# **GraphPIM**: Enabling Instruction-Level PIM Offloading in Graph Computing Frameworks
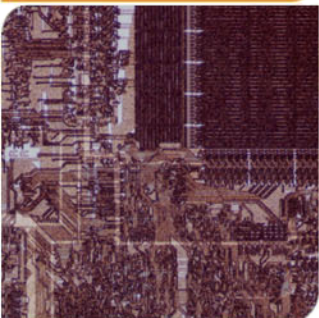
**Lifeng Nai**, Ramyad Hadidi, Jaewoong Sim*,

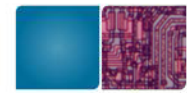Hyojong Kim, Pranith Kumar, Hyesoon Kim

Georgia Tech,   *Intel Labs

Graph computing: processing big network data

▸ Social network, knowledge network, bioinformatics, etc.

Graph computing is inefficient on conventional architectures
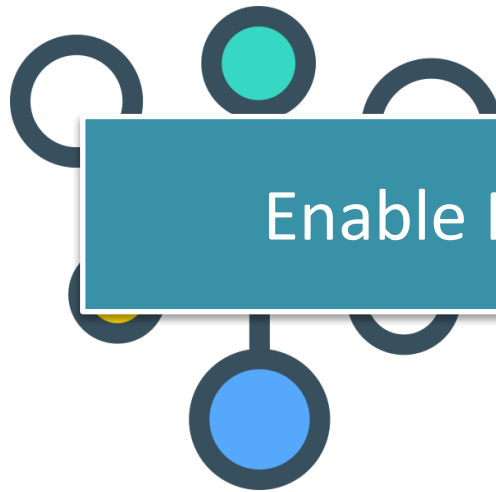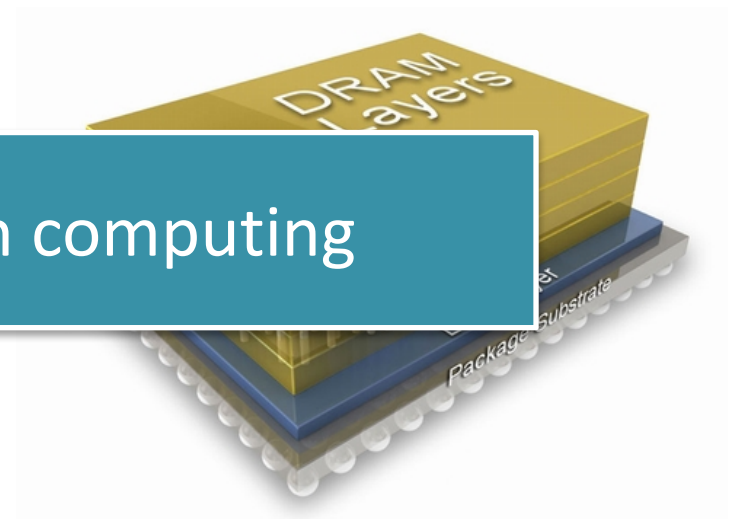
▸ Inefficiency in memory subsystems

## Processing-in-memory (PIM)

▶ PIM has the potential of helping graph computing performance

▶ PIM is being realized in real products: Hybrid memory cube (HMC) 2.0
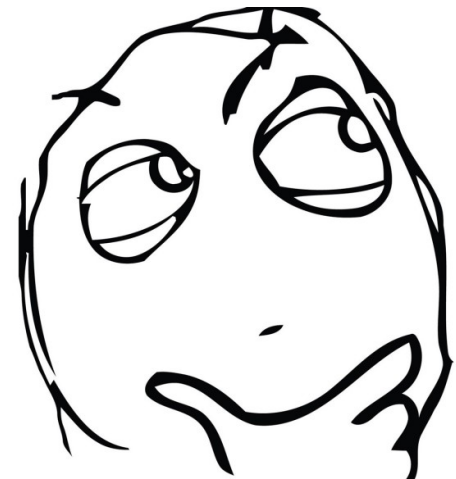


Enable PIM for graph computing
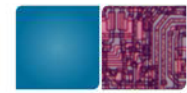
What are the benefits of PIM for graph computing?

- Known benefits of PIM
  - Bandwidth savings, latency reduction, more computation power
- But, they are not good enough
- **We explore something more!**

How to enable PIM for graph in a practical way?

- Minor hardware/software change
- No programmer burden

GraphPIM: a **PIM-enabled** graph **framework**

We identify a new benefit of PIM offloading

We determine PIM offloading targets

We enable PIM without user-application change

**Georgia Tech** | **comparch**

## More computation power

▸ Extra computation units in memory

## Bandwidth savings

▸ Pulling data vs. pushing computation command

|  | Request | Response | Total |
|---|---|---|---|
| 64-byte READ | 1 FLIT (addr) | 5 FLITs (data) | 6 FLITs |
| 64-byte WRITE | 5 FLITs (addr, data) | 1 FLIT (ack) | 6 FLITs |
| CPU rd-modify-wr (rd-Miss; wr-evict) | 6 FLIT | 6 FLIT | 12 FLITs |
| PIM rd-modify-wr | 2 FLITs (addr, imm) | 1 FLIT (ack) | 3 FLITs |

(FLIT:  16 byte, basic flow unit)

Georgia Tech | comparch

More computation power?
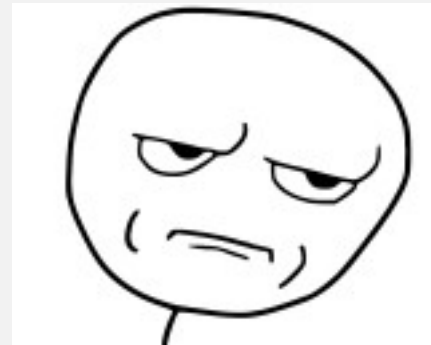
▶ Limited # of FUs in memory

Bandwidth savings?

▶ Not BW saturated

Latency reduction?

▶ Yes, but small

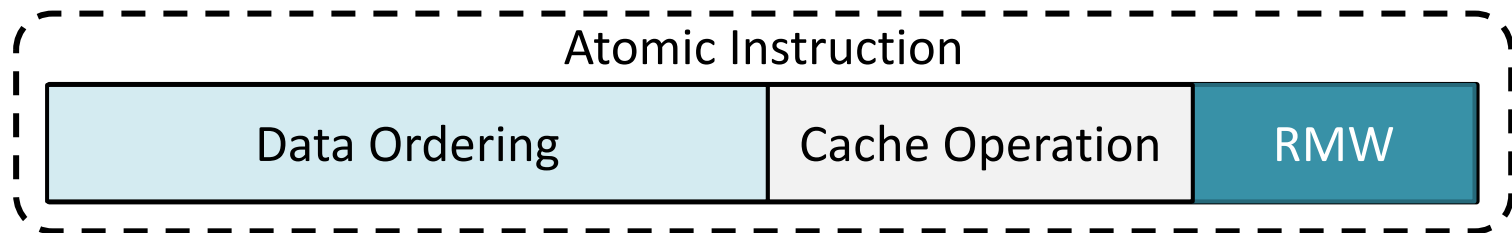## Atomic overhead reduction

▶ Atomic instructions on CPUs have substantial overhead [Schweizer'15]

Atomic Instruction
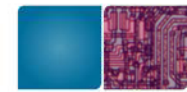
| Data Ordering | Cache Operation | RMW |
|---|---|---|

▸ RMW: read-modify-write

▸ Cache operation: cache-line invalidation, coherence traffic etc.

