

# 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 | 👤 abdul

## 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 architect and develop the Wolfram Language compiler.
- Architected 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+)

### Kernel Developer – WOLFRAM RESEARCH

April. 2010 - Dec. 2018

- Architected 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

---

## DLSpec: A Deep Learning Task Exchange Specification

ABDUL DAKKAK, CHENG LI, JINJUN XIONG, WEN-MEI HWU

USENIX OpML

2020

## Benanza: Automatic $\mu$ Benchmark Generation to Compute “Lower-bound” Latency and Inform Optimizations of Deep Learning Models on GPUs

ABDUL DAKKAK, CHENG LI, JINJUN XIONG, WEN-MEI HWU

IPDPS

2020

## DLBricks: Composable Benchmark Generation to Reduce Deep Learning Benchmarking Effort on CPUs

CHENG LI, ABDUL DAKKAK, JINJUN XIONG, WEN-MEI HWU

ICPE

2020

## XSP: Across-Stack Profiling and Analysis of Machine Learning Models on GPUs

ABDUL DAKKAK, CHENG LI, JINJUN XIONG, WEI WEI, LINGJIE XU, WEN-MEI HWU

IPDPS

2020

## The Design and Implementation of the Wolfram Compiler

ABDUL DAKKAK, TOM WICKHAM-JONES, WEN-MEI HWU

CGO

2020

## MLModelScope: A Distributed Platform for Model Evaluation and Benchmarking at Scale

ABDUL DAKKAK, CHENG LI, JINJUN XIONG, WEN-MEI HWU

ArXiv

2020

## The Design and Implementation of a Scalable DL Benchmarking Platform

ABDUL DAKKAK, CHENG LI, JINJUN XIONG, WEN-MEI HWU

ArXiv

2019

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

ABDUL DAKKAK, CHENG LI, SIMON GARCIA DE GONZALO, JINJUN XIONG, WEN-MEI HWU

Cloud

2019

## Accelerating Reduction and Scan Using Tensor Core Units

ABDUL DAKKAK, CHENG LI, JINJUN XIONG, ISAAC GELADO, WEN-MEI HWU

ICS

2019

## Frustrated with Replicating Claims of a Shared Model? A Solution

ABDUL DAKKAK, CHENG LI, JINJUN XIONG, WEN-MEI HWU

ArXiv

2019

## Evaluating Characteristics of CUDA Communication Primitives on High-Bandwidth Interconnects

CARL PEARSON, ABDUL DAKKAK, SARAH HASHASH, CHENG LI, I CHUNG, JINJUN XIONG, WEN-MEI HWU

ICPE

2019

## Accelerating Reduction Using Tensor Core Units

ABDUL DAKKAK, CHENG LI, JINJUN XIONG, WEN-MEI HWU

HPCaML

2019

## SCOPE: C3SR Systems Characterization and Benchmarking Framework

CARL PEARSON, ABDUL DAKKAK, CHENG LI, SARAH HASHASH, JINJUN XIONG, WEN-MEI HWU

ArXiv

2019

## Challenges and Pitfalls of Reproducing Machine Learning Artifacts

CHENG LI, ABDUL DAKKAK, JINJUN XIONG, WEN-MEI HWU

ArXiv

2019

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

ABDUL DAKKAK, CHENG LI, SIMON GARCIA DE GONZALO, JINJUN XIONG, WEN-MEI HWU

SysML NIPS

2018

## Thoughts on Massively-Parallel Heterogeneous Computing for Solving Large Problems

WEN-MEI HWU, MERT HIDAYETOGLU, WENG CHO CHEW, CARL PEARSON, SIMON GARCIA, SITAO HUANG, ABDUL DAKKAK

CEM

2017

## RAI: A Scalable Project Submission System for Parallel Programming Courses

ABDUL DAKKAK, CARL PEARSON, CHENG LI, WEN-MEI HWU

IPDPSW

2017

## WebGPU: A Scalable Online Development Platform for GPU Programming Courses

ABDUL DAKKAK, CARL PEARSON, WEN-MEI HWU

IPDPSW

2016

## A Programming System for Future Proofing Performance Critical Libraries

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

