

Rethinking NLP Evaluation for LLM Assistants in Education: A Human-Centered Evaluation Framework

Anonymous ACL submission

Abstract

Large language models (LLMs) are increasingly adopted as educational assistants and feedback providers; however, dominant NLP evaluation practices continue to emphasize technical metrics such as accuracy, fluency, and automated judgment. We argue that these approaches are misaligned with educational contexts because they overlook human values, learner agency, and pedagogical goals. This paper presents an argument-driven critique of NLP evaluation practices as applied to educational LLMs and introduces a literature-informed, human-centered, and sociotechnical evaluation framework with guiding questions to support its use. We highlight criteria such as explainability, consistency, and refinement, which are critical to human-AI interactive instructional effectiveness but absent from existing NLP benchmarks, and argue for more human-aligned, co-creative, and long-term evaluation practices for educational LLM systems.

1 Introduction

Large language models (LLMs) are being rapidly adopted in educational settings, yet their evaluation remains governed by general-purpose NLP benchmarks. While these benchmarks support scalable and standardized comparisons, they are often misaligned with the pedagogical, contextual, and ethical requirements of educational use. This paper asks which evaluative criteria are overlooked by dominant NLP benchmarks when educational language models are used as collaborative, teachable, and value-aligned agents. We critically examine current NLP evaluation practices in the context of educational LLMs and argue that existing metrics insufficiently capture learning processes, learner agency, educational context, and stakeholder values. As a starting point for future work, we introduce a human-centered, value-aligned evaluation framework and distill early takeaways to inform

the design, benchmarking, and assessment of educational language models in future NLP research.

2 NLP Evaluation critiques

2.1 Dominant NLP Evaluation Practices

NLP evaluation has traditionally emphasized task correctness through metrics such as accuracy, precision, recall, and F-score measures (Hutchinson et al., 2022). While these metrics capture specific error types, most of them ignore contextual factors. Widely used automated metrics such as BLEU (Papineni et al., 2002) and ROUGE (Lin, 2004) are well-suited to short, constrained tasks with a single correct answer, but break down for longer, open-ended generations where multiple valid responses may exist. Even more semantically informed metrics, such as BERTScore improve alignment at the embedding level, yet still struggle to capture reasoning quality, domain-specific correctness, and human preferences. This highlights a fundamental mismatch between dominant NLP evaluation metrics and the evaluative demands of complex generative tasks (Zhang et al., 2019).

LLM-as-a-Judge frameworks have recently adopted prominence due to their low cost, ease of use, and reported correlations with human stylistic preferences. However, prior work demonstrates that LLM judges exhibit systematic biases that can distort evaluation outcomes, often relying on output accuracy or text similarity (Li et al., 2025; Szymanski et al., 2025). As a result, LLM-as-a-Judge evaluations risk presenting an illusion of holistic assessment while remaining grounded in opaque evaluative criteria. Benchmarks such as MT-Bench (Zheng et al., 2023), and Arena-Hard-Auto (Li et al., 2024) exemplify this trend, emphasizing conversational ability but offering limited insight into model behavior, or contextual appropriateness. Benchmark suites such as GLUE and SuperGLUE (Wang et al., 2018) further exemplify dom-

inant evaluation practices in NLP by aggregating diverse natural language understanding tasks and emphasizing sample-efficient learning and cross-task transfer (Wang et al., 2018, 2019). However, uneven data availability across tasks, genre mismatches between training and test sets, and reliance on privately held test data limit transparency, raising concerns about whether benchmark performance reliably translates to real-world deployment settings.

