

PEI: Processing-in-Memory-Enabled Instruction

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

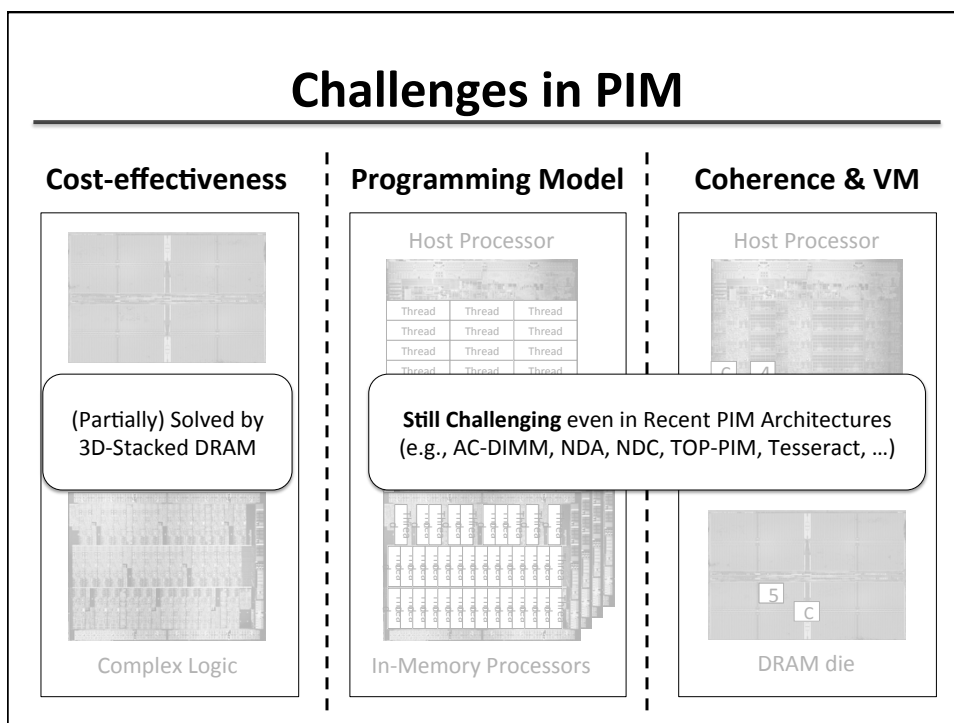
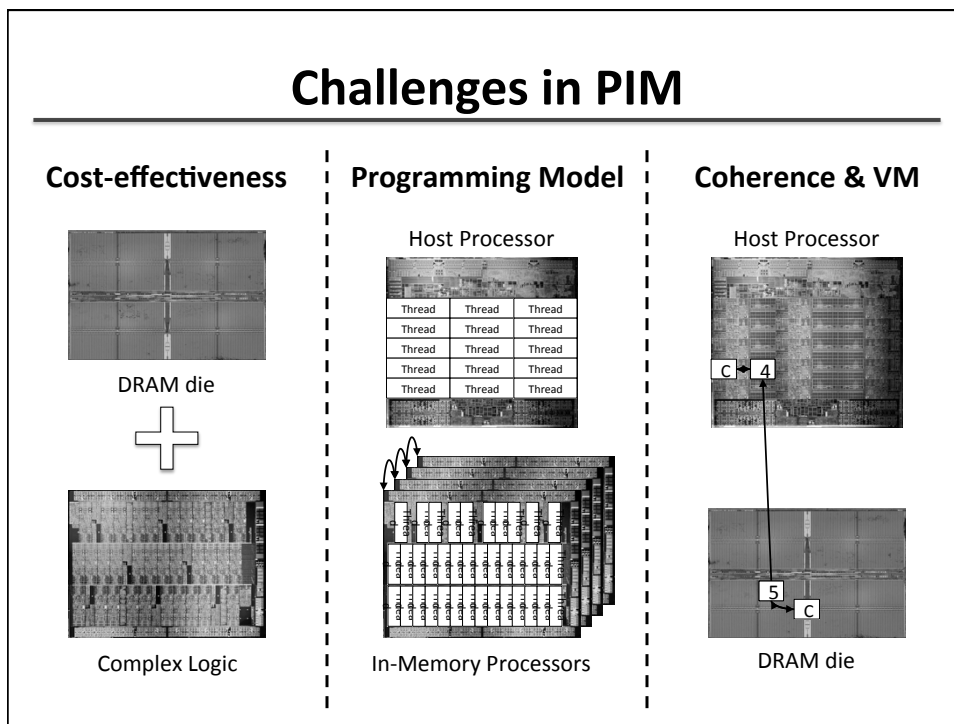
Seoul National University

⁺Carnegie Mellon University

Processing-in-Memory (PIM)

- Move computations 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?



