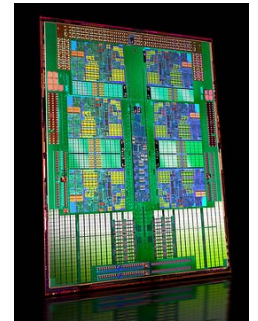# COHESION: A Hybrid Memory Model for Accelerators

**John H. Kelm**, Daniel R. Johnson, William Tuohy, Steven S. Lumetta and Sanjay J. Patel

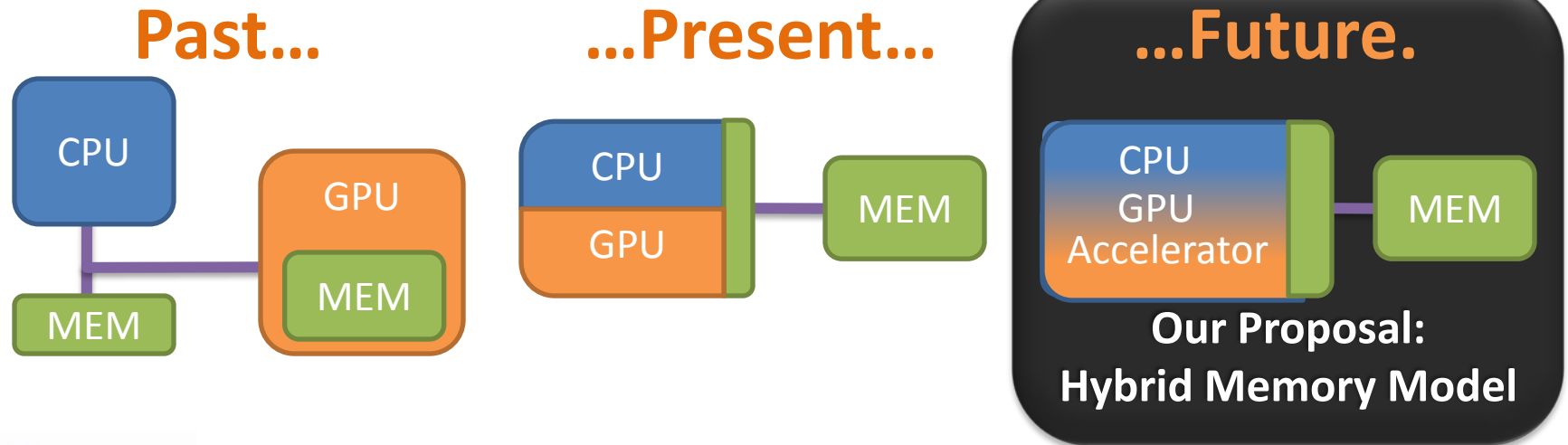University of Illinois at Urbana-Champaign

# Chip Multiprocessors Today

- General-purpose + accelerators (e.g., GPUs)

- General-purpose CMP Challenges:

  1. Programmability
  2. Power/perf density of ILP-centric cores
  3. **Scalability of HW coherence, strict memory models**

- Accelerator Challenges:

  1. Inflexible programming/execution models
  2. Hard to scale irregular parallel apps
  3. **Lack of conventional memory model**

Rigel

ILLINOIS
UNIVERSITY OF ILLINOIS AT URBANA-CHAMPAIGN

John H. Kelm

2

# Chip Multiprocessors Tomorrow

- Industry Trend: Integration over time
- Hybrids: Accelerators + CPUs together on die
- More core/compute heterogeneity but…

  …more homogeneity in memory model



**Past…**

**…Present…**

**…Future.**

CPU

GPU

MEM

MEM

CPU

GPU

MEM

CPU
GPU
Accelerator

MEM

**Our Proposal:
Hybrid Memory Model**

John H. Kelm

# CMP Memory Model Choices

## Conventional Multicore CPU

- Ex: Intel i7, Sun Niagara
- Optimized for:
  - Minimal latency
  - Tightly coupled sharing
  - Fine-grained synchronization
  - Minimal programmer effort
- Provides:
  - **Single address space**
  - **Hardware caching**
  - Strong ordering
  - HW-managed coherence

## Contemporary Accelerator

- Ex: NVIDIA GPU, IBM Cell
- Optimized for:
  - **Maximum throughput**
  - **Loosely coupled sharing**
  - **Coarse-grained synchronization**
  - Short silicon design cycle
- Provides:
  - Multiple address spaces
  - Scratchpad memories
  - Relaxed ordering
  - SW-managed coherence