Overall, evaluating whether an NLP model is appropriate for real-world use requires more than benchmark scores and accuracy, it requires attention to context and responsibility. Prior work highlights the need to avoid the "portability trap", the assumption that what works well in one context will generalize across others by accounting for how models interact with specific human, social, and technical environments (Selbst et al., 2019; Martin Jr et al., 2020). Such evaluations should incorporate multiple metrics, distributional analyses, and qualitative assessments of impact (Mitchell et al., 2019). They should also be reflexive about social and intentional factors in dataset construction and model development (Scheurman et al., 2021; Miceli et al., 2021; Mantelero, 2018; McGregor et al., 2019), explicitly document data provenance and known biases (Andrus et al., 2021), and consider broader obligations such as ethical alignment, institutional accountability, and potential privacy leakage (Mantelero, 2018; McGregor et al., 2019; Yeom et al., 2018; Raji et al., 2021). Together, these considerations highlight that evaluating NLP systems for responsible deployment requires assessing not only performance, but also harms, benefits, and obligations within the application ecosystem (Hutchinson et al., 2022; Raji et al., 2021).

2.2 NLP Evaluation in Educational Contexts

Despite the growing adoption of large language models (LLMs) in educational settings, including reducing instructional workload through automated grading and content generation (Ali et al., 2021), supporting personalized learning interactions (Vincent-Lancrin and Van der Vlies, 2020), supporting digital and AI literacy (Eshet, 2004; Bundy, 2017), and fostering higher-order cognitive skills such as reasoning, creativity, and computational thinking (Baer, 1998; Amabile, 2018; Grover and Pea, 2013), educational contexts have received limited attention in dominant NLP evaluation research, which have primarily been devel-

oped for general-purpose language tasks. As a result, general NLP evaluation approaches often fail in the educational context because they prioritize surface-level correctness and fluency over deep-level learning processes. Accurate or well-formed language does not guarantee conceptual understanding, appropriate scaffolding (Vygotsky and Cole, 1978; Papert and Harel, 1991), or productive struggle (Kapur, 2008), nor does it adequately reflect learners' needs, engagement, or cognitive growth (Shaffer and Resnick, 1999; Shneiderman, 2022; Friedman et al., 2002). These limitations also lead to education-specific language challenges, such as domain-specific terms, ambiguity, sarcasm, and imbalanced data, which make human feedback and learning signals more difficult to interpret. As emphasized in prior work on educational measurement, meaningful evaluation must consider incorporating formative, contextual, and human-centered outcomes that capture how AI systems support learning in authentic, situated environments (Shute and Ventura, 2013; Thomas et al., 2024; Anderson et al., 1990; Shaffer and Resnick, 1999; Friedman et al., 2002).

2.3 Conceptual Human-centered Evaluation Framework

In this section, we revisit our research question and propose an initial NLP evaluation framework, along with a set of guiding questions grounded in prior literature. The framework targets observable linguistic behaviors and focuses on how instructional qualities, particularly those that emerge when educational scenarios act as stress tests, such as explanation transparency, behavioral consistency, and iterative refinement, are reflected in generated language across interaction turns. By defining each component and synthesizing the analysis into preliminary implications for future evaluation and development of NLP-based generative tutoring agents, we highlight limitations of automated evaluation practices, including LLM-as-a-judge approaches, as well as sociotechnical considerations that are often overlooked by dominant NLP benchmarks. This framing provides a starting point for iterative refinement through future research.

2.3.1 Evaluation Framework Components for NLP

We focus on a subset of our framework's components that most directly surface misalignments between current NLP benchmarks and the evaluation

