

LREC-COLING 2024

**The Fourth Workshop on
Human Evaluation of NLP Systems
(HumEval 2024)**

Workshop Proceedings

Editors

Simone Balloccu, Anya Belz, Rudali Huidrom, Ehud Reiter,
João Sedoc and Craig Thomson

21 May, 2024
Torino, Italia

**Proceedings of the Fourth Workshop on Human Evaluation of NLP Systems
(HumEval 2024)**

Copyright ELRA Language Resources Association (ELRA), 2024
These proceedings are licensed under a Creative Commons
Attribution-NonCommercial 4.0 International License (CC BY-NC 4.0)

ISBN 978-2-493814-41-8
ISSN 2951-2093 (COLING); 2522-2686 (LREC)

Jointly organized by the ELRA Language Resources Association
and the International Committee on Computational Linguistics

Preface

Welcome to HumEval 2024!

We are pleased to present the proceedings of the fourth workshop on Human Evaluation of NLP Systems (HumEval) which is taking place as part of the 2024 Joint International Conference on Computational Linguistics, Language Resources and Evaluation (LREC-COLING 2024).

Human evaluation is vital in NLP, and it is often considered as the most reliable form of evaluation. It ranges from large-scale crowd-sourced evaluations to the much smaller experiments routinely encountered in conference papers. With this workshop we wish to create a forum for current human evaluation research, a space for researchers working with human evaluations to exchange ideas and begin to address the issues that human evaluation in NLP currently faces, including aspects of experimental design, reporting standards, meta-evaluation and reproducibility.

We are truly grateful to the authors of the submitted papers that showed interest in human evaluation research. The HumEval workshop accepted 8 submissions. The accepted papers cover a broad range of NLP areas where human evaluation is used: machine translation, natural language generation, text simplification, conversational search. Several papers are addressing reproducibility of human evaluations. The workshop once again hosted the results session of the ReprONLP Shared Task on Reproducibility of Evaluations in NLP which consisted of the presentation of overall results by the organizers and 18 oral and poster presentations by participants.

This workshop would not have been possible without the hard work of the program committee. We would like to express our gratitude to them for writing detailed and thoughtful reviews in a very constrained span of time. We also thank our invited speakers, Mark Diaz and Sheila Castilho, for their contribution to our program. We are grateful for the help from the LREC-COLING 2024 workshop organizers, and to all the people involved in setting up the infrastructure.

You can find more details about the workshop on its website: <https://humeval.github.io/>.

Simone, Craig, João, Anya, Ehud, Rudali

Organizing Committee

Organizers

Simone Balloccu, ADAPT Centre, Dublin City University, Ireland
Anyá Belz, ADAPT Centre, Dublin City University, Ireland
Rudali Huidrom, ADAPT Centre, Dublin City University, Ireland
Ehud Reiter, University of Aberdeen, UK
João Sedoc, New York University, US
Craig Thomson, University of Aberdeen, UK

Program Committee

Gavin Abercrombie, Heriot-Watt University, UK
Jose Maria Alonso-Moral, University of Santiago de Compostela, ES
Mohammad Arvan, University of Illinois, Chicago, US
Anouck Braggaar, Tilburg University, NL
Daniel Braun, University of Twente, NL
Javier Corbelle, University of Santiago de Compostela, ES
Tanvi Dinkar, Heriot-Watt University, UK
Ondrej Dusek, Karls University, CZ
Steffen Eger, University of Mannheim, DE
Manuela Hürlimann, Zurich University of Applied Sciences, CH
Mateusz Lango, Poznań University of Technology, PL
Yiru Li, Groningen University, NL
Michela Lorandi, Dublin City University / ADAPT, IE
Saad Mahamood, trivago, DE
Zola Mahlaza, University of Cape Town, ZA
Gonzalo Mendez, University of Madrid, ES
Margot Mieskes, University of Applied Sciences, Darmstadt, DE
Jie Ruan, Peking University, CN
Patricia Schmidtova, Karls University, CZ
Raj Shah, Georgia Tech, US
Barkavi Sundararajan, University of Aberdeen, UK
Supryadi, Tianjin University, CN
Chris van der Lee, Tilburg University, NL Xiaojun Wan, Peking University, CN
Deyi Xiong, Tianjin University, CN
Emiel van Miltenburg, Tilburg University, NL
Chuang Liu, Tianjin University, CN

