HumEval 2023

# Proceedings of the 3rd Workshop on Human Evaluation of NLP Systems

*associated with*
**The 14th International Conference on**
**Recent Advances in Natural Language Processing'2023**

7 September 2023
Varna, Bulgaria

3RD WORKSHOP ON HUMAN EVALUATION OF NLP SYSTEMS
ASSOCIATED WITH THE 14TH INTERNATIONAL CONFERENCE
RECENT ADVANCES IN
NATURAL LANGUAGE PROCESSING'2023

**PROCEEDINGS**

7 September 2023
Varna, Bulgaria

# Preface

Welcome to HumEval 2023!

We are pleased to present the third workshop on Human Evaluation of NLP Systems (HumEval) which is taking place as part of the Conference on Recent Advances in Natural Language Processing (RANLP 2023).

Human evaluation is vital in NLP, and it is often considered as the most reliable form of evaluation. It ranges from the 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 15 submissions. The accepted papers cover a broad range of NLP areas where human evaluation is used: machine translation, natural language generation, summarisation, text-to-speech. Several papers are addressing reproducibility of human evaluations.

This workshop would not have been possible without the hard work of the programme 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 speaker, Elizabeth Clark, for her contribution to our program. We are grateful for the help from the RANLP organisers, especially Galia Angelova and Ivelina Nikolova, and we are grateful to all the people involved in setting up the infrastructure.

You can find more details about the worskhop on its website: `https://humeval.github.io/`.

Anya, Ehud, Craig, Maja, Joao, Simone, Rudali

# Table of Contents

# Conference Programme

**No Day Set (continued)**