182	needs of educational LLMs. These components	233
183	emphasize model behaviors that are critical to in-	234
184	structional effectiveness.	
185	Explainability Explainability refers to an	
186	agent’s ability to present its reasoning and decision-	
187	making processes in clear, contextually meaning-	
188	ful, and human-understandable terms (Nauta et al.,	
189	2023; Arrieta et al., 2020; Guidotti and Ruggieri,	
190	2019). In educational contexts, explanations are	
191	core instructional acts rather than auxiliary outputs.	
192	However, current NLP evaluations rarely assess	
193	whether explanations are faithful, interpretable,	
194	or supportive of human understanding. Relevant	
195	evaluation sub-components include reasoning clar-	
196	ity, traceability between inputs and outputs, and	
197	alignment between explanations and task outcomes	
198	(Rosenfeld, 2021; Silva et al., 2023; Ribera and	
199	Lapedriza, 2019; Bommasani et al., 2021; Brans-	
200	ford et al., 2000; Abbas et al., 2022; Amershi et al.,	
201	2019; Guidotti et al., 2018; Liao et al., 2020).	
202	Consistency Consistency captures the stability	
203	and trustworthiness of system behavior under simi-	
204	lar conditions, including alignment across prompts,	
205	languages, situations, and evaluators (Nauta et al.,	
206	2023; Nielsen, 1995). While educational trust de-	
207	pends on predictable feedback, most benchmarks	
208	report only aggregate scores, masking behavioral	
209	variance. Indicators of consistency include out-	
210	put stability across runs, sensitivity to prompt per-	
211	turbations, where minor rephrasings of semanti-	
212	cally equivalent prompts lead to divergent outputs,	
213	and cross-evaluator agreement (Carvalho et al.,	
214	2019; Robnik-Šikonja and Bohanec, 2018; Van-	
215	den Abeele et al., 2012).	
216	Refinement Refinement illustrated an agent’s	
217	capacity to support iterative human–AI interaction,	
218	including user-guided error correction, clarification	
219	of vague or biased feedback, and traceable revi-	
220	sions over time (Pan et al., 2024; Wang et al., 2024;	
221	Nielsen, 1995). Existing NLP evaluations primar-	
222	ily evaluate static, single-turn outputs, implicitly	
223	assuming quality is fixed at generation time. In con-	
224	trast, educational use is inherently iterative. Evalu-	
225	ating refinement behaviors such as responsiveness	
226	to feedback, coherence across revisions, time-to-	
227	refine, value alignment over revisions, these met-	
228	rics or indicators are largely absent from current	
229	evaluation components(Guo et al., 2024; Wang	
230	et al., 2024; Pan et al., 2024; Vasconcelos et al.,	
231	2022; Bryant et al., 2017; Hong et al., 2024; Ju-	
232	renka et al., 2024).	
	2.3.2 Reflecting Evaluation Through	
	Questions	
	We re-frame the evaluation of NLP systems in ed-	
	ucational contexts through a small set of guiding	
	questions grounded in a human-centered frame-	
	work(Jacovi and Goldberg, 2020; Shneiderman,	
	2022). These questions present how evaluators	
	can move beyond outcome-focused judgments to-	
	ward more transparent, process-oriented evalua-	
	tion. Practitioners may select, adapt, or extend	
	these questions based on their roles, contexts, and	
	decision-making needs.	
	Explainability	
	• Does the AI explain its decisions or recom-	
	mendations through its generated language in	
	a clear and contextually meaningful way?	
	• Are sources, assumptions, or reasoning made	
	visible to the user?	
	• To what extent would a user trust the system’s	
	explanation, and why?	
	• How could a user verify the explanation if	
	needed?	
	Consistency	
	• Does the system produce linguistically con-	
	sistent outputs across similar prompts, condi-	
	tions, or learning contexts?	
	• Are evaluations stable across sessions or hu-	
	man evaluators?	
	Refinement	
	• Does the AI support iterative revision of its	
	generated responses across multiple attempts	
	rather than providing a single final answer?	
	• Do users understand why a change or refine-	
	ment was suggested?	
	• Does the refinement process remain aligned	
	with the learner’s context and goals?	
	These questions are intended as starting points	
	for reflective evaluation, helping surface pedagog-	
	ical, social, and human-centered considerations	
	that support discussion, comparison, and reflection	
	across stakeholders.	

