# LUCA CARLONI

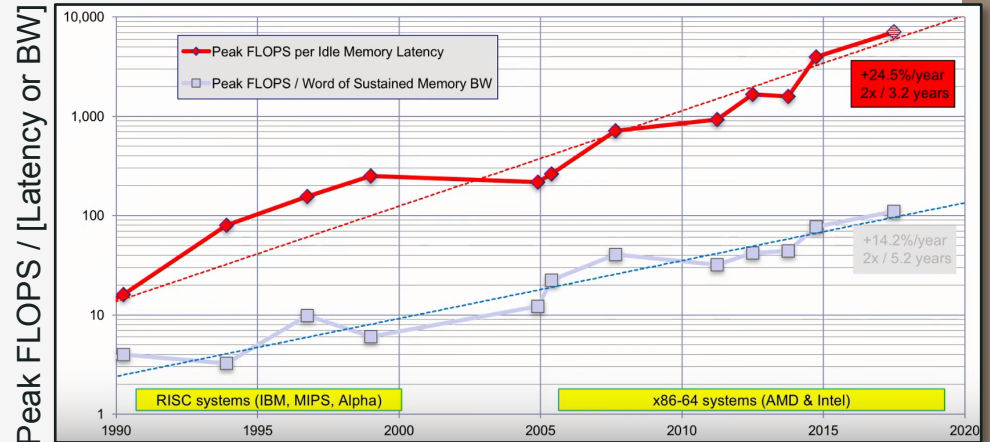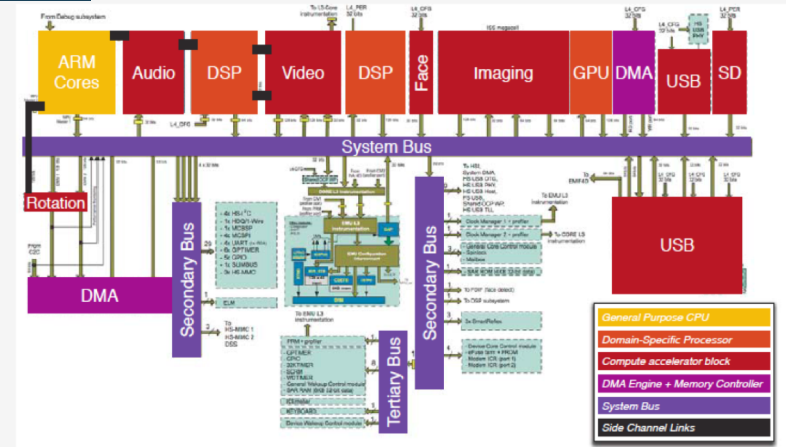DEPARTMENT OF COMPUTER SCIENCE
COLUMBIA UNIVERSITY

# THE DATA SUPPLY CHALLENGE

**Modern computer systems are increasingly heterogeneous**

- Accelerator-oriented parallelism to meet aggressive performance and power targets

**As accelerators have sped up compute portions, the main challenge is data supply**

- Key bottlenecks lie in memory and communication overheads associated with supplying specialized accelerators with data

- Different apps have distinct data supply needs



John McCalpin, SC'16 Keynote

# DECADES: A VERTICALLY-INTEGRATED APPROACH

## Language and Compiler Support

- Enhance data locality
- Optimize spatial mapping of threads
- Enable in-memory computing

## Very Coarse-Grained Reconfigurable Tile-Based Architecture

- Coarser than CGRA → VCGRTA
- 3 classes of reconfigurable tiles
- Reconfigurable interconnection network
- Reconfigurable in-memory computing

## Multi-Tiered Demonstration Strategy

- Scalable full-system simulation
- Multi-FPGA emulation infrastructure
- 225-tile DECADES chip prototype

# DECADES: A VERTICALLY-INTEGRATED APPROACH

## Language and Compiler Support
(M. Martonosi)

- Enhance data locality
- Optimize spatial mapping of threads
- Enable in-memory computing