Invited Speakers

Mark Diaz, Google Research
Sheila Castilho, ADAPT/DCU

Table of Contents

<i>Quality and Quantity of Machine Translation References for Automatic Metrics</i> Vilém Zouhar and Ondřej Bojar	1
<i>Exploratory Study on the Impact of English Bias of Generative Large Language Models in Dutch and French</i> Ayla Rigouts Terryn and Miryam de Lhoneux	12
<i>Adding Argumentation into Human Evaluation of Long Document Abstractive Summarization: A Case Study on Legal Opinions</i> Mohamed Elaraby, Huihui Xu, Morgan Gray, Kevin Ashley and Diane Litman	28
<i>A Gold Standard with Silver Linings: Scaling Up Annotation for Distinguishing Bosnian, Croatian, Montenegrin and Serbian</i> Aleksandra Miletic and Filip Miletic	36
<i>Insights of a Usability Study for KBQA Interactive Semantic Parsing: Generation Yields Benefits over Templates but External Validity Remains Challenging</i> Ashley Lewis, Lingbo Mo, Marie-Catherine de Marneffe, Huan Sun and Michael White	47
<i>Extrinsic evaluation of question generation methods with user journey logs</i> Elie Antoine, Eléonore Besnehard, Frederic Bechet, Geraldine Damnati, Eric Kergosien and Arnaud Laborderie	63
<i>Towards Holistic Human Evaluation of Automatic Text Simplification</i> Luisa Carrer, Andreas Säuberli, Martin Kappus and Sarah Ebling	71
<i>Decoding the Metrics Maze: Navigating the Landscape of Conversational Question Answering System Evaluation in Procedural Tasks</i> Alexander Frummet and David Elswiler	81
<i>The 2024 ReprONLP Shared Task on Reproducibility of Evaluations in NLP: Overview and Results</i> Anya Belz and Craig Thomson	91
<i>Once Upon a Replication: It is Humans' Turn to Evaluate AI's Understanding of Children's Stories for QA Generation</i> Andra-Maria Florescu, Marius Micluta-Campeanu and Liviu P. Dinu	106
<i>Exploring Reproducibility of Human-Labelled Data for Code-Mixed Sentiment Analysis</i> Sachin Sasidharan Nair, Tanvi Dinkar and Gavin Abercrombie	114
<i>Reproducing the Metric-Based Evaluation of a Set of Controllable Text Generation Techniques</i> Michela Lorandi and Anya Belz	125
<i>ReproHum: #0033-03: How Reproducible Are Fluency Ratings of Generated Text? A Reproduction of August et al. 2022</i> Emiel van Miltenburg, Anouck Braggaa, Nadine Braun, Martijn Goudbeek, Emiel Krahmer, Chris van der Lee, Steffen Pauws and Frédéric Tomas	132
<i>ReproHum #0927-03: DExpert Evaluation? Reproducing Human Judgements of the Fluency of Generated Text</i> Tanvi Dinkar, Gavin Abercrombie and Verena Rieser	145

