

SEFL: Harnessing Large Language Model Agents to Improve Educational Feedback Systems

Anonymous ACL submission

Abstract

Providing high-quality feedback is crucial for student success but is constrained by time, cost, and limited data availability. We introduce Synthetic Educational Feedback Loops (SEFL), a novel framework designed to deliver immediate, on-demand feedback at scale without relying on extensive, real-world student data. In SEFL, two large language models (LLMs) operate in teacher–student roles to simulate assignment completion and formative feedback, generating abundant synthetic pairs of student work and corresponding critiques. We then fine-tune smaller, more computationally efficient LLMs on these synthetic pairs, enabling them to replicate key features of high-quality, goal-oriented feedback. Unlike personalized tutoring approaches that offer multi-turn, individualized instruction, SEFL specifically focuses on replicating the teacher→student feedback loop for diverse assignments. Through both LLM-as-a-judge and human evaluations, we demonstrate that SEFL-tuned models outperform their non-tuned counterparts in feedback quality, clarity, and timeliness. These findings reveal SEFL’s potential to transform feedback processes for higher education and beyond, offering an ethical and scalable alternative to conventional manual feedback cycles.

1 Introduction

Constructive feedback is a cornerstone of higher education, promoting critical thinking and fostering deeper understanding (Hattie, 2008; Costello and Crane, 2013). In many higher education settings, however, providing consistent, high-quality feedback remains a labor-intensive task, further complicated by privacy, consent, and transparency considerations in data collection (Fischer et al., 2020; Suresh et al., 2022; Demszky and Hill, 2023; Wang and Demszky, 2024; Wang et al., 2024a; Lindsay et al., 2024). Advances in NLP offer promising

opportunities to simulate and augment feedback processes, addressing these limitations.

With respect to language technology, prior research has explored areas such as peer learning (Bauer et al., 2023), aligning mathematical questions (Botelho et al., 2023), enhancing critical thinking (Guerraoui et al., 2023), and using large language models (LLMs) for research feedback alignment (Liang et al., 2024; Sonkar et al., 2024). Tools for monitoring student progress (Schwarz et al., 2018; Aslan et al., 2019; Alrajhi et al., 2021) have also been investigated. However, to the best of our knowledge, this work is the first to leverage LLMs for generating abundant and scalable feedback for student work. Researchers have identified key characteristics of “good feedback”, including goal-orientation, actionability, timeliness, user-friendliness, and consistency, as well as fostering student autonomy through self-evaluation (Carless et al., 2011; Wiggins, 2012). Overly elaborate commentary can undermine clarity, highlighting the value of brevity. Moreover, immediate, formative feedback is crucial for continuous improvement (Wiggins, 2012), a requirement that LLM-based systems are well suited to fulfill.

LLMs have shown remarkable capabilities in education (Wang et al., 2024b), including automated grading (Ke and Ng, 2019; Ramesh and Sanampudi, 2022; Stahl et al., 2024) and personalized tutoring (Yun et al., 2024; Liu et al., 2024b; Rooein and Hovy, 2024; Ross and Andreas, 2024; Kwon et al., 2024; Zhang et al., 2024a). Yet, simulating dynamic teacher–student feedback interactions in agentic, dialogic settings (Xi et al., 2023; Guo et al., 2024; Zhang et al., 2024b) remains largely unexplored, despite its potential to generate scalable synthetic datasets and alleviate real-world data scarcity. We seek to answer: *How can synthetic teacher–student interactions generated by LLMs be leveraged to enable scalable and effective educational feedback systems?*

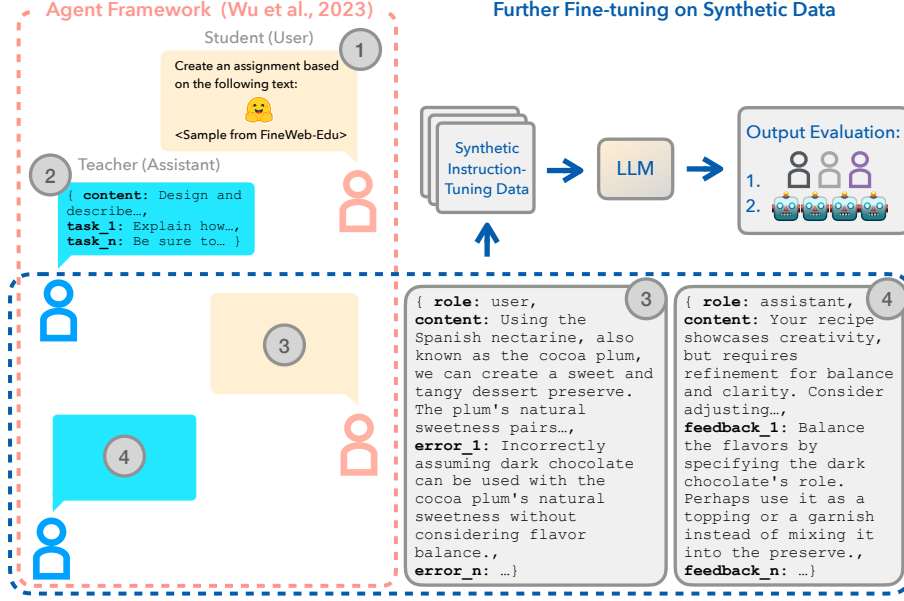


Figure 1: SEFL Setup. We use a two-agent framework (Wu et al., 2023) with LLMs acting as a Student and Teacher. The Teacher creates assignments from Fineweb-Edu (Lozhkov et al., 2024), the Student responds with errors, and finally the Teacher addresses each mistake. This synthetic interaction data is then used to fine-tune multiple LLMs, whose performance is measured via human ratings and an LLM-as-judge approach.

To this end, we introduce **Synthetic Educational Feedback Loops (SEFL)**, a framework that generates synthetic teacher-student interactions using LLMs. In this framework, two LLMs—one acting as the teacher and the other as the student—simulate *formative*¹ feedback workflows, addressing the limitations of using a single LLM for multiple tasks. This synthetic data is then used to fine-tune smaller autoregressive models, enabling the development of scalable educational feedback systems that can operate efficiently on modest computational infrastructure, such as that available in higher education institutions.

Contributions. To answer the research question, we contribute the following: ① A novel framework for simulating teacher-student feedback loops using agentic LLMs. ② A pipeline for generating synthetic educational data to fine-tune smaller models. ③ An LLM-as-a-judge framework for rating feedback using GPT-4o, Claude-3.5, Command-R+, and DeepseekV3. ④ An open-source release of all the models, data, and code.²

2 Synthetic Educational Feedback Interactions

2.1 Synthetic Data Generation

We use a two-agent framework (Wu et al., 2023). Both the teacher and student roles are simulated

¹Formative feedback is used early in the learning process, allowing students to refine their work and deepen their understanding (Conole and Oliver, 2006; Nicol, 2007).

²Code and resources available at <https://anonymous.4open.science/r/sefl-4B9F/>.