Rigel · ILLINOIS UNIVERSITY OF ILLINOIS AT URBANA-CHAMPAIGN

**John H. Kelm**

4

# Roadmap

- Motivation and context
- Problem statement
- COHESION design
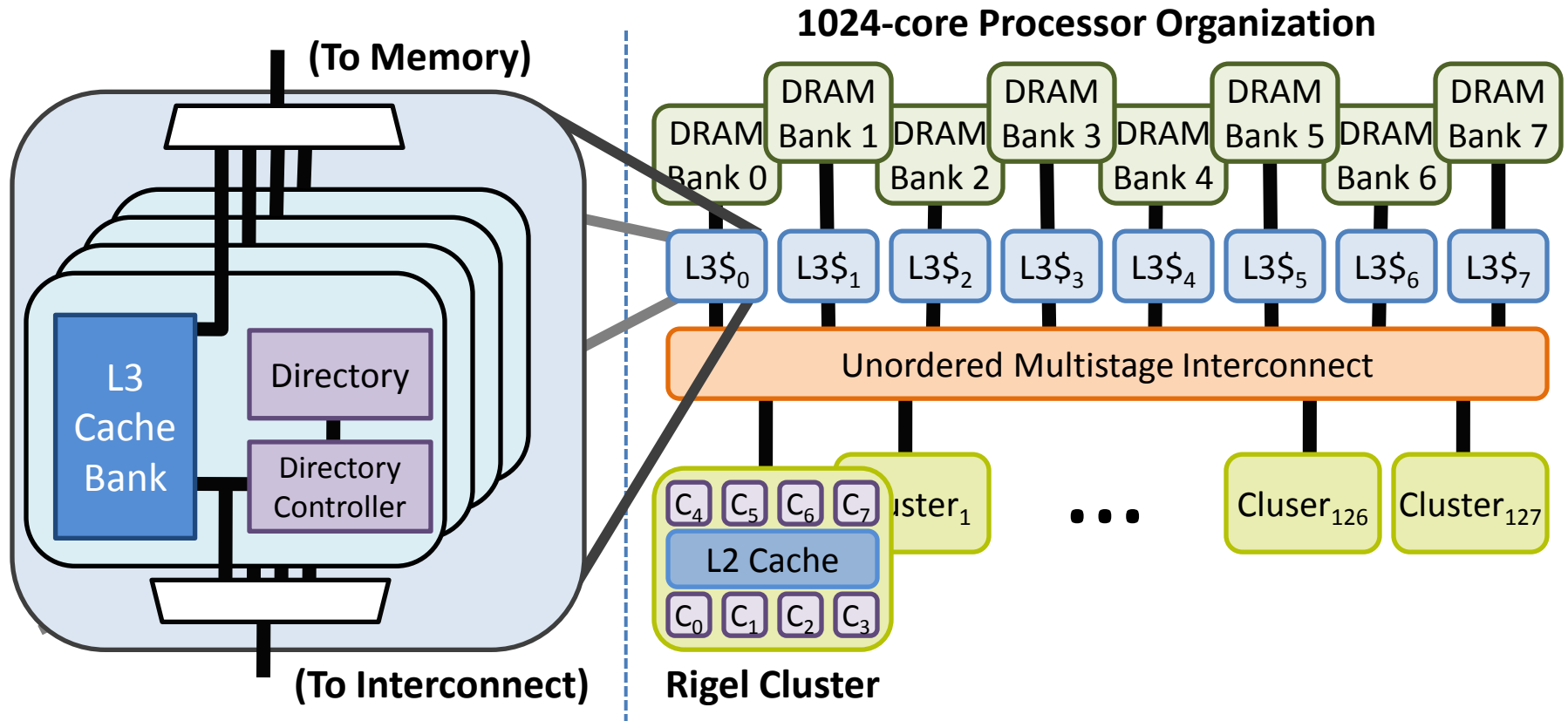- Use cases and programming examples

Addressed in this talk:

1. **Opportunity**: Is combining protocols worthwhile?
2. **Feasibility**: How does one implement hybrid memory models?
3. **Tradeoffs**: What are the tradeoffs in $HW_{cc}$ v. $SW_{cc}$?
4. **Benefit**: What does hybrid coherence get you?

