

Eval4NLP 2023

**The 4th Workshop on Evaluation and Comparison of NLP
Systems**

Proceedings of the Workshop

November 1, 2023

The Eval4NLP organizers gratefully acknowledge the support from the following sponsors.

Sponsors

The Federal Ministry of Education and Research (BMBF) via the grant “Metrics4NLG”



©2023 The Asian Federation of Natural Language Processing and The Association for Computational Linguistics

Order copies of this and other ACL proceedings from:

Association for Computational Linguistics (ACL)
209 N. Eighth Street
Stroudsburg, PA 18360
USA
Tel: +1-570-476-8006
Fax: +1-570-476-0860
acl@aclweb.org

ISBN 979-8-89176-021-9

Introduction

Welcome to the Fourth Workshop on Evaluation and Comparison of NLP Systems (Eval4NLP 2023).

The current year has brought astonishing achievements in NLP. Generative large language models (LLMs) like ChatGPT and GPT4 demonstrate wide capabilities in understanding and performing tasks from in-context descriptions without fine-tuning, bringing world-wide attention to the risks and opportunities that arise from current and ongoing research. Further, the release of open-source models like LLaMA and Falcon LLM, better quantization techniques for inference and training, as well as the adaptation of efficient fine-tuning techniques such as LORA accelerate the research progress by allowing hardware and runtime efficiency. Given the ever growing speed of research, fair evaluations and comparisons are of fundamental importance to the NLP community in order to properly track progress. This concerns the creation of benchmark datasets that cover typical use cases and blind spots of existing systems, the designing of metrics for evaluating the performance of NLP systems on different dimensions, and the reporting of evaluation results in an unbiased manner.

We believe that new insights and methodology, particularly in the last 2-3 years, have led to much renewed interest in the workshop topic. The first workshop in the series, Eval4NLP'20, was the first workshop to take a broad and unifying perspective on the subject matter. The second (Eval4NLP'21) and third (Eval4NLP'22) workshop extended this perspective. We believe the fourth workshop continues the tradition of being a reputed platform for presenting and discussing latest advances in NLP evaluation methods and resources.

This year we especially encouraged the submission of works that consider the evaluation of LLMs and their generated content as well as works that leverage LLMs in their evaluation strategies. In fact, to encourage research in this direction, we ran a successful shared task this year on prompting LLMs as explainable metrics. Participants were given a set of open-source LLMs and were tasked with designing prompts and score retrieval strategies for automatically scoring machine translation and automatic text summarization outputs without using a reference text.

Our workshop and shared task attracted a lot of attention from the research community. Among the 15 submissions, 9 were accepted for presentation after thorough consideration by the program committee. In addition, there were 9 teams that participated in the shared task. This year's program covers a wide range of topics, including creating a benchmark dataset for identifying and quantifying sexism in language models, evaluation metrics for named entity recognition, probing techniques for large language models, and much more.

We would like to thank all of the authors for their contributions, the program committee for their thoughtful reviews, the keynote speaker for sharing their perspective, and all the attendees for their participation. We believe that all of these will contribute to a lively and successful workshop. Looking forward to meeting you all at Eval4NLP 2023!

Eval4NLP 2023 Organizing Committee,
Daniel Deutsch, Rotem Dror, Steffen Eger, Yang Gao, Christoph Leiter, Juri Opitz, Andreas Rücklé

Organizing Committee

Organizing Committee

Daniel Deutsch, Google
Rotem Dror, University of Haifa
Steffen Eger, Bielefeld University
Yang Gao, Google Research
Christoph Leiter, Bielefeld University
Juri Opitz, Heidelberg University
Andreas Rücklé, Amazon

Program Committee

Reviewers

Omri Abend
Jonas Belouadi
Yanran Chen
Daniel Deutsch
Li Dong
Zi-Yi Dou
Rotem Dror
Steffen Eger
George Foster
Anette Frank
Markus Freitag
Yang Gao
Claire Gardent
Juraj Juraska
Ji-Ung Lee
Christoph Leiter
Lucy Lin
Juri Opitz
Ines Rehbein
Ehud Reiter
Leonardo Ribeiro
Ori Shapira
Julius Steen
Benyou Wang
Ran Zhang
Shiyue Zhang
Wei Zhao