A New Direction of PIM

- Objectives
 - Reduce the implementation overhead of PIM units
 - Provide an intuitive programming model for PIM
 - Full support for cache coherence and virtual memory
- Our solution: simple PIM operations as ISA extension
 - Simple: low-overhead implementation
 - ISA extension: intuitive programming model

Potential of ISA Extension as the PIM Interface

- Example: Parallel PageRank (PR) computation

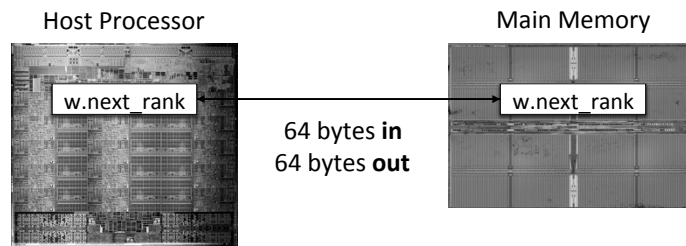
```
for (v: graph.vertices) {  
  for (w: v.successors) {  
    w.next_rank += weight * v.rank;  
  }  
}  
for (v: graph.vertices) {  
  v.rank = v.next_rank; v.next_rank = alpha;  
}
```

Potential of ISA Extension as the PIM Interface

```

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

```



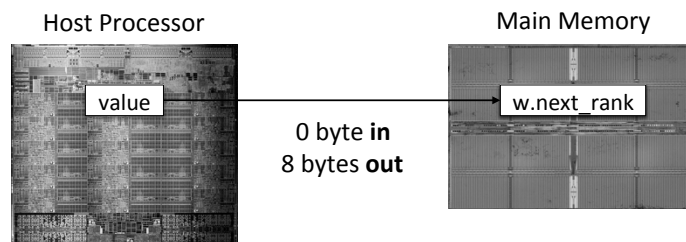
Conventional Architecture

Potential of ISA Extension as the PIM Interface

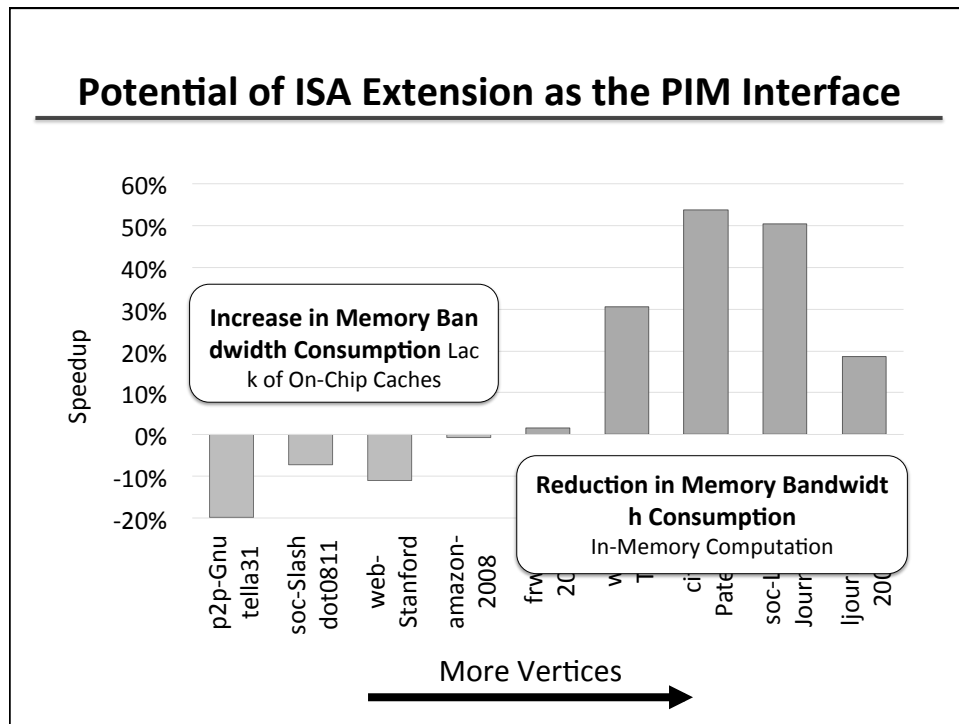
```

for (v: graph.vertices) {
  for (w: v.successors) {
    PIM_add(&w.next_rank, weight * v.rank);
  }
}