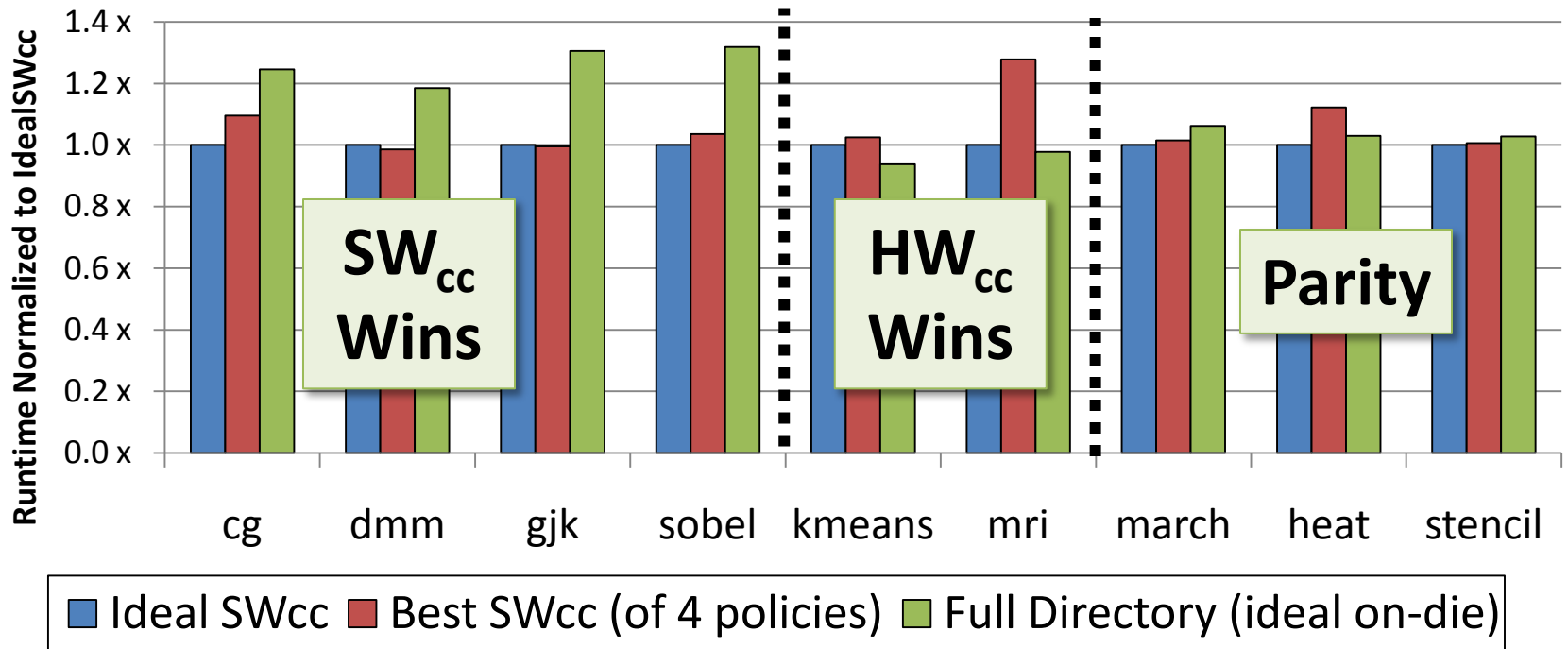# Problem: Scalable Coherence

- Available architectures:
  - **Accelerators**: 100s of cores, TFLOPS, no coherence
  - **CMPs**: <10s of cores, GFLOPS, HW coherence
  - Multiple memory models on-die
- What devs want in heterogeneous CMPs:
  - Hardware caches (locality)
  - Single address space (marshalling)
  - Minimal changes to current practices
- Accelerator scalability w/CMP memory model
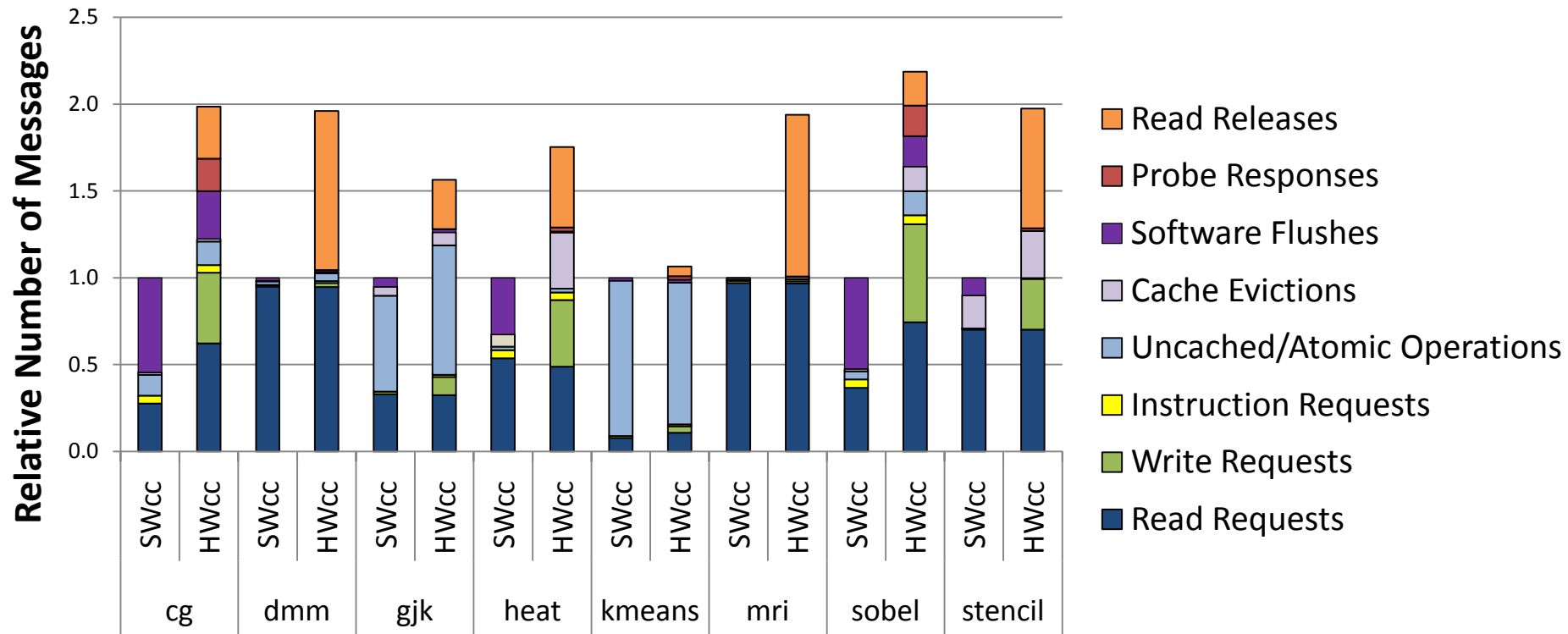
# Baseline Architecture



**1024-core Processor Organization**

(To Memory)

L3 Cache Bank | Directory | Directory Controller

(To Interconnect)

Rigel Cluster

DRAM Bank 0, DRAM Bank 1, DRAM Bank 2, DRAM Bank 3, DRAM Bank 4, DRAM Bank 5, DRAM Bank 6, DRAM Bank 7

$L3\$_0$, $L3\$_1$, $L3\$_2$, $L3\$_3$, $L3\$_4$, $L3\$_5$, $L3\$_6$, $L3\$_7$

Unordered Multistage Interconnect

$C_4$ $C_5$ $C_6$ $C_7$ / L2 Cache / $C_0$ $C_1$ $C_2$ $C_3$ — $Cluster_1$ ... $Cluser_{126}$ $Cluster_{127}$

- Variant of the Rigel Architecture [Kelm et al. ISCA'09]
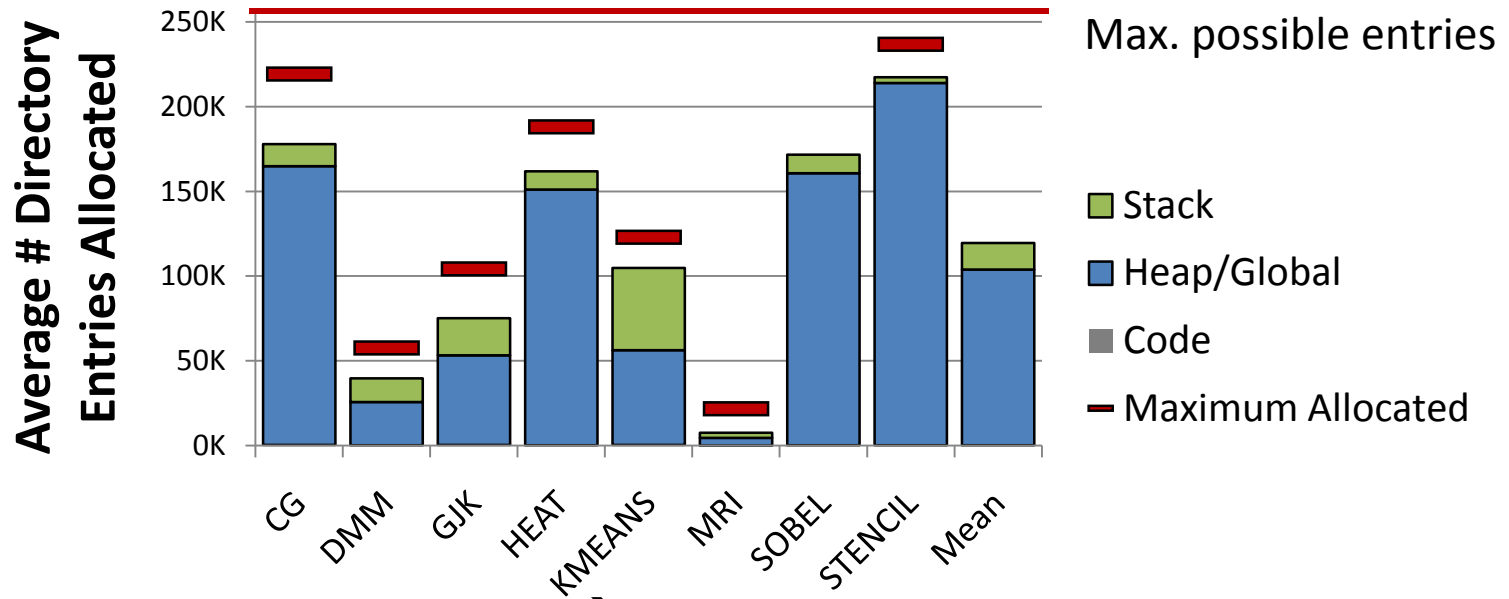- 1024-core CMP, HW caches, single address space, MIMD

