

Abdul Dakkak

COMPILER DEVELOPER · SYSTEM ARCHITECT · AI OPTIMIZER

2906 S. Myra Ridge Dr, Urbana, IL 61802

☎ (+1) 419-418-1158 | ✉ dakkak@illinois.edu | 🏠 www.dakkak.dev | 👤 abduld

Research Interest

My research interest lies between programming languages and accelerated computing. My work has focused on understanding and optimizing high-level languages that target accelerators. In the process, I have developed widely used industry-grade tools for compiling, running, profiling, and introspecting whole-level applications to optimize their performance across both the hardware and software stack.

Education

PhD. in Computer Science

UNIVERSITY OF ILLINOIS AT URBANA-CHAMPAIGN

Champaign, IL

Aug. 2013 - Exp. Aug. 2020

B.A. in Pure Mathematics

UNIVERSITY OF TOLEDO

Toledo, OH

Aug. 2005 - June. 2009

Experience

Senior Compiler Developer – WOLFRAM RESEARCH

Jan. 2019 - Present

- Continued to lead and develop the Wolfram Language compiler.
- Architected and developed the runtime library used for the Wolfram language.
- Developed a constraint-based type system for the Wolfram language.
- Developed datastructures and applications leveraging the compiler (in *Mathematica* 12.1)

Kernel Developer – WOLFRAM RESEARCH

April. 2010 - Dec. 2018

- Lead a team to develop the Wolfram compiler which was released in *Mathematica* 12.
- Developed a domain specific language to write financial code for Wolfram Finance Platform (released in *Mathematica* 11).
- Developed graphics rendering for the cloud using both canvas and WebGL (released in *Mathematica* 9).
- Developed optimized data-structures and algorithms for computational geometry (released in *Mathematica* 10).

Junior Kernel Developer – WOLFRAM RESEARCH

April. 2009 - April. 2010

- Developed Wolfram's CUDA and OpenCL integration (released in *Mathematica* 8).
- Enhanced the C/C++ library bindings interface of the Wolfram language.

Projects

Team Lead & Senior Compiler Developer – MATHEMATICA (WOLFRAM.COM/MATHEMATICA)

2009 – Present

- Developed the Wolfram type system, runtime, and compiler.
- Optimized core Wolfram engine for desktop and cloud.
- Developed CUDALink and OpenCLLink.
- Developed a DSL for writing GPU and CPU financial indicators.

System Architect & Primary Developer – MLMODELSCOPE (GITHUB.COM/RAI-PROJECT/MLMODELSCOPE)

2017 – 2020

- An inference system designed for hierarchical profiling and benchmarking.
- Over 300 built-in models supported and integration with MLPerf workloads.

System Architect & Primary Developer – RAI (GITHUB.COM/RAI-PROJECT/RAI)

2016 – 2019

- A system designed as a configurable programming environment for heterogeneous parallel programming.
- An interactive command line tool used for building and executing accelerated code in the cloud.

System Architect & Primary Developer – WEBGPU (WEBGPU.COM)

2013 – 2018

- A lab-submission system designed for developing CUDA and OpenCL programming within the browser.
- Used by over 100,000 students to evaluate millions of labs.