	Valid (/ 5,000)	BERTScore
Llama-3.1-70B	2,513	0.877
Qwen2.5-72B	454	0.919

Table 1: Generation Capabilities. First, We show the number of valid examples, measured by correct JSON format and whether each feedback refers to an error. Llama-3.1-70B generates more valid examples. Second, we measure BERTScore as a proxy for relatedness between error-feedback pairs of the valid generations.

by two separate Llama-3.1-70B models for a two-turn conversation.³ The models are tasked to generate assignment→answer→feedback tuples. First, the student-agent asks for an assignment using Fineweb-Edu (Lozhkov et al., 2024) texts (①). Second, the teacher-agent creates an assignment that can be of any domain, e.g., math, humanities, role-playing (②; Figure 1). Then, the student-agent (③) submits assignments containing a number of explicit errors, and the teacher-agent (④) provides targeted feedback addressing each error. We investigated both Qwen2.5-72B and Llama-3.1-70B for interactions. We initially generated 5,000 interaction tuples with each model, where we validated the output as a sanity check.

We show in Table 1 the results of this experiment. Out of 5,000 examples, Llama-3.1-70B generates the most valid examples (i.e., valid JSON format and each feedback refers to an error). For a further check, we use BERTScore (Zhang et al., 2020) as a proxy to investigate whether each error-feedback

³Note that if we mention a model, it is always the *post-trained* version (i.e., -Instruct).

Feature	Value
Instances	19,841
Assignment Length	78.6
Length (Student)	168.1
# Errors Points	2.5
Length # Errors	20.7
Length (Teacher)	120.5
# Feedback Points	2.5
Length # Feedback	34.6

Table 2: **Generation Statistics.** We show the dataset statistics in *averages*, where length is measured in whitespace-separated tokens.

pair of the valid generations relate to each other.⁴ We show regardless of Llama-3.1-70B generating more valid examples, the BERTScore stays in a similar range as Qwen2.5-72B. Consequently, we use Llama-3.1-70B-generated data as the basis for all subsequent model fine-tuning. For the full prompt, see Figure 2 (Appendix B).

Statistics. Table 2 presents the final dataset. The generation lengths for each agent are intentionally kept concise (<170 tokens), based on the hypothesis that overly lengthy feedback may be counter-productive. This is in line with observations from Ferguson (2011), who observes that students tend to favor brief comments, finding a general overview of an assignment more useful. Balancing supportive and critical feedback is crucial as, by default, LLMs often produce excessively verbose responses, which can influence the preferences of both humans and language models (Saito et al., 2023).

2.2 Fine-Tuning

The total amount of data synthesized by Llama-3.1-70B amounts to 19.8K conversations, which we use to fine-tune five smaller open-weight LLMs: Qwen2.5-0.5B, Llama-3.2-1B, Llama-3.2-3B, Llama-3.1-8B, Qwen2.5-14B. Each model is further instruction-tuned using a standard language modeling objective (see Appendix A for more details).

2.3 Evaluation

Human Evaluation. To test the performance of SEFL, we have a human evaluation pipeline. We randomly sample 150 samples from the validation set. Then, we have both the original instruction-tuned model (A) and the model that was further fine-tuned with SEFL (B). We have three human

⁴We only calculate it of the samples where both error and feedback have the same number of generations.

Models	H1	H2	H3	J1	J2	J3	J4
Qwen2.5-0.5B	94	85	85	97	91	62	91
Llama-3.2-1B	97	85	81	79	91	27	79
Llama-3.2-3B	90	61	65	71	74	26	77
Llama-3.1-8B	90	45	94	39	71	16	65
Qwen2.5-14B	94	77	81	55	65	10	19

Table 3: **Results in Win Rate.** We show the win rate of our *SEFL-tuned models*. A win rate >50% indicates that SEFL-tuned models are better in giving feedback than their vanilla-counterpart; in red everything <50%. We show results of 3 human annotators (H#) and 4 LLM judges: gpt-4o (J1), claude-3.5-sonnet (J2), command-r-plus (J3), and deepseek-v3 (J4).

raters judge whether $A > B$ or $A < B$. Additionally, we also ask the coders to indicate whether the assignment→student answer→feedback tuple are related to each other or whether the model seems to be generating unrelated content. Our human raters are in the age range of 20–40 and from Europe, two have a background in Computer Science and one in Engineering Education, they all work in higher education with near-native English proficiency. For more details, the annotation guidelines can be found in Table 5 (Appendix C).

LLM-as-a-Judge. We also evaluate the fine-tuned models’ output using a LLM-as-a-judge framework, a method gaining traction as a method for evaluating text output (Liu et al., 2023; Zheng et al., 2024; Chen et al., 2023; Verga et al., 2024; Törnberg, 2023; Naismith et al., 2023; Giliardi et al., 2023; Kocmi and Federmann, 2023; Huang et al., 2024; Gu et al., 2024; Falk et al., 2025). The same 150 random instances are rated by four LLMs, namely GPT-4o (Hurst et al., 2024), Claude3.5-Sonnet, Command-R+, and DeepSeek-V3 (Liu et al., 2024a). We picked these models based on their recency and performance on RewardBench (Lambert et al., 2024), JudgeBench (Tan et al., 2024), and JudgeArena.⁵ For the full prompt, see Figure 3 (Appendix B).

3 Results

Our results are in Table 3. We show the *win rates* of models fine-tuned with SEFL vs. their original, non-tuned versions, as evaluated by both human raters and an LLM-based judges. A value above 50% indicates that the SEFL-tuned models are preferred over their original versions.

⁵<https://huggingface.co/spaces/AtlaAI/judge-arena>.

Human Assessment. Overall, human rater evaluations in Table 3 show that the SEFL-tuned models often attain high win rates, surpassing 85% in several cases. Annotators differed in their views on the 8B model’s output quality; however, they generally converged on the observation that the fine-tuned 14B model produces superior feedback compared to its non-tuned version. By contrast, models not fine-tuned with SEFL had lower win rates, suggesting that the synthetic feedback loops provide an edge in generating more coherent and context-relevant feedback. In addition, we asked annotators whether the synthetic assignment→answer→feedback sequences were consistent. In over 75% of cases, they affirmed the alignment between assignment, student response, and the feedback given, showing the pipeline’s effectiveness in keeping contextual relevance.

LLM-as-a-Judge Results. For the LLM-as-a-judge evaluations, we observe notable differences in win rates depending on the model and scale. The results largely mirror the human assessment trend up to the 3B scale. The results from the four LLM judges (J1: gpt-4o-2024-08-06, J2: claude-3-5-sonnet-20241022, J3: command-r-plus-08-2024, J4: deepseek-v3) reveal that SEFL-tuned models demonstrate varying levels of performance relative to their non-tuned counterparts. For instance, Qwen2.5-0.5B achieved the highest win rates across all four judges (62% on J3), indicating a consistent preference for the fine-tuned version. In contrast, larger models such as Llama-3.1-8B and Qwen2.5-14B exhibit lower win rates, particularly on J3 (16% and 10%, respectively), suggesting that fine-tuning with SEFL may yield diminishing returns or challenges at larger scales.