▸ Data ordering: write buffer draining, pipeline freeze etc.

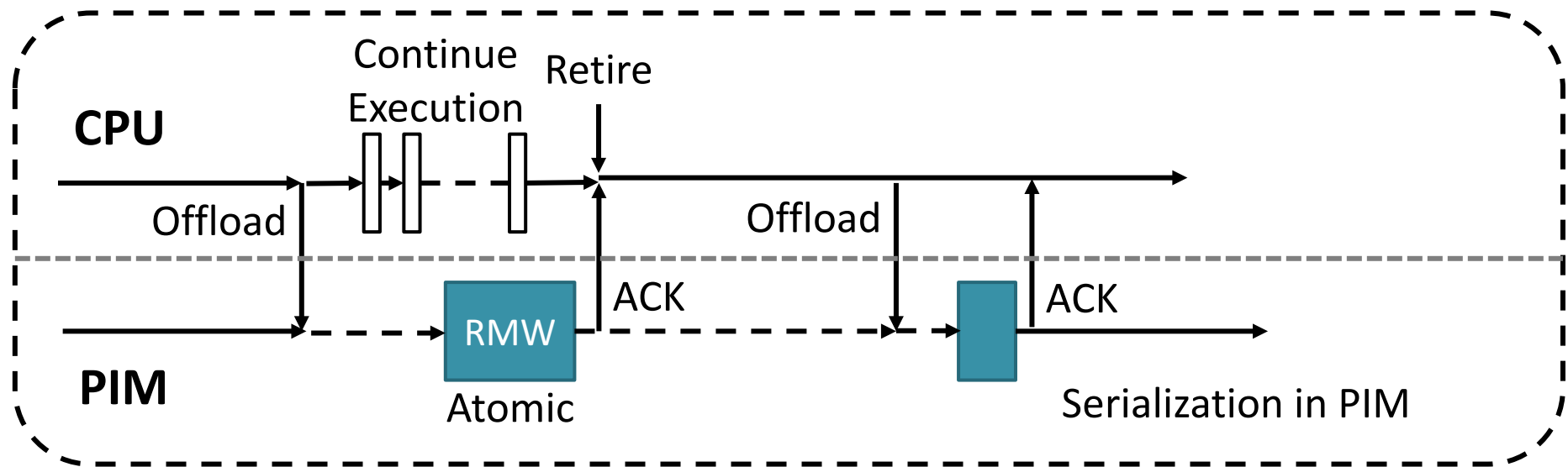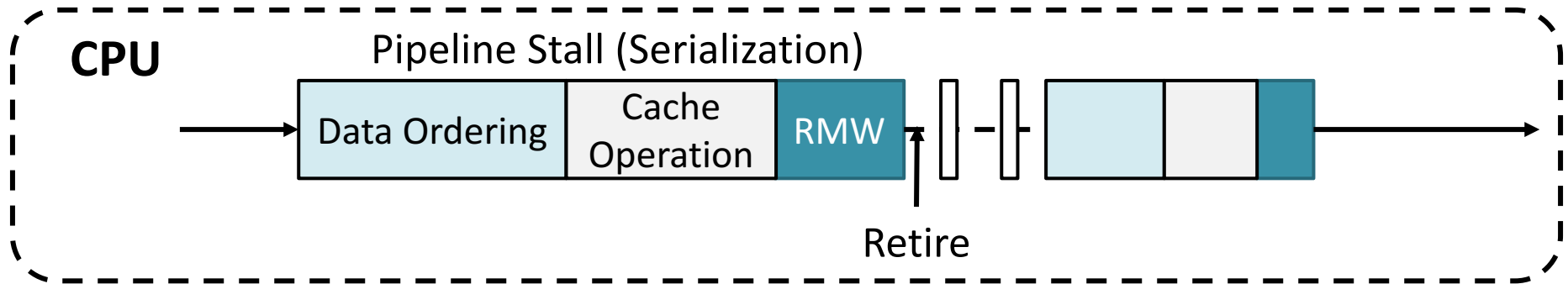▶ Because of the characteristics of graph programming model, PIM offloading can avoid the atomic overhead

*H. Schweizer et al., "Evaluating the Cost of Atomic Operations on Modern Architectures," PACT'15*

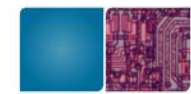**Georgia Tech**  **comparch**
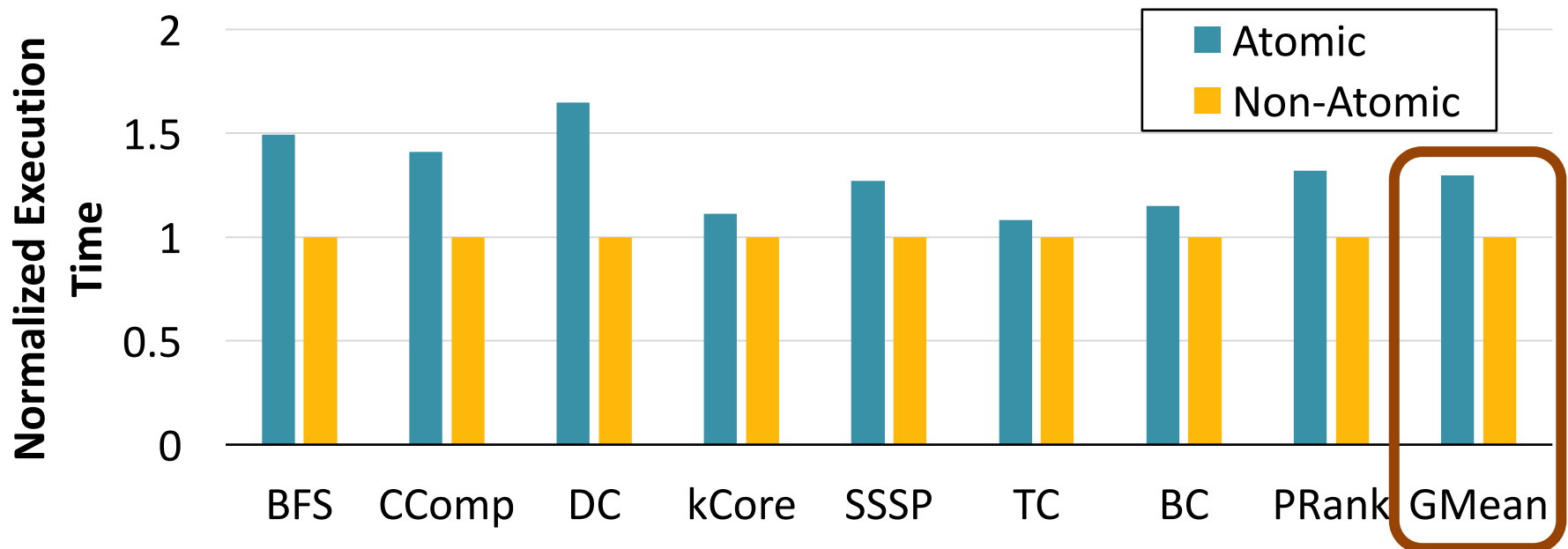
Atomic overhead experiments on a Xeon E5 machine

▶ Atomic RMW → regular load + compute + store

Atomic instructions incur 30% performance degradation



(Non-Atomic: artificial experiment, not precise estimation)

More computation power?

- Limited # of FUs in memory

Bandwidth savings?

- Not BW saturated

Latency reduction?