UNIVERSITY OF ILLINOIS AT URBANA-CHAMPAIGN

# Opportunity: HW$_{cc}$ v. SW$_{cc}$ Shootout



- Note: Lower bars are better
- **Question**: Can we leverage both HW+SW protocols?

**John H. Kelm**

8

# Opportunity: Network Traffic Reduction



- **SW$_{cc}$** w/baseline arch (**left**), HW$_{cc}$ ww/D$_{IR_{FULL}}$ (**right**)
- **SW$_{cc}$**: Fewer L2 messages in network, some flush overhead
- **HW$_{cc}$**: Extraneous msgs for unshared data (Wr$_{Request}$, Rd$_{Release}$)

# Opportunity: Reduce Directory Utilization



- Not all entries used → Wasted die area

- For many, 256K maximum never reached (red line)

- Observations:

  1. Use $SW_{cc}$ when possible to reduce network traffic

  2. Build smaller sparse directory for common case

**John H. Kelm**

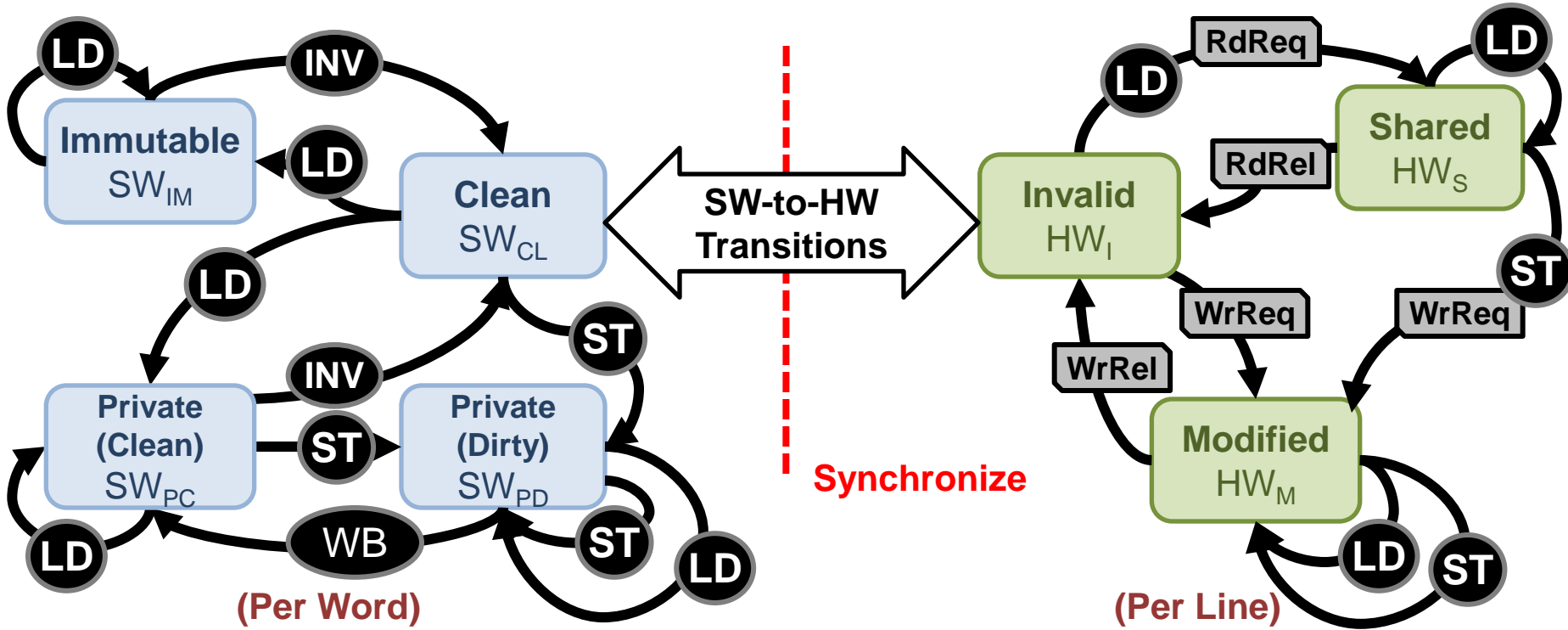# COHESION: Toward a Hybrid Memory Model

- Support for coherence domain transitions
    1. Protocol for safe migration $SW_{cc} \Leftrightarrow HW_{cc}$
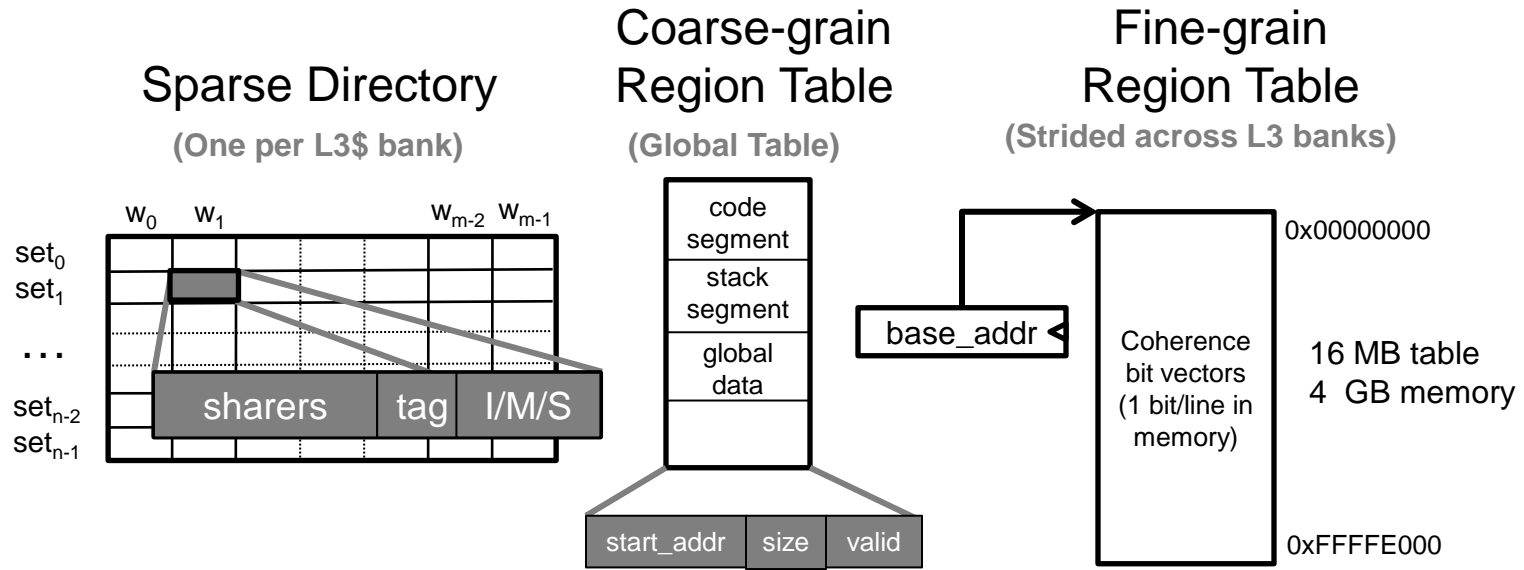    2. Minor architecture extensions

- Motivation:

  ↑ $HW_{cc}$: Supports arbitrary sharing, no SW overhead

  ↓ $HW_{cc}$: Area + message overhead

  ↑ $SW_{cc}$: Removes HW overheads + design complexity

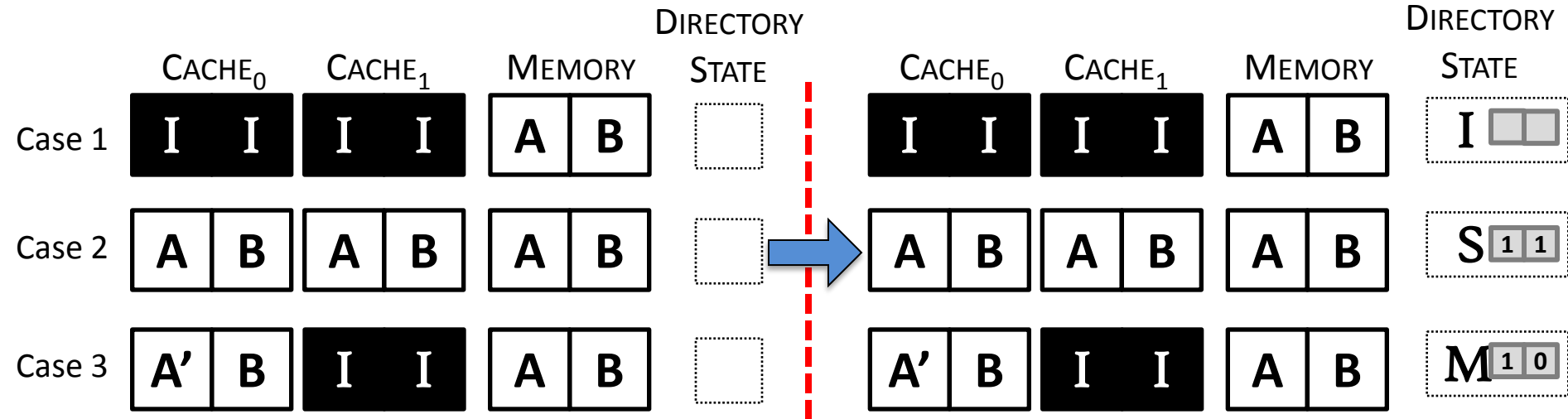  ↓ $SW_{cc}$: Flush overhead + coherence burden on SW

# Protocol Synthesis



- Create a bridge between $SW_{cc}$ and $HW_{cc}$
- Leverage existing $HW_{cc}$ and $SW_{cc}$ techniques

# COHESION Architecture

**Sparse Directory**
*(One per L3$ bank)*

**Coarse-grain Region Table**
*(Global Table)*

**Fine-grain Region Table**
*(Strided across L3 banks)*

$w_0$  $w_1$  $w_{m-2}$  $w_{m-1}$

$set_0$
$set_1$
…
$set_{n-2}$
$set_{n-1}$

sharers | tag | I/M/S

code segment
stack segment
global data

start_addr | size | valid

base_addr

Coherence bit vectors (1 bit/line in memory)

0x00000000

16 MB table
4 GB memory

0xFFFFE000

- Extension to baseline directory protocol
  - Addition 1: Region table/bit vector in memory
  - Addition 2: One bit/line in the L2 cache (not shown)
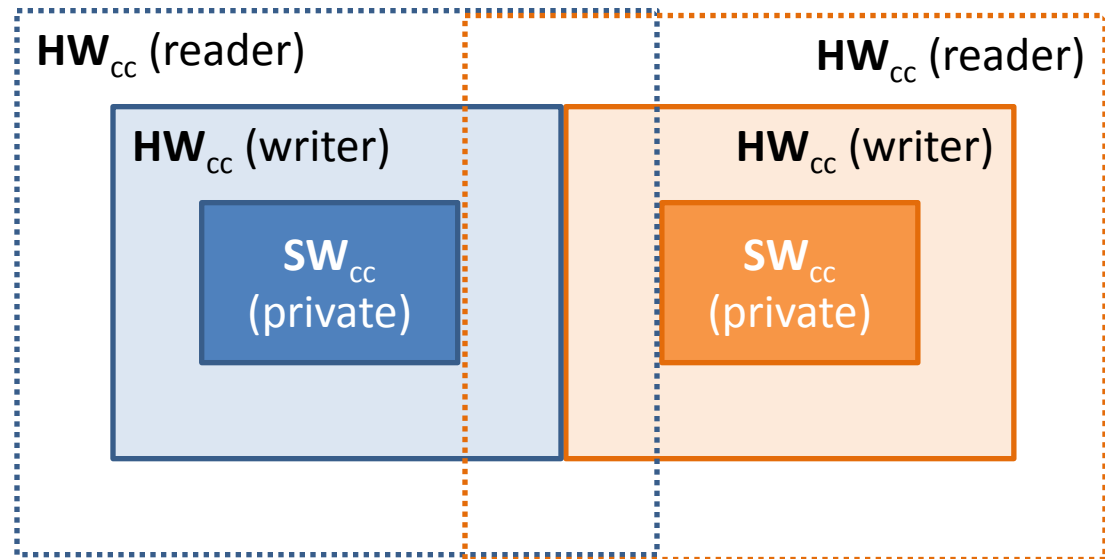- SW writes table → COHESION controller exec's transition