274	2.3.3 Illustrative Application: Evaluating an LLM assistant Response	3.2 Implication 2: Evaluation as a Sociotechnical Process	322
275			323
276	To show how our framework presents differences	Dominant NLP evaluation practices focus on	324
277	missed by current NLP metrics, consider an LLM	single-turn correctness and surface-level similarity	325
278	programming assistant responding to the prompt:	metrics, such as BLEU, ROUGE, Exact Match, and	326
279	"Why does the <code>fitMedia</code> call in this Python code	embedding-based overlap scores. These metrics	327
280	produce no output?" One response provides a di-	work well for short tasks with one correct answer,	328
281	rect answer by identifying a misspelled sound	but they do not capture how educational interac-	329
282	name, while another also explains how fixing er-	tions involve over time through dialogue and adap-	330
283	rors allows the loop to run and connects the fix	tive scaffolding. In educational settings, correct-	331
284	to curriculum-aligned debugging. Standard auto-	ness alone does not indicate whether a model helps	332
285	mated metrics such as BLEU, embedding-based	learners understand concepts, supports productive	333
286	similarity, and LLM-as-a-judge score, rate both	struggle, or aligns with instructional goals. These	334
287	responses similarly as fluent and correct. In con-	limitations highlight evaluation as a sociotechnical	335
288	trast, our human-centered evaluation framework	process shaped by human values, institutional ex-	336
289	distinguishes the second response by its stronger	pectations, and interactional context. Framing eval-	337
290	explainability and closer alignment with instruc-	uation around stakeholder-aligned questions makes	338
291	tional debugging practices in educational stress-test	these assumptions explicit and surfaces concrete	339
292	scenarios.	aspects of model behavior, positioning evaluation	340
		as an integral practice in how educational language	341
293	3 Implications for NLP Evaluation Research	models are assessed and governed.	342
294		Together, these implications suggest that future	343
295	Our analysis highlights key limitations in dom-	NLP benchmarks and evaluation methods should	344
296	inant NLP evaluation practices, particularly the	incorporate interaction-level metrics, multi-turn	345
297	common assumptions about correctness, fluency,	consistency analyses, and hybrid human–AI eval-	346
298	and automated judgment when applied to educa-	uation approaches, rather than relying solely on	347
299	tional LLMs. We distill two early implications that	single-turn automated judgments.	348
300	clarify why current evaluation practices are insuf-		
301	ficient for human-centered, instructional language	4 Limitations	349
302	models.		
303	3.1 Implication 1: Limitation of LLM-as-a-Judge for Evaluation	This work presents a conceptual analysis and does	350
304		not involve human subjects, data collection, or sys-	351
305	LLM-as-a-Judge frameworks offer scalability and	tem deployment at the current stage. The proposed	352
306	convenience and have shown promise in correlat-	framework has not yet been empirically validated	353
307	ing human preference judgments. However, our	or evaluated with relevant stakeholders, and the	354
308	analysis identifies important limitations in their ap-	current discussion focuses on LLMs in educational	355
309	plicability to pedagogical evaluation. Educational	contexts.	356
310	quality depends on domain expertise, sensitivity to	5 Conclusion and Future Work	357
311	learner context, and value-based judgment, which		
312	cannot be reliably inferred from surface-level tex-	This paper argues that current NLP evaluation prac-	358
313	tual comparison alone. In instructional settings,	tices exhibit gaps when applied to educational	359
314	LLM judges tend to prioritize fluent responses	LLMs and emphasizes human-centered criteria	360
315	while overlooking misleading explanations, inap-	such as explainability, consistency, and refinement	361
316	propriate levels of assistance, or misalignment with	in NLP evaluation. At this stage, the work presents	362
317	learning objectives. These patterns indicate that	a conceptual analysis without empirical validation,	363
318	LLM-based judgment captures only a partial view	and the discussion is limited to LLMs in educa-	364
319	of educational quality and highlights the need for	tional contexts. Looking ahead, we plan to it-	365
320	hybrid evaluation approaches that combine auto-	eratively refine the framework through feedback	366
321	mated judgment with human expertise.	from interdisciplinary stakeholders in education	367
		and to conduct empirical and cross-domain studies	368
		to examine how its evaluation principles generalize	369
		within and beyond educational domains.	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
411
412
413
414
415
416
417
418
419
420
421
422

5.1 References

Acknowledgments

References

Abdallah MH Abbas, Khairil Imran Ghauth, and Choo-Yee Ting. 2022. User experience design using machine learning: a systematic review. *IEEE Access*, 10:51501–51514.

Safinah Ali, Hae Won Park, and Cynthia Breazeal. 2021. A social robot’s influence on children’s figural creativity during gameplay. *International Journal of Child-Computer Interaction*, 28:100234.