- Yes, but small

Atomic overhead reduction?

- Yes and significant!
- Main source of PIM benefit for graph

Code snippet: Breadth-first search (BFS)

```
1  F ← {source}
2  while F is not empty
3     F' ← {∅}
4     for each u ∈ F in parallel
5        d ← u.depth + 1
6        for each v ∈ neighbor(u)
7           ret←CAS(v.depth, inf, d)
8           if
...
11       endfor
...
15  endwhile
```

**F**:  frontier vertex set of current step
**F'**: frontier vertex set of next step
**u.depth**:      depth value of vertex u
**neighbor(u)**: neighbor vertices of u
**CAS**(v.depth, inf, d): atomic compare and swap operation

Cache Unfriendly + Atomic

*line 4-5, 8-10*:  accessing **meta data**

*line 6*:  accessing **graph structure**

Cache Friendly

Offload **atomic** operations on graph **property**
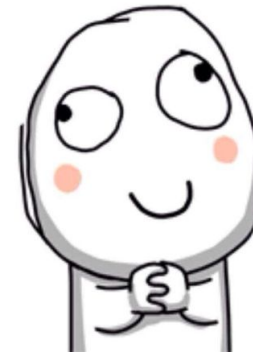
How to indicate offloading targets?

## Option #1: Mark instructions

▸ New instructions in ISA

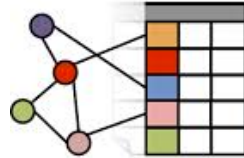▸ Requires changes in user-level applications

## Option #2: Mark memory regions

▸ Special memory region for offloading data
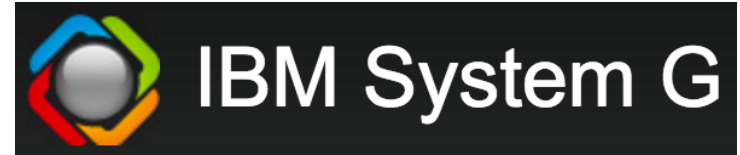
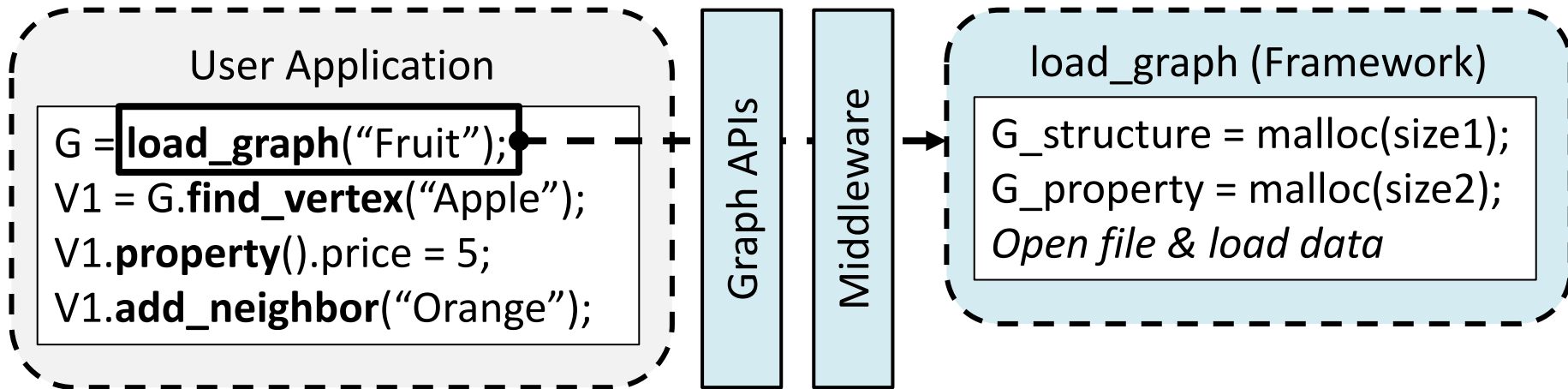▸ Can be transparent to application programmers

**Georgia Tech** comparch

Graph computing is **framework**-based



▶ User application is designed on top of framework interfaces

▶ Data is managed within the framework

| User Application | Graph APIs | Middleware | load_graph (Framework) |
|---|---|---|---|
| G = **load_graph**("Fruit");<br>V1 = G.**find_vertex**("Apple");<br>V1.**property**().price = 5;<br>V1.**add_neighbor**("Orange"); | | | G_structure = malloc(size1);<br>G_property = malloc(size2);<br>*Open file & load data* |

# ENABLE PIM IN GRAPH FRAMEWORK

## PIM memory region (PMR)

▶ Uncacheable memory region in virtual memory space

▶ Utilizing existing uncacheable (UC) support in X86

## Framework change

▶ malloc() → pmr_malloc()
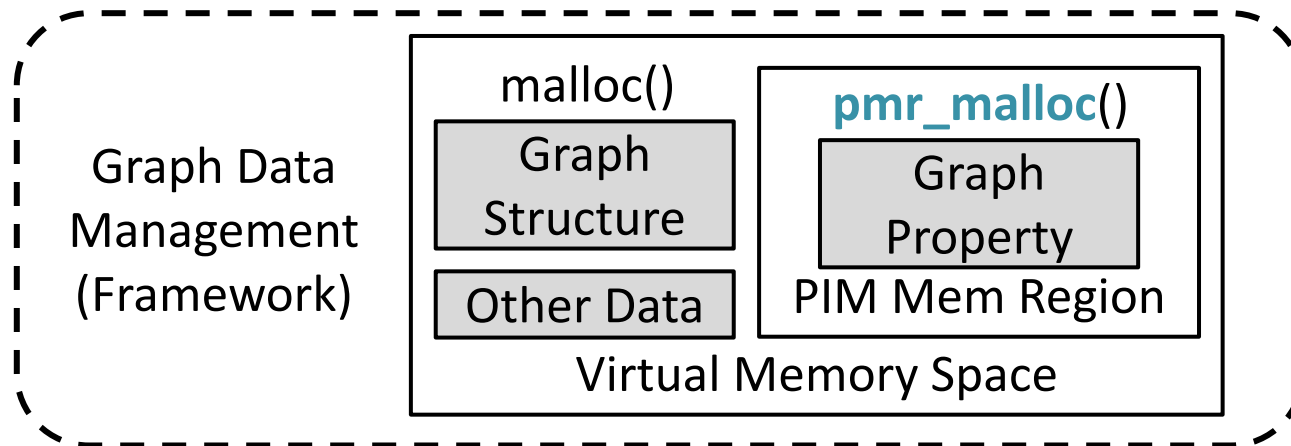
▶ pmr_malloc(): customized malloc function that allocates mem objects in PMR

Graph Data Management (Framework)

malloc()

