# PIM-Enabled Instructions: A Low-Overhead, Locality-Aware PIM Architecture

Junwhan Ahn, Sungjoo Yoo, Onur Mutlu[+], and Kiyoung Choi

Seoul National University          [+]Carnegie Mellon University
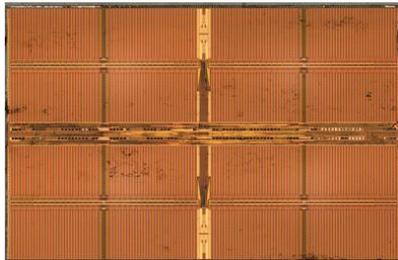
# **Processing-in-Memory**

- Move computation to memory
  - Higher memory bandwidth
  - Lower memory latency
  - Better energy efficiency (e.g., off-chip links vs. TSVs)

- Originally studied in 1990s
  - Also known as processor-in-memory
  - e.g., DIVA, EXECUBE, FlexRAM, IRAM, Active Pages, …
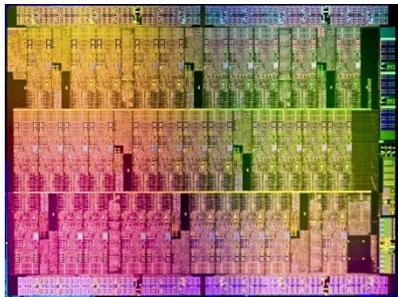  - *Not commercialized in the end*

Why was PIM unsuccessful in its first attempt?

# Challenges in Processing-in-Memory
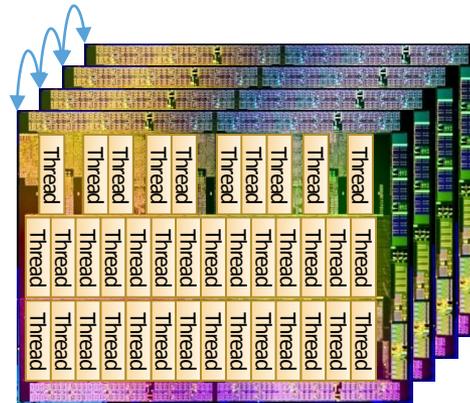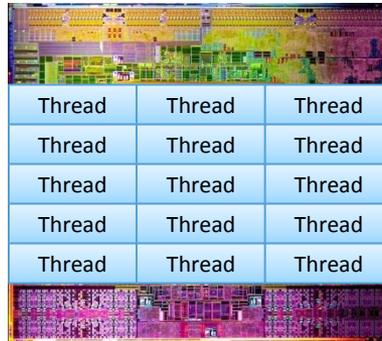
| Cost-effectiveness | Programming Model | Coherence & VM |
|---|---|---|

**Cost-effectiveness**

DRAM die



Complex Logic

**Programming Model**

Host Processor

| Thread | Thread | Thread |
|---|---|---|
| Thread | Thread | Thread |
| Thread | Thread | Thread |
| Thread | Thread | Thread |
| Thread | Thread | Thread |

In-Memory Processors

**Coherence & VM**

Host Processor

C  4

5  C

DRAM die

# Challenges in Processing-in-Memory

| Cost-effectiveness | Programming Model | Coherence & VM |
|---|---|---|

Host Processor

Host Processor

| Thread | Thread | Thread |
|---|---|---|
| Thread | Thread | Thread |
| Thread | Thread | Thread |
| Thread | Thread | Thread |

(Partially) Solved by
3D-Stacked DRAM

**Still Challenging** even in Recent PIM Architectures
(e.g., AC-DIMM, NDA, NDC, TOP-PIM, Tesseract, …)

Complex Logic

In-Memory Processors

DRAM die

# New Direction of PIM

- Objectives
  - Provide an intuitive programming model for PIM
  - Full support for cache coherence and virtual memory
  - Reduce the implementation overhead of PIM units

- Our solution: simple PIM operations as ISA extension
  - Simple: low-overhead implementation
  - PIM operations as host processor instructions: intuitive
  - Conventional PIM : Simple PIM ≈ GPGPU : SSE/AVX

# Potential of ISA Extension as PIM Interface

- Example: Parallel PageRank computation

```
for (v: graph.vertices) {
    value = weight * v.rank;
    for (w: v.successors) {
        w.next_rank += value;
    }
}
for (v: graph.vertices) {
    v.rank = v.next_rank; v.next_rank = alpha;
}
```

