# **Rigel**:
# An Architecture and Scalable Programming Interface for a 1000-core Accelerator

**John H. Kelm**, Daniel R. Johnson,

Matthew R. Johnson, Neal C. Crago, William Tuohy,
Aqeel Mahesri[*], Steven S. Lumetta,
Matthew I. Frank[†], Sanjay J. Patel

*The author is now with NVIDIA.
† The author is now with Intel.

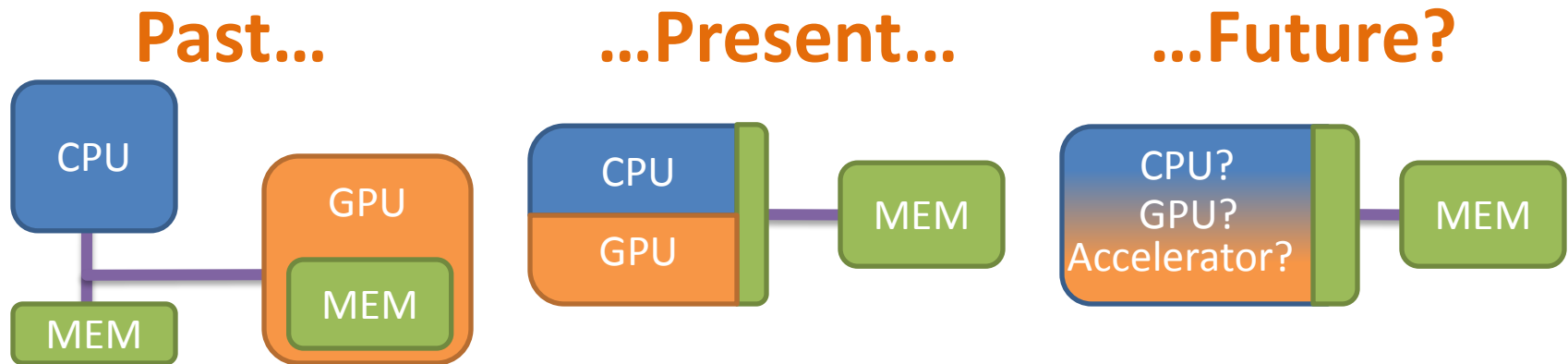# Accelerated Computing: **Today**

> **Programmable accelerator:** HW entity designed to provide advantages for a class of apps including: higher performance, lower power, or lower unit cost relative to a general-purpose CPU.

- Contemporary Accelerators: GPUs, Cell, Larrabee
- **Challenges**:
    1. Inflexible programming models
    2. Lack of conventional memory model
    3. Hard to scale irregular parallel apps

**Effect on Development**: Unattractive time to solution

John H. Kelm

2

# Accelerated Computing: **Tomorrow**

- Why research accelerators?
  - Insight into future general-purpose CMPs
  - **Challenges**: Performance vs. programmer effort
- **Accelerator Trend**: Integration over time