Graph Structure

Other Data

**pmr_malloc**()

Graph Property

PIM Mem Region

Virtual Memory Space
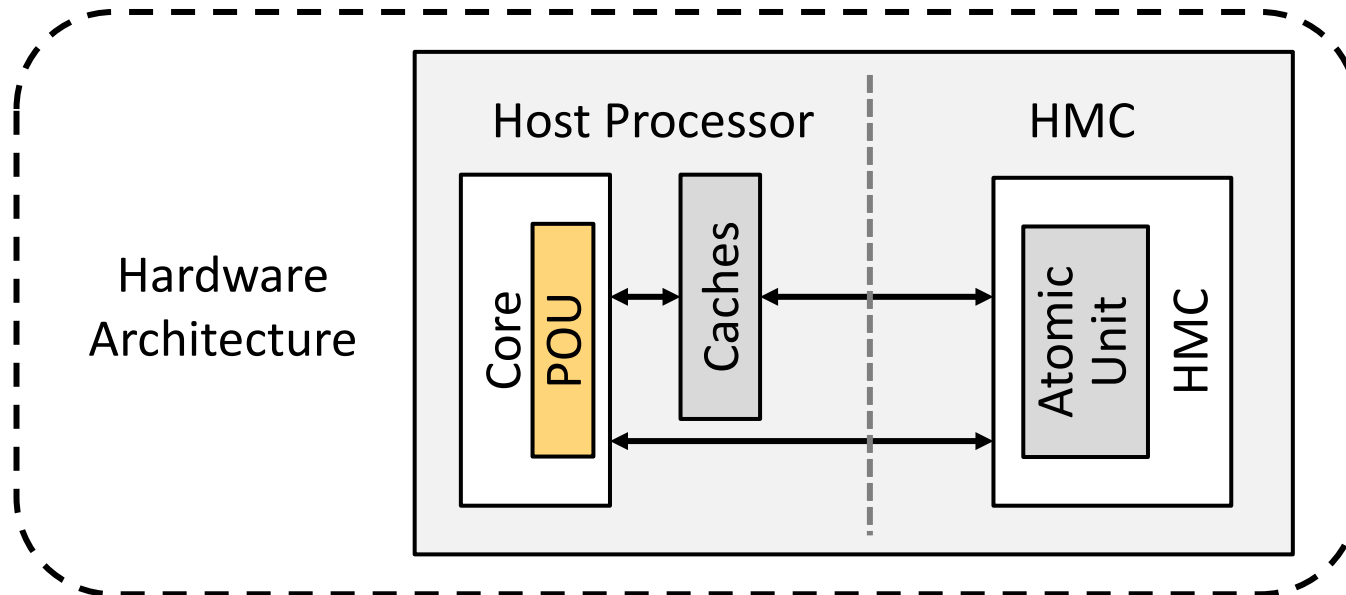
## PIM offloading unit (POU)

▶ Identifies **atomic** instructions that are accessing **PIM Memory Region**

▶ Offloads them as PIM instructions

Software changes:

▸ No **user application** change

▸ Minor change in framework: malloc() → pmr_malloc()

Hardware changes:

▸ PIM memory region: utilizes existing uncacheable (UC) support

▸ PIM offloading unit (POU): identifies offloading targets

No burden on programmers + Minor HW/SW change

## Simulation Environment

▶ SST (framework) + MacSim (CPU) + VaultSim (HMC)

## Benchmark

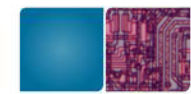▶ GraphBIG benchmark suite [Nai'15] (https://github.com/graphbig)

▶ LDBC dataset from *Linked Data Benchmarking Council* (LDBC)

## Configuration

▶ 16 OoO cores, 2GHz, 4-issue

▶ 32KB L1/256KB L2/16MB shared L3

▶ HMC 2.0 spec, 8GB, 32 vaults, 512 banks, 4 links, 120GB/s per link

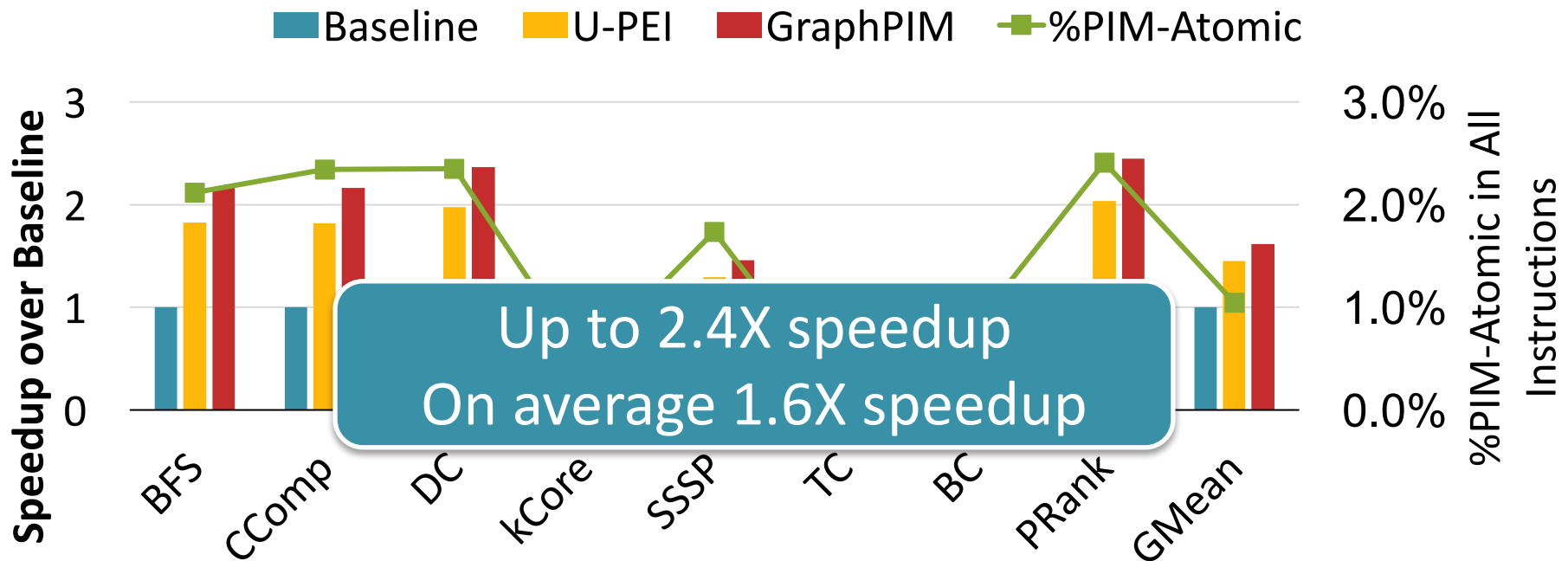*L. Nai et al. "GraphBIG: Understanding Graph Computing in the Context of Industrial Solutions," SC'15*
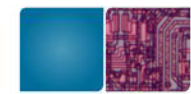
**Georgia Tech** comparch