# Potential of ISA Extension as PIM Interface

```
for (v: graph.vertices) {
    value = weight * v.rank;
    for (w: v.successors) {
        w.next_rank += value;
    }
}
```

Host Processor

Main Memory

w.next_rank  ←———————————→  w.next_rank

64 bytes **in**
64 bytes **out**

**Conventional Architecture**

# Potential of ISA Extension as PIM Interface

```
for (v: graph.vertices) {
    value = weight * v.rank;
    for (w: v.successors) {
        __pim_add(&w.next_rank, value);
    }
}
```

pim.add r1, (r2)

Host Processor

Main Memory

value

w.next_rank

8 bytes **in**
0 bytes **out**

**In-Memory Addition**

# Potential of ISA Extension as PIM Interface



**Increase in Memory Bandwidth Consumption**
Lack of On-Chip Caches

**Reduction in Memory Bandwidth Consumption**
In-Memory Computation

Speedup

60%
50%
40%
30%
20%
10%
0%
-10%
-20%

p2p-Gnu tella31

soc-Slash dot0811

web-Stanford

amazon-2008

frw 20

w T

ci Pate

soc-L Journ

ljour 200

More Vertices

# Overview

1.  How should simple PIM operations be interfaced to conventional systems?

    – Expose PIM operations as *cache-coherent, virtually-addressed host processor instructions*

    – No changes to the existing sequential programming model

2.  What is the most efficient way of exploiting such simple PIM operations?

    – Dynamically determine the location of PIM execution based on data locality without software hints

# PIM-Enabled Instructions

```
for (v: graph.vertices) {
    value = weight * v.rank;
    for (w: v.successors) {
        w.next_rank += value;
    }
}
```
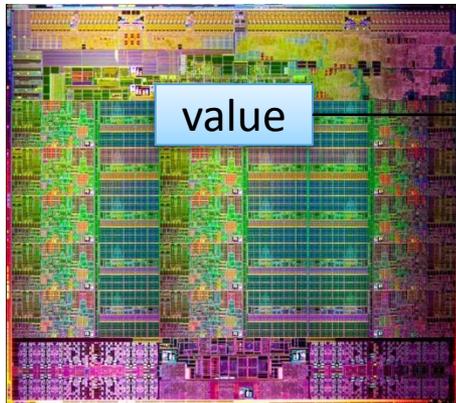
# PIM-Enabled Instructions

```
for (v: graph.vertices) {
    value = weight * v.rank;
    for (w: v.successors) {
        __pim_add(&w.next_rank, value);
    }
}
```

pim.add r1, (r2)

- Executed either in memory or in the host processor
- Cache-coherent, virtually-addressed
- Atomic between different PEIs
- *Not* atomic with normal instructions (use *pfence*)

# PIM-Enabled Instructions

```
for (v: graph.vertices) {
    value = weight * v.rank;
    for (w: v.successors) {
        __pim_add(&w.next_rank, value);
    }
}
pfence();
```

pim.add r1, (r2)

pfence

- Executed either in memory or in the host processor
- Cache-coherent, virtually-addressed
- Atomic between different PEIs
- *Not* atomic with normal instructions (use *pfence*)

# PIM-Enabled Instructions

- Key to practicality: single-cache-block restriction
  - Each PEI can access *at most one last-level cache block*
  - Similar restrictions exist in atomic instructions

- Benefits
  - **Localization**: each PEI is bounded to one memory module
  - **Interoperability**: easier support for cache coherence and virtual memory
  - **Simplified locality monitoring**: data locality of PEIs can be identified by LLC tag checks or similar methods

# Architecture



Proposed PEI Architecture

# Memory-side PEI Execution



pim.add y, &x

# Memory-side PEI Execution

Host Processor

y Of-Order Core

y PCU

**Address Translation for PEIs**

- Done by the host processor TLB (similar to normal instructions)
- *No modifications to existing HW/OS*
- *No need for in-memory TLBs*

Directory

Locality Monitor

Cr...

PCU

DRAM Controller

pim.add y, &x

# Memory-side PEI Execution



Host Processor

| y | Of-Order Core |

| y | PCU |

Wait until x is writable

L1 Cache

L2 Cache

Last-Level Cache

HMC Controller

PMU

PIM Directory

Locality Monitor

HMC

Crossbar Network

PCU

DRAM Controller

PCU

x | RAM Controller

⋮

PCU

DRAM Controller

pim.add y, &x

Host P...

Out-...

XOR-Hash

Address

(Inexact, but *Conservative*)

Reader-writer lock #0
Reader-writer lock #1
Reader-writer lock #2

⋮

Reader-writer lock #N-1

y  P...