**John H. Kelm**

Rigel

ILLINOIS
UNIVERSITY OF ILLINOIS AT URBANA-CHAMPAIGN

# Example Software → Hardware Transitions



- App. initiates transitions between $SW_{cc}$ and $HW_{cc}$
- COHESION controller probes L2's to reconstruct state
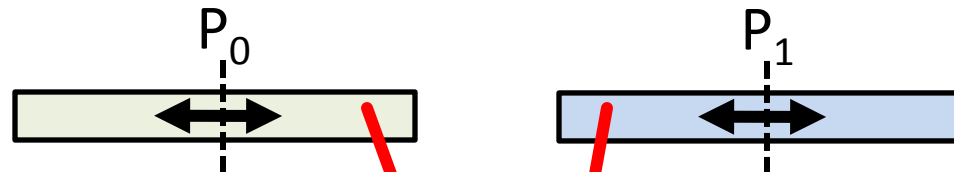- *See paper for other cases and $HW_{cc}$ → $SW_{cc}$*

# Static COHESION Example (1 of 3)

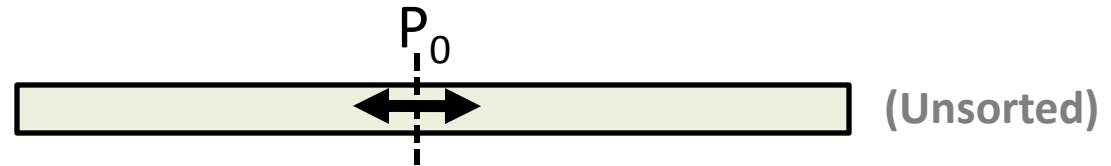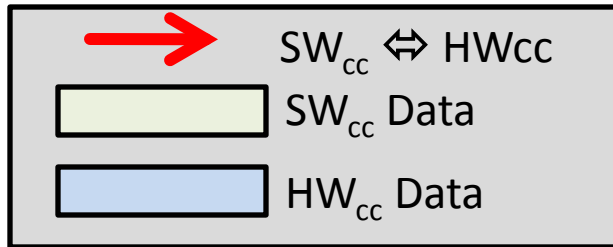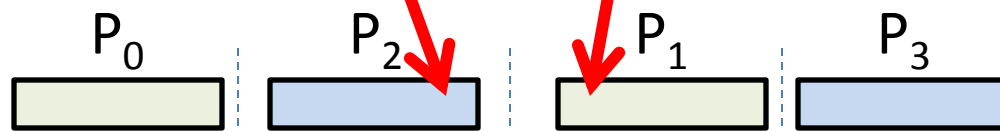Data regions for two grid blocks from a 2D stencil computation

**HW**$_{cc}$ (reader)  **HW**$_{cc}$ (reader)

**HW**$_{cc}$ (writer)  **HW**$_{cc}$ (writer)

**SW**$_{cc}$ (private)  **SW**$_{cc}$ (private)

- COHESION provides static partitioning of data
- (Large) read-only/private regions SW$_{cc}$
- (Small) shared regions HW$_{cc}$

# Dynamic CoHESION Example (2 of 3)
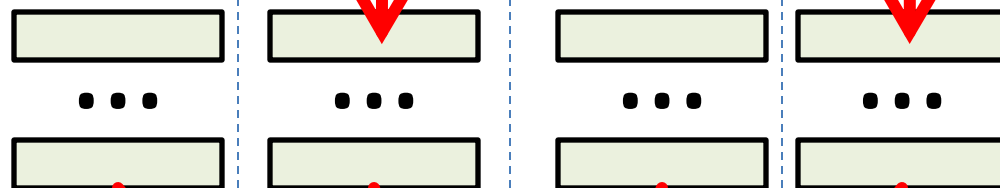
## Parallel Sort (on four cores)

Legend:
- SW$_{cc}$ ⇔ HWcc (red arrow)
- SW$_{cc}$ Data
- HW$_{cc}$ Data

$P_0$

(Unsorted)

$P_0$  $P_1$

1. Parallel Quicksort

$P_0$  $P_2$  $P_1$  $P_3$

2. Sequential Selection
   Sort (Phase 0)

3. Sequential Selection
   Sort (Phases 1-N)

...  ...  ...  ...

4. Result Visible to All

(Sorted)