Baseline: No PIM offloading

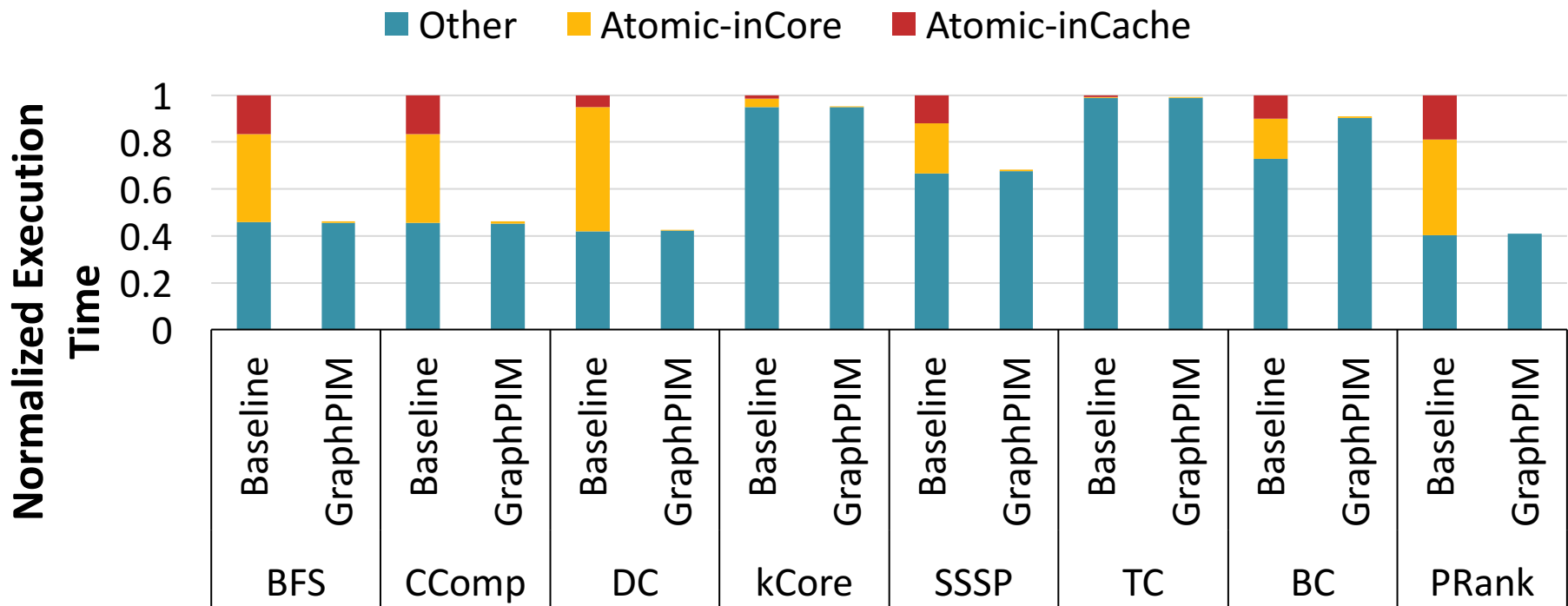U-PEI: Performance upper-bound of PIM-enabled instructions [Ahn'15]



Up to 2.4X speedup
On average 1.6X speedup

*J. Ahn et al. "PIM-Enabled Instructions: A Low-Overhead, Locality-Aware PIM Architecture," ISCA'15*

Breakdown of normalized execution time

▸ Atomic-inCore: **atomic overhead** of offloading targets (atomic inst.)

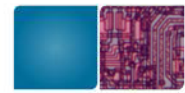▸ Atomic-inCache: **cache-checking overhead** of offloading targets

Graph computing is inefficiency on conventional architectures

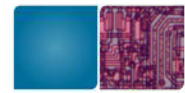GraphPIM enables PIM in graph computing frameworks

▶ Explores a new benefit of PIM offloading: **atomic overhead** reduction

▶ Identifies **atomic operations** on **graph property** as the offloading target

▶ Requires **no user-application change** and only **minor change** in framework and architecture

# THANK YOU!

# BACKUP SLIDES

# DYNAMIC CACHE BEHAVIOR?

GraphPIM marks memory region statically

▶ Cons: cannot be **adaptive** to the working set sizes

▶ But, property accesses to graphs have very high cache misses regardless of graph inputs except for really small graph sizes [JPDC'16, SC'15]

▶ Pros: **coherence** support between memory and processor-cache is not required

**Georgia Tech** comparch

PIM offloading for atomic instructions works fine because…

▸ The programming model of graph applications naturally avoids consistency issues: *all PIM writes are done before reads*

▸ Graph applications require only atomicity from atomic instructions

▸ But, atomic instructions in CPUs don't allow to specify atomicity without fence

▸ We also have a follow-up work discussing the consistency issue for PIM instructions in the context of **general applications** [MEMSYS'17]

**Georgia Tech** **comparch**

Graph applications with BSP model naturally avoids consistency issues

▶ Barriers ensures all PIM writes are done before reads

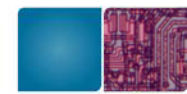| Program Phases | Operation |
|---|---|
| loop:<br>    foreach vertex in task queue:<br>        **read property**<br>        fetch neighbor list<br>        foreach neighbor:<br>            **update neighbor property**<br>            update next-iter task queue<br>    barrier | <br><br>*// Reads*<br><br><br>*// HMC Inst.* |

## GraphPIM

▸ Considers the **separation** of framework and user application

▸ Proposes a **full-stack** solution: SW framework + HW architecture
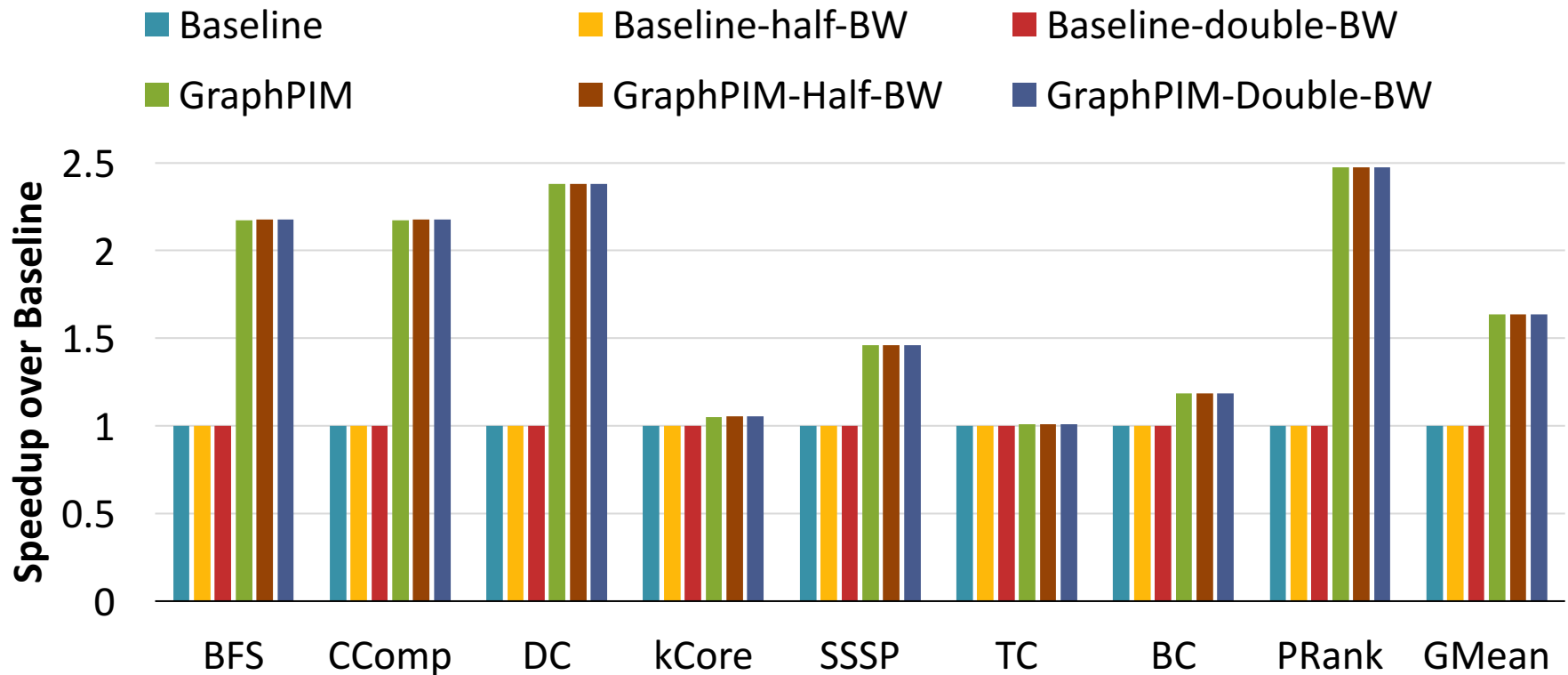
▸ Requires **no** application programmers' efforts

Users can easily enable/disable GraphPIM by switching between different framework libraries.