PMU

PIM
Directory

Locality
Monitor

Wait until x is writable

HMC Contro...

Crossbar Net...

PCU

PCU

RAM
...troller

...AM
Controller

DRAM
Controller

⋮

pim.add y, &x

# Memory-side PEI Execution



Host Processor

Out-Of-Order Core

y PCU

L1 Cache

L2 Cache

Last-Level Cache

HMC Controller

PMU

PIM Directory

Locality Monitor

Wait until x is writable
Check the data locality of x

HMC

Crossbar Network

PCU

DRAM Controller

PCU

x DRAM Controller

PCU

DRAM Controller

pim.add y, &x

## Partial Tag Array

| Tag | Tag | Tag | ... | Tag |
|-----|-----|-----|-----|-----|
| Tag | Tag | Tag | ... | Tag |

. . .

| Tag | Tag | Tag | ... | Tag |

Address

**Hit**: High locality

**Miss**: Low locality

Updated on
- Each LLC access
- Each issue of a PIM operation to memory

Hos

Ou

y

PIM
Directory

HM

Cross

DRAM
Controller

M
ller

M
ler

PCU

Locality
Monitor

Wait until x is writable
Check the data locality of x

pim.add y, &x

# Memory-side PEI Execution

**Host Processor**

Out-Of-Order Core

y PCU

L1 Cache

L2 Cache

Last-Level Cache

HMC Controller

**Low** locality

Wait until x is writable

Check the data locality of x

PMU

PIM Directory

Locality Monitor

**HMC**

Crossbar Network

PCU

DRAM Controller

PCU

x DRAM Controller

PCU

DRAM Controller

pim.add y, &x

# Memory-side PEI Execution

Host Processor

HMC

Out-Of-Order Core

y PCU

L1 Cache

L2 Cache

Last-Level Cache

**Low** locality

PMU

PIM Directory

Locality Monitor

- Back-invalidation for cache coherence
- *No modifications to existing cache coherence protocols*

