

# M SAIFUL BARI

sbmaruf.github.io

sbmaruf@gmail.com

## EDUCATION

---

**Nanyang Technological University**

*Jan 2019 - Aug 2023*

Doctor of Philosophy (Ph.D.)

Natural Language Processing

Computer Science and Engineering

## SKILLS

---

**Computer Languages**

Python, C/C++, Bash

**Tools & Library**

Pytorch, DeepSpeed, fairseq, Huggingface

**Problem Solving**

Deep Learning, Algorithms & Data Structure

## EXPERIENCE

---

**Senior Research Scientist**

Aug 2023 - Present

*Saudi Authority for Data and Artificial Intelligence*

- Working on the Pre-training, Instruction Tuning and the Alignment of the Foundation Models.

**Applied Scientist Intern**

July 2022 - Nov 2022

*Amazon Development Service*

- Worked on parameter efficient multitask inference (PEMI) training in Large Language Model.

**BLOOM LLM training**

Sep 2021 - Mar 2022

*International Research Effort*

- Worked on large-scale LLM training in Architecture WG and Prompt Engineering WG.

**Applied Scientist Intern (Part Time)**

July 2021 - Jan 2022

*Amazon Web Service*

- Working on prompt tuning for multi-lingual models.

**Applied Scientist Intern**

Aug 2020 - Oct 2020

*Amazon Web Service*

- Worked on Cross-lingual Few-shot Adaptation.

**Research Assistant**

Sep 2017 - Aug 2018

*Nanyang Technological University*

- Research on MT, NER and Adversarial Training.

**Aubichol Intelligent Technologies**

Sep 2017 - Aug 2018

*Product Development Intern*

- Help an early *Sports Analytic* start-up to build their MVP.

**Software Engineering Intern**

Nov 2015 - Dec 2015

*XeonBD*

- The course of the internship goes through Cloud Computing and kernel virtualization.

## HONORS & AWARDS

---

<b>Scholarship</b> <i>NTU Research Scholarship, Fully funded Ph.D. scholarship for 4 years.</i>	2019
<b>Scholarship</b> <i>OIC Scholarship for undergraduate study, Islamic University of Technology</i>	2012
<b>Champion</b> <i>IUT Computer Programming Contest</i>	2014
<b>Honorable Mention</b> <i>Human Expedition on Mars Timeline 2018</i>	2014
<b>Champion</b> <i>IUT Computer Programming Contest</i>	2015
<b>2<sup>nd</sup>/100 in Inter University Programming Contest</b> <i>Daffodill International University ACM ICPC world finals warmup contest 2016</i>	2016
<b>6<sup>th</sup>/100+ in Inter University Programming Contest</b> <i>NSU Cybnauts National Programming Contest</i>	2016

## PUBLICATION

---

**Summary:** ACL-19, AACL-20, EMNLP-20, 2\*ACL-21, EMNLP-21, ICLR-22, ACL-22, EMNLP-22, 3\*ACL-23, 3 Preprint, 1 Book Chapter

1. **M Saiful Bari\***, Mohammad Abdullah Matin\* Khan, Xuan Long Do, Weishi Wang, Md Rizwan Parvez, and Shafiq Joty. **xCODEEVAL: A Large Scale Multilingual Multitask Benchmark for Code Understanding, Generation, Translation and Retrieval** (under review at thirty-seventh conference on neural information processing systems, neurips'23), 2023
2. **Bari, M Saiful\***, Laskar Tahmid\*, Rahman Mizanur, Md Amran Hossen Bhuiyan, Shafiq Joty, and Jimmy Huang. **A Systematic Study of ChatGPT on Benchmark Datasets**. In *Findings of the 61th Annual Meeting of the Association for Computational Linguistics, ACL*, 2023
3. **M Saiful Bari**, Aston Zhang, Shuai Zheng, Xingjian Shi, Yi Zhu, Shafiq Joty, and Mu Li. **SPT: Semi-Parametric Prompt Tuning for Multitask Prompted Learning** (under review at the 2023 conference on empirical methods in natural language processing, emnlp'23), 2022
4. Niklas Muennighoff, Thomas Wang, Lintang Sutawika, Adam Roberts, Stella Biderman, Teven Le Scao, **M Saiful Bari**, Sheng Shen, Zheng-Xin Yong, Hailey Schoelkopf, Xiangru Tang, Dragomir Radev, Alham Fikri Aji, Khalid Almubarak, Samuel Albanie, Zaid Alyafeai, Albert Webson, Edward Raff, and Colin Raffel. **Crosslingual Generalization through Multitask Finetuning**. In *Proceedings of the 61th Annual Meeting of the Association for Computational Linguistics, ACL*, 2023
5. Zheng-Xin Yong, Hailey Schoelkopf, Niklas Muennighoff, Alham Fikri Aji, David Ifeoluwa Adelani, Khalid Almubarak, **M Saiful Bari**, Lintang Sutawika, Jungo Kasai, Ahmed Baruwa, Genta Indra Winata, Stella Biderman, Dragomir Radev, and Vassilina Nikoulina. **BLOOM+1: Adding Language Support to BLOOM for Zero-Shot Prompting**. In *Proceedings of the 61th Annual Meeting of the Association for Computational Linguistics, ACL*, 2023