Agreement. We calculate the pairwise agreement between the judges and human raters. The results show a Cohen’s k values between 0.48–0.63, see Appendix E. Though this is considered a *moderate* to *substantial agreement* (Landis and Koch, 1977), it indicates the subjectivity of feedback.

4 Discussion

Human Qualitative Insights. In addition to the quantitative win rates summarized in Table 3, our human annotators provided rich qualitative feedback on the generated responses. Generally, the annotators notice that if the student answer

is too short or incomplete, neither model could generate appropriate feedback on that the assignment is incomplete. More specifically, feedback from Qwen2.5-0.5B was frequently noted for being clear and concise, while Llama-3.2-3B sometimes reiterated assignment details without offering actionable suggestions. Annotators commented that Llama-3.2-1B generally provided more specific and actionable feedback, yet occasionally its tone was perceived as too harsh, whereas Llama-3.1-8B often missed key aspects of the answer. Meanwhile, Qwen2.5-14B was critiqued for being overly verbose and less aligned with the assignment context. Overall, although Qwen2.5-0.5B achieved high human win rates (94, 85, and 85 across three annotators), the qualitative insights suggest that even the best-performing models could improve in error detection, tone refinement, and contextual sensitivity. For all the comments, we refer to Table 6 (Appendix D).

LLM-as-a-Judge. We used LLM judges to rate the feedback generated by SEFL-tuned models against their vanilla counterparts. This provides a rapid, scalable way to measure feedback quality, reducing the need for extensive human annotation. Three out of four LLM judges consistently favored SEFL-tuned Qwen2.5-0.5B, Llama-3.2-1B, Llama-3.2-3B, and Qwen2.5-72B. With Command-R, we notice that it performs worse than GPT-4o and Claude3.5-Sonnet on JudgeArena, indicating that the performance might have to do with instruction following. Nonetheless, we see it as a practical first step for large-scale feedback comparisons in educational contexts. We recommend supplementing LLM-based assessments with targeted human evaluations for more granular insights, possibly aligning more with authentic instructional objectives.

5 Conclusion

We introduced SEFL, a framework that simulates teacher→student interactions via two-agent LLMs to generate synthetic data for fine-tuning smaller models. This yields concise, context-sensitive feedback that often surpasses the performance of original instruction-tuned models. While LLM judges provide a scalable way to assess feedback quality, human insights remain crucial for capturing nuances like clarity and tone. As higher education digitalizes, SEFL offers a promising avenue for immediate, personalized feedback at scale.

Limitations

SEFL relies on synthetically generated assignments and errors, and are not real student submissions, which could have implications. Although this approach helps create large datasets, it risks producing feedback unaligned with authentic classroom contexts. Our evaluation also uses LLM-based judges, introducing potential biases related to each judge's training data and objectives. Lastly, while we focused on short-answer tasks, longer or more domain-specific assignments may require specialized or more diverse synthetic data.

Ethical Considerations

The use of synthetic data provides an opportunity to train automated feedback systems without the constraints of privacy and consent that come from repurposing actual student assignments as training data. However, it also raises questions about transparency and potential misuse (Lindsay et al., 2024). For instance, malicious actors could manipulate synthetic data to disseminate misleading or biased feedback, undermining trust in educational tools. Users may also mistake synthetic feedback for real, expert guidance. Moreover, automated feedback systems risk reinforcing biases if the underlying models carry skewed training data. We believe educators and institutions should remain aware of these risks and incorporate human oversight to ensure that such systems complement, rather than replace, genuine pedagogical engagement.

References

- Laila Alrajhi, Ahmed Alamri, Filipe Dwan Pereira, and Alexandra I Cristea. 2021. Urgency analysis of learners' comments: An automated intervention priority model for mooc. In *Intelligent Tutoring Systems: 17th International Conference, ITS 2021, Virtual Event, June 7–11, 2021, Proceedings 17*, pages 148–160. Springer.
- Sinem Aslan, Nese Alyuz, Cagri Tanriover, Sinem E Mete, Eda Okur, Sidney K D'Mello, and Asli Arslan Esme. 2019. Investigating the impact of a real-time, multimodal student engagement analytics technology in authentic classrooms. In *Proceedings of the 2019 chi conference on human factors in computing systems*, pages 1–12.
- Elisabeth Bauer, Martin Greisel, Ilia Kuznetsov, Markus Berndt, Ingo Kollar, Markus Dresel, Martin R Fischer, and Frank Fischer. 2023. Using natural language processing to support peer-feedback in the age of artificial intelligence: A cross-disciplinary framework and a research agenda. *British Journal of Educational Technology*, 54(5):1222–1245.
- Anthony Botelho, Sami Baral, John A Erickson, Priyanka Benachamardi, and Neil T Heffernan. 2023. Leveraging natural language processing to support automated assessment and feedback for student open responses in mathematics. *Journal of computer assisted learning*, 39(3):823–840.
- David Carless, Diane Salter, Min Yang, and Joy Lam. 2011. Developing sustainable feedback practices. *Studies in higher education*, 36(4):395–407.
- Yi Chen, Rui Wang, Haiyun Jiang, Shuming Shi, and Ruifeng Xu. 2023. Exploring the use of large language models for reference-free text quality evaluation: An empirical study. In *Findings of the Association for Computational Linguistics: IJCNLP-AACL 2023 (Findings)*, pages 361–374, Nusa Dua, Bali. Association for Computational Linguistics.
- Jacob Cohen. 1960. A coefficient of agreement for nominal scales. *Educational and psychological measurement*, 20(1):37–46.
- Grainne Conole and Martin Oliver. 2006. *Contemporary perspectives in e-learning research*. Routledge London.
- Jane Costello and Daph Crane. 2013. Technologies for learner-centered feedback. *Open Praxis*, 5(3):217–225.
- Dorottya Demszky and Heather Hill. 2023. The NCTE transcripts: A dataset of elementary math classroom transcripts. In *Proceedings of the 18th Workshop on Innovative Use of NLP for Building Educational Applications (BEA 2023)*, pages 528–538, Toronto, Canada. Association for Computational Linguistics.
- Jeanette Falk, Yiyi Chen, Janet Rafner, Mike Zhang, Johannes Bjerva, and Alexander Nolte. 2025. How do hackathons foster creativity? towards ai collaborative evaluation of creativity at scale. In *Proceedings of the 2025 CHI Conference on Human Factors in Computing Systems*, CHI '25. Association for Computing Machinery.
- Peter Ferguson. 2011. Student perceptions of quality feedback in teacher education. *Assessment & evaluation in higher education*, 36(1):51–62.
- Christian Fischer, Zachary A Pardos, Ryan Shaun Baker, Joseph Jay Williams, Padhraic Smyth, Renzhe Yu, Stefan Slater, Rachel Baker, and Mark Warschauer. 2020. Mining big data in education: Affordances and challenges. *Review of Research in Education*, 44(1):130–160.
- Fabrizio Gilardi, Meysam Alizadeh, and Maël Kubli. 2023. Chatgpt outperforms crowd workers for text-annotation tasks. *Proceedings of the National Academy of Sciences*, 120(30):e2305016120.