pim.add y, &x

# Memory-side PEI Execution



pim.add y, &x

# Memory-side PEI Execution

**Host Processor**

Out-Of-Order Core

...he ...he ...vel ...e

Directory

Locality Monitor

H

**Completely Localized PIM Memory Accesses without Special Data Mapping**

**HMC**

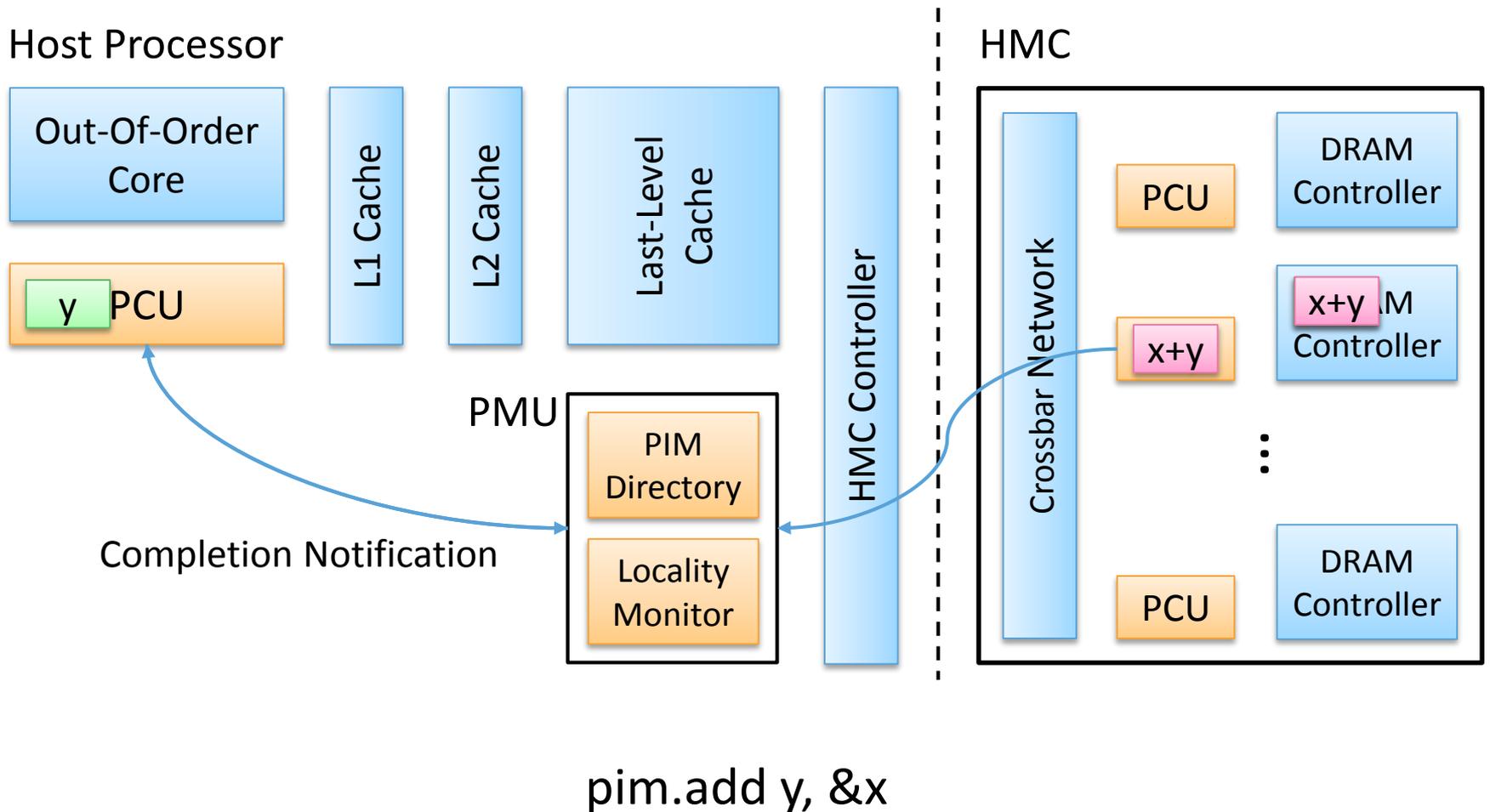Crossbar Network

PCU

DRAM Controller

x+y

x+y M Controller

PCU

DRAM Controller

pim.add y, &x

# Memory-side PEI Execution



Host Processor | HMC

Out-Of-Order Core

y  PCU

L1 Cache

L2 Cache

Last-Level Cache

HMC Controller

PMU

PIM Directory

Locality Monitor

Completion Notification

Crossbar Network

PCU

x+y

x+y AM Controller

PCU

DRAM Controller

DRAM Controller

pim.add y, &x

# Host-side PEI Execution



Host Processor

| y | Of-Order Core |

| y | PCU |

L1 Cache

L2 Cache

| x | Last-Level Cache |

Wait until x is writable

Check the data locality of x

PMU

PIM Directory

Locality Monitor

HMC Controller

HMC

Crossbar Network

PCU

PCU

PCU

DRAM Controller

DRAM Controller

DRAM Controller

pim.add y, &x

# Host-side PEI Execution



High locality

Wait until x is writable
Check the data locality of x

pim.add y, &x

# Host-side PEI Execution

**Host Processor**

Out-Of-Order Core

x+y

y  P  x+y

L1 Cache

x+y

**HMC**

x

L2 Cache

x

**PMU**

PIM Directory

Locality Monitor

HMC Controller

Crossbar Network

PCU

DRAM Controller

PCU

DRAM Controller

No Cache Coherence Issues

pim.add y, &x

# Host-side PEI Execution



pim.add y, &x

# Mechanism Summary

- ## Atomicity of PEIs
  - PIM directory implements reader-writer locks

- ## Locality-aware PEI execution
  - Locality monitor simulates cache replacement behavior

- ## Cache coherence for PEIs
  - Memory-side: back-invalidation/back-writeback
  - Host-side: no need for consideration

- ## Virtual memory for PEIs
  - Host processor performs address translation before issuing a PEI

