

## Yong Zheng-Xin

---

CONTACT INFORMATION	personal website: <a href="https://yongzx.github.io">yongzx.github.io</a> email: <a href="mailto:contact.yong@brown.edu">contact.yong@brown.edu</a>	
EDUCATION	<b>Brown University</b> , Providence, RI Ph.D. Student, Computer Science Advisor: Prof. Stephen H. Bach 07/2021 - Present	
	<b>Minerva University</b> , 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	
WORK EXPERIENCE	<b>Meta</b> — Fundamental AI Research (FAIR) <i>Research Scientist Intern</i> 06/2024 - Present <ul style="list-style-type: none"><li>Measuring intrinsic toxicity in speech models.</li><li>Minimizing accent biases in Massively Multilingual Speech models.</li><li>Mentors: Jean Maillard, Michael Auli &amp; Marta Costa-Jussà</li></ul>	
	<b>Meta</b> — GenAI Safety Alignment <i>Research Collaborator</i> 07/2024 - 10/2024 <ul style="list-style-type: none"><li>Safety research on multilingual LLMs' jailbreaks.</li><li>Deliverable: <b>Mechanistic explanations for cross-lingual finetuning attacks</b></li><li>Mentor: Jianfeng Chi</li></ul>	
	<b>Cohere For AI</b> — Aya Responsible Release Team <i>Research Collaborator</i> 05/2023 - 02/2024 <ul style="list-style-type: none"><li>Safety red-teaming for multilingual LLM Aya-101.</li><li>Deliverable: <b>Aya Model (co-first author for safety)</b></li><li>Mentors: Julia Kreutzer &amp; Sara Hooker</li></ul>	
	<b>BigScience</b> — Multilingual Modeling Group <i>Research Lead/Collaborator</i> 08/2021 - 12/2022 <ul style="list-style-type: none"><li>Led language adaptation research of BLOOM to low-resource languages.</li><li>Deliverables: <b>BLOOM+1 (Research Lead)</b>, <b>PromptSource</b>, <b>To</b>, <b>BLOOMZ</b>, <b>BLOOM</b></li><li>Mentor: Vassilina Nikoulina</li></ul>	
	<b>Google Summer of Code</b> — FrameNet Project <i>Research Intern</i> 06/2019 - 12/2020 <ul style="list-style-type: none"><li>Expanded FrameNet and investigated cross-lingual semantic frame alignment.</li><li>Deliverables: <b>SDEC-AD</b>, <b>Frame Shift Prediction</b></li><li>Mentors: Tiago T. Torrent, Oliver Czulo &amp; Collin F. Baker</li></ul>	
SELECTED PUBLICATIONS (* INDICATES CO-FIRST AUTHORSHIP)	[16] Towards Understanding the Fragility of Multilingual LLMs against Fine-Tuning Attacks Samuele Poppi, <b>Zheng-Xin Yong</b> , Yifei He, Bobbie Chern, Han Zhao, Aobo Yang, Jianfeng Chi <i>Preprint</i> .	
	[15] Preference Tuning For Toxicity Mitigation Generalizes Across Languages Xiaochen Li*, <b>Zheng-Xin Yong*</b> , Stephen H. Bach	

*EMNLP 2024 Findings.*

- [14] 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.*

**Best Paper Award.** Featured on: [The New York Times](#) | [The Washington Post](#).

- [13] Low-Resource Languages Jailbreak GPT-4  
**Zheng-Xin Yong**, Cristina Menghini, Stephen Bach  
*NeurIPS 2023 Socially Responsible Language Modeling Workshop*  
**Best Paper Award.** Featured on: [New Scientist](#) | [The Register](#).