## Graph on CPUs are not very sensitive to BW changes

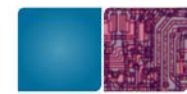▸ Speedup over baseline system with different HMC link bandwidth



- Baseline
- Baseline-half-BW
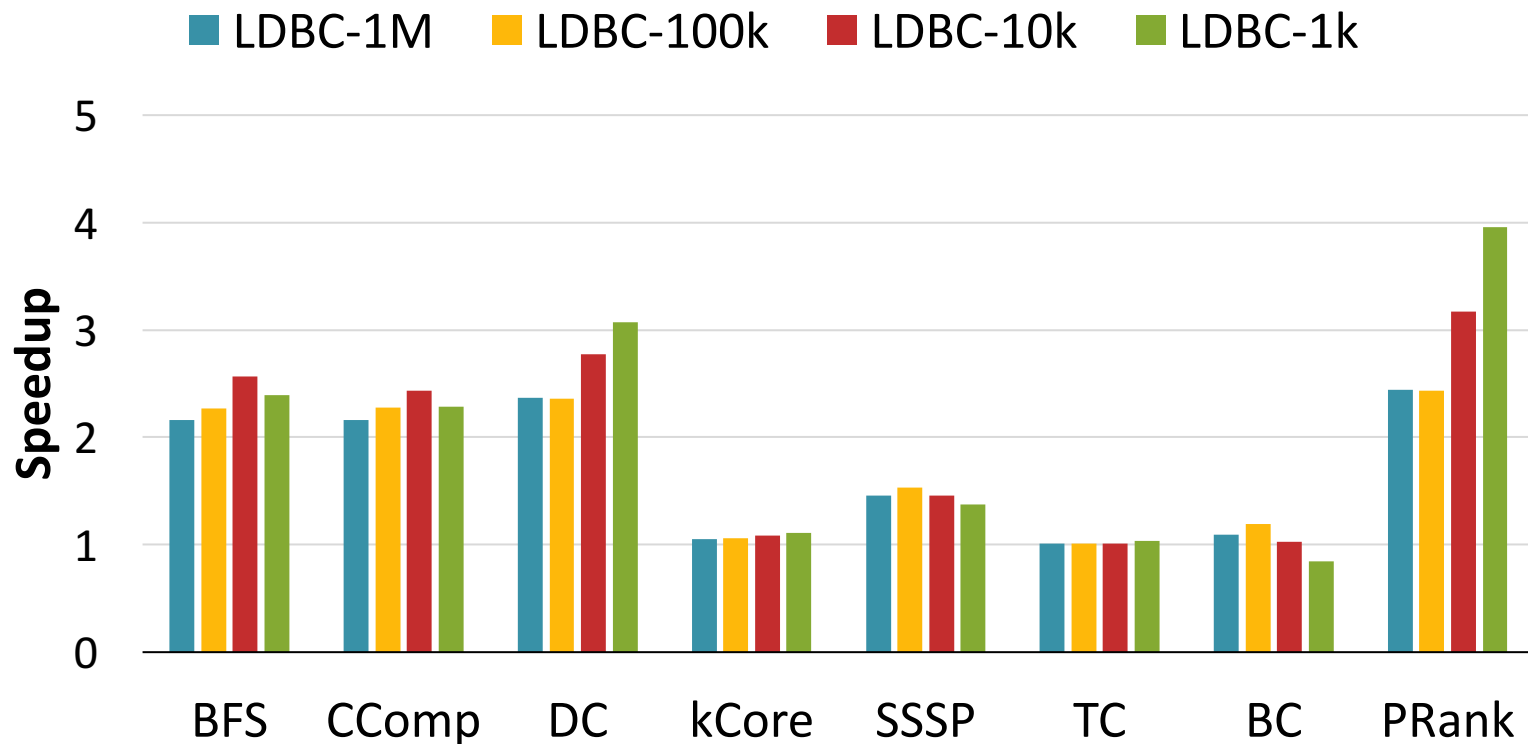- Baseline-double-BW
- GraphPIM
- GraphPIM-Half-BW
- GraphPIM-Double-BW

# APPLICABILITY?

| Category | Workload | Applicable? | | Offloading Target | PIM Inst. |
|---|---|---|---|---|---|
| Graph Traversal | Breadth-first search | ✔ | | lock cmpxchg | CAS if equal |
| | Degree centrality | ✔ | | lock addw | Singed add |
| | Betweenness centrality | ✗ | ✔ | (Floating point add) | (FP add) |
| | Shortest path | ✔ | | lock cmpxchg | CAS if equal |
| | K-core decomposition | ✔ | | lock subw | Singed add |
| | Connected component | ✔ | | lock cmpxchg | CAS if equal |
| | Page rank | ✗ | ✔ | (Floating point add) | (FP add) |
| Dynamic Graph | Graph construction | ✗ | | (Complex operation) | |
| | Graph update | ✗ | | (Complex operation) | |
| | Topology morphing | ✗ | | (Complex operation) | |
| Rich Property | Triangle count | ✔ | | lock add | Singed add |
| | Gibbs inference | ✗ | | (Compute intensive) | |

Georgia Tech comparch

GraphPIM speedup over baseline with different dataset sizes

Normalized uncore energy consumption



On average, GraphPIM saves **37%** of uncore energy because of reduction in **cache accesses** and **memory bandwidth**