Teresa M Amabile. 2018. *Creativity in context: Update to the social psychology of creativity*. Routledge.

Saleema Amershi, Dan Weld, Mihaela Vorvoreanu, Adam Fourney, Besmira Nushi, Penny Collisson, Jina Suh, Shamsi Iqbal, Paul N Bennett, Kori Inkpen, and 1 others. 2019. Guidelines for human-ai interaction. In *Proceedings of the 2019 chi conference on human factors in computing systems*, pages 1–13.

John R Anderson, C Franklin Boyle, Albert T Corbett, and Matthew W Lewis. 1990. Cognitive modeling and intelligent tutoring. *Artificial intelligence*, 42(1):7–49.

McKane Andrus, Elena Spitzer, Jeffrey Brown, and Alice Xiang. 2021. What we can’t measure, we can’t understand: Challenges to demographic data procurement in the pursuit of fairness. In *Proceedings of the 2021 ACM conference on fairness, accountability, and transparency*, pages 249–260.

Alejandro Barredo Arrieta, Natalia Díaz-Rodríguez, Javier Del Ser, Adrien Bannetot, Siham Tabik, Alberto Barbado, Salvador García, Sergio Gil-López, Daniel Molina, Richard Benjamins, and 1 others. 2020. Explainable artificial intelligence (xai): Concepts, taxonomies, opportunities and challenges toward responsible ai. *Information fusion*, 58:82–115.

John Baer. 1998. The case for domain specificity of creativity. *Creativity research journal*, 11(2):173–177.

Rishi Bommasani, Drew A Hudson, Ehsan Adeli, Russ Altman, Simran Arora, Sydney von Arx, Michael S Bernstein, Jeannette Bohg, Antoine Bosselut, Emma Brunskill, and 1 others. 2021. On the opportunities and risks of foundation models. *arXiv preprint arXiv:2108.07258*.

John D Bransford, Ann L Brown, Rodney R Cocking, and 1 others. 2000. *How people learn*, volume 11. Washington, DC: National academy press.

CJ Bryant, Mariano Felice, and Edward Briscoe. 2017. Automatic annotation and evaluation of error types for grammatical error correction. Association for Computational Linguistics.

Alan Bundy. 2017. Preparing for the future of artificial intelligence. 423
424

Diogo V Carvalho, Eduardo M Pereira, and Jaime S Cardoso. 2019. Machine learning interpretability: A survey on methods and metrics. *Electronics*, 8(8):832. 425
426
427

Yoram Eshet. 2004. Digital literacy: A conceptual framework for survival skills in the digital era. *Journal of educational multimedia and hypermedia*, 13(1):93–106. 428
429
430
431

Batya Friedman, Peter Kahn, and Alan Borning. 2002. Value sensitive design: Theory and methods. *University of Washington technical report*, 2(8):1–8. 432
433
434

Shuchi Grover and Roy Pea. 2013. Computational thinking in k–12: A review of the state of the field. *Educational researcher*, 42(1):38–43. 435
436
437

Riccardo Guidotti, Anna Monreale, Salvatore Ruggieri, Franco Turini, Fosca Giannotti, and Dino Pedreschi. 2018. A survey of methods for explaining black box models. *ACM computing surveys (CSUR)*, 51(5):1–42. 438
439
440
441
442

Riccardo Guidotti and Salvatore Ruggieri. 2019. On the stability of interpretable models. In *2019 international joint conference on neural networks (IJCNN)*, pages 1–8. IEEE. 443
444
445
446

Shuchen Guo, Ehsan Latif, Yifan Zhou, Xuan Huang, and Xiaoming Zhai. 2024. Using generative ai and multi-agents to provide automatic feedback. *arXiv preprint arXiv:2411.07407*. 447
448
449
450

Shengxin Hong, Chang Cai, Sixuan Du, Haiyue Feng, Siyuan Liu, and Xiuyi Fan. 2024. " my grade is wrong!": A contestable ai framework for interactive feedback in evaluating student essays. *arXiv preprint arXiv:2409.07453*. 451
452
453
454
455