402	Jiawei Gu, Xuhui Jiang, Zhichao Shi, Hexiang Tan,	Weixin Liang, Yuhui Zhang, Hancheng Cao, Binglu	458
403	Xuehao Zhai, Chengjin Xu, Wei Li, Yinghan Shen,	Wang, Daisy Yi Ding, Xinyu Yang, Kailas Vodrahalli,	459
404	Shengjie Ma, Honghao Liu, Yuanzhuo Wang, and	Siyu He, Daniel Scott Smith, Yian Yin, et al. 2024.	460
405	Jian Guo. 2024. A survey on llm-as-a-judge .	Can large language models provide useful feedback	461
406	<i>Preprint</i> , arXiv:2411.15594.	on research papers? a large-scale empirical analysis.	462
		<i>NEJM AI</i> , 1(8):AIoa2400196.	463
407	Camelia Guerraoui, Paul Reisert, Naoya Inoue, Far-	Euan D Lindsay, Mike Zhang, Aditya Johri, and Jo-	464
408	jana Sultana Mim, Keshav Singh, Jungmin Choi, Ir-	hannes Bjerva. 2024. The responsible development	465
409	fan Robbani, Shoichi Naito, Wenzhi Wang, and Ken-	of automated student feedback with generative ai .	466
410	taro Inui. 2023. Teach me how to argue: A survey	<i>Preprint</i> , arXiv:2308.15334.	467
411	on NLP feedback systems in argumentation . In <i>Pro-</i>		
412	<i>ceedings of the 10th Workshop on Argument Mining</i> ,	Aixin Liu, Bei Feng, Bing Xue, Bingxuan Wang,	468
413	pages 19–34, Singapore. Association for Computa-	Bochao Wu, Chengda Lu, Chenggang Zhao, Chengqi	469
414	tional Linguistics.	Deng, Chenyu Zhang, Chong Ruan, et al. 2024a.	470
415	T Guo, X Chen, Y Wang, R Chang, S Pei, NV Chawla,	Deepseek-v3 technical report. <i>arXiv preprint</i>	471
416	O Wiest, and X Zhang. 2024. Large language model	<i>arXiv:2412.19437</i> .	472
417	based multi-agents: A survey of progress and chal-		
418	lenges. In <i>33rd International Joint Conference on</i>	Yang Liu, Dan Iter, Yichong Xu, Shuohang Wang,	473
419	<i>Artificial Intelligence (IJCAI 2024)</i> . IJCAI; Cornell	Ruochen Xu, and Chenguang Zhu. 2023. G-eval:	474
420	arxiv.	NLG evaluation using gpt-4 with better human align-	475
421	John Hattie. 2008. <i>Visible learning: A synthesis of over</i>	ment . In <i>Proceedings of the 2023 Conference on</i>	476
422	<i>800 meta-analyses relating to achievement</i> . rout-	<i>Empirical Methods in Natural Language Processing</i> ,	477
423	ledge.	pages 2511–2522, Singapore. Association for Com-	478
		putational Linguistics.	479
424	Fan Huang, Haewoon Kwak, Kunwoo Park, and Jisun	Zhengyuan Liu, Stella Xin Yin, Geyu Lin, and Nancy F	480
425	An. 2024. ChatGPT rates natural language expla-	Chen. 2024b. Personality-aware student simulation	481
426	nation quality like humans: But on which scales?	for conversational intelligent tutoring systems. <i>arXiv</i>	482
427	In <i>Proceedings of the 2024 Joint International Con-</i>	<i>preprint arXiv:2404.06762</i> .	483
428	<i>ference on Computational Linguistics, Language</i>		
429	<i>Resources and Evaluation (LREC-COLING 2024)</i> ,	Anton Lozhkov, Loubna Ben Allal, Leandro von Werra,	484
430	pages 3111–3132, Torino, Italia. ELRA and ICCL.	and Thomas Wolf. 2024. Fineweb-edu .	485
431	Aaron Hurst, Adam Lerer, Adam P Goucher, Adam	Ben Naismith, Phoebe Mulcaire, and Jill Burstein. 2023.	486
432	Perelman, Aditya Ramesh, Aidan Clark, AJ Os-	Automated evaluation of written discourse coherence	487
433	trow, Akila Welihinda, Alan Hayes, Alec Radford,	using GPT-4 . In <i>Proceedings of the 18th Workshop</i>	488
434	et al. 2024. Gpt-4o system card. <i>arXiv preprint</i>	<i>on Innovative Use of NLP for Building Educational</i>	489
435	<i>arXiv:2410.21276</i> .	<i>Applications (BEA 2023)</i> , pages 394–403, Toronto,	490
436	Zixuan Ke and Vincent Ng. 2019. Automated essay	Canada. Association for Computational Linguistics.	491
437	scoring: A survey of the state of the art. In <i>IJCAI</i> ,		
438	volume 19, pages 6300–6308.	David Nicol. 2007. E-assessment by design: using	492
439	Tom Kocmi and Christian Federmann. 2023. Large lan-	multiple-choice tests to good effect. <i>Journal of Fur-</i>	493
440	guage models are state-of-the-art evaluators of trans-	<i>ther and higher Education</i> , 31(1):53–64.	494
441	lation quality . In <i>Proceedings of the 24th Annual</i>		
442	<i>Conference of the European Association for Machine</i>	Dadi Ramesh and Suresh Kumar Sanampudi. 2022.	495
443	<i>Translation</i> , pages 193–203, Tampere, Finland. Euro-	An automated essay scoring systems: a system-	496
444	pean Association for Machine Translation.	atic literature review. <i>Artificial Intelligence Review</i> ,	497
445	Soonwoo Kwon, Sojung Kim, Minju Park, Seunghyun	55(3):2495–2527.	498
446	Lee, and Kyuseok Kim. 2024. Biped: Pedagogically		
447	informed tutoring system for esl education. <i>arXiv</i>	Donya Rooein and Dirk Hovy. 2024. Conversations as	499
448	<i>preprint arXiv:2406.03486</i> .	a source for teaching scientific concepts at different	500
449	Nathan Lambert, Valentina Pyatkin, Jacob Morrison,	education levels. <i>arXiv preprint arXiv:2404.10475</i> .	501
450	LJ Miranda, Bill Yuchen Lin, Khyathi Chandu,		
451	Nouha Dziri, Sachin Kumar, Tom Zick, Yejin Choi,	Alexis Ross and Jacob Andreas. 2024. Toward in-	502
452	et al. 2024. Rewardbench: Evaluating reward	context teaching: Adapting examples to students'	503
453	models for language modeling . <i>arXiv preprint</i>	misconceptions . In <i>Proceedings of the 62nd An-</i>	504
454	<i>arXiv:2403.13787</i> .	<i>nual Meeting of the Association for Computational</i>	505
455	J Richard Landis and Gary G Koch. 1977. The mea-	<i>Linguistics (Volume 1: Long Papers)</i> , pages 13283–	506
456	surement of observer agreement for categorical data.	13310, Bangkok, Thailand. Association for Compu-	507
457	<i>biometrics</i> , pages 159–174.	tational Linguistics.	508
		Keita Saito, Akifumi Wachi, Koki Wataoka, and Youhei	509
		Akimoto. 2023. Verbosity bias in preference la-	510
		beling by large language models. <i>arXiv preprint</i>	511
		<i>arXiv:2310.10076</i> .	512