<i>ReproHum #0927-3: Reproducing The Human Evaluation Of The DExperts Controlled Text Generation Method</i>	
Javier González Corbelle, Ainhoa Vivel Couso, Jose Maria Alonso-Moral and Alberto Bugarín-Diz	153
<i>ReproHum #1018-09: Reproducing Human Evaluations of Redundancy Errors in Data-To-Text Systems</i>	
Filip Klubička and John D. Kelleher	163
<i>ReproHum#0043: Human Evaluation Reproducing Language Model as an Annotator: Exploring Dialogue Summarization on AMI Dataset</i>	
Vivian Fresen, Mei-Shin Wu-Urbanek and Steffen Eger	199
<i>ReproHum #0712-01: Human Evaluation Reproduction Report for “Hierarchical Sketch Induction for Paraphrase Generation”</i>	
Mohammad Arvan and Natalie Parde	210
<i>ReproHum #0712-01: Reproducing Human Evaluation of Meaning Preservation in Paraphrase Generation</i>	
Lewis N. Watson and Dimitra Gkatzia	221
<i>ReproHum #0043-4: Evaluating Summarization Models: investigating the impact of education and language proficiency on reproducibility</i>	
Mateusz Lango, Patricia Schmidtova, Simone Balloccu and Ondrej Dusek	229
<i>ReproHum #0033-3: Comparable Relative Results with Lower Absolute Values in a Reproduction Study</i>	
Yiru Li, Huiyuan Lai, Antonio Toral and Malvina Nissim	238
<i>ReproHum #0124-03: Reproducing Human Evaluations of end-to-end approaches for Referring Expression Generation</i>	
Saad Mahamood	250
<i>ReproHum #0087-01: Human Evaluation Reproduction Report for Generating Fact Checking Explanations</i>	
Tyler Loakman and Chenghua Lin	255
<i>ReproHum #0892-01: The painful route to consistent results: A reproduction study of human evaluation in NLG</i>	
Irene Mondella, Huiyuan Lai and Malvina Nissim	261
<i>ReproHum #0087-01: A Reproduction Study of the Human Evaluation of the Coverage of Fact Checking Explanations</i>	
Mingqi Gao, Jie Ruan and Xiaojun Wan	269
<i>ReproHum #0866-04: Another Evaluation of Readers’ Reactions to News Headlines</i>	
Zola Mahlaza, Toky Hajatiana Raboanary, Kyle Seakgwa and C. Maria Keet	274

Conference Program

May 21, 2024

9:00–9:10 **Opening Remarks**

9:10–10:30 **Oral Session 1**

Quality and Quantity of Machine Translation References for Automatic Metrics
Vilém Zouhar and Ondřej Bojar

Exploratory Study on the Impact of English Bias of Generative Large Language Models in Dutch and French
Ayla Rigouts Terryn and Miryam de Lhoneux

Adding Argumentation into Human Evaluation of Long Document Abstractive Summarization: A Case Study on Legal Opinions
Mohamed Elaraby, Huihui Xu, Morgan Gray, Kevin Ashley and Diane Litman

A Gold Standard with Silver Linings: Scaling Up Annotation for Distinguishing Bosnian, Croatian, Montenegrin and Serbian
Aleksandra Miletić and Filip Miletić

10:30–11:00 **Coffee Break**

11:00–11:45 **Invited Talk 1**

Beyond Performance: The Evolving Landscape of Human Evaluation
Sheila Castilho

11:45–13:00 **ReproNLP Shared Task Session 1**

13:00–14:00 **Lunch**

May 21, 2024 (continued)

14:00–14:45 Oral Session 2

Insights of a Usability Study for KBQA Interactive Semantic Parsing: Generation Yields Benefits over Templates but External Validity Remains Challenging

Ashley Lewis, Lingbo Mo, Marie-Catherine de Marneffe, Huan Sun and Michael White

Extrinsic evaluation of question generation methods with user journey logs

Elie Antoine, Eléonore Besnehard, Frederic Bechet, Geraldine Damnati, Eric Kergosien and Arnaud Laborderie

Towards Holistic Human Evaluation of Automatic Text Simplification

Luisa Carrer, Andreas Säuberli, Martin Kappus and Sarah Ebling

14:45–16:00 ReprONLP Shared Task Session 2

16:00–16:30 Coffee Break

16:30–17:15 Invited Talk 2

All That Agrees Is Not Gold: Evaluating Ground Truth and Conversational Safety

Mark Diaz

17:15–18:00 Oral Session 3

Decoding the Metrics Maze: Navigating the Landscape of Conversational Question Answering System Evaluation in Procedural Tasks

Alexander Frummet and David Elswailer

18:00–18:05 Closing Remarks