Ben Hutchinson, Negar Rostamzadeh, Christina Greer, Katherine Heller, and Vinodkumar Prabhakaran. 2022. Evaluation gaps in machine learning practice. In *Proceedings of the 2022 ACM conference on fairness, accountability, and transparency*, pages 1859–1876. 456
457
458
459
460
461

Alon Jacovi and Yoav Goldberg. 2020. Towards faithfully interpretable nlp systems: How should we define and evaluate faithfulness? *arXiv preprint arXiv:2004.03685*. 462
463
464
465

Irina Jurenka, Markus Kunesch, Kevin R McKee, Daniel Gillick, Shaojian Zhu, Sara Wiltberger, Shubham Milind Phal, Katherine Hermann, Daniel Kasenberg, Avishkar Bhoopchand, and 1 others. 2024. Towards responsible development of generative ai for education: An evaluation-driven approach. *arXiv preprint arXiv:2407.12687*. 466
467
468
469
470
471
472

Manu Kapur. 2008. Productive failure. *Cognition and instruction*, 26(3):379–424. 473
474

475	Dawei Li, Bohan Jiang, Liangjie Huang, Alimohammad	Jiayi Pan, Yichi Zhang, Nicholas Tomlin, Yifei Zhou,	530
476	Beigi, Chengshuai Zhao, Zhen Tan, Amrita Bhat-	Sergey Levine, and Alane Suhr. 2024. Autonomous	531
477	tacharjee, Yuxuan Jiang, Canyu Chen, Tianhao Wu,	evaluation and refinement of digital agents. <i>arXiv</i>	532
478	and 1 others. 2025. From generation to judgment:	<i>preprint arXiv:2404.06474</i> .	533
479	Opportunities and challenges of llm-as-a-judge. In		
480	<i>Proceedings of the 2025 Conference on Empirical</i>	Seymour Papert and Idit Harel. 1991. Situating con-	534
481	<i>Methods in Natural Language Processing</i> , pages	structionism. <i>constructionism</i> , 36(2):1–11.	535
482	2757–2791.		
483	Tianle Li, Wei-Lin Chiang, Evan Frick, Lisa Dunlap,	Kishore Papineni, Salim Roukos, Todd Ward, and Wei-	536
484	Tianhao Wu, Banghua Zhu, Joseph E Gonzalez, and	Jing Zhu. 2002. Bleu: a method for automatic evalu-	537
485	Ion Stoica. 2024. From crowdsourced data to high-	ation of machine translation. In <i>Proceedings of the</i>	538
486	quality benchmarks: Arena-hard and benchbuilder	<i>40th annual meeting of the Association for Computa-</i>	539
487	pipeline. <i>arXiv preprint arXiv:2406.11939</i> .	<i>tional Linguistics</i> , pages 311–318.	540
488	Q Vera Liao, Daniel Gruen, and Sarah Miller. 2020.	Inioluwa Deborah Raji, Emily M Bender, Amandalynne	541
489	Questioning the ai: informing design practices for ex-	Paullada, Emily Denton, and Alex Hanna. 2021. Ai	542
490	plainable ai user experiences. In <i>Proceedings of the</i>	and the everything in the whole wide world bench-	543
491	<i>2020 CHI conference on human factors in computing</i>	mark. <i>arXiv preprint arXiv:2111.15366</i> .	544
492	<i>systems</i> , pages 1–15.		
493	Chin-Yew Lin. 2004. Rouge: A package for automatic	Mireia Ribera and Agata Lapedriza. 2019. Can we	545
494	evaluation of summaries. In <i>Text summarization</i>	do better explanations? a proposal of user-centered	546
495	<i>branches out</i> , pages 74–81.	explainable ai. CEUR Workshop Proceedings.	547
496	Alessandro Mantelero. 2018. Ai and big data: A	Marko Robnik-Šikonja and Marko Bohanec. 2018.	548
497	blueprint for a human rights, social and ethical im-	Perturbation-based explanations of prediction models.	549
498	act assessment. <i>Computer Law & Security Review</i> ,	In <i>Human and Machine Learning: Visible, Explain-</i>	550