**Past…**          **…Present…**          **…Future?**

# Accelerated Computing: **Metrics**

**Challenges lead to:**

- **FLOPS/$ (area)**
- **FLOPS/Watt (power)**
- **FLOPS/Programmer Effort**

- Enable new platforms
- Open new markets
- Enable new apps

# **Context**: Project Orion

## Applications, Programming Environments, and Architecture for 1000-core Parallelism

**1000-core Architecture**

**Programming Environments**

Illinois Image Formation and Processing
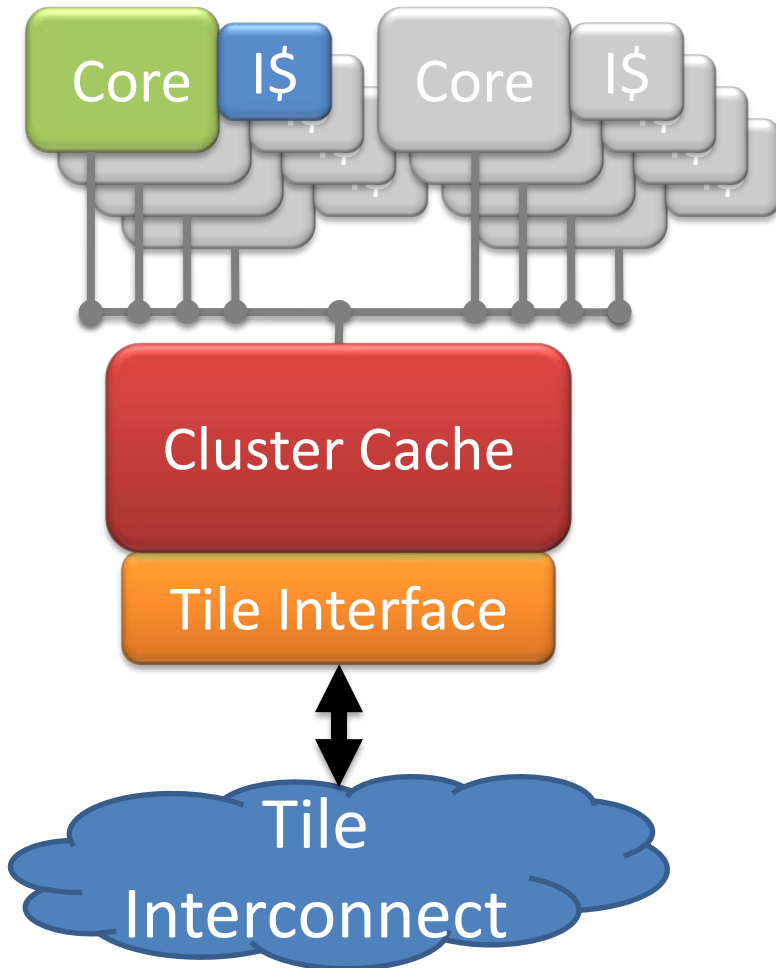
**Applications**

John H. Kelm

# Rigel Design Goals

- **What**: Future programming models
  - Apps and models may not exist yet
  - We have ideas (**visual computing**), but who knows?
  - Flexible design → easier to retarget
- **How:** Focus on scalability, programmer effort
  - Room to play: Raised HW/SW interface
  - Focusing design effort: **Five Elements**

John H. Kelm

# Outline

- Motivation

- Rigel architecture

- **Elements** in context of Rigel architecture

- Evaluation:
  - Area and power
  - Scalability
  - SW Task management

- Future work and conclusions
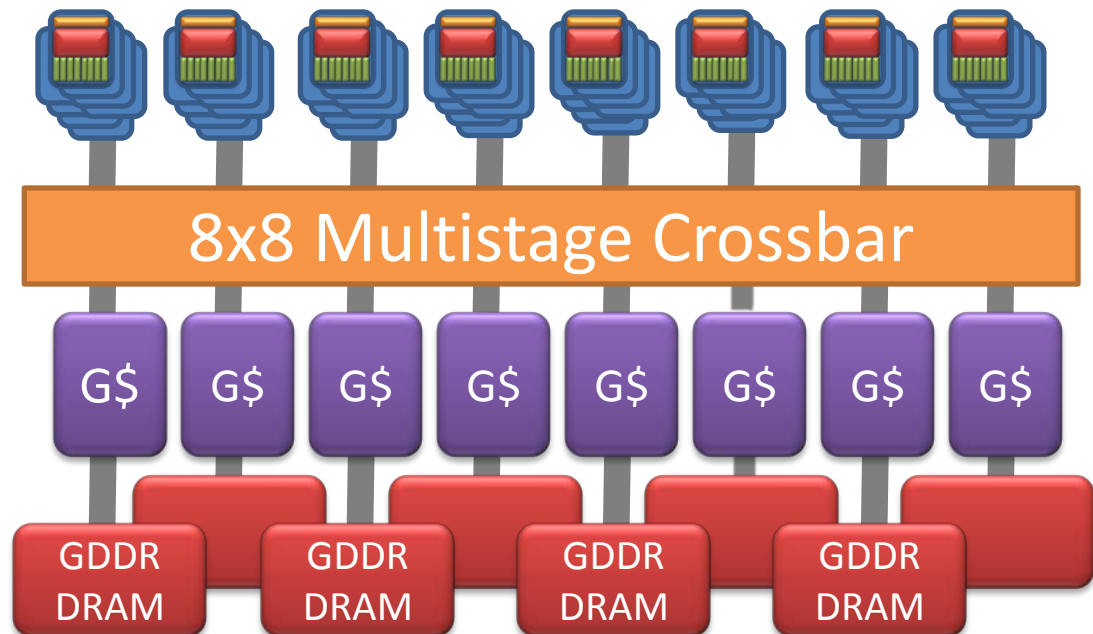
# Rigel Architecture: **Cluster View**

