

Yong Zheng-Xin

CONTACT INFORMATION	personal website: yongzx.github.io email: contact.yong@brown.edu	
EDUCATION	Brown University , Providence, RI Ph.D. Student, Computer Science Advisor: Prof. Stephen H. Bach Minerva University , San Francisco, CA B.Sc., Computer Science (Major) and Business (Minor) Advisor: Prof. Patrick D Watson Major GPA: 4.0/4, Cumulative GPA: 4.0/4	07/2021 - Present 09/2017 - 05/2021
RELEVANT EMPLOYMENT	Brown University , Ph.D. Student <i>Supervisor: Stephen Bach</i> Meta AI , Research Scientist Intern <i>Fundamental AI Research (FAIR)</i> • Mitigating accent bias in Massively Multilingual Speech models. • <i>Supervisors: Jean Maillard & Michael Auli</i> <i>GenAI Safety Alignment Team</i> • Multilingual safety research collaboration. • <i>Supervisor: Jianfeng Chi</i> Cohere For AI , Research Collaborator <i>Aya Responsible Release</i> • Safety red-teaming for multilingual LLM Aya-101. • Deliverable: Aya Model (co-first author) • <i>Supervisors: Julia Kreutzer & Sara Hooker</i> BigScience , Research Collaborator/Lead <i>Multilingual Modeling Group</i> • Led language adaptation research of BLOOM to low-resource languages. • Helped with data collection for the earliest instruction-following models, namely To (English) and BLOOMZ (multilingual) . • Deliverables: BLOOM+1 (Research Lead) , To , PromptSource , BLOOMZ , BLOOM • <i>Supervisor: Vassilina Nikoulina</i> FrameNet Project , Research Intern <i>Google Summer of Code 2019 and 2020</i> • Expanded FrameNet through semi-supervised learning and anomaly detection. • Investigated cross-lingual alignment of semantic frames in FrameNet graph. • Deliverables: SDEC-AD , Frame Shift Prediction • <i>Supervisors: Tiago T. Torrent, Oliver Czulo & Collin F. Baker</i>	07/2021 - Present 06/2024 - 12/2024 07/2024 - 11/2024 05/2023 - 02/2024 08/2021 - 12/2022 06/2019 - 12/2020
AWARDS	Best Paper Award at ACL 2024 Best Paper Award at NeurIPS 2023 Socially Responsible Language Modeling (SoLaR) Workshop 3rd Best App Overall at FirstNet Public Safety Hackathon Best Use of ESRI at FirstNet Public Safety Hackathon	2024 2023 2018 2018

Grand Prize Winner at #100Hacks Hackathon for Puerto Rico	2017
Grand Prize Winner at NCSV Innovate NUS Hackathon	2017
Open Category Grand Prize at SIA Airlines AppChallenge	2017
Global Finalist at Google Science Fair	2016
Outstanding A-Level Cambridge Learner Award for Perfect Score in Maths	2016

FEATURED
PUBLICATIONS
(* INDICATES
CO-FIRST
AUTHORSHIP)

Multilingual AI Safety

- [1] Preference Tuning For Toxicity Mitigation Generalizes Across Languages
Xiaochen Li*, **Zheng-Xin Yong***, Stephen H. Bach
EMNLP 2024 Findings.
- [2] Aya Model: An Instruction Finetuned Open-Access Multilingual Language Model
Ahmet Üstün*, Viraat Aryabumi*, **Zheng-Xin Yong***, Wei-Yin Ko*, Daniel D'souza*,
Gbemileke Onilude, Neel Bhandari, Shivalika Singh, Hui-Lee Ooi, Amr Kayid, Fred-
die Vargus, Phil Blunsom, Shayne Longpre, Niklas Muennighoff, Marzieh Fadaee,
Julia Kreutzer, Sara Hooker
ACL 2024. (Role in project: Multilingual safety red-teaming)
Best Paper Award. Work also featured in *The New York Times* and other media.
- [3] Low-Resource Languages Jailbreak GPT-4
Zheng-Xin Yong, Cristina Menghini, Stephen H. Bach
NeurIPS 2023 Socially Responsible Language Modeling Workshop
Best Paper Award. Work also featured in *New Scientist* and other media.

