EACL 2021

# Human Evaluation of NLP Systems (HumEval)

## Proceedings of the Workshop

April 19, 2021
Online

# Introduction

Welcome to HumEval 2021!

We are pleased to present the first workshop on Human Evaluation of NLP Systems (HumEval) that is taking place virtually as part of the 16th Conference of the European Chapter of the Association for Computational Linguistics (EACL 2021).

Human evaluation plays an important role in NLP, 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.

The HumEval workshop accepted 9 submissions as long papers, and 6 as short papers. The accepted papers cover a broad range of NLP areas where human evaluation is used: natural language generation, machine translation, summarisation, dialogue, and word embeddings. There are also papers dealing with evaluation practices and methodology in NLP.

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, Lucia Specia, and Margaret Mitchell, for their contribution to our program. As the workshop is part of EACL, we appreciated help from the EACL Workshop Chairs, Jonathan Berant, and Angeliki Lazaridou, from the EACL Publication Chairs, Valerio Basile, and Tommaso Caselli, and we are grateful to all the people involved in setting up the virtual infrastructure.

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

Anya, Shubham, Yvette, Ehud, Anastasia

# Invited Speaker: Lucia Specia, Imperial College London

## Disagreement in Human Evaluation: Blame the Task not the Annotators

**Abstract**: It is well known that human evaluators are prone to disagreement and that this is a problem for reliability and reproducibility of evaluation experiments. The reasons for disagreement can fall into two broad categories: (1) human evaluator, including under-trained, under-incentivised, lacking expertise, or ill-intended individuals, e.g., cheaters; and (2) task, including ill-definition, poor guidelines, suboptimal setup, or inherent subjectivity. While in an ideal evaluation experiment many of these elements will be controlled for, I argue that task subjectivity is a much harder issue. In this talk I will cover a number of evaluation experiments on tasks with variable degrees of subjectivity, discuss their levels of disagreement along with other issues, and cover a few practical approaches do address them. I hope this will lead to an open discussion on possible strategies and directions to alleviate this problem.

# Invited Speaker: Margaret Mitchell

## The Ins and Outs of Ethics-Informed Evaluation

**Abstract**: The modern train/test paradigm in Artificial Intelligence (AI) and Machine Learning (ML) narrows what we can understand about AI models, and skews our understanding of models' robustness in different environments. In this talk, I will work through the different factors involved in ethics-informed AI evaluation, including connections to ML training and ML fairness, and present an overarching evaluation protocol that addresses a multitude of considerations in developing ethical AI.

# Table of Contents

# Workshop Program

**Monday, April 19, 2021**

9:00–9:10      *Opening*
               Anya Belz

**9:10–10:00   *Invited Talk: Lucia Specia***

**10:00–11:00  Oral Session 1: NLG**

10:00–10:20    *It's Commonsense, isn't it? Demystifying Human Evaluations in Commonsense-Enhanced NLG Systems*
               Miruna-Adriana Clinciu, Dimitra Gkatzia and Saad Mahamood

10:20–10:40    *Estimating Subjective Crowd-Evaluations as an Additional Objective to Improve Natural Language Generation*
               Jakob Nyberg, Maike Paetzel and Ramesh Manuvinakurike

10:40–11:00    *Trading Off Diversity and Quality in Natural Language Generation*
               Hugh Zhang, Daniel Duckworth, Daphne Ippolito and Arvind Neelakantan

**11:00–11:30  *Break***

**11:30–12:10  Oral Session 2: MT**

11:30–11:50    *Towards Document-Level Human MT Evaluation: On the Issues of Annotator Agreement, Effort and Misevaluation*
               Sheila Castilho

11:50–12:10    *Is This Translation Error Critical?: Classification-Based Human and Automatic Machine Translation Evaluation Focusing on Critical Errors*
               Katsuhito Sudoh, Kosuke Takahashi and Satoshi Nakamura