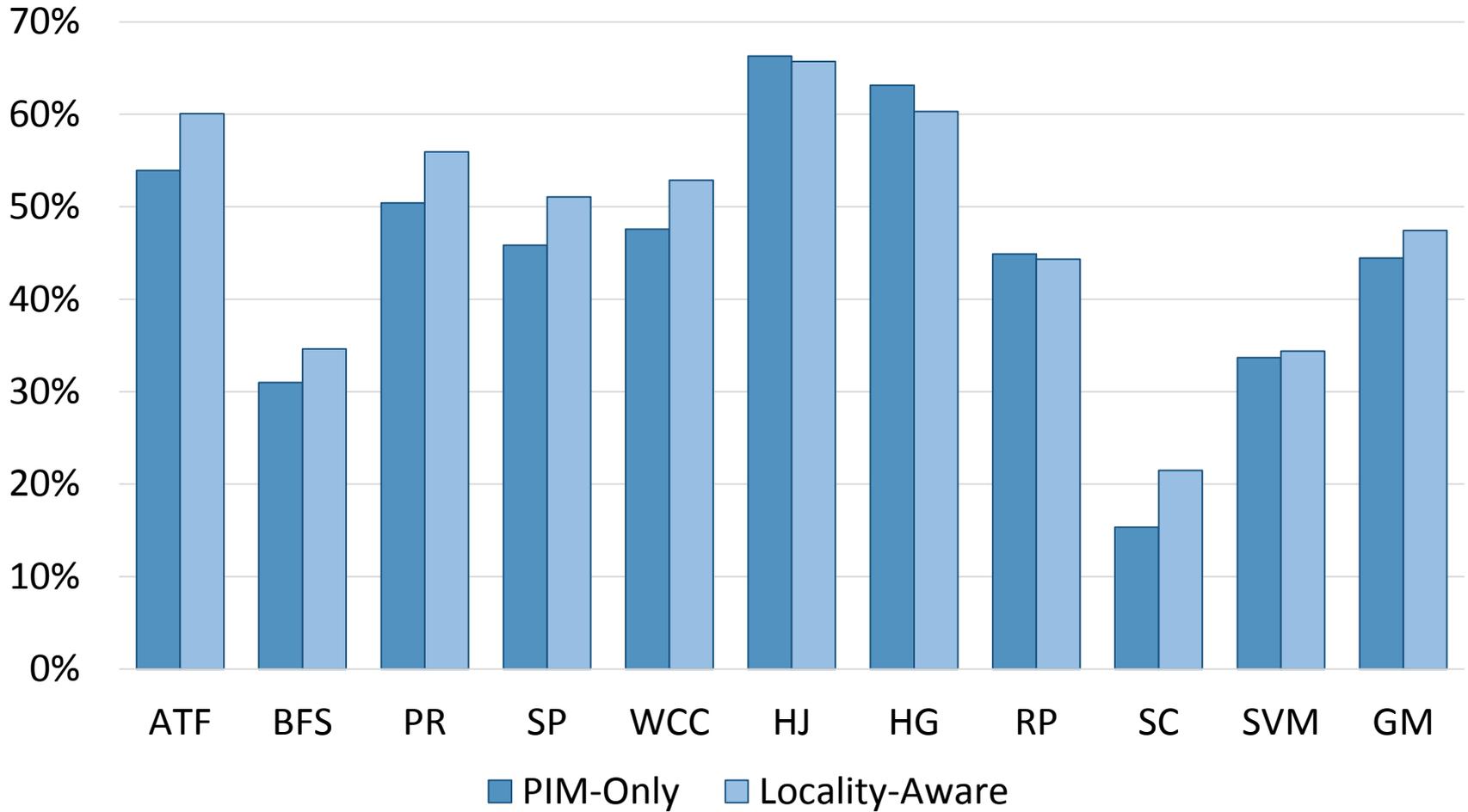# Simulation Configuration

- In-house x86-64 simulator based on Pin
  - 16 out-of-order cores, 4GHz, 4-issue
  - 32KB private L1 I/D-cache, 256KB private L2 cache
  - 16MB shared 16-way L3 cache, 64B blocks
  - 32GB main memory with 8 daisy-chained HMCs (80GB/s)

- PCU
  - 1-issue computation logic, 4-entry operand buffer
  - 16 host-side PCUs at 4GHz, 128 memory-side PCUs at 2GHz

- PMU
  - PIM directory: 2048 entries (3.25KB)
  - Locality monitor: similar to LLC tag array (512KB)

# Target Applications

- Ten emerging data-intensive workloads
  - Large-scale graph processing
    - Average teenage followers, BFS, PageRank, single-source shortest path, weakly connected components
  - In-memory data analytics
    - Hash join, histogram, radix partitioning
  - Machine learning and data mining
    - Streamcluster, SVM-RFE