## Very Coarse-Grained Reconfigurable Tile-Based Architecture
(L. Carloni)

- Coarser than CGRA → VCGRTA
- 3 classes of reconfigurable tiles
- Reconfigurable interconnection network
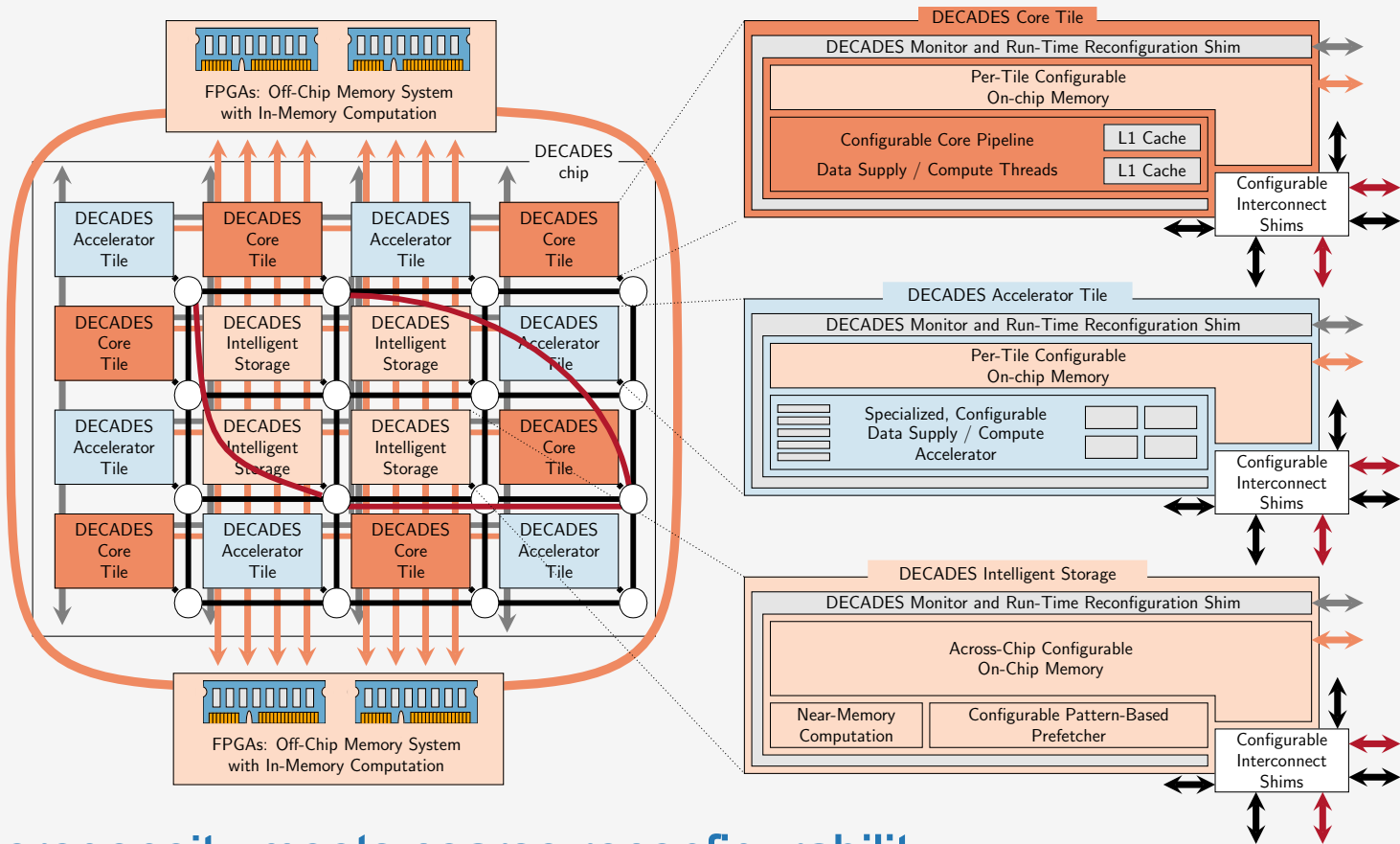- Reconfigurable in-memory computing