Core  I$  Core  I$

Cluster Cache

Tile Interface

Tile Interconnect

- Basic building block
- Eight 32b RISC cores
- Per-core SP FPUs
- 64 kB shared cache
- Cache line buffer

Rigel

ILLINOIS
UNIVERSITY OF ILLINOIS AT URBANA-CHAMPAIGN

John H. Kelm

8

# Rigel Architecture: Full Chip View

- Cluster caches not HW coherent (8 MB total)
- G$ fronts mem. controllers (4 MB total)
- Uniform cache access

**8 Tiles Per Chip**
**16 Clusters per Tile**

**Global Interconnect**

**Global Cache Banks**

**8 GDDR Channels**

8x8 Multistage Crossbar

G$ G$ G$ G$ G$ G$ G$ G$

GDDR DRAM  GDDR DRAM  GDDR DRAM  GDDR DRAM

John H. Kelm

UNIVERSITY OF ILLINOIS AT URBANA-CHAMPAIGN

# Design Elements

- **Challenges** in accelerator computing
- **FLOPS/dev. effort** → Difficult to quantify
- Guiding our 1000-core architecture
- Room to Play: Raising the HW/SW interface

# Design Elements

1. **Execution Model**: ISA, SIMD vs. MIMD, VLIW vs. OoOE, MT
2. **Memory Model**: Caches vs. scratchpad, ordering, coherence
3. **Work Distribution**:  Scheduling, spectrum of SW/HW choices
4. **Synchronization**:  Scalability, influence on prog. model
5. **Locality Management**
   – Moving data costs perf. and power
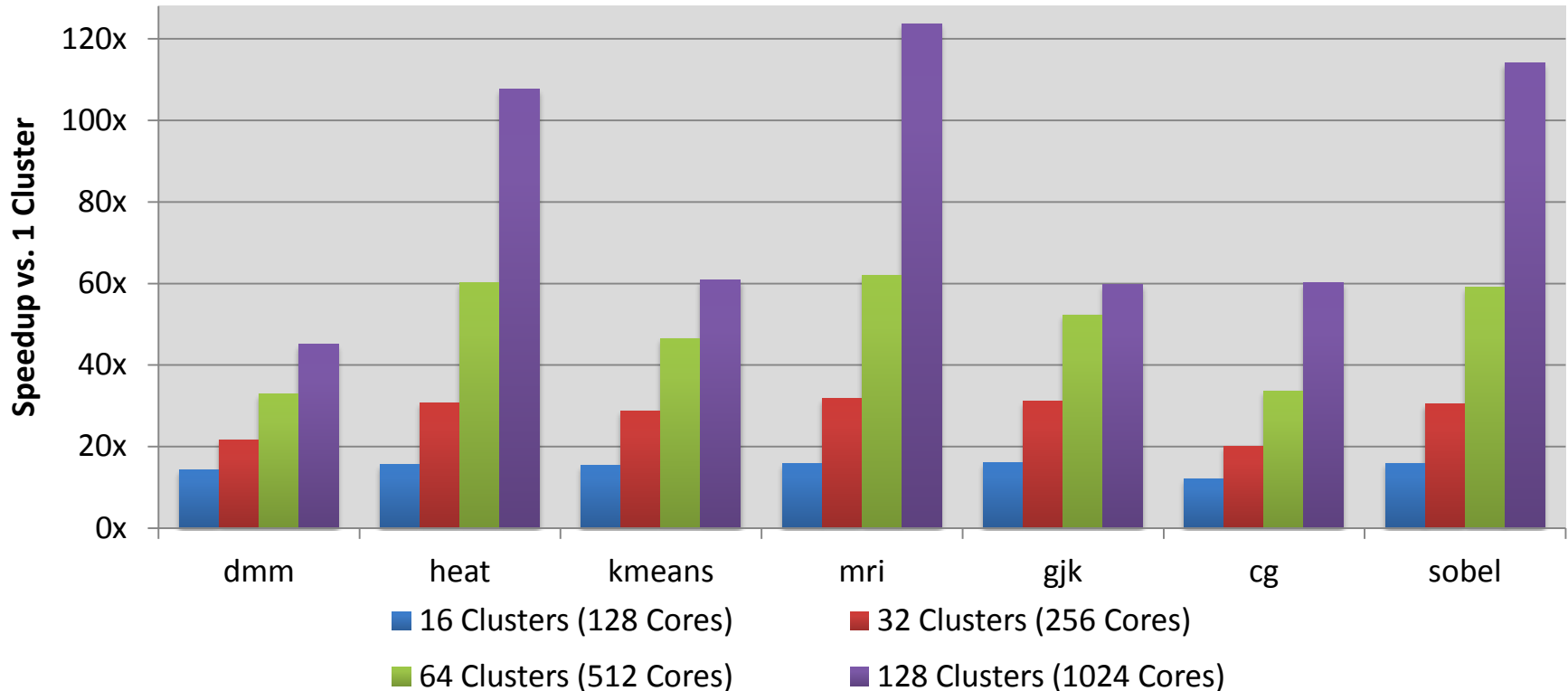   – Balance: dev. effort, compiler, runtime, HW

