# Research Infrastructures for Hardware Accelerators

# Synthesis Lectures on Computer Architecture

# Research Infrastructures for Hardware Accelerators

Yakun Sophia Shao and David Brooks

Harvard University

## ABSTRACT

Hardware acceleration in the form of customized datapath and control circuitry tuned to specific applications has gained popularity for its promise to utilize transistors more efficiently. Historically, the computer architecture community has focused on general-purpose processors, and extensive research infrastructure has been developed to support research efforts in this domain. Envisioning future computing systems with a diverse set of general-purpose cores and accelerators, computer architects must add accelerator-related research infrastructures to their toolboxes to explore future heterogeneous systems. This book serves as a primer for the field, as an overview of the vast literature on accelerator architectures and their design flows, and as a resource guidebook for researchers working in related areas.

## KEYWORDS

accelerators, specialized architecture, SoC, high-level synthesis, simulators, design space exploration, workload characterization, benchmarks

# Contents

# Preface

Specialized architectures have been a growing topic in both academic research and commercial development for the past decade. As traditional technology scaling slows, specialization becomes a viable solution for computer architects to continue performance growth and energy efficiency improvements without relying on technological advances.

This book aims to present a high-level overview of the state-of-the-art accelerator research in both industry and academia, with a special emphasis on research infrastructure available for accelerator-related research. This book begins by describing the technology trends that have led accelerator research to prominence. In Chapter 2, we present a taxonomy of accelerator research and practice, with the goal of introducing the reader to the flavor of accelerator designs that have been proposed in recent years. Chapter 3 presents the standard accelerator design flow from RTL generation, simulation, and synthesis. Recent advances in high-level synthesis (HLS) tools provide a promising path for accelerator development in the future, and we describe the capabilities of commercial tools like Xilinx's Vivado HLS and their limitations. Chapter 4 discusses pre-RTL modeling approaches to facilitate the rapid exploration of the design space of accelerators as well as the interaction between accelerators and the rest of the system. Chapter 5 focuses on workload characterization approaches in the context of accelerators and Chapter 6 discusses benchmarking. We end this book with a discussion on the challenges and opportunities of accelerator architectures and design tools in Chapter 7.

Yakun Sophia Shao and David Brooks
October 2015

# Acknowledgments