## Multi-Tiered Demonstration Strategy
(D. Wentzlaff)

- Scalable full-system simulation
- Multi-FPGA emulation infrastructure
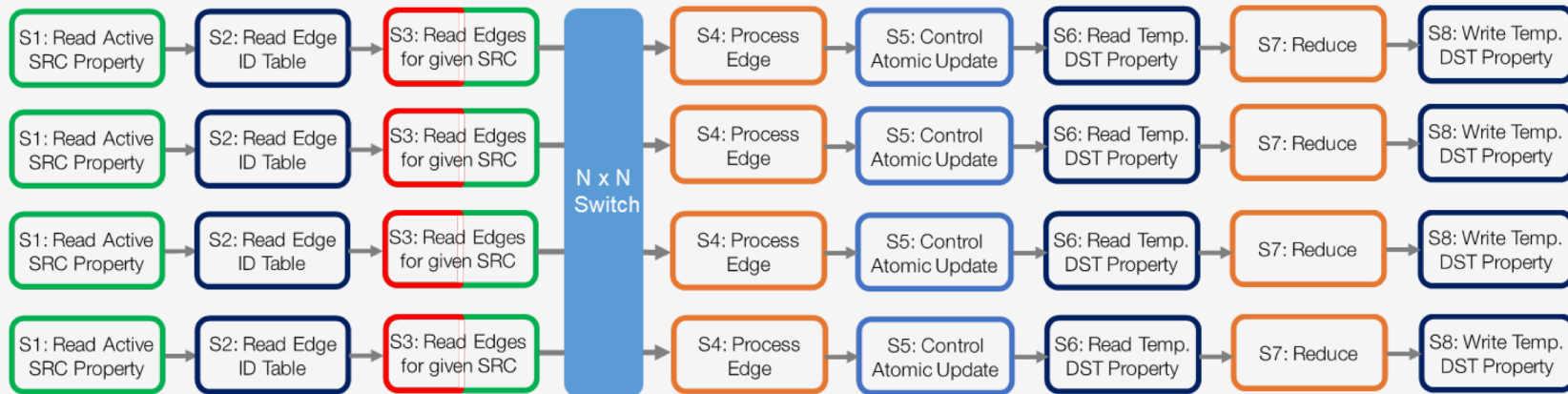- 225-tile DECADES chip prototype

# DECADES PLATFORM ARCHITECTURE

Heterogeneity meets coarse reconfigurability
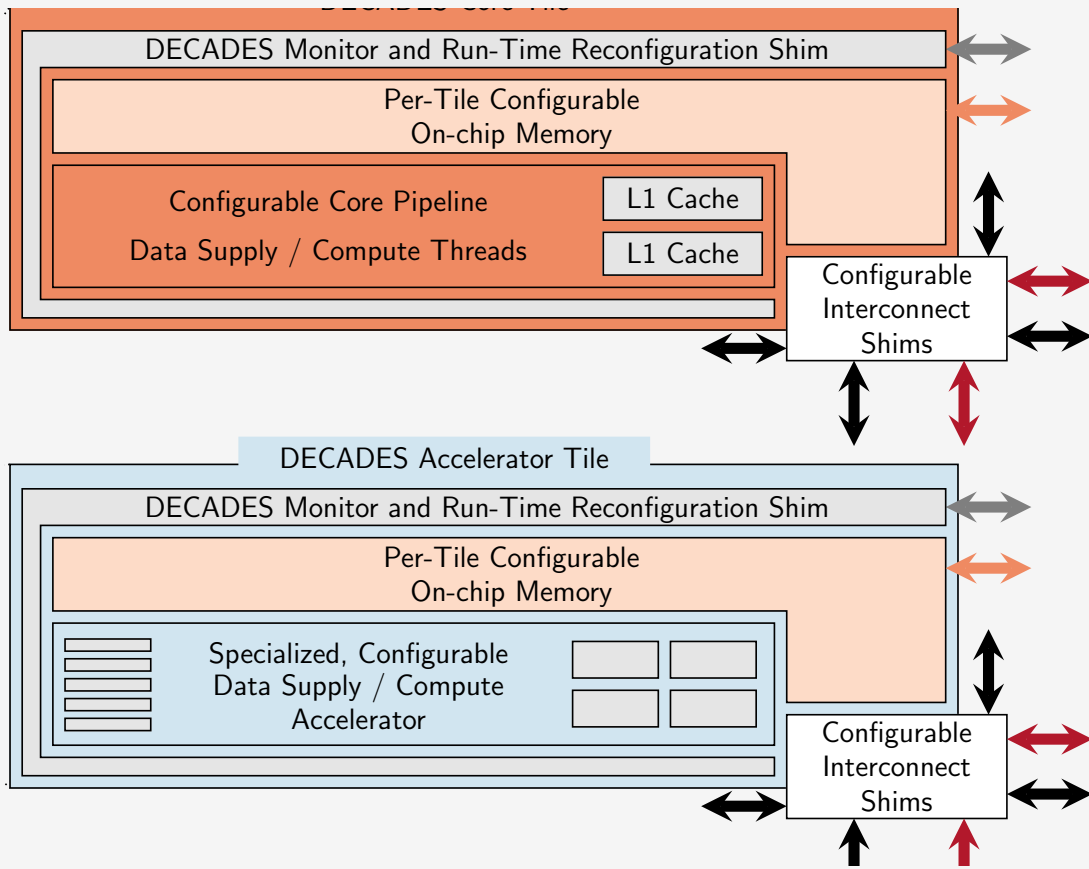
# PRIOR WORK: GRAPHICIONADO