Table of Contents

| | |
|--|-----|
| <i>WRF: Weighted Rouge-F1 Metric for Entity Recognition</i> Lukas Jonathan Weber, Krishnan Jothi Ramalingam, Matthias Beyer and Axel Zimmermann . . . | 1 |
| <i>Assessing Distractors in Multiple-Choice Tests</i> Vatsal Raina, Adian Liusie and Mark Gales | 12 |
| <i>Delving into Evaluation Metrics for Generation: A Thorough Assessment of How Metrics Generalize to Rephrasing Across Languages</i> Yixuan Wang, Qingyan Chen and Duygu Ataman | 23 |
| <i>EduQuick: A Dataset Toward Evaluating Summarization of Informal Educational Content for Social Media</i> Zahra Kolagar, Sebastian Steindl and Alessandra Zarcone | 32 |
| <i>Zero-shot Probing of Pretrained Language Models for Geography Knowledge</i> Nitin Ramrakhiani, Vasudeva Varma, Girish Keshav Palshikar and Sachin Pawar | 49 |
| <i>Transformers Go for the LOLs: Generating (Humourous) Titles from Scientific Abstracts End-to-End</i> Yanran Chen and Steffen Eger | 62 |
| <i>Summary Cycles: Exploring the Impact of Prompt Engineering on Large Language Models' Interaction with Interaction Log Information</i> Jeremy E Block, Yu-Peng Chen, Abhilash Budharapu, Lisa Anthony and Bonnie J Dorr | 85 |
| <i>Large Language Models As Annotators: A Preliminary Evaluation For Annotating Low-Resource Language Content</i> Savita Bhat and Vasudeva Varma | 100 |
| <i>Can a Prediction's Rank Offer a More Accurate Quantification of Bias? A Case Study Measuring Sexism in Debaised Language Models</i> Jad Doughman, Shady Shehata, Leen Al Qadi, Youssef Nafea and Fakhri Karray | 108 |
| <i>The Eval4NLP 2023 Shared Task on Prompting Large Language Models as Explainable Metrics</i> Christoph Leiter, Juri Opitz, Daniel Deutsch, Yang Gao, Rotem Dror and Steffen Eger | 117 |
| <i>HIT-MI&T Lab's Submission to Eval4NLP 2023 Shared Task</i> Rui Zhang, Fuhai Song, Hui Huang, Jinghao Yuan, Muyun Yang and Tiejun Zhao | 139 |
| <i>Understanding Large Language Model Based Metrics for Text Summarization</i> Abhishek Pradhan and Ketan Kumar Todi | 149 |
| <i>LTRC_IITB's 2023 Submission for Prompting Large Language Models as Explainable Metrics Task</i> Pavan Baswani, Ananya Mukherjee and Manish Shrivastava | 156 |
| <i>Which is better? Exploring Prompting Strategy For LLM-based Metrics</i> JoongHoon Kim, Sangmin Lee, Seung Hun Han, Saeran Park, Jiyeon Lee, Kiyoon Jeong and Pilsung Kang | 164 |
| <i>Characterised LLMs Affect its Evaluation of Summary and Translation</i> Yuan Lu and Yu-Ting Lin | 184 |
| <i>Reference-Free Summarization Evaluation with Large Language Models</i> Abbas Akkasi, Kathleen C. Fraser and Majid Komeili | 193 |

| | |
|---|-----|
| <i>Little Giants: Exploring the Potential of Small LLMs as Evaluation Metrics in Summarization in the Eval4NLP 2023 Shared Task</i> | |
| Neema Kotonya, Saran Krishnasamy, Joel R. Tetreault and Alejandro Jaimes | 202 |
| <i>Exploring Prompting Large Language Models as Explainable Metrics</i> | |
| Ghazaleh Mahmoudi | 219 |
| <i>Team NLLG submission for Eval4NLP 2023 Shared Task: Retrieval-Augmented In-Context Learning for NLG Evaluation</i> | |
| Daniil Larionov, Vasiliy Viskov, George Kokush, Alexander Panchenko and Steffen Eger . . . | 228 |