Rigel

ILLINOIS
UNIVERSITY OF ILLINOIS AT URBANA-CHAMPAIGN
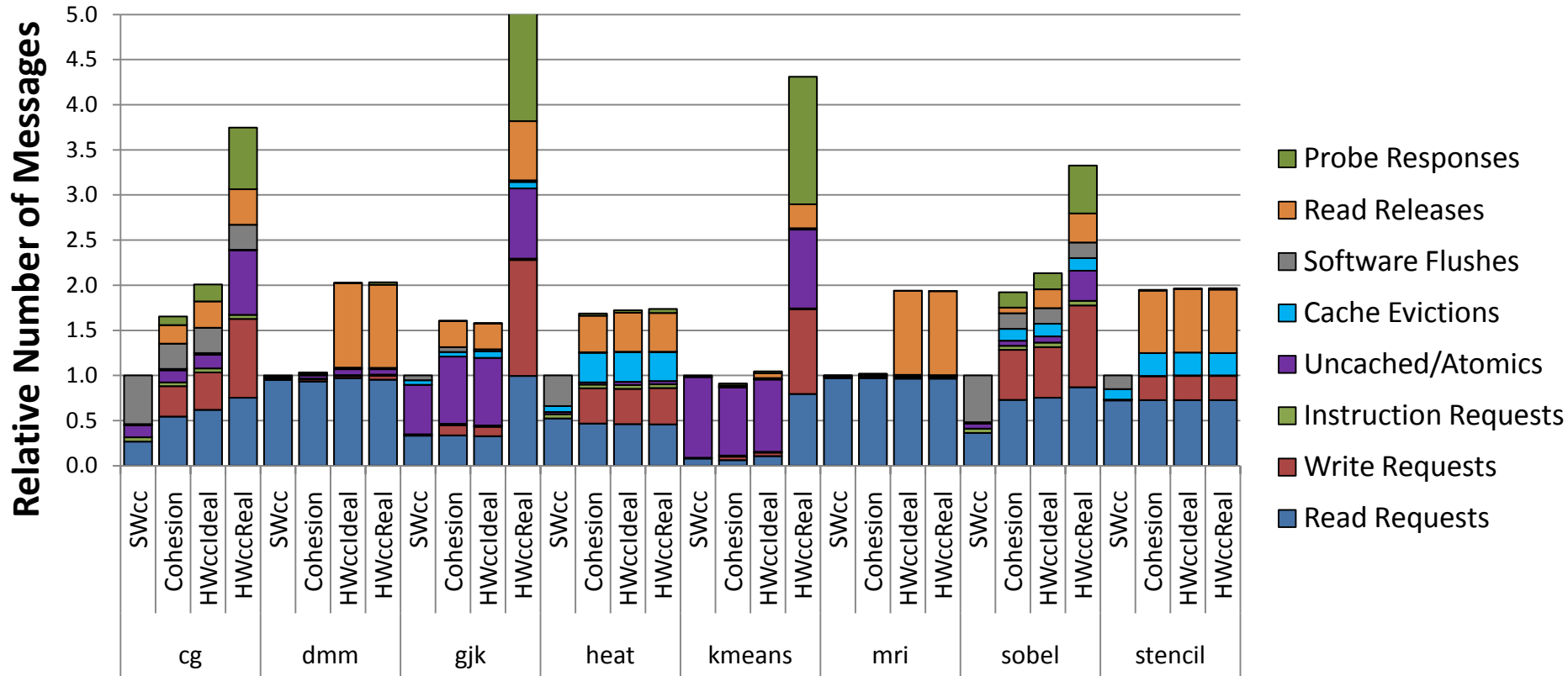
# System SW COHESION Example (3 of 3)
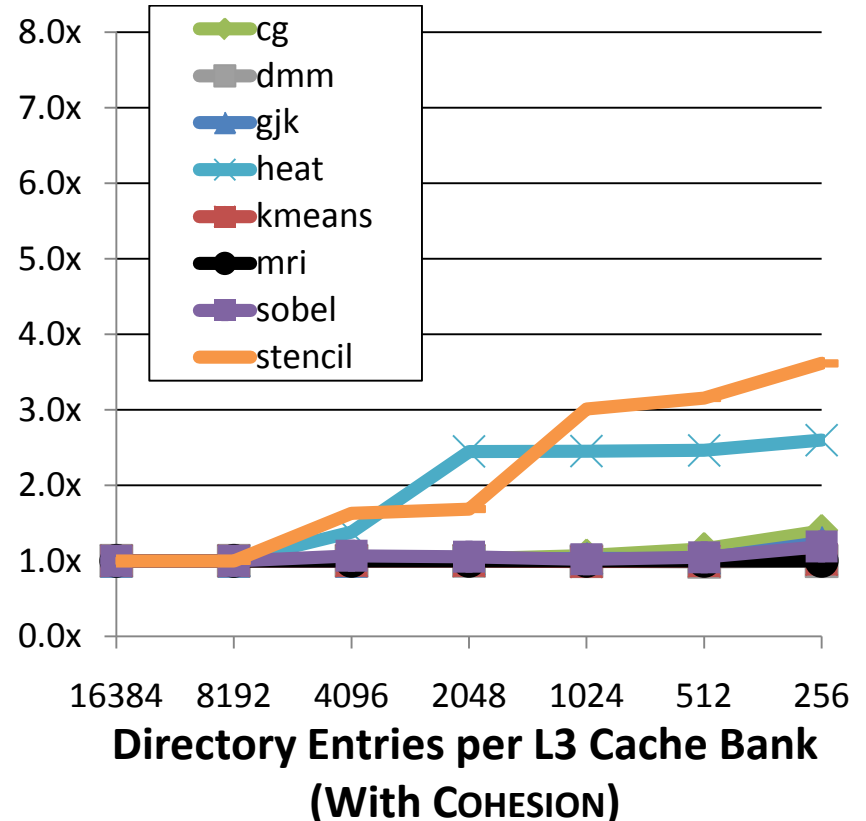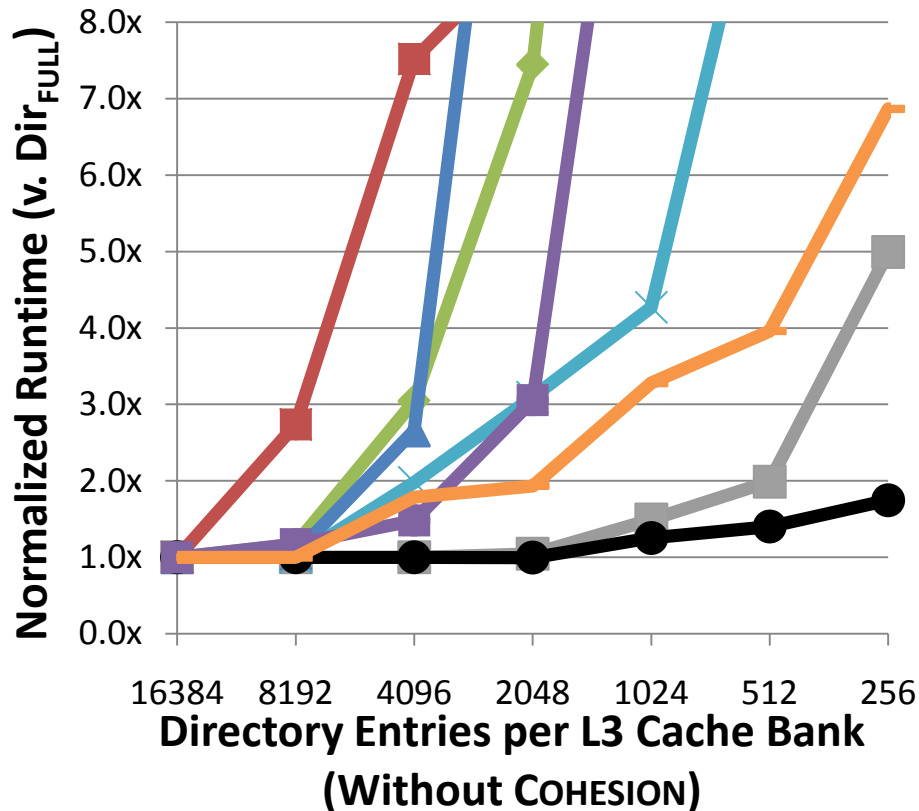
- Problem: Supporting multitasking w/$SW_{cc}$
- OS process creation workflow
    1. Runtime allocates proc's memory $HW_{cc}$
    2. Start new process
    3. Process runs, migrates, $SW_{cc} \Leftrightarrow HW_{cc}$ transitions
    4. Exit process
    5. Runtime makes allocated memory $HW_{cc}$
- COHESION enables: Migration, isolation, cleanup