- Application-Specific Memory Hierarchy for Bandwidth-Bound Graph Analytics
  - Customized memory hierarchy to minimize off-chip memory access traffic [3x reduction]
  - Ease the design/use of accelerators
  - Dataflow pipeline based on high-level abstraction eases the programming and enables hardware reuse for different graph applications
  - Specialized HW accelerator for graph analytics successfully achieves ~3x speedup and 50x+ energy saving compared to state-of-the-art software framework on 32-core CPU

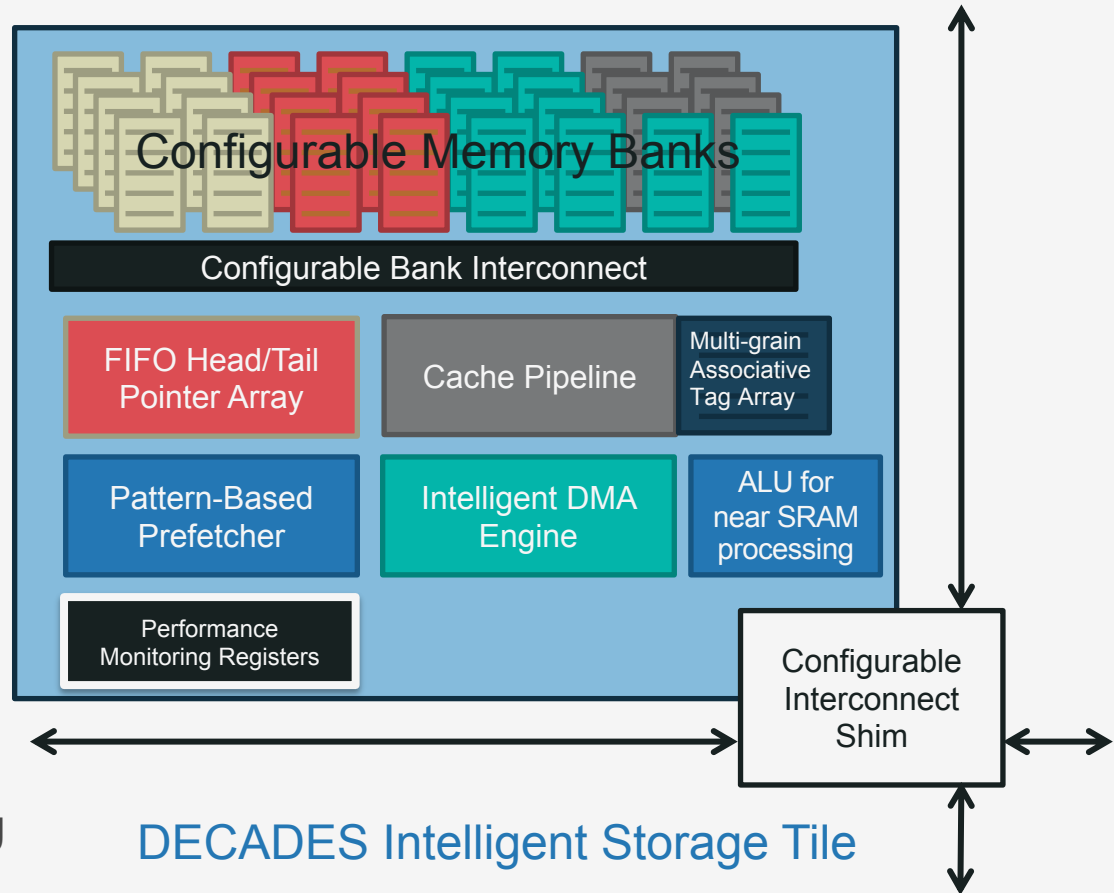[Ham et al., MICRO-49, 2016. IEEE Micro Top Picks Honorable Mention]

