

# Andrew Z. Luo

MACHINE LEARNING · SYSTEMS · SOFTWARE ENGINEERING

☎ 425-241-9772 | ✉ andrew.zhao.luo@gmail.com | 🌐 andrew-zhao-luo

## Summary

Engineer with interest in machine learning systems, compilers, and performance. Has experience in writing compilers to make it easier to make machine learning systems run fast. Proficient in **Python, C/C++**, working knowledge of **CUDA**, and experience with **MLIR**. Enjoys solving problems, quickly iterating, and taking ownership in ambiguous environments.

## Experience

### Senior Compiler Engineer

Seattle, WA

MODULAR AI, CORE TECHNOLOGIES

Jun. 2023 - Now

- Top contributor in the past year (by commits) on a team of 10 to a MLIR based compiler for deep learning models.
- Major contributor to shape propagation, quantization, and programmability systems of compiler.
- Lead efforts around inter-operability with Mojo with a group of 7 engineers in creating a Triton-lang like experience.

### Staff Engineer

Seattle, WA

OCTOML, APPLIED COMPILER ENGINEERING - PERFORMANCE

Mar. 2021 - Jun. 2023

- **Official committer to Apache TVM (9.2k stars on GitHub)**, an open source autotuning compiler written in C++ and Python.
- Contribute quantization, mixed precision, ONNX support in TVM. Write kernels and benchmarking tools to improve iteration.
- Implementing features in SaaS product such as compiler cache, adding new tuning algorithms, and lowering resource usage.

### Machine Learning Engineer

Seattle, WA

APPLE, AI/ML MACHINE INTELLIGENCE NEURAL DESIGN

Jan. 2020 - Mar. 2021

- Using quantization, sparsity, and hardware-specific knowledge to train models for Siri, Homepod, and future products
- Developing in-house solutions for training vision models and deploying/benchmarking on FPGA and ASIC environments
- Developing demos show-casing technologies to internal stake-holders.

XNOR.AI, MACHINE LEARNING TEAM (ACQUIRED BY APPLE)

Aug. 2019 - Jan 2020

- Training performant computer vision models that can run on bespoke and edge hardware. Part time until Jan 2019.
- Creating demos showcasing XNOR's technologies to key executives at major tech companies.

## Education

### University of Washington

Seattle, WA

DOUBLE MAJOR IN COMPUTER ENGINEERING AND BIOENGINEERING

Sep. 2015 - Jun. 2019

- **Coursework:** Machine Learning, Probability and Statistics, Real Analysis, Operating Systems, Compilers, Embedded Systems
- **GPA:** 3.95, *Summa Cum Laude*

## Projects and Potpourri

PAST INTERNSHIPS

Summers 2015-2018

- **Sift:** I rewrote HBase snapshot system, **saving over \$1.5 million in S3 costs a year** and integrated data store with BigQuery.
- **Facebook:** I implemented statistical models to predict ad reach demographics for customers with multi-million yearly spend
- **The Institute For Systems Biology:** I helped scientists analyze large datasets of gene expression data.

PROJECT: FPGA IMAGE CONVOLUTION PHOTOBOOTH

2019

- Implemented streaming algorithm to run kernel convolutions on streamed images, implemented in FPGA on Altera Cyclone V
- Integrated with camera and VGA, creating a variety of filters like Sobel edge detector, Gaussian blur, and image sharpening

PUBLICATION

2017

*Automatic Characterization of User Errors in Spirometry.* **Andrew Luo**, Eric Whitmire, James Stout, Drew Martenson, Shwetak Patel. *IEEE EMBC 2017 (Oral Presentation + Paper)*

PATENT APPLICATION

2020

*Compressed Neural Network Models.* US Patent App. 16/788261. James Gabriel et al. and **Andrew Luo**. Filed 13 August 2020.