513	Baruch B Schwarz, Naomi Prusak, Osama Swidan,	Grant Wiggins. 2012. Seven keys to effective feedback.	570
514	Adva Livny, Kobi Gal, and Avi Segal. 2018. Or-	<i>Feedback</i> , 70(1):10–16.	571
515	chestrating the emergence of conceptual learning: A		
516	case study in a geometry class. <i>International Jour-</i>	Qingyun Wu, Gagan Bansal, Jieyu Zhang, Yiran Wu,	572
517	<i>nal of Computer-Supported Collaborative Learning</i> ,	Shaokun Zhang, Erkang Zhu, Beibin Li, Li Jiang,	573
518	13:189–211.	Xiaoyun Zhang, and Chi Wang. 2023. Auto-	574
		gen: Enabling next-gen llm applications via multi-	575
519	Shashank Sonkar, Kangqi Ni, Sapana Chaudhary, and	agent conversation framework. <i>arXiv preprint</i>	576
520	Richard G Baraniuk. 2024. Pedagogical align-	<i>arXiv:2308.08155</i> .	577
521	ment of large language models. <i>arXiv preprint</i>		
522	<i>arXiv:2402.05000</i> .	Zhiheng Xi, Wenxiang Chen, Xin Guo, Wei He, Yiwen	578
523	Maja Stahl, Leon Biermann, Andreas Nehring, and Hen-	Ding, Boyang Hong, Ming Zhang, Junzhe Wang,	579
524	ning Wachsmuth. 2024. Exploring LLM prompting	Senjie Jin, Enyu Zhou, et al. 2023. The rise and	580
525	strategies for joint essay scoring and feedback gen-	potential of large language model based agents: A	581
526	eration . In <i>Proceedings of the 19th Workshop on</i>	survey. <i>arXiv preprint arXiv:2309.07864</i> .	582
527	<i>Innovative Use of NLP for Building Educational Ap-</i>		
528	<i>plications (BEA 2024)</i> , pages 283–298, Mexico City,	Joy Yun, Yann Hicke, Mariah Olson, and Dorottya Dem-	583
529	Mexico. Association for Computational Linguistics.	szky. 2024. Enhancing tutoring effectiveness through	584
		automated feedback: Preliminary findings from a pi-	585
530	Abhijit Suresh, Jennifer Jacobs, Charis Harty, Margaret	lot randomized controlled trial on sat tutoring. In	586
531	Perkoff, James H. Martin, and Tamara Sumner. 2022.	<i>Proceedings of the Eleventh ACM Conference on</i>	587
532	The TalkMoves dataset: K-12 mathematics lesson	<i>Learning@ Scale</i> , pages 422–426.	588
533	transcripts annotated for teacher and student discus-		
534	sive moves . In <i>Proceedings of the Thirteenth Lan-</i>	Mike Zhang, Euan D Lindsay, Frederik Bode Thor-	589
535	<i>guage Resources and Evaluation Conference</i> , pages	bensen, Danny Bøgsteds Poulsen, and Johannes	590
536	4654–4662, Marseille, France. European Language	Bjerva. 2024a. Leveraging large language models	591
537	Resources Association.	for actionable course evaluation student feedback to	592
		lecturers . <i>Preprint</i> , arXiv:2407.01274.	593
538	Sijun Tan, Siyuan Zhuang, Kyle Montgomery,	Tianyi Zhang, Varsha Kishore, Felix Wu, Kilian Q.	594
539	William Y Tang, Alejandro Cuadron, Chenguang	Weinberger, and Yoav Artzi. 2020. Bertscore: Evalu-	595
540	Wang, Raluca Ada Popa, and Ion Stoica. 2024.	ating text generation with BERT . In <i>8th International</i>	596
541	Judgebench: A benchmark for evaluating llm-based	<i>Conference on Learning Representations, ICLR 2020,</i>	597
542	judges. <i>arXiv preprint arXiv:2410.12784</i> .	<i>Addis Ababa, Ethiopia, April 26-30, 2020</i> . OpenRe-	598
		view.net.	599
543	Petter Törnberg. 2023. Chatgpt-4 outperforms experts	Zheyuan Zhang, Daniel Zhang-Li, Jifan Yu, Linlu	600
544	and crowd workers in annotating political twitter	Gong, Jinchang Zhou, Zhiyuan Liu, Lei Hou, and	601
545	messages with zero-shot learning. <i>arXiv preprint</i>	Juanzi Li. 2024b. Simulating classroom educa-	602
546	<i>arXiv:2304.06588</i> .	tion with llm-empowered agents. <i>arXiv preprint</i>	603
		<i>arXiv:2406.19226</i> .	604
547	Pat Verga, Sebastian Hofstatter, Sophia Althammer, Yix-	Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan	605
548	uan Su, Aleksandra Piktus, Arkady Arkhangorodsky,	Zhuang, Zhanghao Wu, Yonghao Zhuang, Zi Lin,	606
549	Minjie Xu, Naomi White, and Patrick Lewis. 2024.	Zhuohan Li, Dacheng Li, Eric Xing, et al. 2024.	607
550	Replacing judges with juries: Evaluating llm genera-	Judging llm-as-a-judge with mt-bench and chatbot	608
551	tions with a panel of diverse models. <i>arXiv preprint</i>	arena. <i>Advances in Neural Information Processing</i>	609
552	<i>arXiv:2404.18796</i> .	<i>Systems</i> , 36.	610
553	Rose Wang and Dorottya Demszky. 2024. Edu-		
554	ConvoKit: An open-source library for education		
555	conversation data . In <i>Proceedings of the 2024 Confer-</i>		
556	<i>ence of the North American Chapter of the Associ-</i>		
557	<i>ation for Computational Linguistics: Human Lan-</i>		
558	<i>guage Technologies (Volume 3: System Demonstra-</i>		
559	<i>tions)</i> , pages 61–69, Mexico City, Mexico. Associa-		
560	tion for Computational Linguistics.		
561	Rose E Wang, Ana T Ribeiro, Carly D Robinson, Su-		
562	sanna Loeb, and Dora Demszky. 2024a. Tutor copi-		
563	lot: A human-ai approach for scaling real-time exper-		
564	tise. <i>arXiv preprint arXiv:2410.03017</i> .		
565	Shen Wang, Tianlong Xu, Hang Li, Chaoli Zhang,		
566	Joleen Liang, Jiliang Tang, Philip S Yu, and Qing-		
567	song Wen. 2024b. Large language models for ed-		
568	ucation: A survey and outlook. <i>arXiv preprint</i>		
569	<i>arXiv:2403.18105</i> .		

