

Additional Research Questions to Be Included in the Camera-Ready Version

ANONYMOUS AUTHOR(S)

1 RQ6: Open-Source Large Language Model

The evaluation in the previous sections suggests that ChatGPT-based scores have a stronger correlation with human scores than lexical similarity-based metrics. In this section, we investigate whether we can replace the closed-source ChatGPT with open-source LLMs while keeping a strong correlation with human scores. To this end, in this section, we replace ChatGPT with DeepSeek-Coder [2], a well-known open-source large language model. Notably, we reuse all of the prompts for ChatGPT and employ the default settings of DeepSeek-Coder.

To have a quantitative and objective assessment of the correlation between DeepSeek-Coder-based scores and human scores, we compute the Spearman Rank Correlation between them as shown in Table 1. The Spearman Rank Correlation Coefficient is 0.49, much larger than those of the lexical similarity-based scores (e.g., 0.21 of METEOR, 0.22 of BLEU, and 0.25 of ROUGE). It is even slightly greater than that (0.47) of ChatGPT-based scores. The results may suggest that ChatGPT could be safely replaced with open-source LLMs like DeepSeek-Coder while keeping the strong correlation between LLM-based scores and human scores. It also suggests that the proposed approach can work well with various LLMs.

To visualize the consistency between DeepSeek-Coder-based scores and human scores, we present a heatmap of their overlap in Fig. 1. From this figure, we observe that the data distribution is primarily concentrated along the diagonal, suggesting that DeepSeek-Coder-based scores are often aligned with human scores. We also use the Kolmogorov-Smirnov (K-S) statistic [1] to evaluate whether the distributions of DeepSeek-Coder-based scores significantly differ across reviews of different quality levels. Results are presented in Table 1. From this table, we observe that the K-S statistic of DeepSeek-Coder-based scores is very close to that of ChatGPT-based scores, suggesting that DeepSeek-Coder-based scores are comparable to ChatGPT-based scores in distinguishing high-quality reviews.

We conclude based on the preceding analysis that ChatGPT could be safely replaced with the state-of-the-art open-source LLMs concerning the proposed LLM-based scoring approach.

References

- [1] Giovanni Fasano and Alberto Franceschini. 1987. A multidimensional version of the Kolmogorov-Smirnov test. *Monthly Notices of the Royal Astronomical Society* 225, 1 (1987), 155–170.
- [2] Daya Guo, Qihao Zhu, Dejian Yang, Zhenda Xie, Kai Dong, Wentao Zhang, Guanting Chen, Xiao Bi, Yu Wu, YK Li, et al. 2024. DeepSeek-Coder: When the Large Language Model Meets Programming-The Rise of Code Intelligence. *arXiv preprint arXiv:2401.14196* (2024).

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than the author(s) must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.

Conference acronym 'XX, June 03–05, 2018, Woodstock, NY

© 2018 Copyright held by the owner/author(s). Publication rights licensed to ACM.

ACM ISBN 978-1-4503-XXXX-X/18/06

<https://doi.org/XXXXXXX.XXXXXXX>

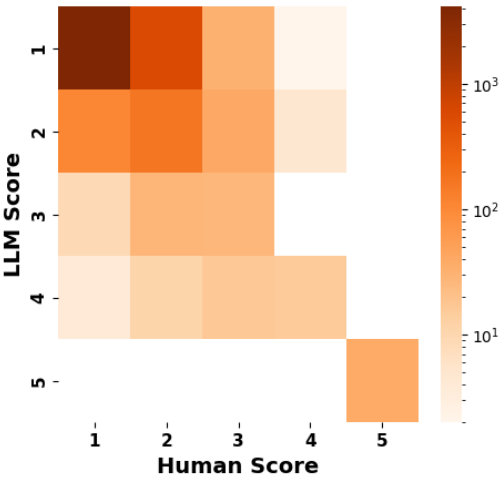


Fig. 1. Heatmap of Overlap Between LLM(deepSeek) Scores and Human Scores

Table 1. Kolmogorov-Smirnov Statistic Across Reviews of Varying Quality Levels

Quality Level	2-Point	3-Point	4-Point	5-Point
1-Point	BLEU: 0.26	BLEU:0.36	BLEU: 0.44	BLEU: 1.00
	Embedding: 0.55	Embedding: 0.71	Embedding: 0.73	Embedding: 1.00
	ChatGPT: 0.55	ChatGPT: 0.71	ChatGPT: 0.73	ChatGPT: 1.00
	DeepSeek: 0.53	DeepSeek: 0.73	DeepSeek: 0.79	DeepSeek: 1.00
2-Point		BLEU:0.16	BLEU: 0.30	BLEU: 0.99
		Embedding: 0.26	Embedding: 0.35	Embedding: 1.00
		ChatGPT: 0.24	ChatGPT: 0.47	ChatGPT: 1.00
		DeepSeek: 0.28	DeepSeek: 0.55	DeepSeek: 1.00
3-Point			BLEU: 0.20	BLEU: 1.00
			Embedding: 0.22	Embedding: 1.00
			ChatGPT: 0.28	LLM: 1.00
			DeepSeek: 0.33	DeepSeek: 1.00
4-Point				BLEU: 0.97
				Embedding: 1.00
				ChatGPT: 1.00
				DeepSeek: 1.00