```



In-Memory Addition



Overview

1. How should simple PIM operations be interfaced to conventional systems?
 - Expose PIM operations as *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) {  
  for (w: v.successors) {  
    w.next_rank += weight * v.rank;  
  }  
}
```

PIM-Enabled Instructions

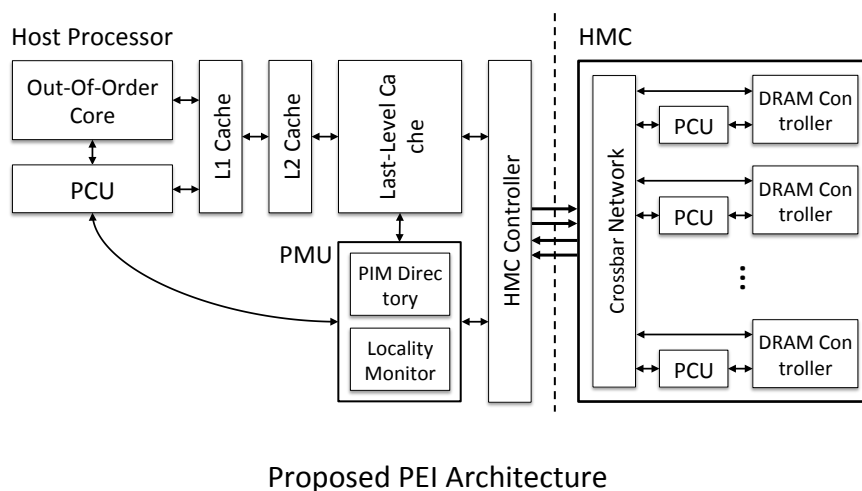
```
for (v: graph.vertices) {  
  for (w: v.successors) {  
    PIM_add(&w.next_rank, weight * v.rank);  
  }  
}  
pfence();
```

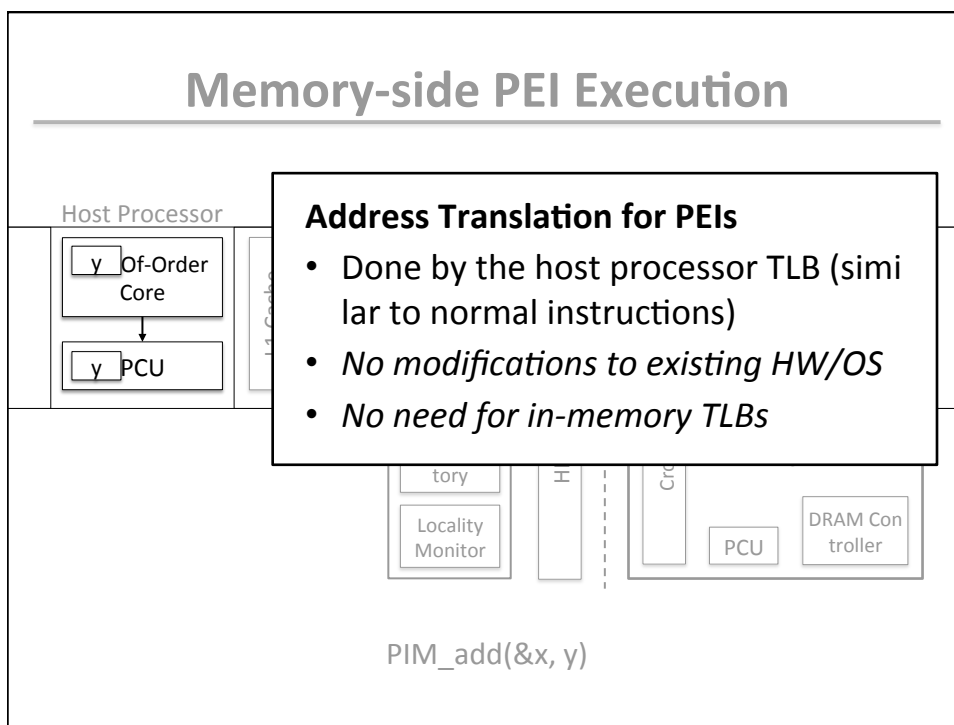
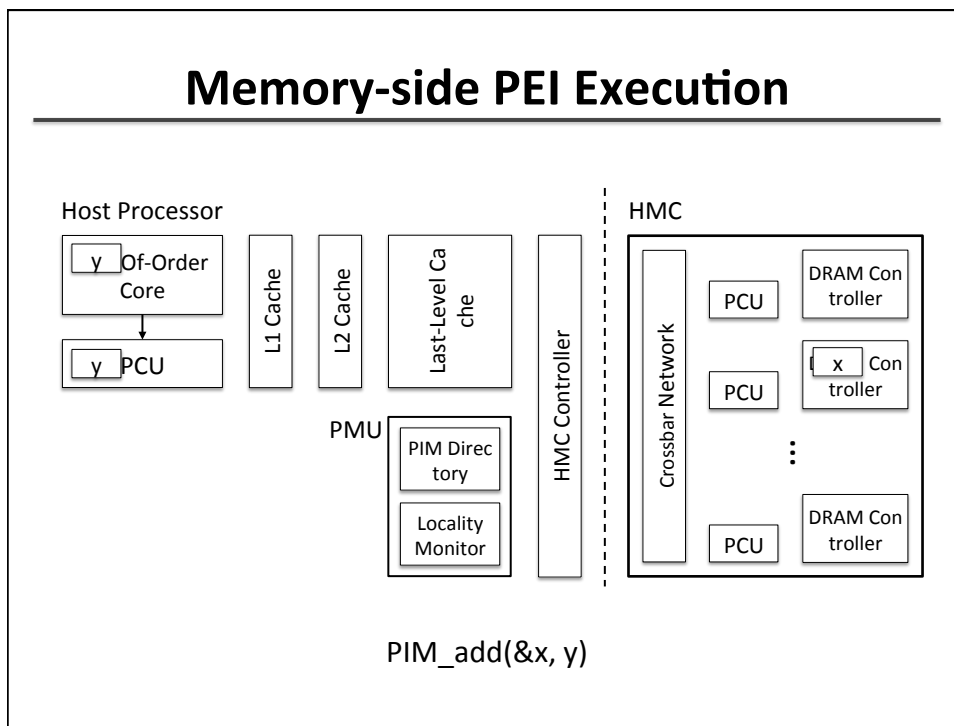
- Executed either in memory or in the host processor