Low-Resource NLP and Synthetic Data Generation

- [4] LexC-Gen: Generating Data for Extremely Low-Resource Languages with Large
Language Models and Bilingual Lexicons
Zheng-Xin Yong, Cristina Menghini, Stephen H. Bach
EMNLP 2024 Findings
- [5] Prompting Multilingual Large Language Models to Generate Code-Mixed Texts:
The Case of South East Asian Languages
Zheng-Xin Yong, Ruochen Zhang, Jessica Zosa Forde, Skyler Wang, Arjun Subramo-
nian, Holy Lovenia, Samuel Cahyawijaya, Genta Indra Winata, Lintang Sutawika,
Jan Christian Blaise Cruz, Yin Lin Tan, Long Phan, Rowena Garcia, Thamar Solorio,
Alham Fikri Aji
EMNLP 2023 CALCS Workshop
Work also featured on *WIRED*.
- [6] BLOOM+1: Adding Language Support to BLOOM for Zero-Shot Prompting
Zheng-Xin Yong, Hailey Schoelkopf, Niklas Muennighoff, Alham Fikri Aji, David
Ifeoluwa Adelani, Khalid Almubarak, M Saiful Bari, Lintang Sutawika, Jungo Ka-
sai, Ahmed Baruwa, Genta Indra Winata, Stella Biderman, Edward Raff, Dragomir
Radev, Vassilina Nikoulina
ACL 2023
- [7] CVQA: Culturally-diverse Multilingual Visual Question Answering Benchmark
David Romero, ..., **Zheng-Xin Yong**, ..., Thamar Solorio, Alham Fikri Aji (75 authors)
NeurIPS 2024 Datasets and Benchmarks
- [8] SEACrowd: A Multilingual Multimodal Data Hub and Benchmark Suite for South-
east Asian Languages
Holy Lovenia, ..., **Zheng-Xin Yong**, Samuel Cahyawijaya (61 authors)
EMNLP 2024

OTHER
PUBLICATIONS

- [9] A Safe Harbor for AI Evaluation and Red Teaming
Shayne Longpre, ..., **Zheng-Xin Yong**, ..., Percy Liang, Peter Henderson (23 authors)
ICML 2024 (Position Paper)
Oral, 1.6% of Submitted Papers (160/9653)
- [10] Representativeness as a Forgotten Lesson for Multilingual and Code-switched Data Collection and Preparation
A. Seza Doğruöz, Sunayana Sitaram, **Zheng-Xin Yong**
EMNLP Findings 2023
- [11] Crosslingual Generalization through Multitask Finetuning
Niklas Muennighoff, ..., **Zheng-Xin Yong**, ..., Edward Raff, Colin Raffel (19 authors)
ACL 2023
- [12] PromptSource: An Integrated Development Environment and Repository for Natural Language Prompts
Stephen H. Bach*, Victor Sanh*, **Zheng-Xin Yong**, ..., Alexander M. Rush (27 authors)
ACL Demo 2023
- [13] The Decades Progress on Code-Switching Research in NLP: A Systematic Survey on Trends and Challenges
Genta Indra Winata, Alham Fikri Aji, **Zheng-Xin Yong**, Tamar Solorio
ACL Findings 2023
- [14] BLOOM: A 176B-Parameter Open-Access Multilingual Language Model
BigScience Workshop (including **Zheng-Xin Yong** and 300+ authors)
Preprint 2023
- [15] Multitask Prompted Training Enables Zero-Shot Task Generalization
Victor Sanh*, ..., **Zheng-Xin Yong**, ..., Alexander M. Rush (41 authors)
ICLR 2022
Spotlight, 5% of submitted papers (176/3391)

SELECTED TALKS	<p>The Alan Turing Institute, London Data Week: <i>Multilingual AI Safety</i>, Jul 2024.</p> <p>Cohere For AI, C4AI Gatherings: <i>Aya Multilingual Safety</i>, Feb 2024.</p> <p>Cohere For AI, Aya Grande Finale: <i>Aya Responsible Release</i>, Feb 2024.</p> <p>Cohere For AI, Closing the Contribution Chapter: <i>Malay Ambassador</i>, Dec 2023.</p>
SELECTED PRESS & MEDIA	<p>GPT-4 gave advice on planning terrorist attacks when asked in Zulu (New Scientist, 2023)</p> <p>ChatGPT Is Cutting Non-English Languages Out of the AI Revolution (WIRED, 2023)</p>
SERVICES	<p>Area Chair: EMNLP 2023 (Multilingualism and Linguistic Diversity)</p> <p>Conference/Workshop Reviewer: COLM 2024 ARR 2021-Present ACL 2021-2023</p> <p>Outreach and Mentorship Service Programs: Deep Learning Indaba Mentorship Service (2022 - 2023) Brown Computer Science exploreCSR (2022) Minerva-Masason Mentoring Program (2019 - 2021)</p>