John H. Kelm

11

# **Element 1**: Execution Model

- **Tradeoff 1**: **MIMD** vs. SIMD [Mahesri MICRO'08]
  - Irregular data parallelism
  - Task parallelism
- **Tradeoff 2:** Latency vs. **Throughput** [Azizi DasCMP'08]
  - Simple in-order cores
- **Tradeoff 3:** **Full RISC ISA** vs. Specialized Cores
  - Complete ISA → conventional code generation
  - Wide range of apps
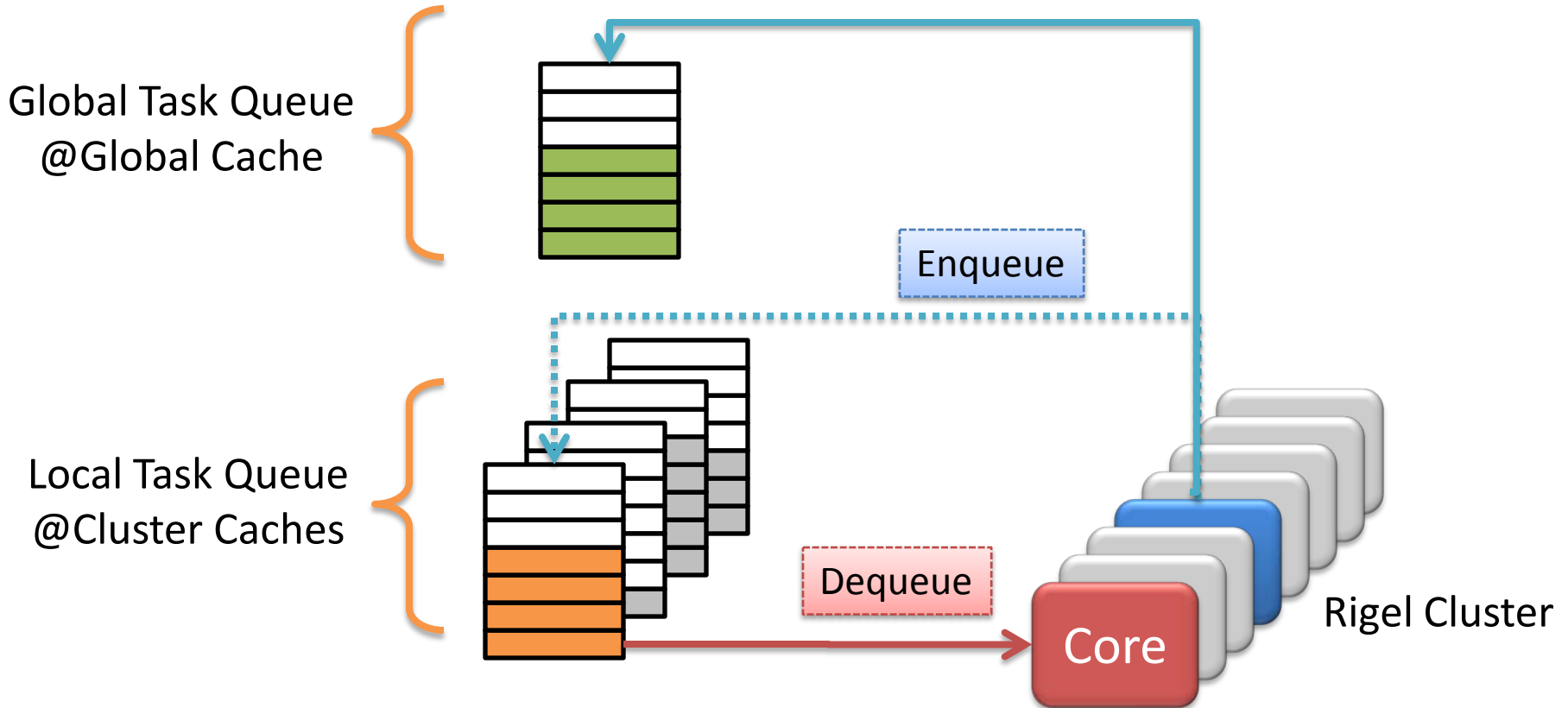
# Element 2: Memory Model

