### Designing Programming Languages for Heterogeneous Hardware Adrian Sampson Cornell

"The complexity for minimum component costs has increased at a rate of roughly a factor of two per year."



## **Gordon Moore** co-founder of Intel 1965

## The size of a single transistor decreases by half every 18 months.



## **Gordon Moore** co-founder of Intel 1965

# A manual cost, and power... The size of a single transistor decreases by half every 18 months.



## **Gordon Moore** co-founder of Intel 1965



10 GHz



year of introduction

10 GHz



**CPU** core frequency













The performance returns from Moore's Law ended in 2015!

The only way forward is to trade off generality for efficiency!

A New Golden Age f A New Control Contr

John L. Hennessy and David A. Patterson







### A Reconfigurable Fubric for Accelerating Large-Scale Datacenter Services

Andrew Putnam Adrian M. Caulfield Eric S. Chung Derek Chiou<sup>1</sup> Kypros Constantinides<sup>2</sup> John Demme<sup>3</sup> Hadi Esmaeilzadeh<sup>4</sup> Jeremy Fowers Gopi Prashanth Gopal Jan Gray Michael Haselman Scott Hauck<sup>5</sup> Stephen Heil Amir Hormati<sup>6</sup> Joo-Young Kim Sitaram Lanka James Larus<sup>7</sup> Eric Peterson Simon Pope Aaron Smith Jason Thong Phillip Yi Xiao Doug Burger

#### Microsoft

Datacenter workloads demand high computational capabilities, flexibility, power efficiency, and low cost. It is challenging to improve all of these factors simultaneously. To advance datacenter capabilities beyond what commodity server designs can provide, we have designed and built a composable, reconfigurable fabric to accelerate portions of large-scale software services. Each instantiation of the fabric consists of a 6x8 2-D torus of high-end Stratix V FPGAs embedded into a half-rack of 48 machines. One FPGA is placed into each server, accessible through PCIe, and wired directly to other FPGAs with pairs of 10 Gb SAS cables.

Abstraction

In this paper, we describe a medium-scale deployment of this fabric on a bed of 1,632 servers, and measure its efficacy in accelerating the Bing web search engine. We describe the requirements and architecture of the system, detail the desirable to reduce management issues and to provide a consistest olatform and applications can rely on. Second, datacenter services evolve extremely rapidly, making non-programmable hardware features impractical. Thus, datacenter providers are faced with a conundrum: they need continued improvements in performance and efficiency, but cannot obtain those improvements from general-purpose systems.

2 AUTHORS!

Reconfigurable chips, such as Field Programmable Gate Arrays (FPGAs), offer the potential for flexible acceleration of many workloads. However, as of this writing, FPGAs have not been widely deployed as compute accelerators in either datacenter infrastructure or in client devices. One challenge traditionally associated with FPGAs is the need to fit the accelerated function into the available reconfigurable area. One could virtualize the FPGA by reconfiguring it at run-time to support more functions than could fit into a single device.