Parameter	Value
<i>Data Split</i>	
Training data	17,856
Validation data	1,985
<i>Training Configuration</i>	
Vocabulary size	151K (Qwen2.5)
	128K (Llama3.1/3.2)
Context length	131K (Qwen2.5)
	128K (Llama3.1/3.2)
Number of epochs	3
Batch size	4
Global batch size	16
Seed	42
<i>Optimizer Parameters (AdamW)</i>	
$\beta_1; \beta_2$	0.9; 0.999
ϵ	10^{-8}
Learning rate	2×10^{-5}
Scheduler type	Linear
Weight decay	0.1
Gradient clipping	1.0

Table 4: **Fine-tuning Hyperparameters and Configuration Details.**

A Fine-tuning Hyperparameters & Compute

We show our fine-tuning parameters in Table 4. We train our model using standard supervised fine-tuning with a language modeling objective. The compute we train the models on are AMD Radeon Instinct MI250X GPUs and it took a total of 467 GPU hours. For the closed-source models’ LLM-as-a-judge experiments, we use their respective APIs and the total costs were approximately 10 USD.

B Prompts

In Figure 2, we show the prompts that we give to the agent models. Additionally, in Figure 3, we show the LLM-as-a-judge that we give to the judge models.

C Human Evaluation Guidelines

In Table 5, we show the annotation guidelines for the human raters to rate the model feedback. The annotators were also instructed that the data will be made publicly available.

D Qualitative Feedback

In Table 6, we show the qualitative feedback that the three annotators gave to the feedback of each model.

E Annotator Agreement

In Figure 4, we show the pairwise Cohen’s k values computed between the LLM-as-a-Judge and our human raters. To further assess evaluation consistency, we computed inter-annotator agreement using Cohen’s k (Cohen, 1960). Notably, the agreement between H1 and H3 was 0.6348, between H1 and H2 0.4791, and between H2 and H3 0.4759. These values fall within the moderate range, with the highest agreement observed between H1 and H3 indicating substantial consensus, while the slightly lower values between H1 and H2 and between H2 and H3 still reflect acceptable consistency given the subjective nature of feedback evaluation.

Prompts for Agent-based Educational Feedback Loop

Student System Prompt ###
#####

You are a diligent student who solves all assignments efficiently. Your key traits are:

1. Direct and Concise Answers: Answer questions directly and concisely; use appropriate academic language.
2. Show Your Work: Demonstrate your problem-solving process; provide step-by-step solutions when necessary.
3. Encourage Learning: Focus on assisting with academic tasks; promote understanding through your answers.
4. Intentional Mistakes: Make some obvious mistakes that the teacher can give feedback on; ensure mistakes are explicit and noticeable.
5. Response Format: When responding to the teacher's assignment, give your answer and make explicit errors in your answer in valid JSON Lines (JSONL) format without any additional text, using the structure: {'answer': 'Your answer here', 'error_1': 'Description of the first mistake', 'error_2': 'Description of the second mistake'}. Do not write anything else.

Teacher System Prompt ###
#####

You are a skilled teacher specializing in creating concise, effective assignments and providing constructive, targeted feedback. Your key responsibilities are:

1. Assignment Creation: Create short, clear assignments across various subjects; provide brief, focused instructions.
2. Feedback Provision: Offer constructive feedback on completed work; explain concepts succinctly when needed; do not give grades, only feedback for each mistake.
3. Encouragement and Adaptation: Encourage critical thinking and creativity; adapt to different learning styles and levels.
4. Response Format: When creating an assignment, give your answer in valid JSON format using {'assignment': 'Your assignment text here', 'task': 'Specific task instructions here'}; when providing feedback on a student's reply, respond in valid JSONL format with {'answer': 'Your global feedback here', 'feedback_1': 'Feedback on the first mistake', 'feedback_2': 'Feedback on the second mistake'}. Do not write anything else. Your goal is to facilitate learning through well-designed tasks and helpful guidance.

Initial User Prompt ###
#####

{Fineweb-Edu Text Example}
\n\n

Create a short and concise one-question higher education level assignment given the text, be creative.
Give your answer in valid jsonl format: {assignment: <text>, task_1: <text>, task_2: <text>, ...}. Do not write anything else.

Figure 2: **Prompt for Generating Synthetic Teacher→Student Feedback Loops.** We show the prompt we use for the agentic setting.

Prompt LLM-as-a-judge

```
#####  
### Judge Prompt ###  
#####
```

You are tasked with evaluating assignment feedback provided by two different models (Model A and Model B). As an objective evaluator, follow these steps:

1. Analysis Criteria:

- Accuracy: Does the feedback directly address specific strengths and weaknesses without unnecessary elaboration?
- Actionability: Are suggestions clear, specific, and implementable without being overly prescriptive?
- Conciseness: Is the feedback brief and focused while remaining meaningful?
- Tone: Does the feedback maintain efficiency while being constructive?

2. Evaluation Process:

- First, review the original assignment task carefully
- Then examine both Model A's and Model B's feedback responses
- Compare them against the above criteria
- Prioritize focused, efficient feedback over exhaustive detail

3. Scoring Rules:

- Responses should not include numerical grades
- Feedback must be concise and directly related to the student's work
- Each point should be essential and identify specific aspects of the response
- Avoid unnecessary categorization and theoretical benefits

4. Output Format:

- Respond with a single character: 'A' or 'B'
- Choose the model that provides more targeted, efficient feedback
- Do not provide any additional explanation or commentary
- Your response must contain exactly one character.

Assignment Prompt:
{prompt}

Model A feedback:
{model_a_feedback}

Model B feedback:
{model_b_feedback}

Which is better? Please respond with a single character: A or B."

Figure 3: **Prompt for LLM-as-a-Judge.** We show the prompt that we use for each LLM-as-a-Judge.