# GRAPHPIM: EVALUATION

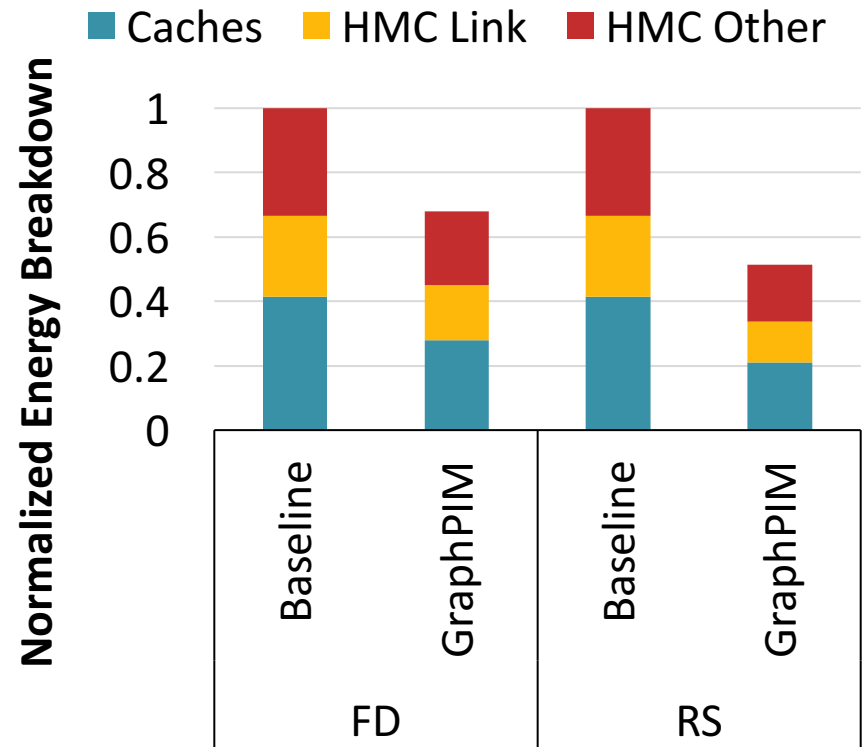# GRAPHPIM: EVALUATION

Normalized bandwidth consumption with request/response breakdown
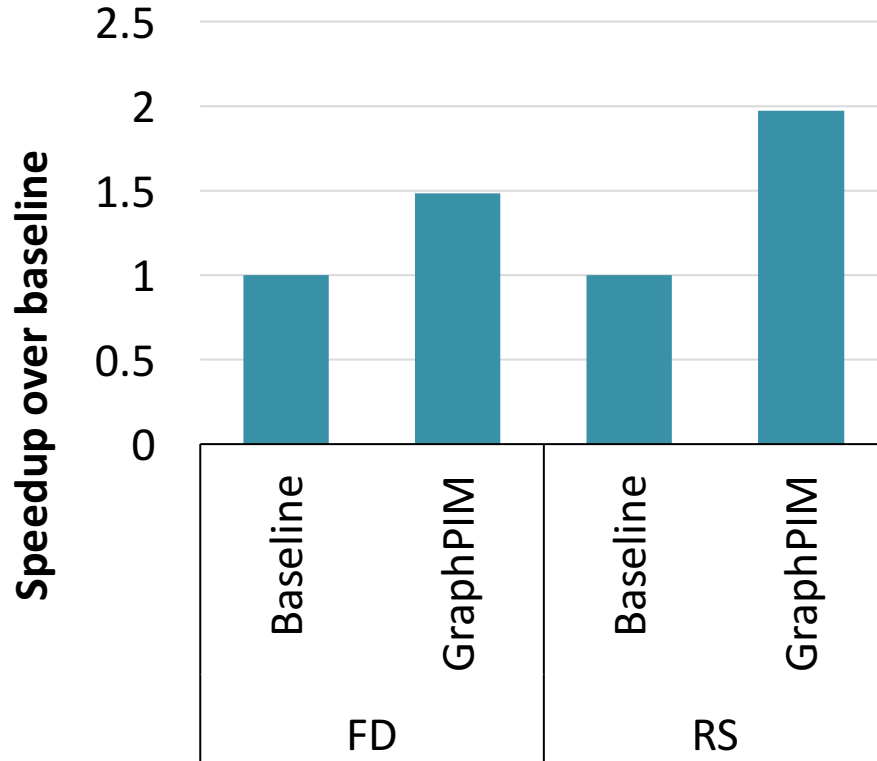
## Performance and energy results of two real-world applications

▸ Based on an analytical model
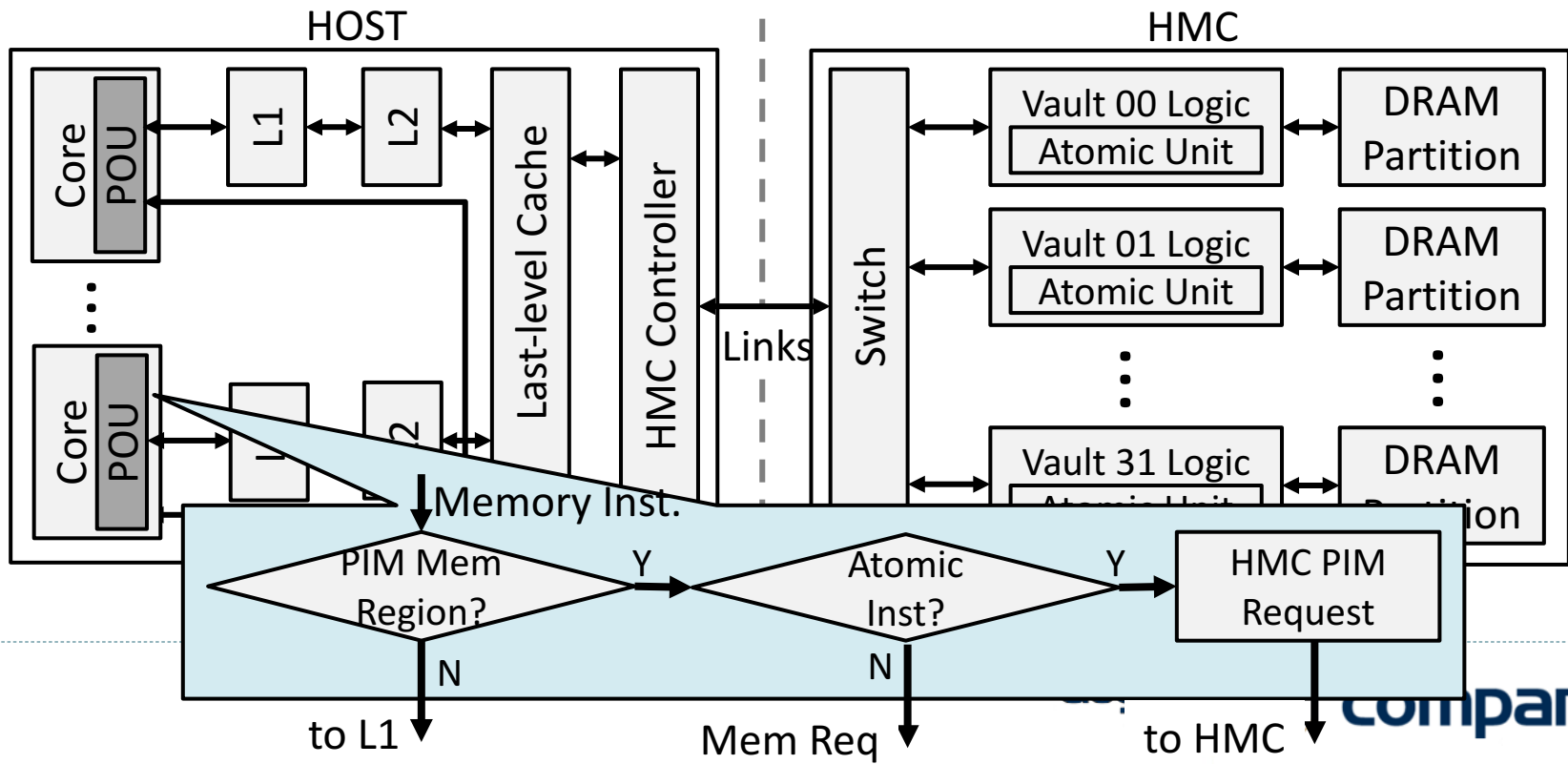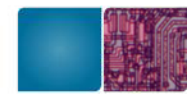
▸ FD: Financial fraud detection; RS: Recommender system