# DECADES CORE AND ACCELERATOR TILES

- Computations mapped onto core tiles or available accelerator tiles

  - Each tile is wrapped in monitor/reconfiguration shim

- Dynamic reconfiguration of Supply-Compute decoupling, power-performance tradeoffs, and interconnect
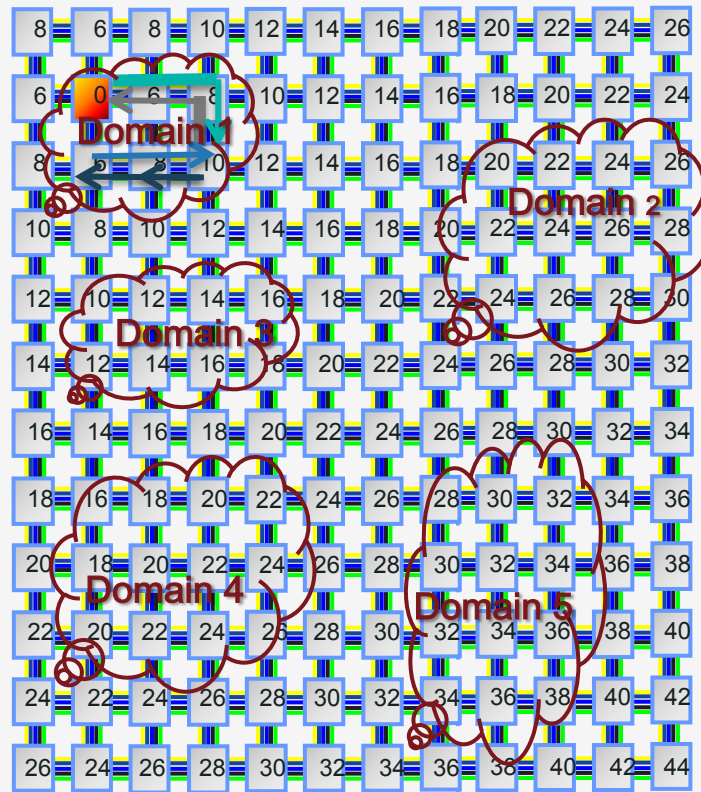
# INTELLIGENT DATA MANAGEMENT

- Specialization #1: Map apps onto mix of compute tiles and intelligent storage (IS) tiles

- Specialization #2: Select and configure appropriate storage features within IS

  - Configurable memory banks + address and prefetching features

  - Simple near-SRAM ALU



Configurable Memory Banks

Configurable Bank Interconnect

FIFO Head/Tail Pointer Array

Cache Pipeline

Multi-grain Associative Tag Array

Pattern-Based Prefetcher

Intelligent DMA Engine

ALU for near SRAM processing

Performance Monitoring Registers

Configurable Interconnect Shim

DECADES Intelligent Storage Tile

- Flexible memory system on top of cache coherent system
  - Enables the exact minimal communication needed
  - Build incoherent coherent domains
- Restriction on application- or page-level
  - Improves performance
    - Shorter network on-chip distances
    - Less interfering memory coherence traffic
  - Reduces energy
    - Fewer on-chip network links need to be transited
    - Less area dedicated to tracking cache line sharers
  - Reduces area
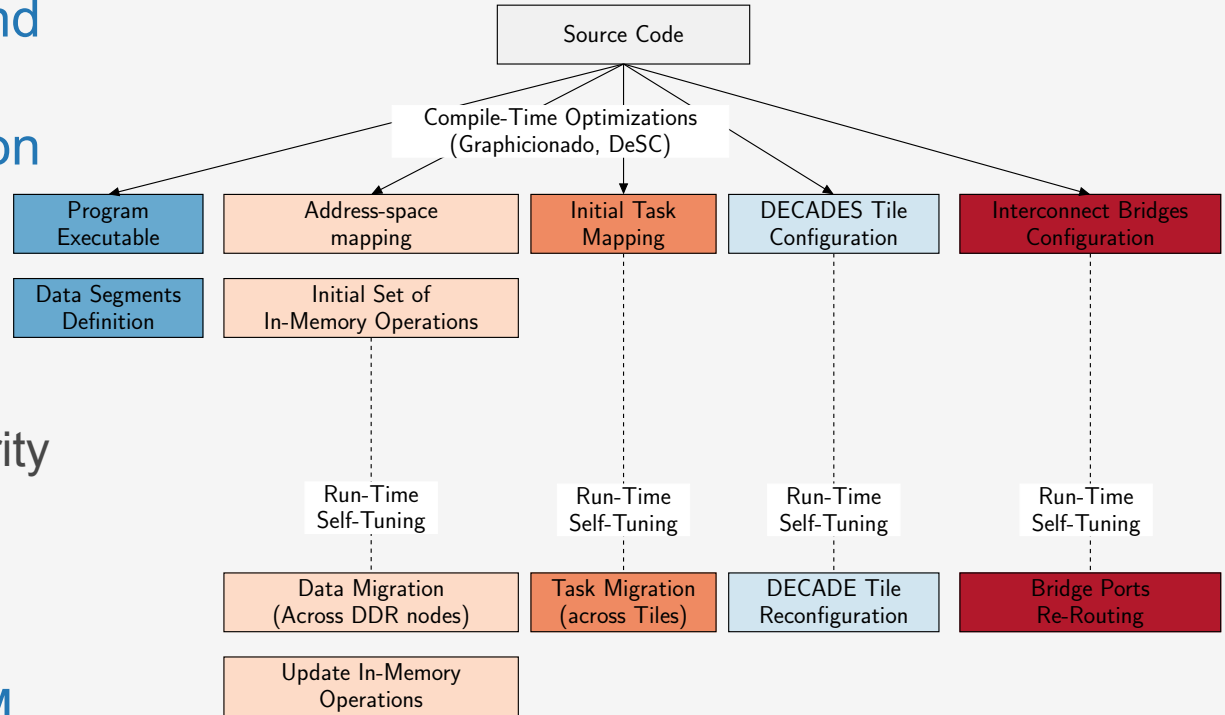    - Track fewer sharers on large configurations



