

Yufan Xu

CONTACT INFORMATION

19608 Pruneridge Ave, *Phone:* (352) 278-6832
Cupertino, *Email:* yf.xu@utah.edu
California 95014 *Profile:* Google Scholar

EDUCATION

University of Utah, Salt lake City, Utah USA (GPA 4.0)
Ph.D, Computer Science, May, 2024
The Ohio State University, Columbus, Ohio USA (GPA 3.71)
Ph.D Student, Computer Science, August, 2017 - August, 2019 (Transfer to Utah)
University of Florida, Gainesville, Florida USA (GPA 3.65)
M.S., Computer Science, May, 2016
Soochow University, Suzhou, Jiangsu China (GPA 3.50)
B.Eng., Software Engineering, May, 2014

PUBLICATION

Accelerated Auto-Tuning of GPU Kernels for Tensor Computations
ICS 24

- Chendi Li*, Yufan Xu*, Sina Mahdipour Saravani, P. Sadayappan

CoNST: Code Generator for Sparse Tensor Networks
TACO 24

- Saurabh Rajе, Yufan Xu, Atanas Rountev, Edward F. Valeev, P. Sadayappan

PEAK: Generating High-Performance Schedules in MLIR
LCPC 23

- Amir Tavakkoli*, Sameeran Joshi*, Shreya Singh, Yufan Xu, P. Sadayappan, Marry Hall

Effective Performance Modeling and Domain-Specific Compiler Optimization of CNNs for GPU
PACT 22

- Yufan Xu, Qiwei Yuan, Erik Curtis Barton, Rui Li, P. Sadayappan, Aravind Sukumaran-Rajam

Training of Deep Learning Pipelines on Memory-Constrained GPUs via Segmented Fused-Tiled Execution
CC 22

- Yufan Xu, Saurabh Rajе, Atanas Rountev, Gerald Sabin, Aravind Sukumaran-Rajam, P. Sadayappan

Efficient Distributed Algorithms for Convolutional Neural Networks
SPAA 21

- Rui Li, Yufan Xu, Aravind Sukumaran-Rajam, Atanas Rountev, P. Sadayappan

Analytical characterization and design space exploration for optimization of CNNs
ASPLOS 21

- Rui Li, Yufan Xu, Aravind Sukumaran-Rajam, Atanas Rountev, P. Sadayappan

Dependence-aware, unbounded sound predictive race detection
OOPSLA 19

- Kaan Genç, Jake Roemer, Yufan Xu, Michael D. Bond

RESEARCH EXPERIENCE	<i>Uber Technologies Inc.</i>	May, 2024 - Now
	<ul style="list-style-type: none"> • Work on memory allocation optimization with PGO in GO compiler • Work on error propagation fixing with GenAI in GO monorepo • Work on ML infrastructure inference performance 	
	<i>University of Utah</i>	August, 2019 - May, 2024
	<ul style="list-style-type: none"> • Worked on search space optimization in TVM Improve consistency, efficiency of TVM internal candidate configuration selection algorithm • Worked on design space exploration for optimizing CNN for GPUs Prune the kernel configuration space by using data-driven analysis and design a hybrid model for configuration quick selection • Worked on memory efficiency for large input on ML system (pytorch) Solve memory constraint issue of a large input image training on a single GPU • Worked on opmin optimization pass for tensor contraction in CCSD benchmark on MLIR Implement a MLIR pass to reduce total number of float operation of high order tensor contraction expressions • Worked on a tile-size optimization problem for affine programs in polyhedral model Build an approximate modeling method for tile size selection 	
	<i>The Ohio State University</i>	August, 2017 - May, 2019
	<ul style="list-style-type: none"> • Worked on data race detection in Java program 	
TEACHING & ADVISING	Course Instructor	
	<i>The Ohio State University</i>	Fall, 2018, Spring, 2019
	<ul style="list-style-type: none"> • Instructor for two semesters of <i>CS1223 Introduction to Computer Programming In Java</i>. • Taught the general concepts of computer programming and programming languages by providing practical experience programming in the Java. 	
	Teaching Assistant	
	<i>University of Utah</i>	Spring, 2020
	<ul style="list-style-type: none"> • Teaching Assistant for <i>CS 6230 Parallel Computing and HPC</i>. • Planned course project, graded assignments and projects. 	
WORKING EXPERIENCE	Uber , Sunnyvale, CA, USA	
	Software Engineer II May, 2024 - Now	
	Uber , Sunnyvale, CA, USA	
	PhD Software Engineer(Intern) May, 2023 - August, 2023	
	LatentAI , Princeton, NJ, USA	
	Compiler Engineer(Intern) May, 2022 - August, 2022	
	T-CETRA , Columbus, OH, USA	
	Software Engineer(Intern) May, 2019 - August, 2019	
SERVICE	Program Committee:	
	<i>CGO</i> '25	
	Artifact Evaluation Committee:	
	<i>ASPLOS</i> '21, '22 ; <i>CGO</i> '23, '24 ; <i>MICRO</i> '23 ; <i>CC</i> '24	