# Feisi Fu

8 Saint Mary St, Photonics Center – Boston, MA, 02215

✉ fufeisi2023@gmail.com   🔗 Homepage
📱 857-869-0444   **in** LinkedIn/feisifu

## SUMMARY

- A passionate **Ph.D.** candidate in **Machine Learning** with a strong **mathematical** and **statistical** background. Contributed to the development of **PyTorch**, with experience **building pipeline for training GPT-like models**, and publications in **top-tier Machine Learning conferences (ICLR, NeurIPS).**
- Have six years of research experience in machine learning:
  1. Familiar with data preprocessing techniques, such as **Pandas** and **Sklearn**;
  2. Build state-of-the-art machine learning models using **PyTorch**, **Huggingface** and **DeepSpeed**, including models such as **ViT**, **Bloom**, and **LLaMA**;
  3. Utilize state-of-the-art strategies for multi-GPU training, such as **DDP**, **ZerO**, and **FDSP**;
  4. **Quantize** a trained neural network to improve the **inference efficiency**;
  5. Visualize results using **Matplotlib** and **Wandb**.

## SKILLS

- Areas of Study: **Trustworthy A.I., Neural Network Post-training, Neural Network Quantization**
- Skills: **Python, C++, CUDA Programming, PyTorch, Tensorflow**

## EXPERIENCE

- **Research Intern at ByteDance** (ML Infra Team)                      *Sep. 2023 – Nov. 2023(expected)*
  1. Extend their **LLM training strategy** to support PyTorch Fully Sharded Data Parallel (FSDP), which can improve the training speed of **Bloom-3b** and **LLaMA-(7b, 13b)** models by about 5% compared to their current DeepSpeed ZerO-3 implementation.
  2. Support **distributed checkpoint** for **PyTorch FSDP**. This improvement significantly reduces checkpoint save/load time, theoretically by up to 1/N (where N is the number of GPUs). In our tests, we observed an 85% efficiency gain when training **LLaMA-13b** on 8*A100s.
  3. Extend their **LLM training pipeline** to support for **RWP** ([PDF]), which enhances the model's robustness to quantization. This enhancement results in significant reductions of up to 40% in the quantization gap when **RWP** is applied during the fine-tuning of **Bloom-1b1** and **LLaMA-7b**.
  **Skills: Python, PyTorch, Quantization, FSDP.**
- **Research Scientist Intern at Meta** (PyTorch ArchOpt Team)                      *Aug. 2022 – Dec. 2022*
  Created ML infrastructure for quantizing activation maps in neural network training, achieving a significant 50% reduction in memory usage with only a 0.01% drop observed on ResNet50 and RoBerta-large models. [Code]
  **Skills:Python, C++, CUDA Programming, Quantization.**

## EDUCATION

- **Boston University**                      **Advisor: Prof. Wenchao Li**
  *Doctorate degree in Systems Engineering, GPA: 4.0/4.0*                      *2018 – Early. 2024*
- **Chinese Academy of Sciences, China**                      **Advisor: Prof. Baohua Fu**
  *Master degree in Mathematics, GPA: 84.5/100*                      *2014 – 2018*

- **Sichuan University, China**
  *Bachelor degree in Applied Mathematics, GPA: 89.2/100*                    *2010 – 2014*

## PAPERS

- **Sound and Complete Neural Network Repair with Minimality and Locality Guarantees**[PDF]
  Feisi Fu, Wenchao Li

  **Accept as a poster paper at International Conference on Learning Representations (ICLR), 2022**
  We present the first neural network repair (post-training) methodology that guarantees the removal of buggy behavior while applying only a localized change in the function space.
  Experiment Performance: 1. Repair Rate 96% (ReTrain) $\rightarrow$ guarantee 100% (ours); 2. Negative Side Effect 22.11% (Fine-Tuning) $\rightarrow$ 0.12% (ours).
- **REGLO: Provable Neural Network Repair for Global Robustness Properties**[PDF]
  Feisi Fu, Zhilu Wang, Jiameng Fan, Yixuan Wang, Chao Huang, Xin Chen, Zhu Qi, Wenchao Li
  **Accept as a workshop paper at Neural Information Processing Systems (NeurIPS), 2022**
  We present REGLO, the first work that enables provable repair of a neural network for global robustness properties.
- **Dormant Neural Trojans**[PDF]
  Feisi Fu, Panagiota Kiourti, Wenchao Li
  **Accept as a long paper at IEEE International Conference on Machine Learning and Applications (ICMLA), 2023.**
  We propose a novel methodology for neural network backdoor attacks, inserting a Trojan that will remain dormant until activated. The dormant Trojan can bypass the most state-of-the-art backdoor detention methods.
- **OVLA: Neural Network Ownership Verification using Latent Watermarks**[PDF]
  Feisi Fu, Wenchao Li
- **A Tool for Neural Network Global Robustness Certification and Training**[PDF]
  Zhilu Wang, Yixuan Wang, Feisi Fu, Ruochen Jiao, Chao Huang, Wenchao Li, Qi Zhu

## COMPETITIONS

- **TrojAI Competition by National Institute of Standards and Technology**
  Feisi Fu, Jiameng Fan, Weichao Zhou, Panagiota Kiourti, Sabbir Ahmad, Wenchao Li
  We train a neural network to analyze the eigenvalues of a given network's weights and detect if such network has a Trojan. Our approach achieves the top 5 ROC-AUC among all approaches.

## REVIEWER FOR JOURNALS & CONFERENCE ARTICLES

- **AAAI Conference on Artificial Intelligence (AAAI), 2023, 2024**
- **Design Automation Conference (DAC), 2020, 2022**
- **Design Automation and Test in Europe (DATE), 2021, 2022**
- **Hybrid Systems: Computation and Control (HSCC), 2020**
- **International Conference on Computer-Aided Design (ICCAD), 2021, 2022**
- **International Conference on Dependable Systems and Networks (DSN), 2021, 2022**
- **Workshop on Trustworthy and Socially Responsible Machine Learning, NeurIPS, 2022**