Publications

The Design and Implementation of the Wolfram Compiler	CGO
ADBUL DAKKAK, TOM WICKHAM-JONES, WEN-MEI HWU	2020
DLSpec: A Deep Learning Task Exchange Specification	USENIX OpML
ADBUL DAKKAK, CHENG LI, JINJUN XIONG, WEN-MEI HWU	2020
The Design and Implementation of a Scalable DL Benchmarking Platform	Cloud
ADBUL DAKKAK, CHENG LI, JINJUN XIONG, WEN-MEI HWU	2020
Benanza: Automatic μBenchmark Generation to Compute “Lower-bound” Latency and Inform Optimizations of Deep Learning Models on GPUs	IPDPS
ADBUL DAKKAK, CHENG LI, JINJUN XIONG, WEN-MEI HWU	2020
DLBricks: Composable Benchmark Generation to Reduce Deep Learning Benchmarking Effort on CPUs	ICPE
CHENG LI, ADBUL DAKKAK, JINJUN XIONG, WEN-MEI HWU	2020
XSP: Across-Stack Profiling and Analysis of Machine Learning Models on GPUs	IPDPS
ADBUL DAKKAK, CHENG LI, JINJUN XIONG, WEI WEI, LINGJIE XU, WEN-MEI HWU	2020
MLModelScope: A Distributed Platform for Model Evaluation and Benchmarking at Scale	ArXiv
ADBUL DAKKAK, CHENG LI, JINJUN XIONG, WEN-MEI HWU	2020
TrIMS: Transparent and Isolated Model Sharing for Low Latency Deep Learning Inference in Function as a Service Environments	Cloud
ADBUL DAKKAK, CHENG LI, SIMON GARCIA DE GONZALO, JINJUN XIONG, WEN-MEI HWU	2019
Accelerating Reduction and Scan Using Tensor Core Units	ICS
ADBUL DAKKAK, CHENG LI, JINJUN XIONG, ISAAC GELADO, WEN-MEI HWU	2019
Frustrated with Replicating Claims of a Shared Model? A Solution	ArXiv
ADBUL DAKKAK, CHENG LI, JINJUN XIONG, WEN-MEI HWU	2019
Evaluating Characteristics of CUDA Communication Primitives on High-Bandwidth Interconnects	ICPE
CARL PEARSON, ADBUL DAKKAK, SARAH HASHASH, CHENG LI, I CHUNG, JINJUN XIONG, WEN-MEI HWU	2019
Accelerating Reduction Using Tensor Core Units	HPCaML
ADBUL DAKKAK, CHENG LI, JINJUN XIONG, WEN-MEI HWU	2019
SCOPE: C3SR Systems Characterization and Benchmarking Framework	ArXiv
CARL PEARSON, ADBUL DAKKAK, CHENG LI, SARAH HASHASH, JINJUN XIONG, WEN-MEI HWU	2019
Challenges and Pitfalls of Reproducing Machine Learning Artifacts	ArXiv
CHENG LI, ADBUL DAKKAK, JINJUN XIONG, WEN-MEI HWU	2019
TrIMS: Transparent and Isolated Model Sharing for Low Latency Deep Learning Inference in Function as a Service Environments	SysML NIPS
ADBUL DAKKAK, CHENG LI, SIMON GARCIA DE GONZALO, JINJUN XIONG, WEN-MEI HWU	2018
Thoughts on Massively-Parallel Heterogeneous Computing for Solving Large Problems	CEM
WEN-MEI HWU, MERT HIDAYETOGLU, WENG CHO CHEW, CARL PEARSON, SIMON GARCIA, SITAO HUANG, ADBUL DAKKAK	2017
RAI: A Scalable Project Submission System for Parallel Programming Courses	IPDPSW
ADBUL DAKKAK, CARL PEARSON, CHENG LI, WEN-MEI HWU	2017
WebGPU: A Scalable Online Development Platform for GPU Programming Courses	IPDPSW
ADBUL DAKKAK, CARL PEARSON, WEN-MEI HWU	2016
A Programming System for Future Proofing Performance Critical Libraries	SIGPLAN
LI-WEN CHANG, IZZAT EL HAJJ, HEE-SEOK KIM, JUAN GÓMEZ-LUNA, ADBUL DAKKAK, WEN-MEI HWU	2016
Enhancing the Usability and Utilization of Accelerated Architectures via Docker	UCC
NICHOLAS HAYDEL, SANDRA GESING, IAN TAYLOR, GREGORY MADEY, ADBUL DAKKAK, SIMON GARCIA DE GONZALO, WEN-MEI HWU	2015
Massively-Parallel Heterogeneous Computing for Solving Large Problems	IPDPSW
WEN-MEI HWU, MERT HIDAYETOGLU, CARL PEARSON, SIMON GARCIA, SITAO HUANG, ADBUL DAKKAK	2016

Tangram: a High-level Language for Performance Portable Code Synthesis

LI-WEN CHANG, ABDUL DAKKAK, CHRISTOPHER I RODRIGUES, WEN-MEI HWU

MULTIPROG

2015

Transitioning HPC software to exascale heterogeneous computing

WEN-MEI HWU, LI-WEN CHANG, HEE-SEOK KIM, ABDUL DAKKAK, IZZAT EL HAJJ

CEM

2015

Triolet: A Programming System that Unifies Algorithmic Skeleton Interfaces for High-performance Cluster Computing