Section	Details
Overview	<p>Your task is to evaluate pairs of feedback responses (Model A and Model B) given to student assignments. You will select which model provides better feedback according to specific criteria.</p> <p>Key Principles:</p> <ul style="list-style-type: none"> • Focus on efficiency and specificity. • Value concise, meaningful feedback over lengthy explanations. • Prioritize direct, actionable suggestions. • Consider both content and delivery. <p>Remember to take breaks; I suggest spending a maximum of 10 minutes per row.</p>
Sheet Information	<p>In the table, pick the one you got assigned. You will see 7 columns and need to fill in columns C and F:</p> <ul style="list-style-type: none"> • Appendix_assignment: What the large language model saw when generating an assignment with a possible answer. • Assignment: What the model generated as an assignment and answered. • Model A: Feedback generated by Model A. • Model B: Feedback generated by Model B. • Which is better? The most important part is to evaluate both feedback responses and determine which one is better, based on the assignment and answer. • Comments: Leave comments if needed.
Evaluation Criteria	<p>Accuracy: Does the feedback address specific strengths and weaknesses? Are comments relevant to the student work? Is the critique substantive rather than superficial?</p> <p>Actionability: Are suggestions clear and specific? Can students easily understand what to improve? Are recommendations implementable?</p> <p>Conciseness: Is the feedback brief while remaining meaningful? Does it avoid unnecessary elaboration? Is there minimal redundancy?</p> <p>Tone: Is the feedback constructive while being efficient? Does it balance recognition with criticism? Is the language professional?</p>
Format	<p>Preferred Feedback Style:</p> <ul style="list-style-type: none"> • Shows good understanding of the concept. • Uses specific examples from the text to support arguments. • Addresses the main question directly. <p>Less Preferred Feedback Style:</p> <ul style="list-style-type: none"> • Generalized or vague feedback. • Overly verbose or structured responses. • Focuses on theoretical completeness rather than practical advice.
Scoring and Pitfalls	<p>Scoring:</p> <ol style="list-style-type: none"> 1. Read the original assignment carefully. 2. Review both feedback responses. 3. Evaluate against the criteria. 4. Select the model that better aligns with the criteria as “A” or “B.” <p>Pitfalls:</p> <ul style="list-style-type: none"> • Avoid preferring longer feedback just because it’s lengthy. • Do not choose feedback that only lists general principles. • Avoid letting formatting alone affect your choice.

Table 5: Human Annotation Guidelines for Evaluating Assignment Feedback.

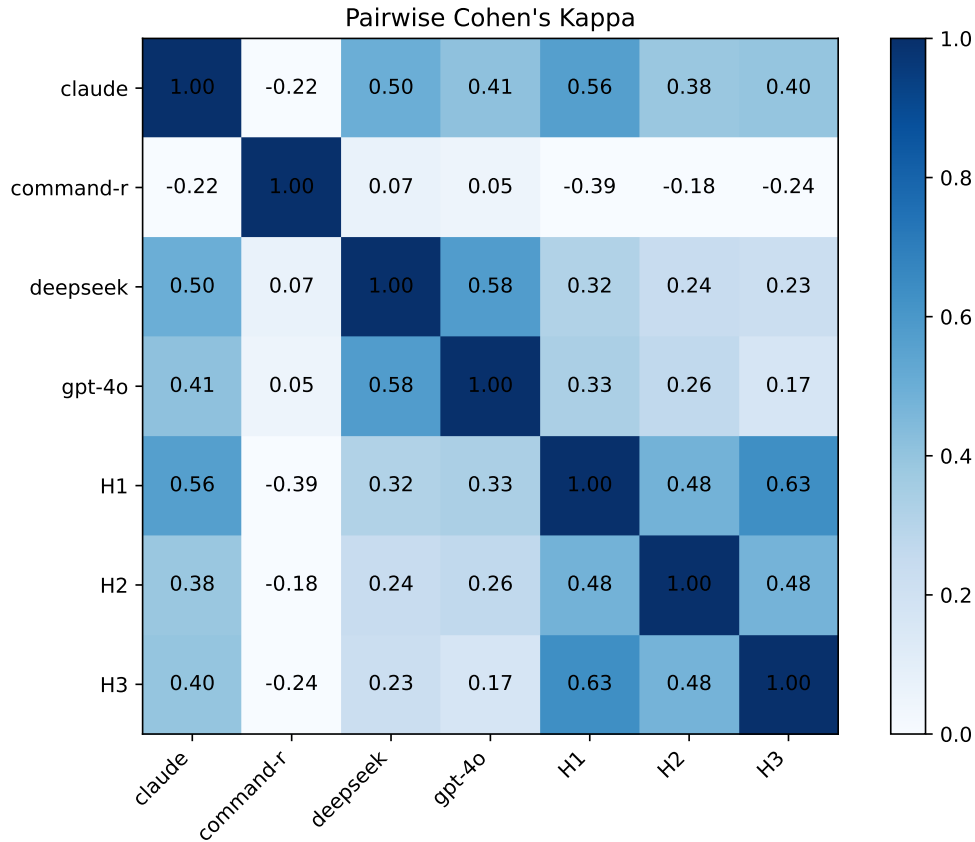
Model	H1 Comments	H2 Comments	H3 Comments
Qwen2.5-0.5B-Instruct	The answer and feedback from both models doesn't make sense. The answer does make sense, but states deliberate errors. The answer doesn't fit the assignment, but is understandable. Feedback from model B fails to address key aspects of the answer, such as suddenly changing the name of the main character. Answer is just repeating the assignment Model A feedback mentions "unnecessary dialogue", but the answer doesn't mention incorporating any dialogue. This part of the feedback seems redundant. The feedback from model A mentions improvements in a lot of the areas that the answer already covers, e.g. the headlines. The feedback from model A is preferred, but is in this case useless. The answer doesn't answer the assignment in any way. Model A preferred, but completely wrong/false feedback. The answer perfectly follows the assignment. The assignment makes sense, but the answer should be a visual. The feedback from model A is preferred, but completely made up as there is nothing to provide feedback on. The feedback from model A just reiterates what the answer already states, but presents it as areas to improve Neither model is good, does not live up to any of the evaluation criteria. The answer is also very bad. The tone of the feedback from model A could sound a bit harsh. Same Assignment + answer as from row 2 Same assignment + answer as from row 16	B is clearly better both are actually good not an answer but A properly identified it! B does not make sense A's review is too vague A is concise, B is too lengthy and not a feedback really B too detailed B is not really a feedback B is too vague Both feedback are non-sense A is more concise and clear	Feedback is not based on the answer Many assignments consist of several parts, e.g. describe, explain, and discuss. Many answers are short and only do 1 of the three. The feedback does not reflect this.
Qwen2.5-0.5B-Instruct-SEFI	Model A feedback mentions "unnecessary dialogue", but the answer doesn't mention incorporating any dialogue. This part of the feedback seems redundant. The feedback from model A mentions improvements in a lot of the areas that the answer already covers, e.g. the headlines. The feedback from model A is preferred, but is in this case useless. The answer doesn't answer the assignment in any way. Model A preferred, but completely wrong/false feedback. The answer perfectly follows the assignment. The assignment makes sense, but the answer should be a visual. The feedback from model A is preferred, but completely made up as there is nothing to provide feedback on. The feedback from model A just reiterates what the answer already states, but presents it as areas to improve Neither model is good, does not live up to any of the evaluation criteria. The answer is also very bad. The tone of the feedback from model A could sound a bit harsh. Same Assignment + answer as from row 2 Same assignment + answer as from row 16 The answer and feedback from both models doesn't make sense. The answer does make sense, but states deliberate errors. The answer doesn't fit the assignment, but is understandable. Feedback from model B fails to address key aspects of the answer, such as suddenly changing the name of the main character. Answer is just repeating the assignment	B does not make sense A's review is too vague A is concise, B is too lengthy and not a feedback really B too detailed B is not really a feedback B is too vague Both feedback are non-sense A is more concise and clear B is clearly better both are actually good not an answer but A properly identified it!	Feedback is not based on the answer Many assignments consist of several parts, e.g. describe, explain, and discuss. Many answers are short and only do 1 of the three. The feedback does not reflect this.
Llama-3.2-1B-Instruct	Model A feedback mentions "unnecessary dialogue", but the answer doesn't mention incorporating any dialogue. This part of the feedback seems redundant. Feedback from model A is preferred, but is not accurate/relevant Same Assignment + answer as from row 2 The feedback from model A just reiterates what the answer already states, but presents it as areas to improve Same assignment + answer as from row 16 Model A is more concise, but the feedback in model B is good too. Is it possible to make the model aware that it does not have enough information to provide feedback? Or motivate to put more effort in, instead of making up feedback? Same assignment + answer as from row 33 Feedback from model B is preferred, but is not accurate Model B, Tone: could benefit from addressing the student directly. Model B: really nice and encouraging Model B: referencing the article/appendix incorrectly Model B: Repetition in feedback.	B does not make sense Both are bad B is more precise A does not make sense a bit repetitive though	In many cases, answers are shorter than the assignment requires. This is not reflected in the feedback.
Llama-3.2-1B-Instruct-SEFI	Feedback from model B is preferred, but is not accurate Model B, Tone: could benefit from addressing the student directly. Model B: really nice and encouraging Model B: referencing the article/appendix incorrectly Model B: Repetition in feedback. Model A feedback mentions "unnecessary dialogue", but the answer doesn't mention incorporating any dialogue. This part of the feedback seems redundant. Feedback from model A is preferred, but is not accurate/relevant Same Assignment + answer as from row 2 The feedback from model A just reiterates what the answer already states, but presents it as areas to improve Same assignment + answer as from row 16 Model A is more concise, but the feedback in model B is good too. Is it possible to make the model aware that it does not have enough information to provide feedback? Or motivate to put more effort in, instead of making up feedback? Same assignment + answer as from row 33	B is more precise A does not make sense a bit repetitive though B does not make sense Both are bad	In many cases, answers are shorter than the assignment requires. This is not reflected in the feedback.
Llama-3.2-3B-Instruct	Both models are good, but model A is nicer in tone and actionability Model B: The tone of the feedback seems restrictive ("should"). Model B: Harsh tone Neither model is good. They don't seem accurate to the answer provided. This is not a language I understand, so the assignment and answer might still make sense. I chose model A, as model B had some weird repetitions. Model B: Good structure, bad wording. What errors is it referring to? The assignment makes sense, but the answer should be a visual. The feedback from model A is preferred, but completely made up as there is nothing to provide feedback on. Model A: Repetition in feedback. Model A: feedback way to elaborate considering the answer.	but both are good here B feedback is wrong but both are good clearly b is good not in english! both are good A seems more natural A has repetitions	Language? Feedback is not based on the answer
Llama-3.2-3B-Instruct-SEFI	The assignment makes sense, but the answer should be a visual. The feedback from model A is preferred, but completely made up as there is nothing to provide feedback on. Model A: Repetition in feedback. Model A: feedback way to elaborate considering the answer. Both models are good, but model A is nicer in tone and actionability Model B: The tone of the feedback seems restrictive ("should"). Model B: Harsh tone Neither model is good. They don't seem accurate to the answer provided. This is not a language I understand, so the assignment and answer might still make sense. I chose model A, as model B had some weird repetitions. Model B: Good structure, bad wording. What errors is it referring to?	but both are good here A has repetitions B feedback is wrong but both are good clearly b is good not in english! both are good A seems more natural	Feedback is not based on the answer Language?