- **Tradeoff 1: Single vs. multiple address space**
- **Tradeoff 2: Hardware caches vs. scratchpads**
  - Hardware exploits locality
  - Software manages global sharing
- **Tradeoff 3: Hierarchical vs. Distributed (NUCA)**
  - Cluster cache/global cache hierarchy
  - ISA provides local/global mem. Operations
  - Non-uniformity → Programmer effort

John H. Kelm

# Some Results: **Scalability**



- Based on cycle-accurate, execution-driven simulation
- Library and run-time system code simulated
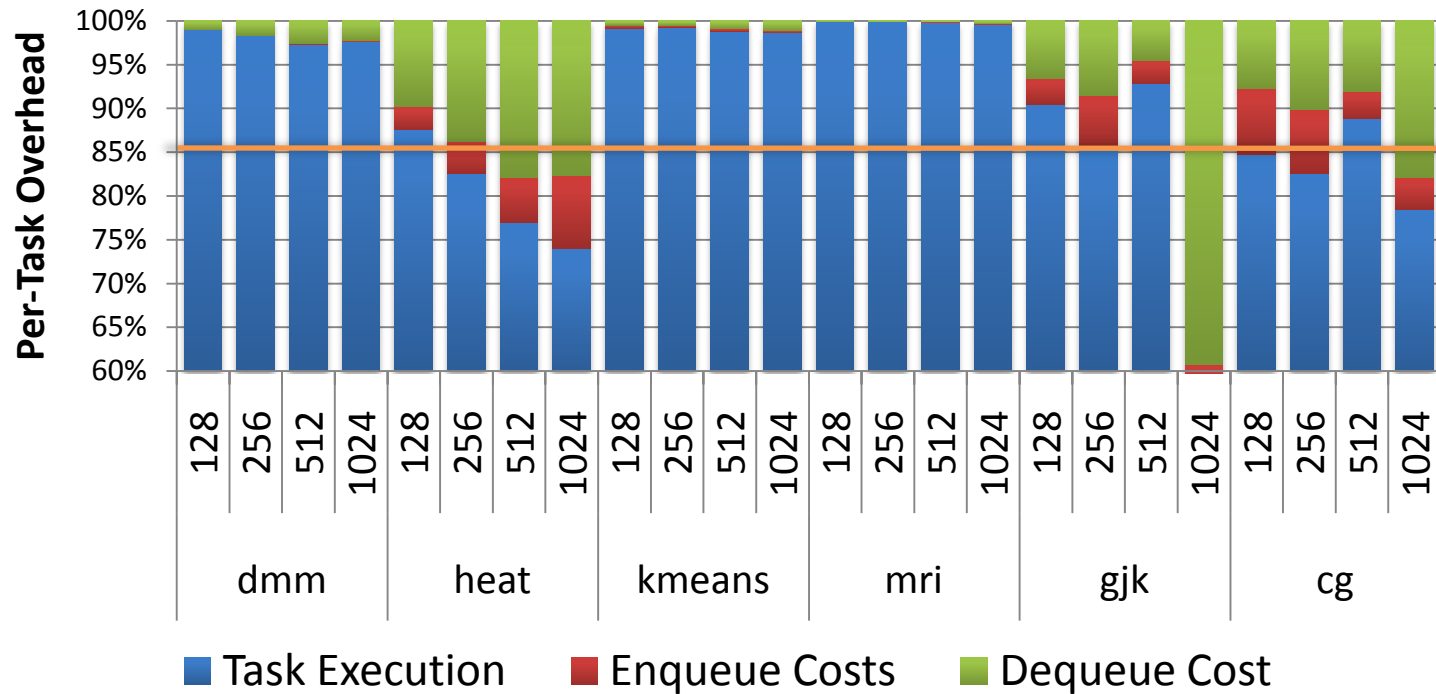- Regular C code + parallel library, standard C compiler