- Cache-coherent, virtually-addressed
- Atomic between different PEIs
- *Not* atomic with normal instructions (use *pfence*)

PIM-Enabled Instructions

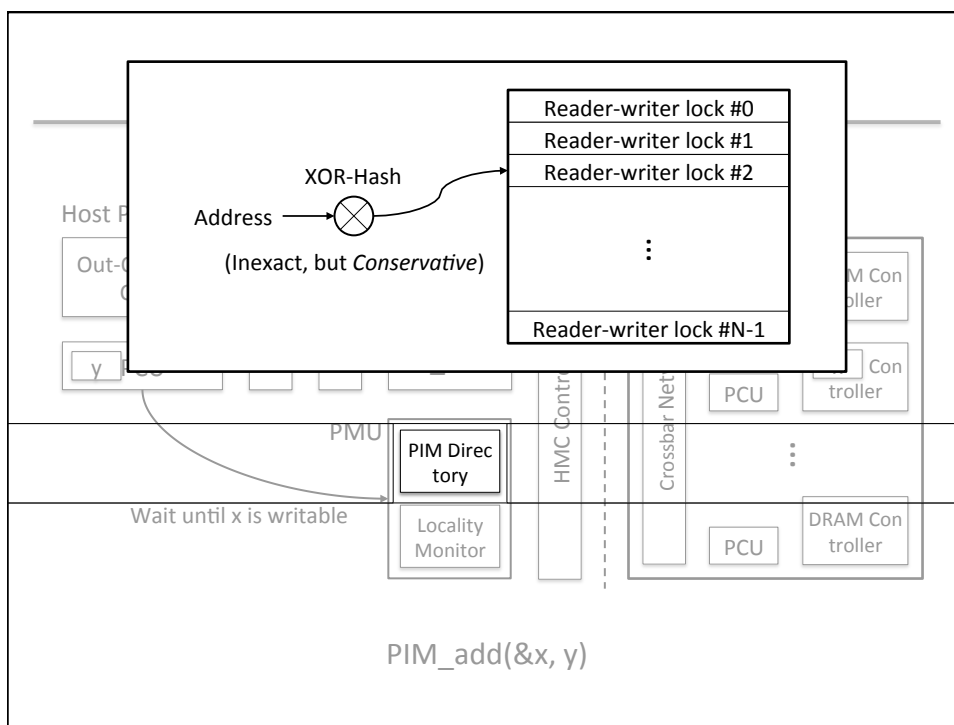
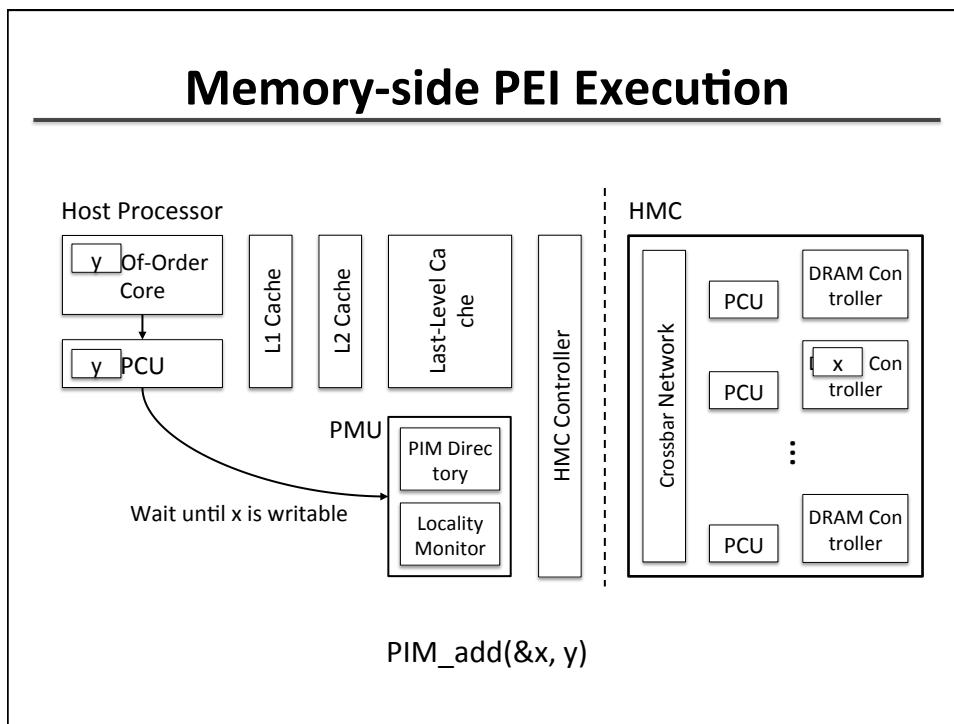
- Single-cache-block restriction
 - Each PEI can access *at most one last-level cache block*