SIGPLAN

2016

## Enhancing the Usability and Utilization of Accelerated Architectures via Docker

NICHOLAS HAYDEL, SANDRA GESING, IAN TAYLOR, GREGORY MADEY, ABDUL DAKKAK, SIMON GARCIA DE GONZALO, WEN-MEI HWU

UCC

2015

## Massively-Parallel Heterogeneous Computing for Solving Large Problems

WEN-MEI HWU, MERT HIDAYETOGLU, CARL PEARSON, SIMON GARCIA, SITAO HUANG, ABDUL DAKKAK

IPDPSW

2016

<b>Tangram: a High-level Language for Performance Portable Code Synthesis</b>	MULTIPROG
LI-WEN CHANG, <b>ABDUL DAKKAK</b> , CHRISTOPHER I RODRIGUES, WEN-MEI HWU	2015
<b>Transitioning HPC software to exascale heterogeneous computing</b>	CEM
WEN-MEI HWU, LI-WEN CHANG, HEE-SEOK KIM, <b>ABDUL DAKKAK</b> , IZZAT EL HAJJ	2015
<b>Triolet: A Programming System that Unifies Algorithmic Skeleton Interfaces for High-performance Cluster Computing</b>	PPoPP
CHRISTOPHER RODRIGUES, THOMAS JABLIN, <b>ABDUL DAKKAK</b> , WEN-MEI HWU	2014
<b>Recovering Missing Depth Information from Microsoft's Kinect</b>	EVA
<b>ABDUL DAKKAK</b> , AMMAR HUSAIN	2012
<b>CUDA &amp; Heterogeneous Programming with the Wolfram Language</b>	Wolfram Whitepaper
<b>ABDUL DAKKAK</b> , ULISES CERVANTES-PIMENTEL	2012
<b>CUDA Programming Using Wolfram Finance Platform</b>	Wolfram Whitepaper
<b>ABDUL DAKKAK</b> , ULISES CERVANTES-PIMENTEL	2011

## Presentations

---

<b>Using Tensor Cores for Accelerating Reduction and Scan</b>	San Jose, CA
GPU TECHNOLOGY CONFERENCE	Mar. 2020
<b>MLPerf-Bench Tutorial</b>	Lausanne, Switzerland
ASPLOS	Mar. 2020
<b>Challenges and Solutions for End-to-End and Across Stack ML Benchmarking Tutorial</b>	Denver, CA
SUPER COMPUTING	Aug. 2019
<b>Challenges and Solutions for End-to-End and Across Stack ML Benchmarkin Tutorial</b>	Orlando, FL
IISWC	Aug. 2019
<b>MLModelScope: Evaluate and Profile ML Models at Scale and Across Stack</b>	Pala Alto, CA
HOT CHIPS	Aug. 2019
<b>MLPerf-Bench Tutorial</b>	Pheonix, AZ
ISCA	Jun. 2019
<b>MLPerf-Bench Tutorial</b>	Providence, RI
ASPLOS	Apr. 2019
<b>MLModelScope</b>	Champaign, IL
MLPERF DECEMBER MEETING 2018	Dec. 2018
<b>Advanced Compilation Techniques</b>	Champaign, IL
WOLFRAM TECHNOLOGY CONFERENCE	Oct. 2018
<b>TensorOps: Accelerating Reduction Using Tensor Core Units</b>	San Jose, CA
NVIDIA GPU TECHNOLOGY CONFERENCE	Mar. 2019
<b>TriMS: Transparent and Isolated Model Sharing for Low Latency Deep Learning Inference</b>	San Jose, CA
NVIDIA GPU TECHNOLOGY CONFERENCE	Mar. 2019
<b>MLModelScope: Evaluate and Measure Machine Learning Models within AI Pipelines</b>	San Jose, CA
NVIDIA GPU TECHNOLOGY CONFERENCE	Mar. 2019
<b>MLModelScope: Evaluate and Measure Machine Learning Models within AI Pipelines</b>	Dallas, TX
SUPER COMPUTING	Nov. 2018
<b>TRIMS: Transparent and Isolated Model Sharing for Low Latency Deep Learning Inference in Function as a Service Environments</b>	Boston, MA
IBM AI SYSTEMS DAY	Oct. 2018
<b>RAI: A Scalable Submission System for GPU Applications</b>	San Jose, CA
NVIDIA GPU TECHNOLOGY CONFERENCE	Mar. 2018
<b>CarML: Common Artifacts for Machine Learning</b>	Denver, CO
SUPER COMPUTING	Nov. 2017

## Teaching Experience

---

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