# Hongyu Zhu

Male, Chinese
Date of Birth: Nov. 21, 1990
E-mail: `serailhydra@gmail.com`

## *Research Interests*

Systems, Machine Learning, Parallel Computing, Distributed Systems

I am currently working on predictive performance for DNN training

## *Education*

- **Ph.D**                                                                     *Now*
  University of Toronto, Toronto, ON, Canada
  Supervised by Prof. Bianca Schroeder (co-supervised by Prof. Gennady Pekhimenko)

- **Master**                                                                   *April 2016*
  McGill University, Montreal, QC, Canada
  GPA: 3.94/4.0

- **Bachelor (prestigious permission)**                                        *June 2013*
  Shanghai Jiao Tong University
  ACM class
  Advisor: Prof. Yong Yu, Prof. Minyi Guo
  GPA: 86.6/100

## *Work Experience*

| | |
|---|---|
| **Jun. 2018 — Sept. 2018** | Research Intern in Microsoft Research, supervised by Amar Phanishayee. Working on performance debugging for DNN computation |
| **Jul. 2014 — Dec. 2014** | Research Intern in New York University in Shanghai (NYU-Shanghai), supervised by Prof. Zheng Zhang. Working on Deep Learning Training Platform (Minerva Project) |
| **Jul. 2012 — Feb. 2013** | Research Intern in Microsoft Research Asia (MSRA), supervised by Dr. Zheng Zhang, the vice Dean of MSRA, and Dr. Zhengping Qian. Working on Streaming DAG Distributed System (TimeStream project) |
| **Jul. 2011 — Jul. 2013** | Research Intern in Embedded Pervasive Computing Center (EPCC), Shanghai Jiao Tong University (SJTU), supervised by Prof. Minyi Guo. Working on GPU communication framework |

## *Projects*

- **Predictive Profiling for DNN computation:** In this project, we aim to enable predictive profiling for DNN computation via answering some what-if questions: how the overall training time will change if I double the number of GPUs, add or remove a layer, use network with higher bandwidth, or speed up some kernels, etc. We build an execution model by representing the entire training job with a dependency graph. We then use the dependency graph as a basis to estimate the training time under new changes, and perform critical path analysis to show optimization potentials. This project is a sub-project of Project Fiddle of Microsoft Research, and it is currently ongoing.

- **Benchmarking and Analyzing DNN Training:** The goal of this project is to understand the performance bottlecks of training modern DNNs. The project consists of several parts: maintaining a diverse benchmark suite with state-of-the-art DNN models, defining performance metrics, and building tools to extract these metrics. We apply our

toolchains to our benchmark suite, gaining some insights for optimizations. A paper has been published in 2018 IEEE International Symposium on Workload Characterization (IISWC18), and the project is still active.

- **A Distributed-system Infrastructure for Real-time Streaming Applications (TimeStream):** TimeStream is a distributed system designed specifically for real-time continuous processing of big streaming data on a large cluster of commodity machines. This project is led by Dr. Zhengping Qian in Dr. Zheng Zhang's group, the vice Dean of MSRA. In this project I participate in designing, analyzing and implementing the fault tolerance protocol, the key part of this project. Besides, I lead the work of implementing two benchmarks which are real applications from industry, and do experiments for scalability and latency. A paper is published on 8th European Conference on Computer Systems (EuroSys 2013).

## *Professional Skills*

- Master in C, C++, Java, C#, Python and Matlab
- Familiar with system design & implementations of main-stream distributed deep learning systems
  - Instrumenting MXNet framework for memory profiler
  - Instrumenting PyTorch framework to discover performance bottlenecks
- Understanding common machine learning and deep learning algorithms
- Adept at algorithms and data structures
  - Competition in Informatics Olympiad about algorithms for 2 years
  - Experience in solving online problem archives (Ural, UVa, etc)

## *Publications*

- Zhu H., Akrout M., Zheng B., Pelegris A., Phanishayee A., Schroeder B., Pekhimenko G. (2018). Benchmarking and Analyzing Deep Neural Network Training. In IEEE International Symposium on Workload Characterization 2018 (IISWC18).
- El-Sayed N, Zhu H, Schroeder B. Learning from Failure Across Multiple Clusters: A Trace-Driven Approach to Understanding, Predicting, and Mitigating Job Terminations[C]//Distributed Computing Systems (ICDCS), 2017 IEEE 37th International Conference on. IEEE, 2017: 1333-1344.
- Qian, Z., He, Y., Su, C., Wu, Z., Zhu, H., Zhang, T., ... , Zhang, Z. Timestream: Reliable stream computation in the cloud. In Proceedings of the 8th ACM European Conference on Computer Systems, ACM, 2013.

## *Invited Talks*

- Holistic Approach to DNN Training Efficiency: Analysis and Optimizations
  Guest Lecture for Advanced Computer Architecture course          Mar 2019
- Benchmarking and Analyzing DNN Training
  Huawei, Markham, ON, Canada          Apr 2018

## *Service*

- Owner of the Deep Speech 2 reference model for MLPerf inference
- Member of Artifact Evaluation Committee of SysML 2019