[Fu et al, MICRO 2015]

- Flexible memory system on top of cache coherent system
  - Enables the exact minimal communication needed
  - Build incoherent coherent domains
- Restriction on application- or page-level
  - Improves performance
    - Shorter network on-chip distances
    - Less interfering memory coherence traffic
  - Reduces energy
    - Fewer on-chip network links need to be transited
    - Less area dedicated to tracking cache line sharers
  - Reduces area
    - Track fewer sharers on large configurations
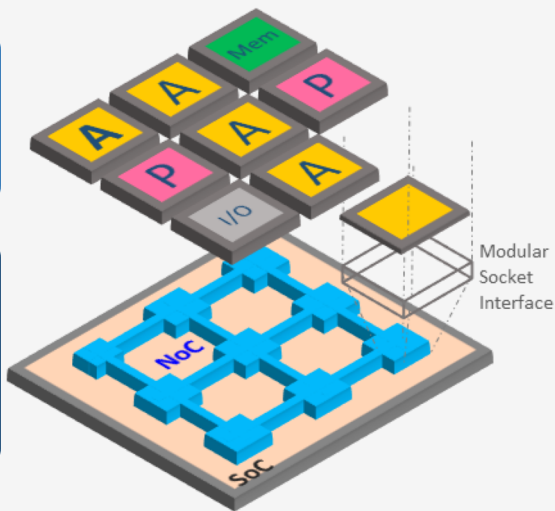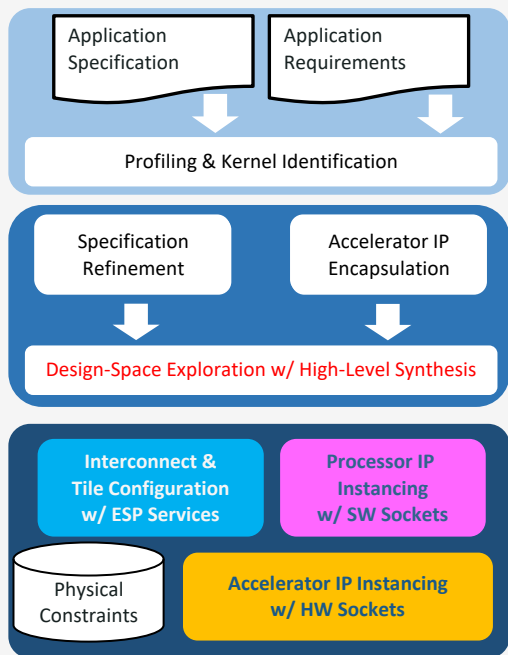


[Fu et al, MICRO 2015]

