

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 07/2021 - Present | |
| | 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 09/2017 - 05/2021 | |
| RELEVANT EMPLOYMENT | Brown University , Ph.D. Student Supervisor: <i>Stephen Bach</i> 07/2021 - Present | |
| | Meta AI , Research Scientist Intern <i>Fundamental AI Research (FAIR)</i> 06/2024 - Present <ul style="list-style-type: none">Multimodal safety research for multilingual speech models.Supervisor: <i>Jean Maillard</i> | |
| | Cohere For AI , Research Collaborator <i>Aya Responsible Release</i> 05/2023 - 02/2024 <ul style="list-style-type: none">Multilingual safety red-teaming for Aya-101 LLMs.Deliverable: Aya Model (co-first author)Supervisors: <i>Sara Hooker & Julia Kreutzer</i> | |
| | BigScience , Research Collaborator/Lead <i>Multilingual Modeling Group</i> 08/2021 - 12/2022 <ul style="list-style-type: none">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, BLOOMSupervisor: <i>Vassilina Nikoulina</i> | |
| | FrameNet Project , Research Intern <i>Google Summer of Code 2019 and 2020</i> 06/2019 - 12/2020 <ul style="list-style-type: none">Expanded FrameNet through semi-supervised learning and anomaly detection.Investigated cross-lingual alignment of semantic frames in FrameNet graph.Deliverables: SDEC-AD, Frame Shift PredictionSupervisors: <i>Collin F. Baker, Tiago T. Torrent, Oliver Czulo</i> | |
| AWARDS | Best Paper Award at NeurIPS 2023 Socially Responsible Language Modeling Workshop 12/2023 | |
| | 3rd Best App Overall at FirstNet Public Safety Hackathon 03/2018 | |
| | Best Use of ESRI at FirstNet Public Safety Hackathon 03/2018 | |
| | Grand Prize Winner at #100Hacks Hackathon for Puerto Rico 12/2017 | |
| | Grand Prize Winner at NCSV Innovate NUS Hackathon 11/2017 | |
| | Open Category Grand Prize at SIA Airlines AppChallenge 10/2017 | |
| | Global Finalist at Google Science Fair 08/2016 | |

FEATURED
PUBLICATIONS
(* INDICATES
CO-AUTHORSHIP)

Responsible AI and AI Safety

- [1] Preference Tuning For Toxicity Mitigation Generalizes Across Languages
Xiaochen Li*, **Zheng-Xin Yong***, Stephen H. Bach
In submission.
- [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)
- [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 Multilingual 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
In submission.
- [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.

OTHER
PUBLICATIONS

- [7] CVQA: Culturally-diverse Multilingual Visual Question Answering Benchmark
David Romero, ..., **Zheng-Xin Yong**, ..., Thamar Solorio, Alham Fikri Aji (75 authors)
In submission
- [8] SEACrowd: A Multilingual Multimodal Data Hub and Benchmark Suite for South-
east Asian Languages
Holy Lovenia, ..., **Zheng-Xin Yong**, Samuel Cahyawijaya (61 authors)
In submission
- [9] A Safe Harbor for AI Evaluation and Red Teaming
Shayne Longpre, ..., **Zheng-Xin Yong**, ..., Percy Liang, Peter Henderson (23 authors)
ICML 2024 Oral (Position Paper).
- [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**, Thamar 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.

SELECTED TALKS **The Alan Turing Institute, London Data Week:** *Multilingual AI Safety*, Jul 2024.
Cohere For AI, C4AI Gatherings: *Aya Multilingual Safety*, Feb 2024.
Cohere For AI, Aya Grande Finale: *Aya Responsible Release*, Feb 2024.
Cohere For AI, Closing the Contribution Chapter: *Malay Ambassador*, Dec 2023.

SELECTED PRESS & MEDIA **GPT-4 gave advice on planning terrorist attacks when asked in Zulu** (New Scientist, 2023)
ChatGPT Is Cutting Non-English Languages Out of the AI Revolution (WIRED, 2023)

SERVICES **Area Chair:**
EMNLP 2023 (Multilingualism and Linguistic Diversity)

Conference Reviewer:
COLM 2024
ARR 2021-Present
ACL 2021-2023

Outreach and Mentorship Service Programs:
Deep Learning Indaba Mentorship Service (2022 - 2023)
Brown Computer Science exploreCSR (2022)
Minerva-Masason Mentoring Program (2019 - 2021)

MENTORING **Jacob Xiaochen Li** (Sep 2023 - Present). Brown University.
Edward Ajayi (Jun 2023 - Nov 2023, *Deep Learning Indaba*). Now pursuing M.Sc. at CMU-Africa with fully-funded MasterCard Foundation Scholarship.
Imam Hasan Araf (Jan 2022 - May 2022, *Brown Computer Science exploreCSR*), Rowan University.
Shutaro Aoyama (Aug 2019 - Jan 2021, *Minerva-Masason Mentoring Program*), Gunma Kokusai Academy (High School). Now undergrad at Columbia University.