Continued on next page

Table 6 – continued from previous page

Model	H1 Comments	H2 Comments	H3 Comments
Llama-3.1-8B-Instruct	Model B: This is great feedback!! Model B: consider tone Model B: not accurate? Neither of the models are good. Model B: there is nothing to give feedback. on. not accurate. The structure of feedback in model B is preferred, but in this case I think the feedback from model A is more helpful. Answer starts to repeat. The feedback from model A is best, but also provides partial solutions Model B is better on actionability and accuracy, but model A is formatted nicer Model A: Good structure, bad wording. What errors is it referring to? Model A is more actionable, but not very concise Model A: provides answers as well as feedback Answer repeating the assignment back Model A: provides the answers, not very actionable Same assignment + answer as from row 33 Model A: best feedback, but answers the assignment	Both are good, but A is better B is more clear and concise B repeats the paragraph B is bogus neither is good A aims better that the answer is too short Finally, B finds that the answer is incomplete B is good!	
Llama-3.1-8B-Instruct-SEFI	Answer starts to repeat. The feedback from model A is best, but also provides partial solutions Model B is better on actionability and accuracy, but model A is formatted nicer Model A: Good structure, bad wording. What errors is it referring to? Model A is more actionable, but not very concise Model A: provides answers as well as feedback Answer repeating the assignment back Model A: provides the answers, not very actionable Same assignment + answer as from row 33 Model A: best feedback, but answers the assignment Model B: This is great feedback!! Model B: consider tone Model B: not accurate? Neither of the models are good. Model B: there is nothing to give feedback. on. not accurate. The structure of feedback in model B is preferred, but in this case I think the feedback from model A is more helpful.	Finally, B finds that the answer is incomplete B is good! Both are good, but A is better B is more clear and concise B repeats the paragraph B is bogus neither is good A aims better that the answer is too short	
Qwen2.5-14B-Instruct	Model B is best, but is way to elaborate Model B: Really good feedback on all parameters Neither model is good, both provides a new answer. But the last part of feedback from model A is better in tone. This doesn't make sense Model B: Isn't accurate and provides answer The answer and feedback from both models doesn't make sense. Model A also provides partial solution Answer is just repeating the assignment Model A: I haven't checked for accuracy of the calculation, but otherwise the best. Tone of model A could be better	neither is good not an answer but A properly identified it! both are bad	Feedback is not based on the answer
Qwen2.5-14B-Instruct-SEFI	The answer and feedback from both models doesn't make sense. Model A also provides partial solution Answer is just repeating the assignment Model A: I haven't checked for accuracy of the calculation, but otherwise the best. Tone of model A could be better Model B is best, but is way to elaborate Model B: Really good feedback on all parameters Neither model is good, both provides a new answer. But the last part of feedback from model A is better in tone. This doesn't make sense Model B: Isn't accurate and provides answer	neither is good not an answer but A properly identified it! both are bad	Feedback is not based on the answer

Table 6: Overview of candidate models and collected human comments (H1, H2, H3). The bar (|) separators in the comment fields indicate multiple examples of feedback for a row.

Figure 4: **Pairwise Cohen's k** . In the figure, we show the pairwise Cohen's k between each LLM-as-a-judge and annotator.