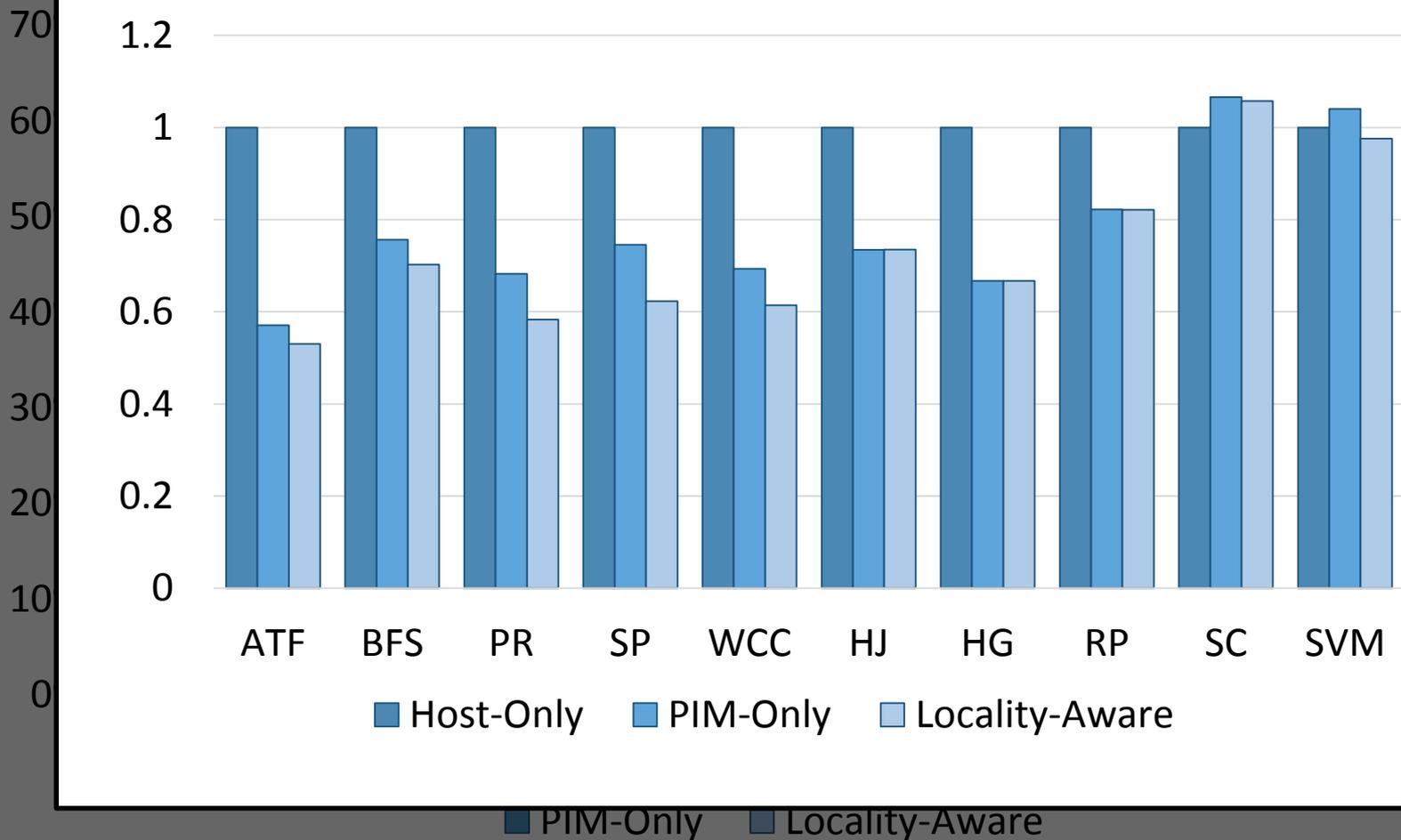- Three input sets (small, medium, large) for each workload to show the impact of data locality
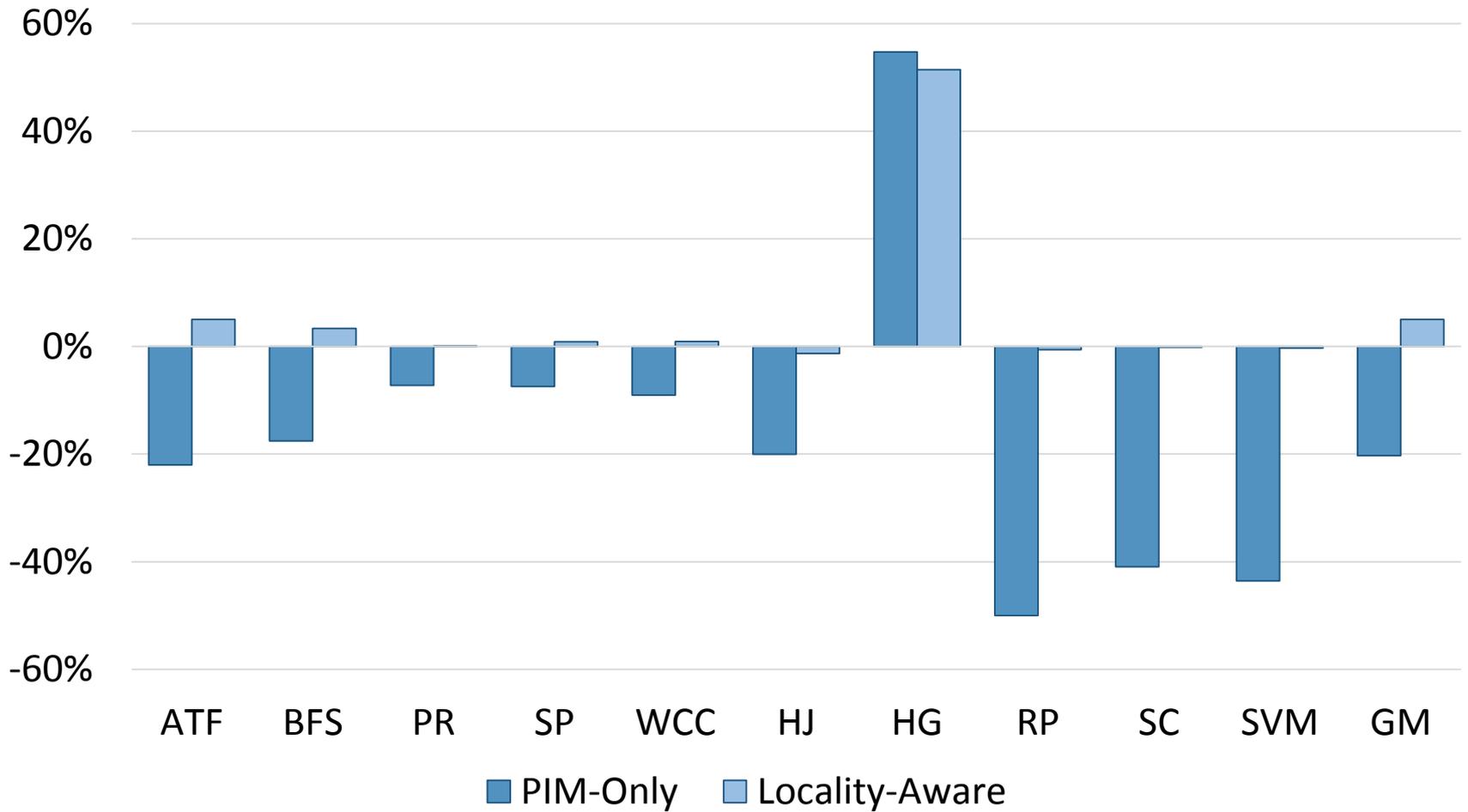
# Speedup