OTHER  
PUBLICATIONS

- [12] 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*
- [11] LexC-Gen: Generating Data for Extremely Low-Resource Languages with Large  
Language Models and Bilingual Lexicons  
**Zheng-Xin Yong**, Cristina Menghini, Stephen Bach  
*EMNLP 2024 Findings*
- [10] 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*
- [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)**
- [8] 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*
- [7] Prompting Multilingual Large Language Models to Generate Code-Mixed Texts:  
The Case of South East Asian Languages  
**Zheng-Xin Yong**, Ruochen Zhang, ..., Alham Fikri Aji (15 authors)  
*EMNLP 2023 CALCS Workshop*  
**Featured on: WIRED.**
- [6] BLOOM+1: Adding Language Support to BLOOM for Zero-Shot Prompting  
**Zheng-Xin Yong**, Hailey Schoelkopf, ..., Vassilina Nikoulina (15 authors)  
*ACL 2023*
- [5] Crosslingual Generalization through Multitask Finetuning  
Niklas Muennighoff, ..., **Zheng-Xin Yong**, ..., Colin Raffel (19 authors)  
*ACL 2023*
- [4] PromptSource: An Integrated Development Environment and Repository for Nat-  
ural Language Prompts  
Stephen Bach\*, Victor Sanh\*, **Zheng-Xin Yong**, ..., Alexander M. Rush (27 authors)  
*ACL Demo 2023*

	<p>[3] The Decades Progress on Code-Switching Research in NLP: A Systematic Survey on Trends and Challenges Genta Indra Winata, Alham Fikri Aji, <b>Zheng-Xin Yong</b>, Thamar Solorio <i>ACL Findings 2023</i></p> <p>[2] BLOOM: A 176B-Parameter Open-Access Multilingual Language Model BigScience Workshop (including <b>Zheng-Xin Yong</b> and 300+ authors) <i>Preprint 2023</i></p> <p>[1] Multitask Prompted Training Enables Zero-Shot Task Generalization Victor Sanh*, ..., <b>Zheng-Xin Yong</b>, ..., Alexander M. Rush (41 authors) <i>ICLR 2022</i> <b>Spotlight, 5% of submitted papers (176/3391)</b></p>	
AWARDS	<p><b>Outstanding Reviewer Award</b> at EMNLP 2024 2024</p> <p><b>Best Paper Award</b> at ACL 2024 2024</p> <p><b>Best Paper Award</b> at NeurIPS 2023 Socially Responsible Language Modeling (SoLaR) Workshop 2023</p> <p><b>3rd Best App Overall</b> at FirstNet Public Safety Hackathon 2018</p> <p><b>Best Use of ESRI</b> at FirstNet Public Safety Hackathon 2018</p> <p><b>Grand Prize Winner</b> at #100Hacks Hackathon for Puerto Rico 2017</p> <p><b>Grand Prize Winner</b> at NCSV Innovate NUS Hackathon 2017</p> <p><b>Open Category Grand Prize</b> at SIA Airlines AppChallenge 2017</p> <p><b>Global Finalist</b> at Google Science Fair 2016</p> <p><b>Outstanding A-Level Cambridge Learner Award</b> for Perfect Score in Maths 2016</p>	
SELECTED TALKS	<p><b>Meta AI: LLM Detoxification Generalizes Across Languages</b>, Aug 2024.</p> <p><b>The Alan Turing Institute, London Data Week: Multilingual AI Safety</b>, Jul 2024.</p> <p><b>Cohere For AI, C4AI Gatherings: Aya Multilingual Safety</b>, Feb 2024.</p> <p><b>Cohere For AI, Aya Grande Finale: Aya Responsible Release</b>, Feb 2024.</p> <p><b>Cohere For AI, Closing the Contribution Chapter: Malay Ambassador</b>, Dec 2023.</p>	
SELECTED PRESS & MEDIA	<p><b>GPT-4 gave advice on planning terrorist attacks when asked in Zulu</b> (New Scientist, 2023)</p> <p><b>ChatGPT Is Cutting Non-English Languages Out of the AI Revolution</b> (WIRED, 2023)</p>	
SERVICES	<p><b>Area Chair/Program Committee:</b> EMNLP 2023 (Multilingualism and Linguistic Diversity)</p> <p><b>Conference/Workshop Reviewer:</b> EMNLP 2024 (<b>Outstanding Reviewer</b>) COLM 2024 ARR 2021-Present (including all *ACL conferences) ACL 2021-2023</p> <p><b>Outreach and Mentorship Service Programs:</b> Deep Learning Indaba Mentorship Service (2022 - 2023) Brown Computer Science exploreCSR (2022) Minerva-Masason Mentoring Program (2019 - 2021)</p>	