

# Yufan Xu

---

## CONTACT INFORMATION

19608 Pruneridge Ave,  
Cupertino,  
California 95014

*Phone:* (352) 278-6832  
*Email:* yf.xu@utah.edu  
*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**

*arxiv 24*

- Saurabh Raje, [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 Raje, 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	<b>Uber Technologies Inc.</b> <span style="float: right;"><b>May, 2024 - Now</b></span> <ul style="list-style-type: none"> <li>• Work on memory allocation optimization with PGO in GO compiler</li> <li>• Work on data race fixing with GenAI in GO monorepo</li> <li>• Work on ML infrastructure inference performance</li> </ul>
	<b>University of Utah</b> <span style="float: right;"><b>August, 2019 - May, 2024</b></span> <ul style="list-style-type: none"> <li>• Worked on search space optimization in TVM Improve consistency, efficiency of TVM internal candidate configuration selection algorithm</li> <li>• 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</li> <li>• 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</li> <li>• 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</li> <li>• Worked on a tile-size optimization problem for affine programs in polyhedral model Build an approximate modeling method for tile size selection</li> </ul>
	<b>The Ohio State University</b> <span style="float: right;"><b>August, 2017 - May, 2019</b></span> <ul style="list-style-type: none"> <li>• Worked on data race detection in Java program</li> </ul>
TEACHING & ADVISING	<b>Course Instructor</b> <span style="float: right;"><b>Fall, 2018, Spring, 2019</b></span> <b>The Ohio State University</b> <ul style="list-style-type: none"> <li>• Instructor for two semesters of <i>CS1223 Introduction to Computer Programming In Java</i>.</li> <li>• Taught the general concepts of computer programming and programming languages by providing practical experience programming in the Java.</li> </ul>
	<b>Teaching Assistant</b> <span style="float: right;"><b>Spring, 2020</b></span> <b>University of Utah</b> <ul style="list-style-type: none"> <li>• Teaching Assistant for <i>CS 6230 Parallel Computing and HPC</i>.</li> <li>• Planned course project, graded assignments and projects.</li> </ul>
	<b>WORKING EXPERIENCE</b> <b>Uber</b> , Sunnyvale, CA, USA Software Engineer II May, 2024 - Now  <b>Uber</b> , Sunnyvale, CA, USA PhD Software Engineer(Intern) May, 2023 - August, 2023  <b>LatentAI</b> , Princeton, NJ, USA Compiler Engineer(Intern) May, 2022 - August, 2022  <b>T-CETRA</b> , 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