# LANGUAGE, COMPILER & RUNTIME SYSTEM

- **Compiler Analysis and Support for memory hierarch specialization**
  - Bandwidth optimizations through cache optimizations and locality/granularity tailoring
  - Latency tolerance through decoupling
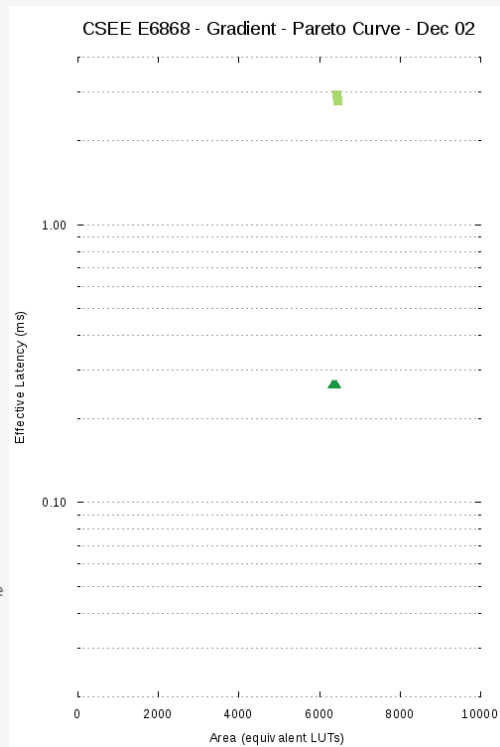- **Build on DeSC LLVM compiler infrastructure**



[Ham/Aragon/Martonosi, MICRO-48, 2015]

# PRIOR WORK: EMBEDDED SCALABLE PLATFORMS

- Flexible Tile-Based Architecture

- System-Level Design Methodology

- SoC Design Productivity



Application Specification
Application Requirements

Profiling & Kernel Identification

Specification Refinement
Accelerator IP Encapsulation

Design-Space Exploration w/ High-Level Synthesis

Interconnect & Tile Configuration w/ ESP Services
Processor IP Instancing w/ SW Sockets

Physical Constraints
Accelerator IP Instancing w/ HW Sockets

CSEE E6868 - Gradient - Pareto Curve - Dec 02

Effective Latency (ms)

Area (equivalent LUTs)

In the span of 1 month:

21: student teams
661: improved designs
32: avg. number of improved designs per team
1.5: avg. number of new designs per team/day
99: Pareto curve changes
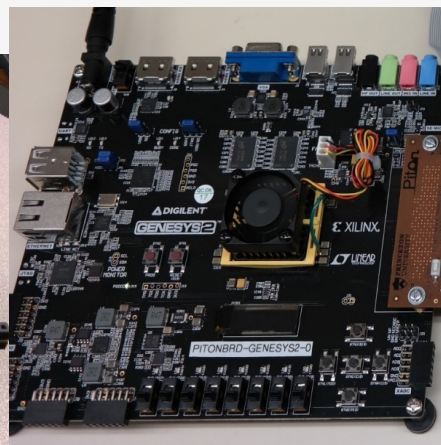11: final number of Pareto-optimal designs
26x: Performance range
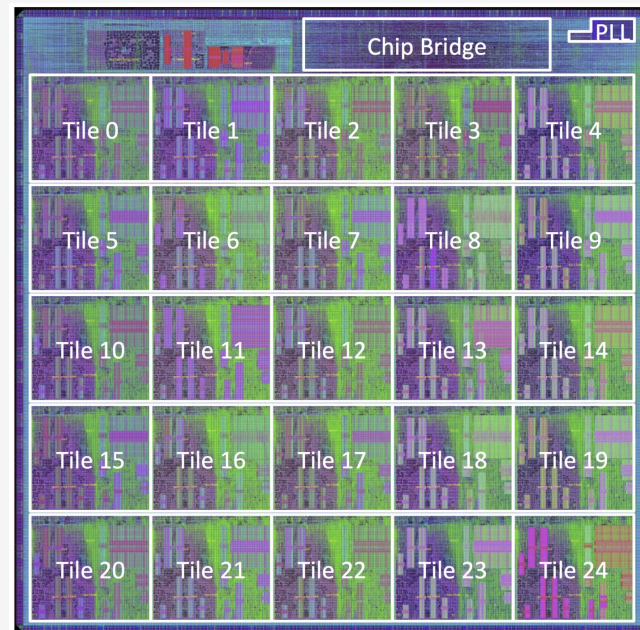10x: Area range

[Carloni, DAC 2016]

# EMULATION & PROTOTYPING

- Take DECADES architecture to FPGA

- Continued design refinement throughout program

- Prototype chip to de-risk architecture



- Multi-FPGA emulation infrastructure

- Recent 25-core manycore system built by our team

# SUMMARY & IMPACT

## Language/Compiler/Runtime:

- Latency: >4X per thread performance benefits from memory data supply decoupling

- Bandwidth: Granularity management and Multiplicative outer-loop parallelism up to bandwidth limit

- Total of 50X over single-thread from software

## Configurable Hardware Platform:

- Hardware speedups from accelerators for address calculation, memory fetch, or compute

