

DO CAPTIONING METRICS REFLECT MUSIC SEMANTIC ALIGNMENT?

Jinwoo Lee¹

Kyogu Lee^{1,2,3}

¹ Department of Intelligence and Information, Seoul National University

² Interdisciplinary Program in Artificial Intelligence, Seoul National University

³ Artificial Intelligence Institute, Seoul National University

{e--jw, kglee}@snu.ac.kr

ABSTRACT

Music captioning has emerged as a promising task, fueled by the advent of advanced language generation models. However, the evaluation of music captioning relies heavily on traditional metrics such as BLEU, METEOR, and ROUGE which were developed for other domains, without proper justification for their use in this new field. We present cases where traditional metrics are vulnerable to syntactic changes, and show they do not correlate well with human judgments. By addressing these issues, we aim to emphasize the need for a critical reevaluation of how music captions are assessed.

1. INTRODUCTION

The advancement of music information retrieval (MIR) parallels significant developments in the music industry, particularly regarding tasks like music captioning [1–5]. Despite the promise of music captioning, the evaluation of generated captions poses significant challenges. Current metrics, borrowed from natural language processing tasks, do not adequately address the unique qualities of musical content. These traditional metrics, such as BLEU [6], METEOR [7], and ROUGE [8], primarily rely on n-gram overlap, focusing on superficial similarities between generated and reference captions [3]. This approach inherently favors syntactic similarity over semantic meaning, which may result in misleading evaluations of caption quality. For example, a paraphrased caption may receive a substantially lower score, even though it conveys the same core idea as the original.

Moreover, these metrics fail to effectively capture the semantics of musical contexts, such as genre, instruments, and mood. This can result in evaluations that may overlook critical aspects of the music being described, leading to poor alignment with how humans perceive and interpret musical content. This work aims to illustrate the limitations of these traditional metrics, and highlight instances

where they score inaccurately high.

2. METHOD

To evaluate the effectiveness of traditional metrics in music captioning, we conduct a human evaluation study aimed at correlating human judgments with the scores generated by these metrics.

We use Amazon Mechanical Turk [9] to recruit 50 participants for a listening test. From the evaluation set of the MusicCaps [10] dataset, we randomly sample 30 audio clips. Each audio sample is paired with three different types of captions:

- Original: The ground truth caption associated with the audio in the dataset.
- Inference: The caption generated by the audio captioning model. For the inference, we use music captioning model by LP-MusicCaps [2].
- Paraphrased: A version of the original caption syntactically altered without changing the meaning.

Participants are asked to evaluate how accurately the captions describe the audio clips. To maintain the integrity of the evaluation and minimize potential bias, each participant is presented with only two of the three caption versions for each audio sample. This design helps prevent participants from recognizing that two of three captions are likely ground truth captions, while the other one is not.

3. RESULTS

3.1 Syntactic & Semantic variations

Following the collection of human evaluation scores, we report the scores of evaluation metrics (including the Mean Opinion Score (MOS)) for original and paraphrased captions, as shown in Figure 1. Notably, evaluation metrics except for FENSE show a significant decrease when comparing the original captions to their paraphrased versions. This indicates that these metrics may be overly sensitive to syntactic changes rather than semantic content. As illustrated in Table 1, we also showcase examples of various caption types along with their corresponding evaluation scores. We also include distorted captions that are

arXiv:2411.11692v1 [cs.LG] 18 Nov 2024

Late-Breaking / Demo Session Extended Abstract, ISMIR 2024 Conference



Table 1: Examples of different caption types and their evaluation scores. We report BLEU (B) as their average.

Caption	Human	N-gram Overlap				Embedding
	MOS	B	M	R	S	FENSE
<i>Original</i> Someone is playing a melody on an e-guitar with a tremolo effect . This song may be playing at home practicing guitar .	4.60	1.00	1.00	1.00	1.00	1.00
<i>Distorted (Semantic X, Syntactic O)</i> Someone is playing a melody on a french horn with a very reverb . This song may be playing at home practicing french horn .	-	0.61	0.38	0.73	0.33	0.56
<i>Paraphrased (Semantic O, Syntactic X)</i> An electric guitar with tremolo effect plays a melody, possibly during home practice.	4.70	0.13	0.25	0.28	0.36	0.83
<i>Original</i> This is a yodeling music piece. There is a female vocalist that is singing happily in the lead. The melody is provided by medium and high pitch woodwinds . In the background, the bass line is played by an upright bass while the rhythm is provided by an acoustic drum. The atmosphere is very lively . This piece could be used in the soundtrack of a comedy movie or a children’s show .	4.64	1.00	1.00	1.00	1.00	1.00
<i>Distorted (Semantic X, Syntactic O)</i> This is a hard rock piece. There is a male vocalist that is singing in energetic mood. The melody is provided by low-pitch electric guitars . Meanwhile, the melody is played by a high-pitched violin while the rhythm is provided by electronic drums. This piece could be used in the soundtrack of a urgent action movie .	-	0.43	0.26	0.59	0.32	0.66
<i>Paraphrased (Semantic O, Syntactic X)</i> A lively yodeling song with a female vocalist, woodwinds, upright bass , and acoustic drums, suitable for a comedy or children’s show .	4.32	0.04	0.15	0.29	0.30	0.72

Table 2: Experimental results of correlation coefficient (Pearson’s r) between evaluation metrics vs. human evaluation (MOS) for 20 entries in the evaluation set of MusicCaps dataset.