499	34(4):754–772.	<i>able, Trustworthy and Transparent</i> , pages 159–175.	551
500	Donald Martin Jr, Vinodkumar Prabhakaran, Jill	Springer.	552
501	Kuhlberg, Andrew Smart, and William S Isaac. 2020.	Avi Rosenfeld. 2021. Better metrics for evaluating ex-	553
502	Extending the machine learning abstraction bound-	plainable artificial intelligence. In <i>Proceedings of the</i>	554
503	ary: A complex systems approach to incorporate	<i>20th international conference on autonomous agents</i>	555
504	societal context. <i>arXiv preprint arXiv:2006.09663</i> .	<i>and multiagent systems</i> , pages 45–50.	556
505	Lorna McGregor, Daragh Murray, and Vivian Ng. 2019.	Morgan Klaus Scheuerman, Alex Hanna, and Remi Den-	557
506	International human rights law as a framework for	ton. 2021. Do datasets have politics? disciplinary	558
507	algorithmic accountability. <i>International & Compar-</i>	values in computer vision dataset development. <i>Pro-</i>	559
508	<i>ative Law Quarterly</i> , 68(2):309–343.	<i>ceedings of the ACM on Human-Computer Interac-</i>	560
509	Milagros Miceli, Tianling Yang, Laurens Naudts, Mar-	<i>tion</i> , 5(CSCW2):1–37.	561
510	tin Schuessler, Diana Serbanescu, and Alex Hanna.	Andrew D Selbst, Danah Boyd, Sorelle A Friedler,	562
511	2021. Documenting computer vision datasets: An	Suresh Venkatasubramanian, and Janet Vertesi. 2019.	563
512	invitation to reflexive data practices. In <i>Proceedings</i>	Fairness and abstraction in sociotechnical systems.	564
513	<i>of the 2021 ACM Conference on Fairness, Account-</i>	In <i>Proceedings of the conference on fairness, ac-</i>	565
514	<i>ability, and Transparency</i> , pages 161–172.	<i>countability, and transparency</i> , pages 59–68.	566
515	Margaret Mitchell, Simone Wu, Andrew Zaldivar,	David Williamson Shaffer and Mitchel Resnick. 1999.	567
516	Parker Barnes, Lucy Vasserman, Ben Hutchinson,	"thick" authenticity: New media and authentic	568
517	Elena Spitzer, Inioluwa Deborah Raji, and Timnit	learning. <i>Journal of interactive learning research</i> ,	569
518	Gebriu. 2019. Model cards for model reporting. In	10(2):195–216.	570
519	<i>Proceedings of the conference on fairness, account-</i>	Ben Shneiderman. 2022. <i>Human-centered AI</i> . Oxford	571
520	<i>ability, and transparency</i> , pages 220–229.	University Press.	572
521	Meike Nauta, Jan Trienes, Shreyasi Pathak, Elisa	Valerie Shute and Matthew Ventura. 2013. <i>Stealth as-</i>	573
522	Nguyen, Michelle Peters, Yasmin Schmitt, Jörg	<i>essment: Measuring and supporting learning in</i>	574
523	Schlötterer, Maurice Van Keulen, and Christin Seifert.	<i>video games</i> . The mit press.	575
524	2023. From anecdotal evidence to quantitative eval-	Andrew Silva, Mariah Schrum, Erin Hedlund-Botti,	576
525	uation methods: A systematic review on evaluating	Nakul Gopalan, and Matthew Gombolay. 2023. Ex-	577
526	explainable ai. <i>ACM Computing Surveys</i> , 55(13s):1–	plainable artificial intelligence: Evaluating the ob-	578
527	42.	jective and subjective impacts of xai on human-	579
528	Jakob Nielsen. 1995. Ten usability heuristics for user	agent interaction. <i>International Journal of Human-</i>	580
529	interface design . Online; accessed July X, 2025.	<i>Computer Interaction</i> , 39(7):1390–1404.	581

