2022

# Challenges & Perspectives in Creating Large Language Models

## Proceedings of the Workshop

May 27, 2022

The organizers gratefully acknowledge the support from the following sponsors.

**Sponsor**

# Introduction

Two years after the appearance of GPT-3, large language models seem to have taken over NLP. Their capabilities, limitations, societal impact and the potential new applications they unlocked have been discussed and debated at length. A handful of replication studies have been published since then, confirming some of the initial findings and discovering new limitations. This workshop aims to gather researchers and practitioners involved in the creation of these models in order to:

1. Share ideas on the next directions of research in this field, including—but not limited to—grounding, multi-modal models, continuous updates and reasoning capabilities.

2. Share best-practices, brainstorm solutions to identified limitations and discuss challenges, such as infrastructure, data, ethical & legal frameworks, evaluation, training efficiency, etc.

This workshop is organized by the BigScience[1] initiative and will also serve as the closing session of this one year-long initiative aimed at developing a multilingual large language model, which is gathering 1.000+ researchers from more than 60 countries and 250 institutions and research labs. Its goal is to investigate the creation of a large scale dataset and model from a very wide diversity of angles.

---

[1] https://bigscience.huggingface.co/

# Organizing Committee

**Organization Committee**

Angela Fan, Meta AI
Matthias Gallé, Naver Labs Europe
Suzana Ilić, HuggingFace
Thomas Wolf, HuggingFace

**Steering Committee**

Yoav Goldberg, Bar Ilan University & Allen Institute for Artificial Intelligence
Percy Lang, Stanford University
Margaret Mitchell, HuggingFace & Ethical AI LLC
Alice Oh, KAIST
Alexander Rush, Cornell University

# Program Committee

**Program Chairs**

Angela Fan, Facebook
Matthias Gallé, Naver Labs Europee
Suzana Ilic, HuggingFace
Thomas Wolf, HuggingFace

# Table of Contents

# Program

**Friday, May 27, 2022**

12:30 - 11:00     *Poster Session*

14:00 - 15:00     *BigScience*

15:00 - 15:20     *Data Governance*

15:20 - 15:40     *Data*

15:40 - 16:00     *Modeling*

16:00 - 16:20     *Prompt Engineering*

16:20 - 16:40     *Evaluation*