## PIM Memory Region (PMR)

▶ A uncacheable memory region in virtual memory space

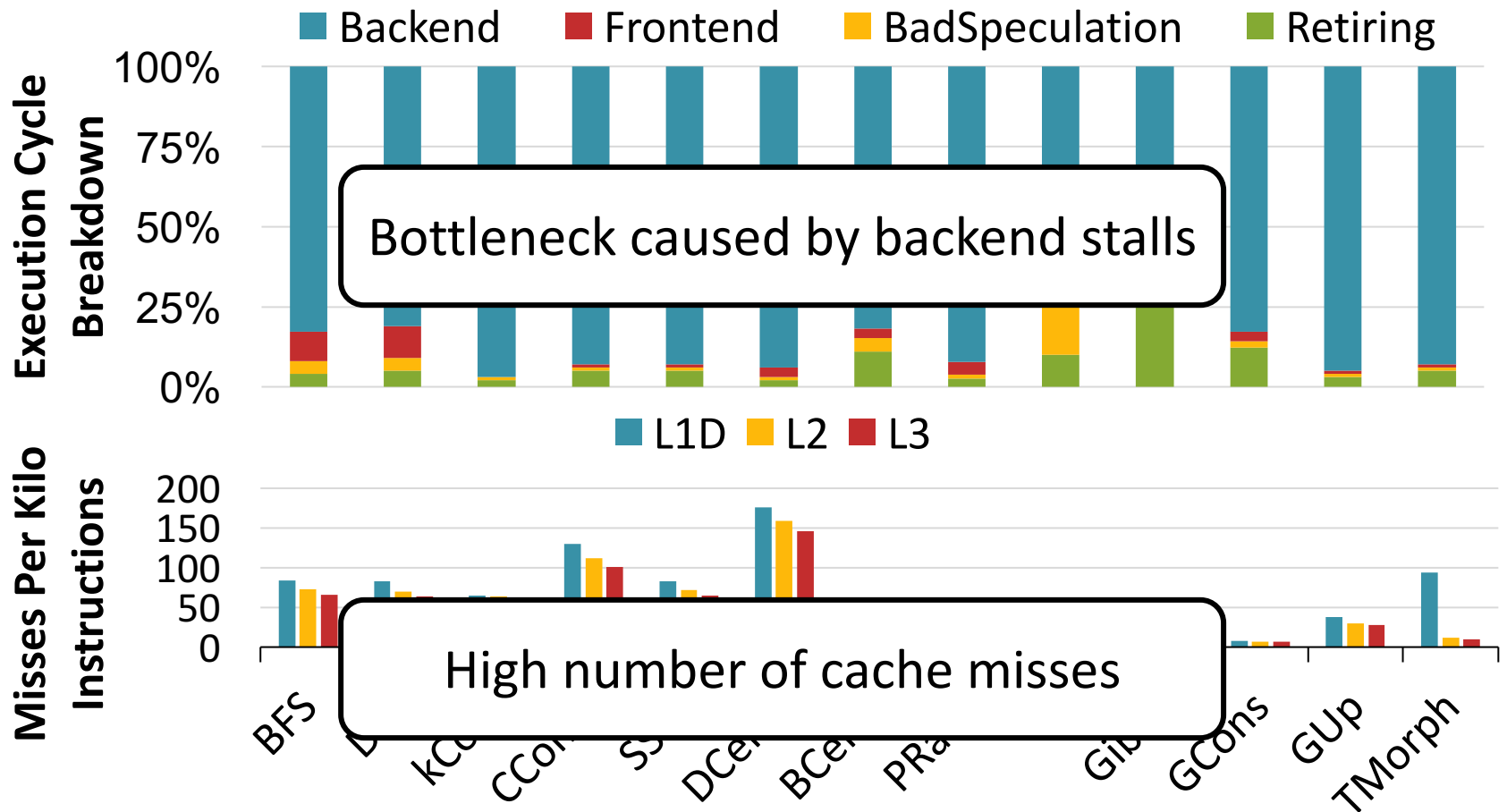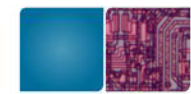▶ Utilizing existing **uncacheable** (UC) support in X86

## PIM Offloading Unit (POU)

## Profiling using HW performance counters

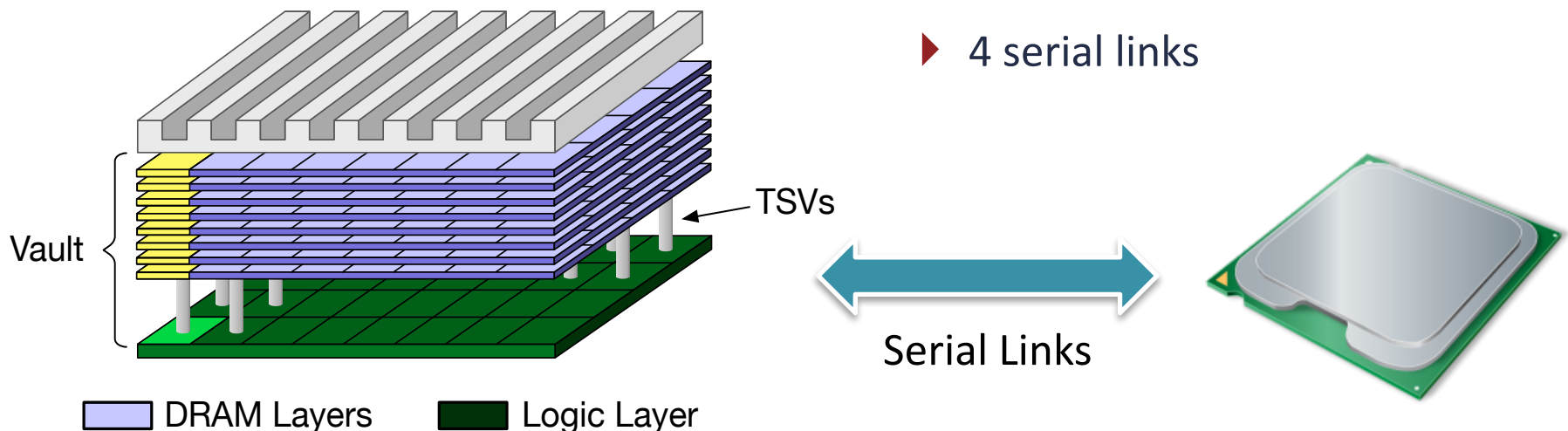▸ Execution cycle breakdown: top-down methodology from Intel

## Hybrid Memory Cube (HMC) 2.0

▸ One of the first industrial PIM proposals

▸ **Instruction-level** PIM offloading

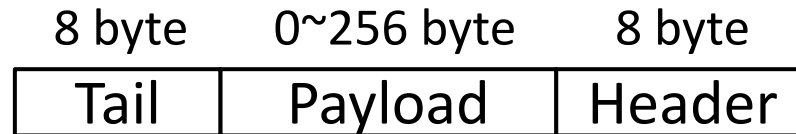▸ 1 logic die + 4/8 DRAM dies

▸ 32 Vaults

▸ 4 serial links

TSVs

Vault

Serial Links

| DRAM Layers | Logic Layer |

## Packet-based protocol

| 8 byte | 0~256 byte | 8 byte |
|--------|------------|--------|
| Tail | Payload | Header |

## Regular READ/WRITE

▸ FLIT: 16-byte; basic flow unit

|  | Request | Response |
|--|---------|----------|
| 64-byte READ | 1 FLIT | 5 FLITs |
| 64-byte WRITE | 5 FLITs | 1 FLIT |

Georgia Tech    comparch

PIM Instruction: **read-modify-write** (RMW) operation

▸ Similar as regular READ/WRITE, just different **CMD** in the **Header**

▸ DRAM bank is locked during the whole RMW for atomicity

PIM-ADD(addr, imm)

ACK

| Tail | addr, imm | Header (PIM-ADD) |

Georgia Tech · comparch