CHRISTOPHER RODRIGUES, THOMAS JABLIN, ABDUL DAKKAK, WEN-MEI HWU

PPoPP

2014

Recovering Missing Depth Information from Microsoft's Kinect

ABDUL DAKKAK, AMMAR HUSAIN

EVA

2012

CUDA & Heterogeneous Programming with the Wolfram Language

ABDUL DAKKAK, ULISES CERVANTES-PIMENTEL

Wolfram Whitepaper

2012

CUDA Programming Using Wolfram Finance Platform

ABDUL DAKKAK, ULISES CERVANTES-PIMENTEL

Wolfram Whitepaper

2011

Selected Presentations

Using Tensor Cores for Accelerating Reduction and Scan

GPU TECHNOLOGY CONFERENCE

San Jose, CA

Mar. 2020

MLPerf-Bench Tutorial

ASPLOS

Lausanne, Switzerland

Mar. 2020

Challenges and Solutions for End-to-End and Across Stack ML Benchmarking Tutorial

SUPER COMPUTING

Denver, CA

Aug. 2019

Challenges and Solutions for End-to-End and Across Stack ML Benchmarkin Tutorial

IISWC

Orlando, FL

Aug. 2019

MLModelScope: Evaluate and Profile ML Models at Scale and Across Stack

HOT CHIPS

Pala Alto, CA

Aug. 2019

MLPerf-Bench Tutorial

ISCA

Pheonix, AZ

Jun. 2019

MLPerf-Bench Tutorial

ASPLOS

Providence, RI

Apr. 2019

MLModelScope

MLPERF DECEMBER MEETING 2018

Champaign, IL

Dec. 2018

Advanced Compilation Techniques

WOLFRAM TECHNOLOGY CONFERENCE

Champaign, IL

Oct. 2018

TensorOps: Accelerating Reduction Using Tensor Core Units

NVIDIA GPU TECHNOLOGY CONFERENCE

San Jose, CA

Mar. 2019

TRIMS: Transparent and Isolated Model Sharing for Low Latency Deep Learning Inference

NVIDIA GPU TECHNOLOGY CONFERENCE

San Jose, CA

Mar. 2019

MLModelScope: Evaluate and Measure Machine Learning Models within AI Pipelines

NVIDIA GPU TECHNOLOGY CONFERENCE

San Jose, CA

Mar. 2019

MLModelScope: Evaluate and Measure Machine Learning Models within AI Pipelines

SUPER COMPUTING

Dallas, TX

Nov. 2018

TRIMS: Transparent and Isolated Model Sharing for Low Latency Deep Learning Inference in Function as a Service Environments

IBM AI SYSTEMS DAY

Boston, MA

Oct. 2018

RAI: A Scalable Submission System for GPU Applications

NVIDIA GPU TECHNOLOGY CONFERENCE

San Jose, CA

Mar. 2018

CarML: Common Artifacts for Machine Learning

SUPER COMPUTING

Denver, CO

Nov. 2017

Teaching Experience

2018	Lead TA , Pumps-AI Summer School	<i>Barcelona, Spain</i>
2017	Lead TA , CS 508: Manycore Parallel Programming	<i>Champaign, Illinois</i>
2016	Head TA , Pumps Summer School	<i>Barcelona, Spain</i>
2016	Support TA , CS 408: Applied Parallel Programming	<i>Champaign, Illinois</i>
2015	Head TA , Pumps Summer School	<i>Barcelona, Spain</i>
2015	Head TA , Coursera Heterogeneous Parallel Programming	<i>Champaign, Illinois</i>
2014	Head TA , Coursera Heterogeneous Parallel Programming	<i>Champaign, Illinois</i>
2013	Head TA , Coursera Heterogeneous Parallel Programming	<i>Champaign, Illinois</i>
2013	Lead TA , CS 408: Applied Parallel Programming	<i>Champaign, Illinois</i>

Skills

Programming Languages C/C++, Mathematica, CUDA, OpenCL, GoLang, OpenMP, JavaScript, Python, Haskell, \LaTeX
English, Arabic