John H. Kelm

14

# Element 3: Work Distribution

Global Task Queue
@Global Cache

Enqueue

Local Task Queue
@Cluster Caches

Dequeue

Core

Rigel Cluster

- **Tradeoff (Spectrum):** HW vs. **SW Implementation**
  - SW task management: Hierarchical queues
  - Flexible policies + little specialized HW
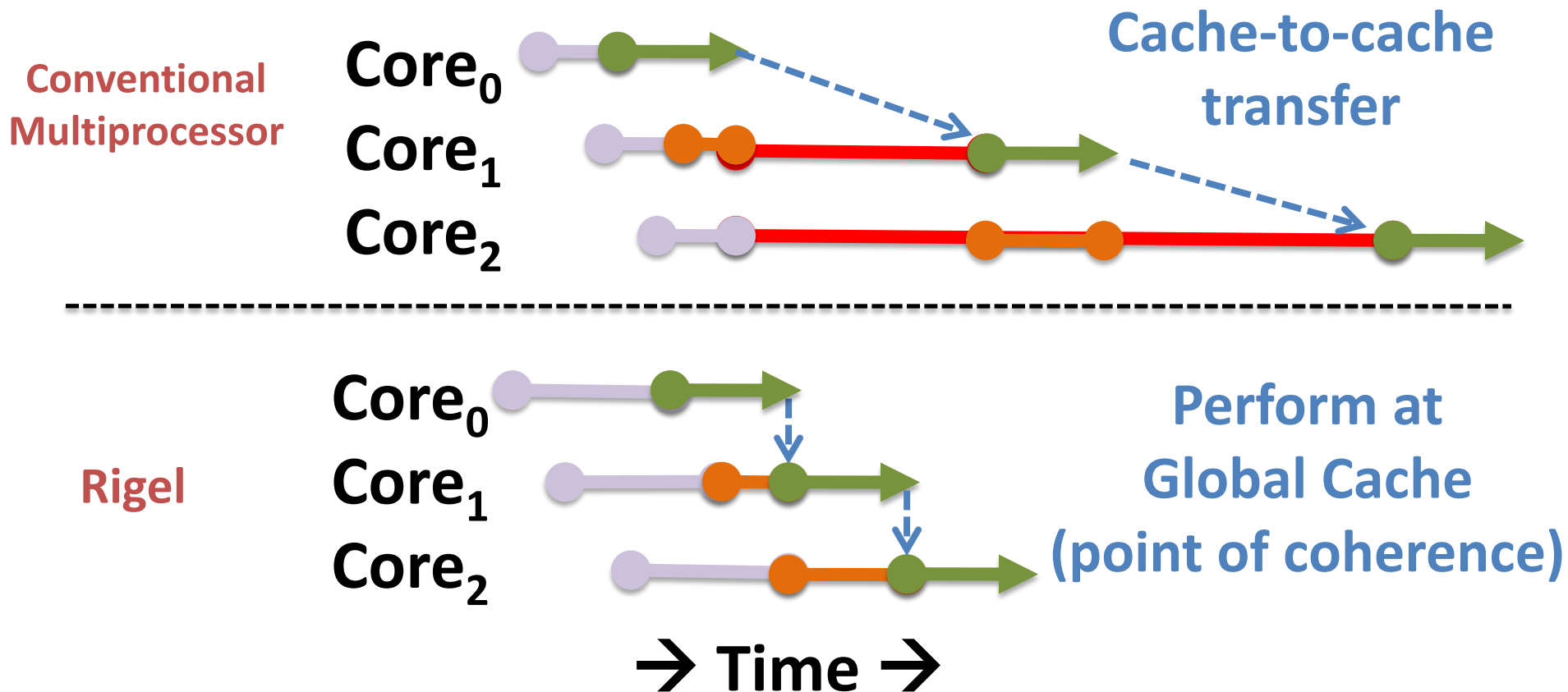
# **Work Distribution:** Rigel Task Model