(Large Inputs, Baseline: Host-Only)
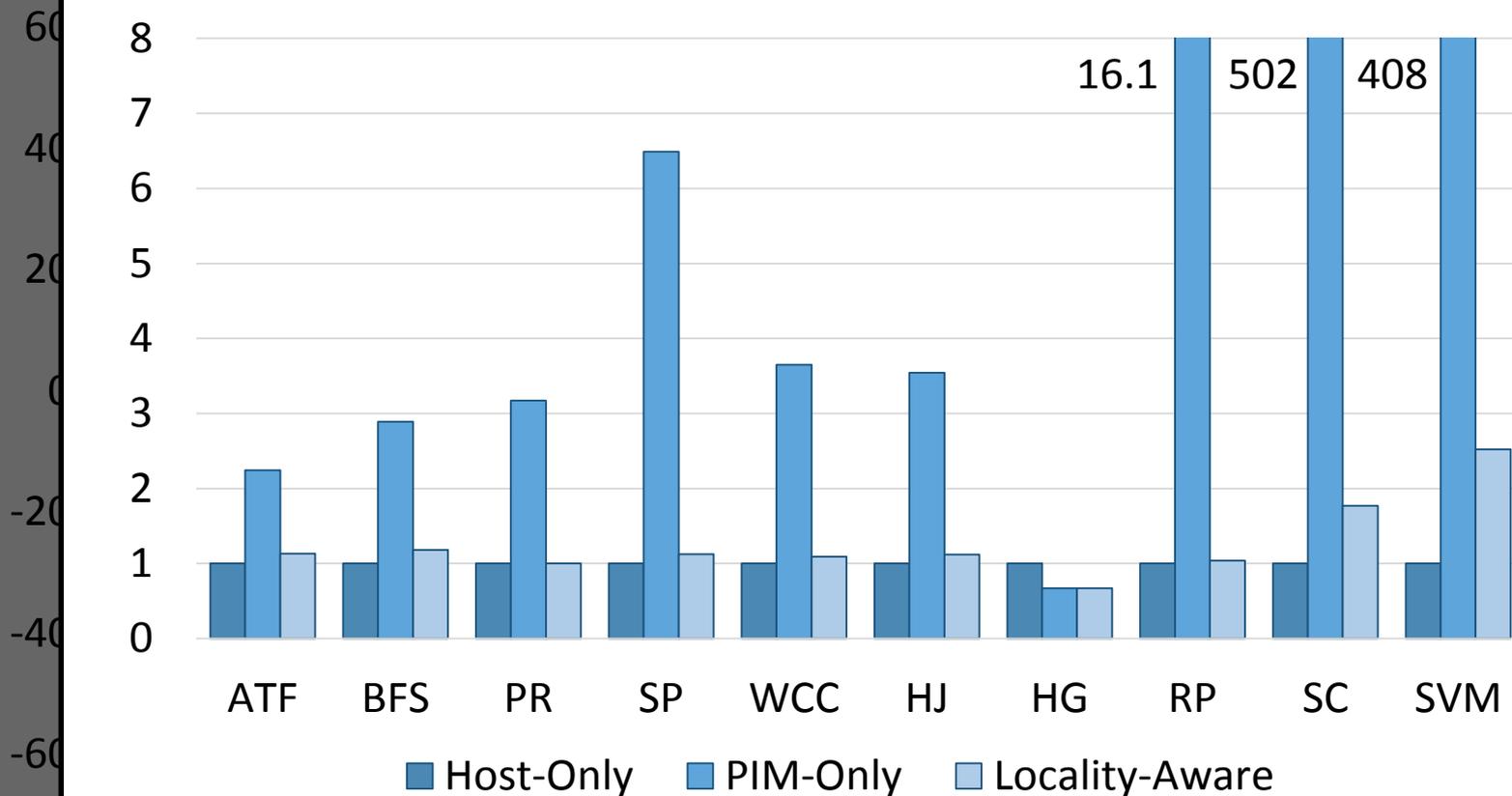
# Speedup