 - Similar restrictions exist in atomic instructions
- Benefits
 - **Localization:** each PEI is bounded to one memory module
 - **Interoperability:** easier support for cache coherence and virtual memory
 - **Simplified locality monitoring:** data locality of PEIs can be identified by LLC tag checks or similar methods

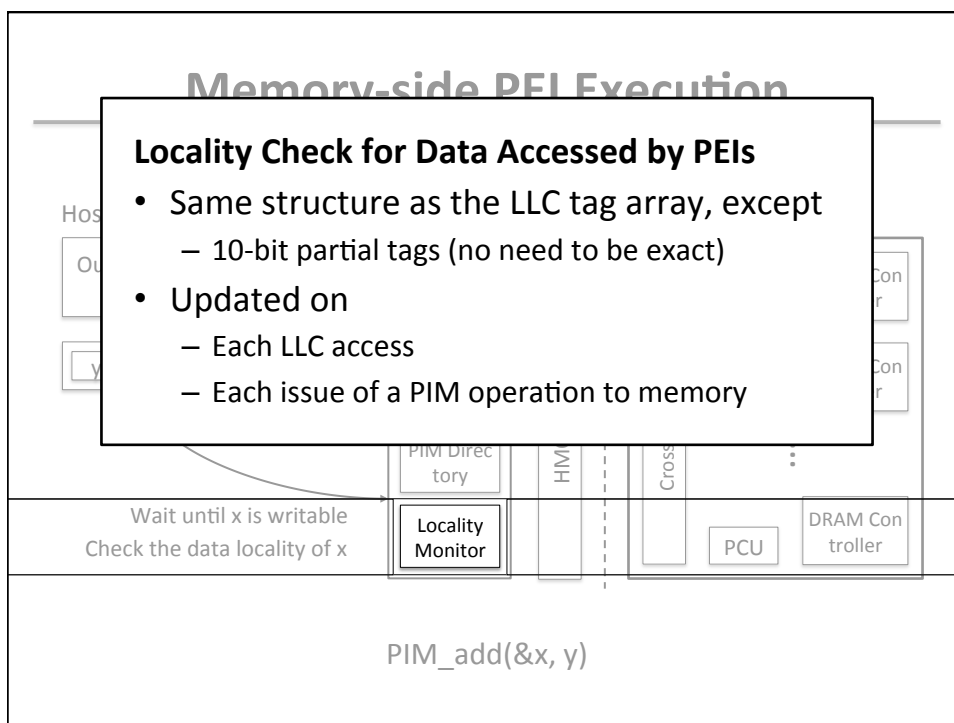
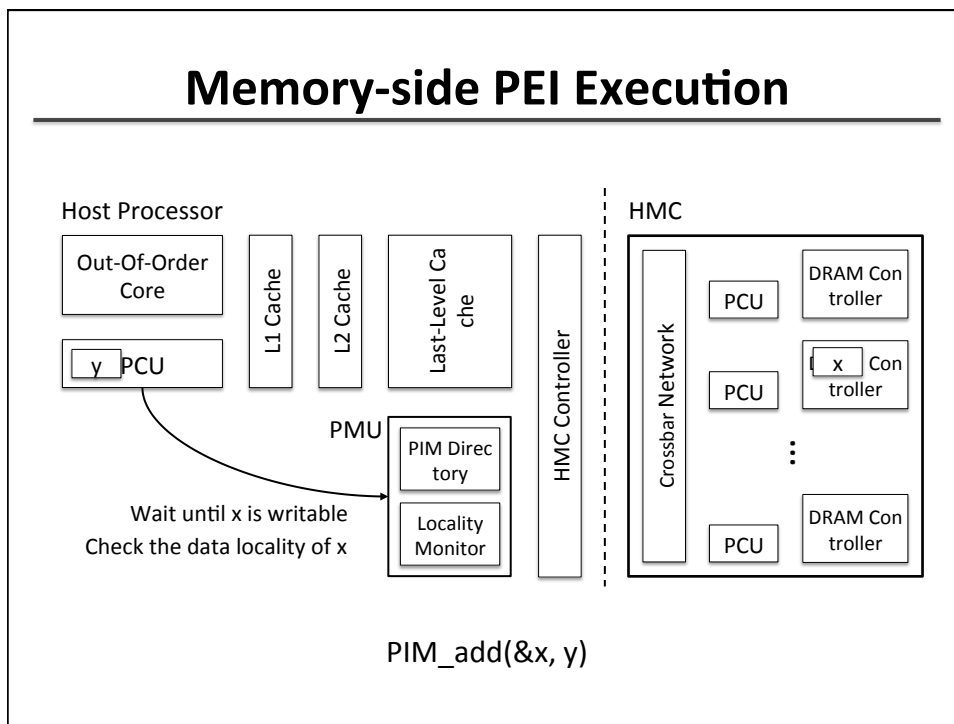
Architecture



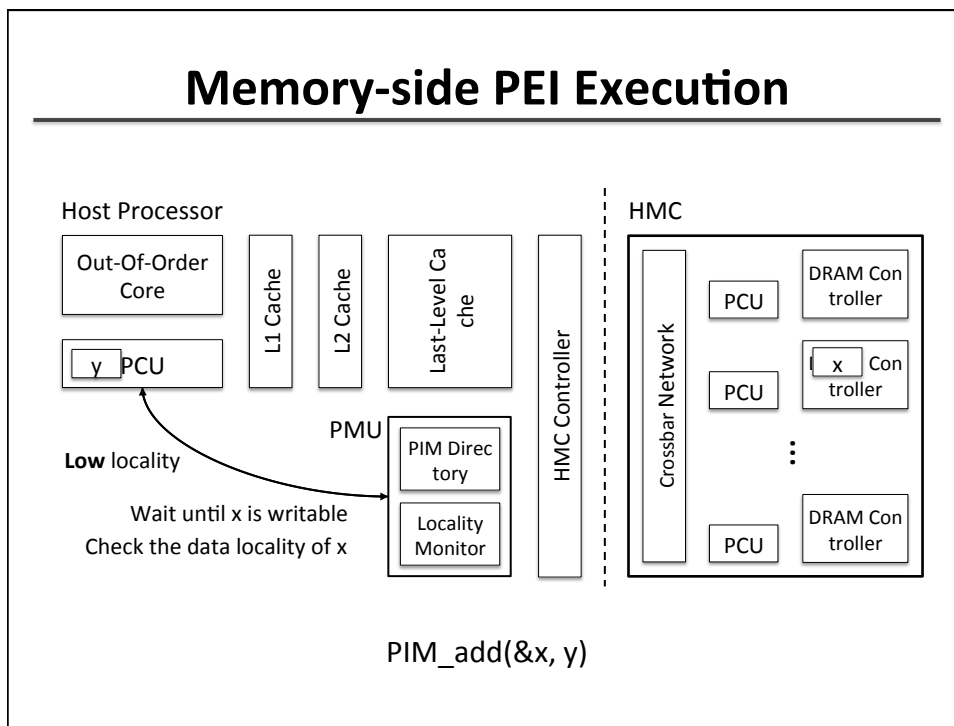


Memory-side PEI Execution

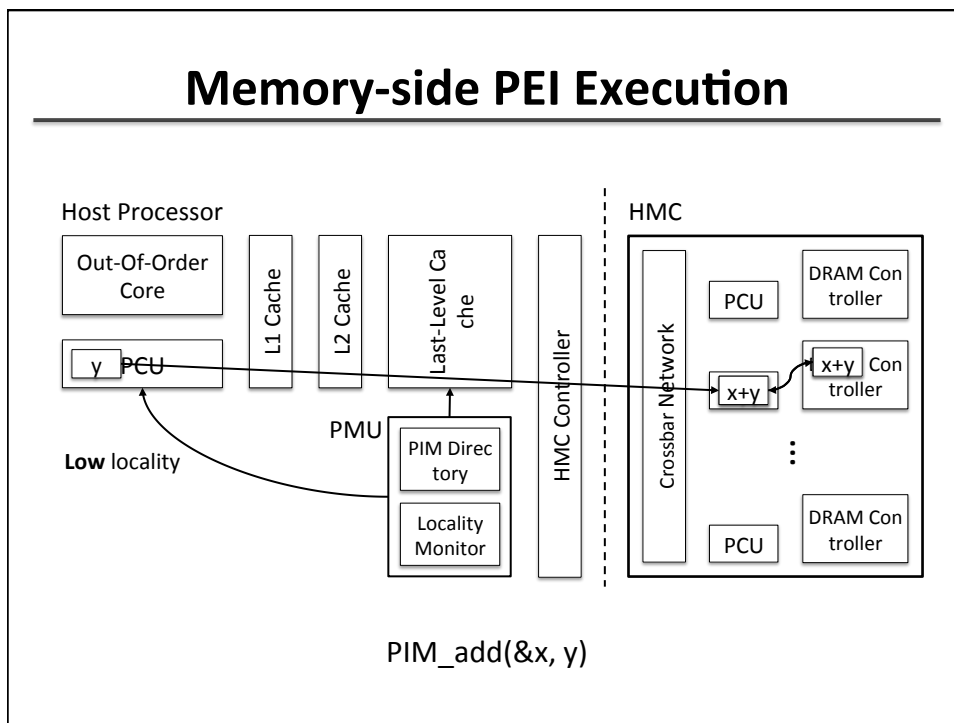




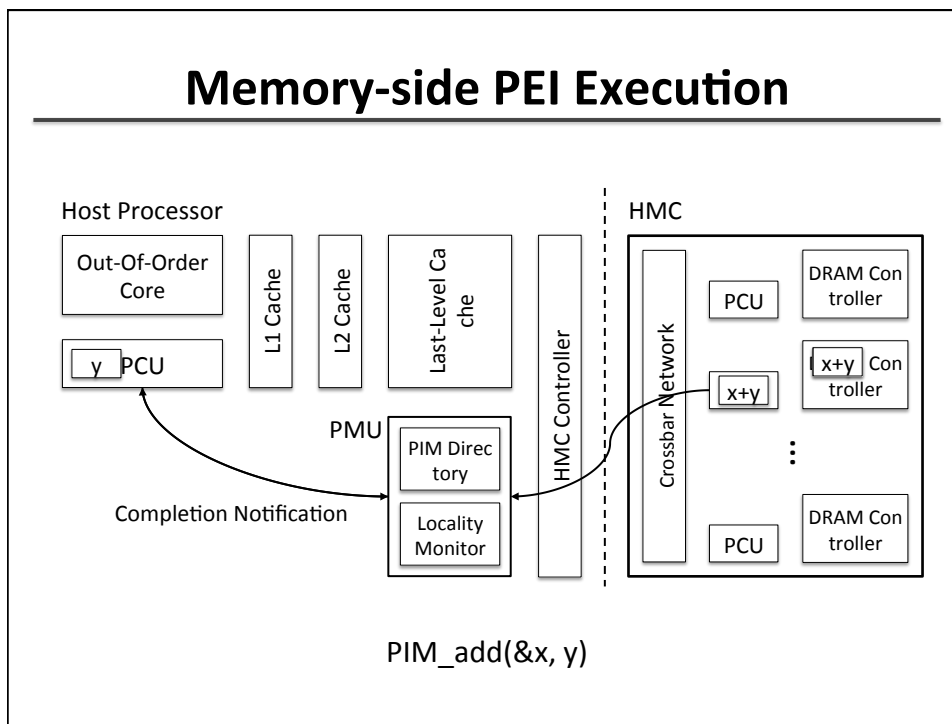
Memory-side PEI Execution



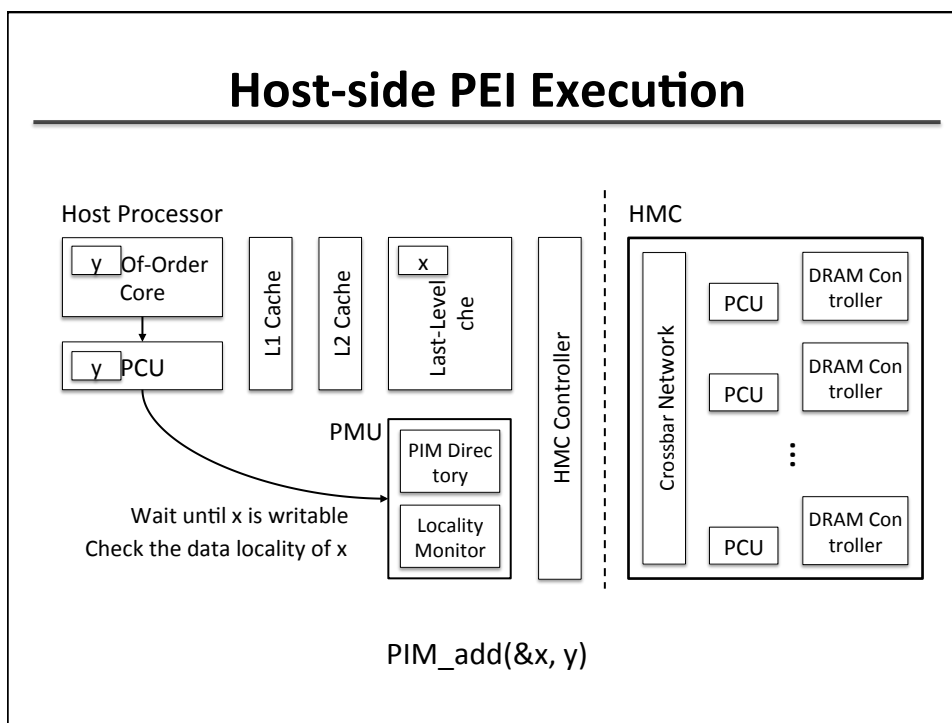