6. Teven Le Scao, Thomas Wang, Daniel Hesslow, Lucile Saulnier, Stas Bekman, **M Saiful Bari**, Stella Biderman, Hady Elsahar, Jason Phang, Ofir Press, Colin Raffel, Victor Sanh, Sheng Shen, Lintang Sutawika, Jaesung Tae, Zheng Xin Yong, Julien Launay, and Iz Beltagy. **What Language Model to Train if You Have One Million GPU Hours?** In *Findings in EMNLP, 2022*, 2022
7. Stephen H. Bach, Victor Sanh, Zheng-Xin Yong, Albert Webson, Colin Raffel, Nihal V. Nayak, Abheesht Sharma, Taewoon Kim, **M Saiful Bari**, Thibault Fevry, Zaid Alyafeai, Manan Dey, Andrea Santilli, Zhiqing Sun, Srulik Ben-David, Canwen Xu, Gunjan Chhablani, Han Wang, Jason Alan Fries, Maged S. Al-shaibani, Shanya Sharma, Urmish Thakker, Khalid Almubarak, Xiangru Tang, Dragomir Radev, Mike Tian-Jian Jiang, and Alexander M. Rush. **PromptSource: An Integrated Development Environment and Repository for Natural Language Prompts.** In *Meeting of the Association for Computational Linguistics (ACL) Demonstration*, 2022
8. Victor Sanh, Albert Webson, Colin Raffel, Stephen Bach, Lintang Sutawika, Zaid Alyafeai, Antoine Chaffin, Arnaud Stiegler, Arun Raja, Manan Dey, **M Saiful Bari**, Canwen Xu, Urmish Thakker, Shanya Sharma, Eliza Szczechla, Taewoon Kim, Gunjan Chhablani, Nihal Nayak, Debajyoti Datta, Jonathan Chang, Mike Tian-Jian Jiang, Han Wang, Matteo Manica, Sheng Shen, Zheng Xin Yong, Harshit Pandey, Rachel Bawden, Thomas Wang, Trishala Neeraj, Jos Rozen, Abheesht Sharma, Andrea Santilli, Thibault Fevry, Jason Alan Fries, Ryan Teehan, Teven Le Scao, Stella Biderman, Leo Gao, Thomas Wolf, and Alexander M Rush. **Multitask Prompted Training Enables Zero-Shot Task Generalization.** In *International Conference on Learning Representations, ICLR, 2022*
9. Teven et al, **BLOOM: A 176B-Parameter Open-Access Multilingual Language Model**, Arxiv preprint, (Under Review at Journal of Machine Learning Research, JMLR), 2022
10. **Bari, M Saiful**, Batool Haider, and Saab Mansour. **Nearest Neighbour Few-Shot Learning for Cross-lingual Classification.** In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 1745–1753, Online and Punta Cana, Dominican Republic, November 2021. Association for Computational Linguistics
11. **Bari, M Saiful**, Tasnim Mohiuddin, and Shafiq Joty. **UXLA: A Robust Unsupervised Data Augmentation Framework for Zero-Resource Cross-Lingual NLP.** In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing, ACL (Volume 1: Long Papers)*, pages 1978–1992, Online, August 2021. Association for Computational Linguistics
12. Tasnim Mohiuddin, **M Saiful Bari**, and Shafiq Joty. **AugVic: Exploiting BiText Vicinity for Low-Resource NMT.** In *Findings of the Association for Computational Linguistics: ACL-IJCNLP 2021*, Online, 2021. Association for Computational Linguistics
13. Tasnim Mohiuddin, **M. Saiful Bari**, and Shafiq R. Joty. **LNMap: Departures from Isomorphic Assumption in Bilingual Lexicon Induction Through Non-Linear Mapping in Latent Space.** In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing*, Virtual, November 2020
14. **M Saiful Bari**, Shafiq Joty, and Prathyusha Jwalapuram. **Zero-Resource Cross-Lingual Named Entity Recognition.** In *Proceedings of the 34th AAAI Conference on Artificial Intelligence*, AAAI '20, New York, USA, 2020. AAAI
15. Xiang Lin, Shafiq Joty, Prathyusha Jwalapuram, and **M Saiful Bari**. **A Unified Linear-Time Framework for Sentence-Level Discourse Parsing.** In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics, ACL '19*, Florence, Italy, 2019. ACL
16. **M Saiful Bari**. *Regression for Data Analytics*, chapter 2, pages 33–54. CRC Press, Boca Raton, September 2018. (in book: Data Analytics: Concepts, Techniques and Applications)

## PROJECTS

---

★ Prompt Engineering :: *Promptsource*

Details can be found here, <https://github.com/bigscience-workshop/promptsource>

★ LLM Pipelining :: *Megatron/DeepSpeed*

Details can be found here, <https://github.com/bigscience-workshop/Megatron-DeepSpeed>

★ Dataset :: *xCodeEval*

Details can be found here, <https://github.com/ntunlp/xCodeEval>

★ Deep learning :: *UXLA*

Details can be found here, <https://github.com/sbmaruf/UXLA>

★ Deep learning :: *Cross-lingual Few Shot Learning*

Details can be found here, <https://github.com/amazon-science/nearest-neighbor-crosslingual-classification>

★ Deep learning :: *Zero-Resource Cross-lingual Named Entity Recognition*

Details of the project can be found here, <https://github.com/ntunlp/Zero-Shot-Cross-Lingual-NER>

★ Deep learning :: *Malay English Machine Translation System*

Details can be found here, <https://sbmaruf.github.io/project/mt-system/>

★ Machine learning system :: *A CBIR System*

Details of the project can be found here, <http://103.82.172.44:8080/xmlui/handle/123456789/93>

★ Algorithms :: *Algorithm-Code-Library*

Details of the project can be found here, <https://github.com/sbmaruf/Algorithms-Code-Library>