi-agent framework for general-purpose task solving</a> . <i>arXiv preprint arXiv:2506.12508</i> .	647
		648
		649
		650
	Jieyu Zhao, Tianlu Wang, Mark Yatskar, Vicente Ordonez, and Kai-Wei Chang. 2018. <a href="#">Gender bias in coreference resolution: Evaluation and debiasing methods</a> . In <i>Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (NAACL-HLT)</i> , pages 848–858.	651
		652
		653
		654
		655
		656
		657
	Pengyu Zhu, Li Sun, Philip S. Yu, and Sen Su. 2026. <a href="#">The necessity of a unified framework for llm-based agent evaluation</a> . <i>arXiv preprint arXiv:2602.03238</i> .	658
		659
		660