- < 5% overhead for most data-parallel workloads
- < 15% for most irregular data-parallel workloads
- Task lengths: 100's-100k instructions
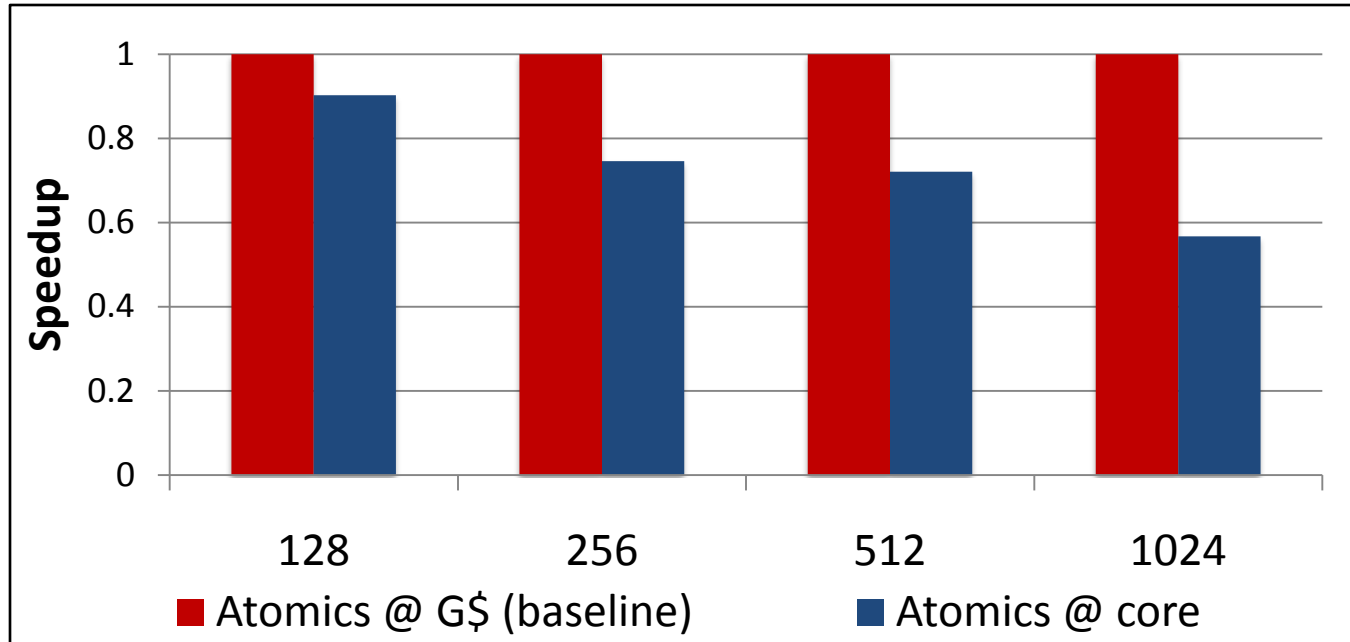
# Element 4: Synchronization

- Uses of coherence mechanisms:

  1. Control synchronization

  2. Data sharing

- Broadcast update

  – Use cases: flags and barriers

  – Reduce contention from polling

  – **Case Study**: 2x speedup for conjugate gradient (CG)

- Atomic primitives (example)