BLEU ₁	BLEU ₂	BLEU ₃	BLEU ₄	METEOR	ROUGE	SPICE	FENSE
0.075	0.074	0.077	0.085	0.096	0.083	0.097	0.091

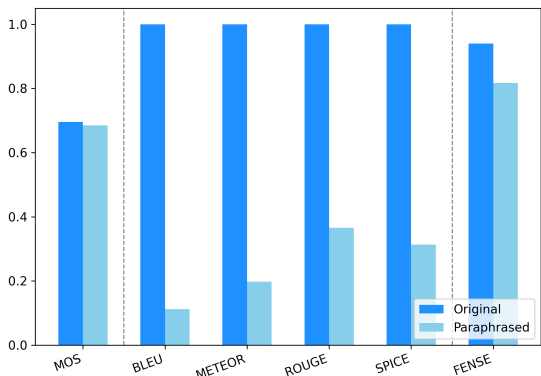


Figure 1: Experimental results on scores across caption type and evaluation metrics from MusicCaps evaluation set. MOS values are rescaled to [0, 1].

semantically altered while preserving their syntactic structure. Notably, most n-gram-based metrics tend to favor

distorted captions over paraphrased captions.

3.2 Correlation with human judgment

To further validate our findings, we compute the correlation coefficient between the scores obtained from human evaluations and each of the evaluation metrics. As presented in Table 2, results indicate that most evaluation metrics, including FENSE [11], do not show significant correlations with human evaluations.

4. CONCLUSION

We demonstrate existing metrics are overly sensitive to syntactic variations, and they lack alignment with actual human evaluations. Given these findings, we conclude that a more nuanced evaluation framework is necessary to address these challenges.

5. REFERENCES

- [1] I. Manco, E. Benetos, E. Quinton, and G. Fazekas, “Muscaps: Generating captions for music audio,” in *2021 International Joint Conference on Neural Networks (IJCNN)*. IEEE, 2021, pp. 1–8.
- [2] S. Doh, K. Choi, J. Lee, and J. Nam, “Lp-musiccaps: Llm-based pseudo music captioning,” *arXiv preprint arXiv:2307.16372*, 2023.
- [3] J. P. Gardner, S. Durand, D. Stoller, and R. M. Bittner, “Lark: A multimodal instruction-following language model for music,” in *Forty-first International Conference on Machine Learning*, 2023.
- [4] Y. Li, R. Yuan, G. Zhang, Y. Ma, X. Chen, H. Yin, C. Lin, A. Ragni, E. Benetos, N. Gyenge *et al.*, “Mert: Acoustic music understanding model with large-scale self-supervised training,” *arXiv preprint arXiv:2306.00107*, 2023.
- [5] S. Doh, M. Lee, D. Jeong, and J. Nam, “Enriching music descriptions with a finetuned-llm and metadata for text-to-music retrieval,” in *ICASSP 2024 - 2024 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, 2024, pp. 826–830.
- [6] K. Papineni, S. Roukos, T. Ward, and W.-J. Zhu, “Bleu: a method for automatic evaluation of machine translation,” in *Proceedings of the 40th annual meeting of the Association for Computational Linguistics*, 2002, pp. 311–318.
- [7] S. Banerjee and A. Lavie, “Meteor: An automatic metric for mt evaluation with improved correlation with human judgments,” in *Proceedings of the acl workshop on intrinsic and extrinsic evaluation measures for machine translation and/or summarization*, 2005, pp. 65–72.
- [8] C.-Y. Lin, “Rouge: A package for automatic evaluation of summaries,” in *Text summarization branches out*, 2004, pp. 74–81.
- [9] “Amazon mechanical turk,” <https://www.mturk.com/>, accessed: 2024-10-05.
- [10] A. Agostinelli, T. I. Denk, Z. Borsos, J. Engel, M. Verzetti, A. Caillon, Q. Huang, A. Jansen, A. Roberts, M. Tagliasacchi *et al.*, “Musiclm: Generating music from text,” *arXiv preprint arXiv:2301.11325*, 2023.
- [11] Z. Zhou, Z. Zhang, X. Xu, Z. Xie, M. Wu, and K. Q. Zhu, “Can audio captions be evaluated with image caption metrics?” in *ICASSP 2022-2022 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2022, pp. 981–985.