CoNLL 2023

# The BabyLM Challenge at the 27th Conference on Computational Natural Language Learning

# Proceedings of the BabyLM Challenge

December 6-7, 2023

Order copies of this and other ACL proceedings from:

# Introduction

Greetings "babies"! and welcome to the proceedings and session of the 2023 BabyLM Challenge, held on December 6, 2023 as part of CoNLL (co-hosted with EMNLP) in Singapore. This challenge aims to bring together researchers interested in developmentally plausible pre-training, sample efficiency, and human language acquisition. Our challenge encourages researchers to "think small" by using training corpora containing 100 million words—approximately the amount of data available to human language learners, but far less data than is typically used for pre-training language models.

We received 31 papers, all of which were accepted on the basis of scientific and technical validity, rather than model performance. We received 162 individual model submissions, the scores of which are hosted online, at `www.https://dynabench.org/babylm`.

We are grateful to the participants for advancing our understanding of how best to train language models on scaled-down and more developmentally plausible corpora.. Their contributions have provided insight into important questions related to cognitive modeling, computational psycholinguistics, and sample-efficient language modeling. We are also grateful to the program committee for their thoughtful reviews of the submissions we received this year. Likewise, we are thankful to the CoNLL organizers for their work in integrating the BabyLM challenge into their program.

– The BabyLM Organization Committee: Alex Warstadt, Aaron Mueller, Leshem Choshen, Ethan Wilcox, Chengxu Zhuang, Adina Williams, Tal Linzen, Ryan Cotterell.

# Organizing Committee

**Organizers**

Alex Warstadt, ETH Zürich, Switzerland
Aaron Mueller, Northeastern University
Leshem Choshen, MIT, IBM
Ethan Wilcox, ETH Zürich, Switzerland
Chengxu Zhuang, Massachusetts Institute of Technology
Juan Ciro, MLCommons
Rafael Mosquera, MLCommons
Bhargavi Paranjabe, University of Washington
Adina Williams, Meta AI (FAIR)
Tal Linzen, New York University
Ryan Cotterell, ETH Zürich, Switzerland

# Table of Contents