582	Annalisa Szymanski, Noah Ziem, Heather A Eicher-Miller, Toby Jia-Jun Li, Meng Jiang, and Ronald A Metoyer. 2025. Limitations of the llm-as-a-judge approach for evaluating llm outputs in expert knowledge tasks. In <i>Proceedings of the 30th International Conference on Intelligent User Interfaces</i> , pages 952–966.		
583			
584			
585			
586			
587			
588			
589	Danielle R Thomas, Jionghao Lin, Erin Gatz, Ashish Gurung, Shivang Gupta, Kole Norberg, Stephen E Fancsali, Vincent Aleven, Lee Branstetter, Emma Brunskill, and 1 others. 2024. Improving student learning with hybrid human-ai tutoring: A three-study quasi-experimental investigation. In <i>Proceedings of the 14th Learning Analytics and Knowledge Conference</i> , pages 404–415.		
590			
591			
592			
593			
594			
595			
596			
597	Vero Vanden Abeele, Erik Hauters, and Bieke Zaman. 2012. Increasing the reliability and validity of quantitative laddering data with ladderux. In <i>CHI'12 Extended Abstracts on Human Factors in Computing Systems</i> , pages 2057–2062.		
598			
599			
600			
601			
602	Helena Vasconcelos, Gagan Bansal, Adam Fournay, Q Vera Liao, and Jennifer Wortman Vaughan. 2022. Generation probabilities are not enough: Improving error highlighting for ai code suggestions. In <i>HCAI Workshop at NeurIPS</i> .		
603			
604			
605			
606			
607	Stéphan Vincent-Lancrin and Reyer Van der Vlies. 2020. Trustworthy artificial intelligence (ai) in education: Promises and challenges. <i>OECD education working papers</i> , (218):0_1–17.		
608			
609			
610			
611	Lev Semenovich Vygotsky and Michael Cole. 1978. <i>Mind in society: Development of higher psychological processes</i> . Harvard university press.		
612			
613			
614	Alex Wang, Yada Pruksachatkun, Nikita Nangia, Amanpreet Singh, Julian Michael, Felix Hill, Omer Levy, and Samuel Bowman. 2019. Superglue: A sticker benchmark for general-purpose language understanding systems. <i>Advances in neural information processing systems</i> , 32.		
615			
616			
617			
618			
619			
620	Alex Wang, Amanpreet Singh, Julian Michael, Felix Hill, Omer Levy, and Samuel Bowman. 2018. Glue: A multi-task benchmark and analysis platform for natural language understanding. In <i>Proceedings of the 2018 EMNLP workshop BlackboxNLP: Analyzing and interpreting neural networks for NLP</i> , pages 353–355.		
621			
622			
623			
624			
625			
626			
627	Han Wang, An Zhang, Nguyen Duy Tai, Jun Sun, Tat-Seng Chua, and 1 others. 2024. Ali-agent: Assessing llms' alignment with human values via agent-based evaluation. <i>Advances in Neural Information Processing Systems</i> , 37:99040–99088.		
628			
629			
630			
631			
632	Samuel Yeom, Irene Giacomelli, Matt Fredrikson, and Somesh Jha. 2018. Privacy risk in machine learning: Analyzing the connection to overfitting. In <i>2018 IEEE 31st computer security foundations symposium (CSF)</i> , pages 268–282. IEEE.		
633			
634			
635			
636			
		Tianyi Zhang, Varsha Kishore, Felix Wu, Kilian Q Weinberger, and Yoav Artzi. 2019. Bertscore: Evaluating text generation with bert. <i>arXiv preprint arXiv:1904.09675</i> .	637 638 639 640
		Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan Zhuang, Zhanghao Wu, Yonghao Zhuang, Zi Lin, Zhuohan Li, Dacheng Li, Eric Xing, and 1 others. 2023. Judging llm-as-a-judge with mt-bench and chatbot arena. <i>Advances in neural information processing systems</i> , 36:46595–46623.	641 642 643 644 645 646
		A Example Appendix	647
		This is an appendix.	648