## Normalized Amount of Off-chip Transfer

Legend: ■ Host-Only ■ PIM-Only ☐ Locality-Aware

Categories: ATF, BFS, PR, SP, WCC, HJ, HG, RP, SC, SVM

■ PIM-Only ■ Locality-Aware

# Speedup

## (Small Inputs, Baseline: Host-Only)

# Speedup



## Normalized Amount of Off-chip Transfer

Values shown at top: 16.1, 502, 408

Categories: ATF, BFS, PR, SP, WCC, HJ, HG, RP, SC, SVM

Legend: Host-Only, PIM-Only, Locality-Aware

PIM-Only    Locality-Aware

# Speedup

## (Medium Inputs, Baseline: Host-Only)

Chart showing speedup percentages for PIM-Only and Locality-Aware across benchmarks: ATF, BFS, PR, SP, WCC, HJ, HG, RP, SC, SVM, GM.

Legend: ■ PIM-Only  ■ Locality-Aware

# Sensitivity to Input Size

# Multiprogrammed Workloads



Host-Only

PIM-Only — Locality-Aware

# Energy Consumption



Legend: Cache, HMC Link, DRAM, Host-side PCU, Memory-side PCU, PMU

Categories: Small, Medium, Large

Bar groups labeled: Host-Only, PIM-Only, Locality-Aware

# Conclusion

- Challenges of PIM architecture design
  - Cost-effective integration of logic and memory
  - Unconventional programming models
  - Lack of interoperability with caches and virtual memory

- PIM-enabled instruction: low-cost PIM abstraction & HW
  - Interfaces PIM operations as ISA extension
  - Simplifies cache coherence and virtual memory support for PIM
  - Locality-aware execution of PIM operations

- Evaluations
  - 47%/32% speedup over Host/PIM-Only in large/small inputs
  - Good adaptivity across randomly generated workloads

# PIM-Enabled Instructions: A Low-Overhead, Locality-Aware PIM Architecture

Junwhan Ahn, Sungjoo Yoo, Onur Mutlu[+], and Kiyoung Choi

Seoul National University        [+]Carnegie Mellon University