# Element 4: Atomic Primitives



**Conventional Multiprocessor**

Core$_0$
Core$_1$
Core$_2$

**Cache-to-cache transfer**

**Rigel**

Core$_0$
Core$_1$
Core$_2$

**Perform at Global Cache (point of coherence)**

→ Time →

**1. Network Latency**  **2. Overlapped Latency**  **3. Exposed Latency**  **4. Operation Execution (atom.inc)**
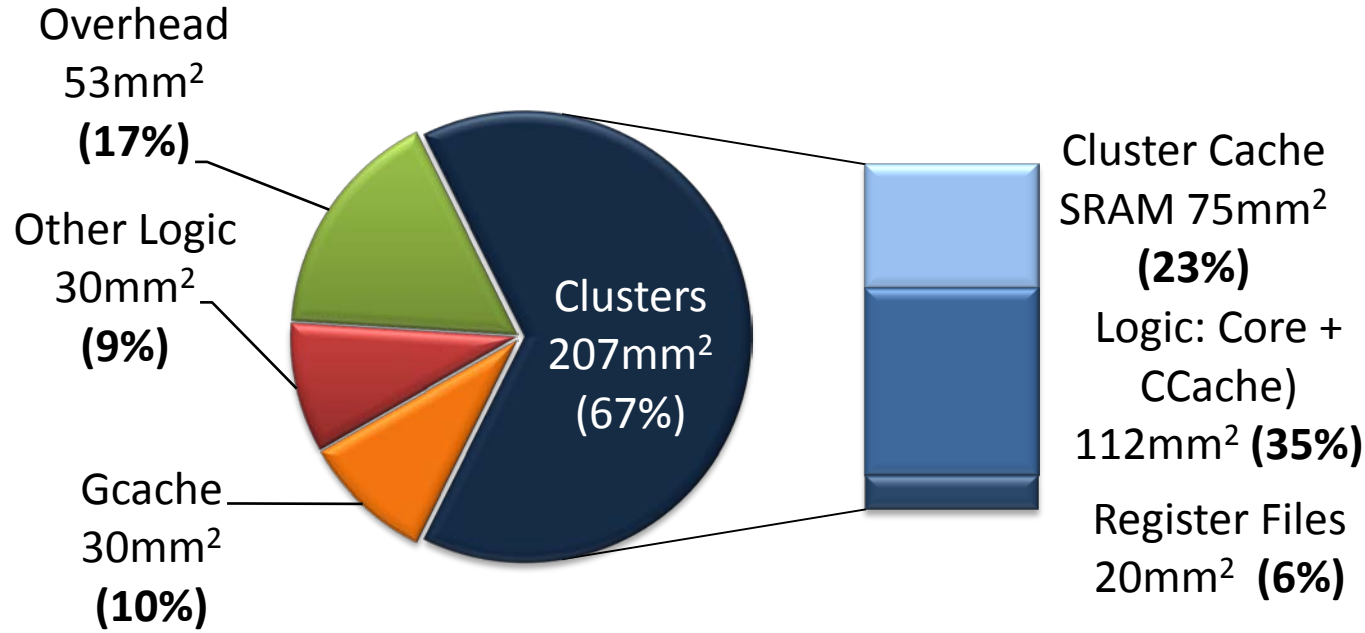
John H. Kelm

18

# Evaluation: Atomic Operations



**K-means Clustering**

- Need global histogramming
- With G$ atomics → Pipelined in network
- Without atomics → Exposed transfer latency

# So, Can We Build It?

Overhead
53mm² **(17%)**

Other Logic
30mm² **(9%)**

Gcache
30mm² **(10%)**

Clusters
207mm²
(67%)

Cluster Cache
SRAM 75mm² **(23%)**

Logic: Core + CCache)
112mm² **(35%)**

Register Files
20mm² **(6%)**

- RTL synthesis results + memory compiler + datasheets
- Targeting 45nm process @ 1.2 GHz
- 320 mm² total die area, <100W average power
- Estimate FLOPS/W and FLOPS/mm² **match or exceed GPUs**

John H. Kelm

20

# Current and Future Work

- RTL implementation
- Coherence and memory model [Kelm et al. PACT'09]
- Other programming models
- Multi-threading (1-4 threads)
- Element Five: Locality Management

# Conclusions

- **FLOPS/Dev. Effort** →Elements can drive design
- Software coherence viable approach
- Task management requires little HW
- 1000-core accelerator is feasible
  - **Area/performance**: 8 GFLOPS/mm² @ ~100W
  - **Programmability**: Task API + MIMD execution