# Network Message Reductions



Legend:
- Probe Responses
- Read Releases
- Software Flushes
- Cache Evictions
- Uncached/Atomics
- Instruction Requests
- Write Requests
- Read Requests

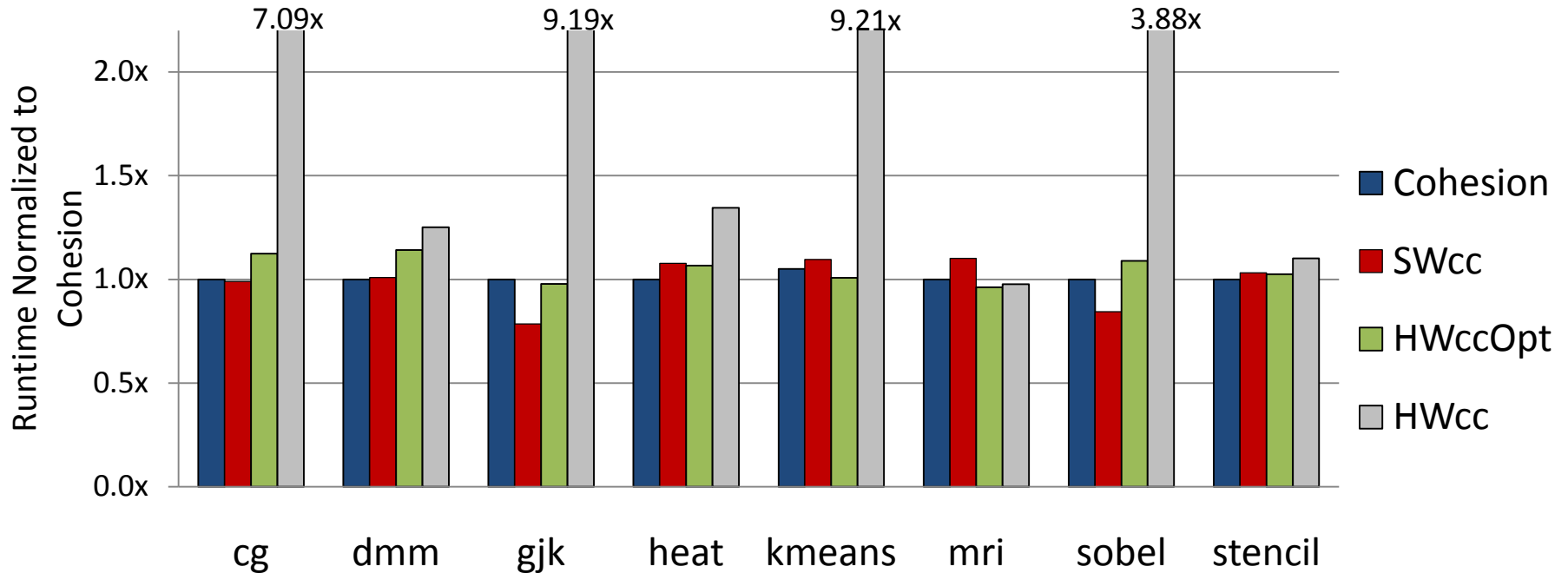Benchmarks (each with SWcc, Cohesion, HWccIdeal, HWccReal): cg, dmm, gjk, heat, kmeans, mri, sobel, stencil

- HW$_{cc}$Real: HW$_{cc}$-only w/sparse directory used by COHESION
- HW$_{cc}$Ideal: Full on-die directory
- **Benefit**: lessens constraints on network design

# Directory Size Sensitivity



**Normalized Runtime (v. $Dir_{FULL}$)**

Legend: cg, dmm, gjk, heat, kmeans, mri, sobel, stencil

**Directory Entries per L3 Cache Bank** (Without COHESION)

**Directory Entries per L3 Cache Bank** (With COHESION)

- Reduces perf. cliffs in sparse directory designs

- **Benefit**: Smaller on-die coherence structures

John H. Kelm

# Runtime: Cohesion, SW$_{cc}$, HW$_{cc}$



- Perf. close to SW$_{cc}$ and full-directory HW$_{cc}$
- Reduce network/directory overhead w/o perf. loss
- Further HW$_{cc}$→SW$_{cc}$ optimizations possible

**John H. Kelm**

# Conclusions

- Why COHESION? CMPs w/multiple mem. models
- Usage scenarios identified
  - System software/migratory tasks w/$SW_{cc}$
  - App uses: Static, dynamic, and host+accel
  - Optimization Path: Piecemeal $HW_{cc} \rightarrow SW_{cc}$
- Hybrid memory model has potential
  - Reduces strain on $HW_{cc}$ implementation
  - Reduces network constraints
  - Competitive performance

**John H. Kelm**