- Fine-grained, low-overhead measurements drive adaptation and module depowering

- 10-20X multiplicative power/performance benefits
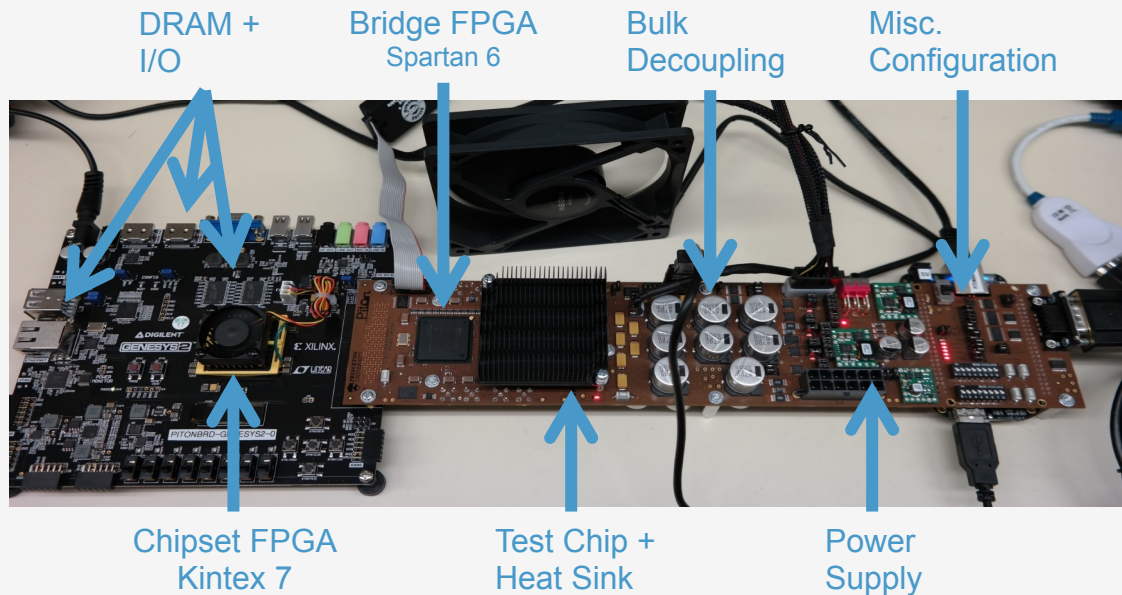
# TECHNOLOGY TRANSFER PLANS

- Outputs:
  - Software ecosystem
  - Chip design
  - FPGA emulation system

- Technology transfer plans:
  - Release of software, hardware, and data where possible
  - Commercialization and licensing
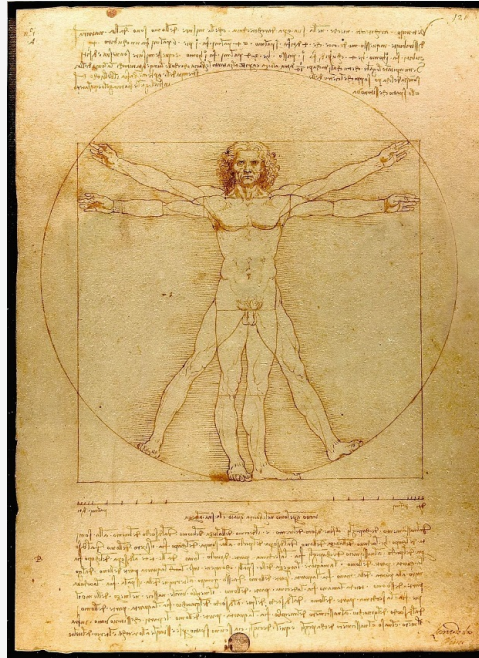
- Leverage extensive past experience:
  - Widely-used open-source software (Wattch, *Check tools, scalable QEMU)
  - Patents licensed to major companies (Power-efficient ALUs)
  - Technology transferred from academia to startups (Tilera)
  - Open-Source Hardware (OpenPiton)



DRAM + I/O

Bridge FPGA
Spartan 6

Bulk Decoupling

Misc. Configuration

Chipset FPGA
Kintex 7

Test Chip + Heat Sink

Power Supply

# TOWARDS A COMPUTER DESIGN RENAISSANCE

- The end of silicon dimensional scaling and
  the rise of heterogeneous reconfigurable computing bring
  an opportunity for a Computer Design Renaissance

- …by supporting the creativity of application developers to realize innovative architectures, chips, systems and products



- … through richly reconfigurable substrates and intelligent compilation and mapping