Memory-side PEI Execution



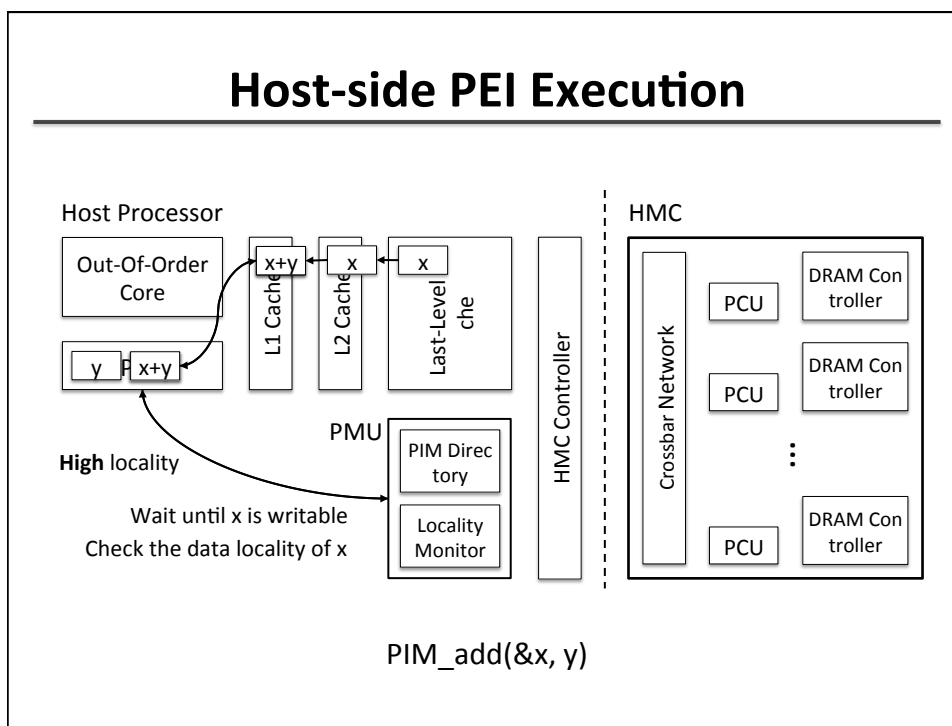
Memory-side PEI Execution



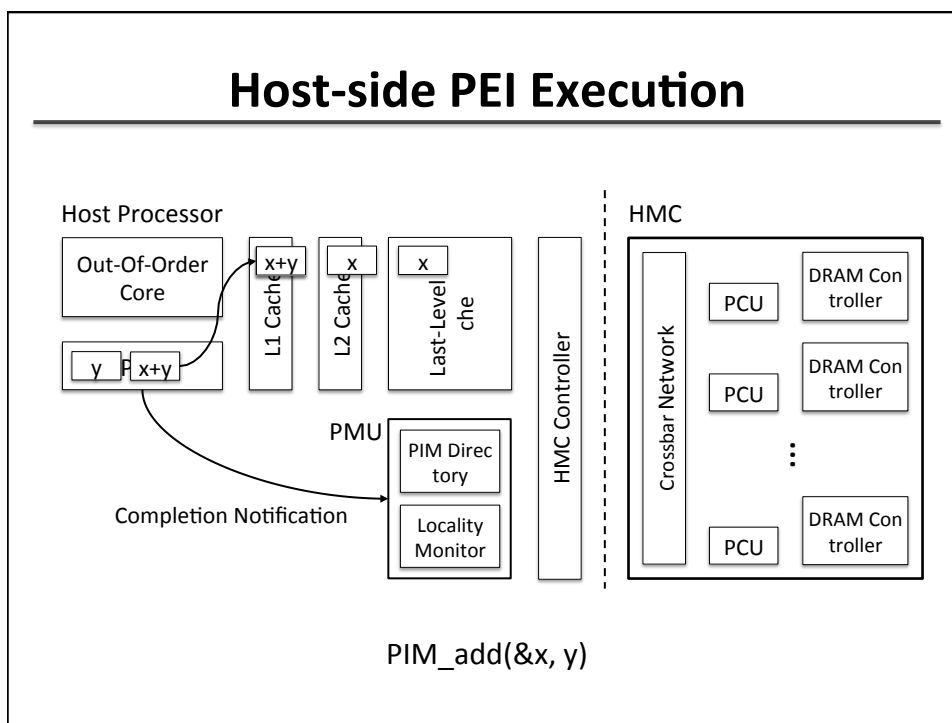
Host-side PEI Execution



Host-side PEI Execution



Host-side PEI Execution

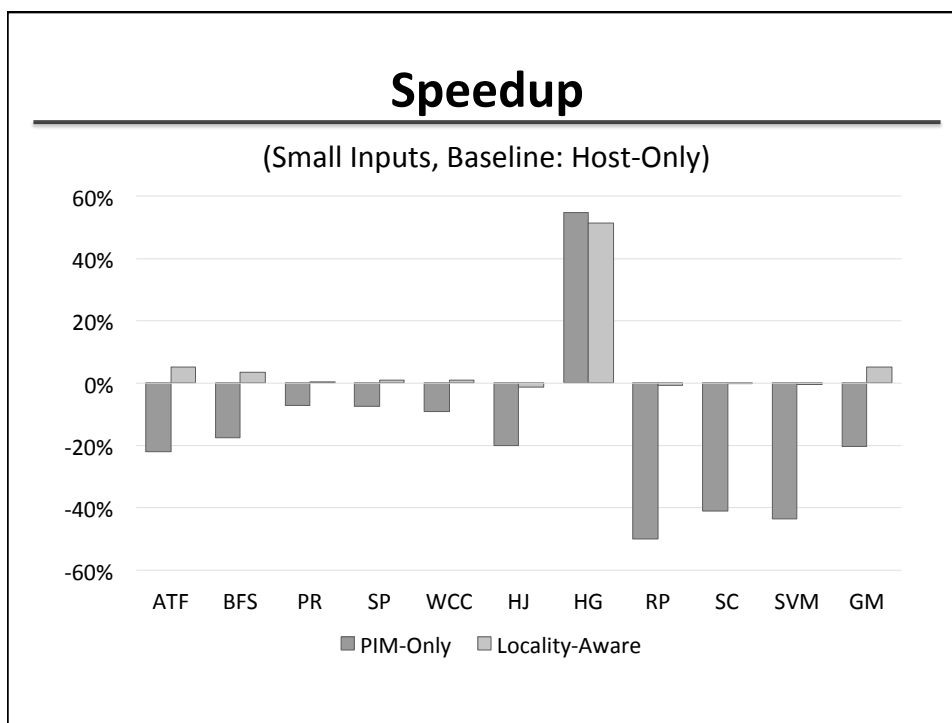
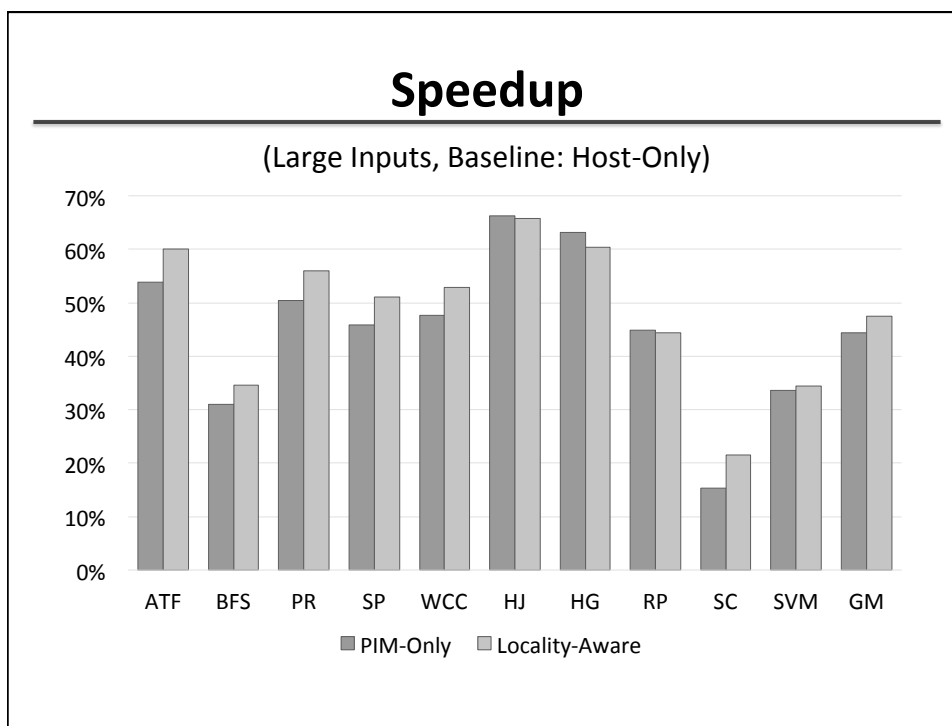


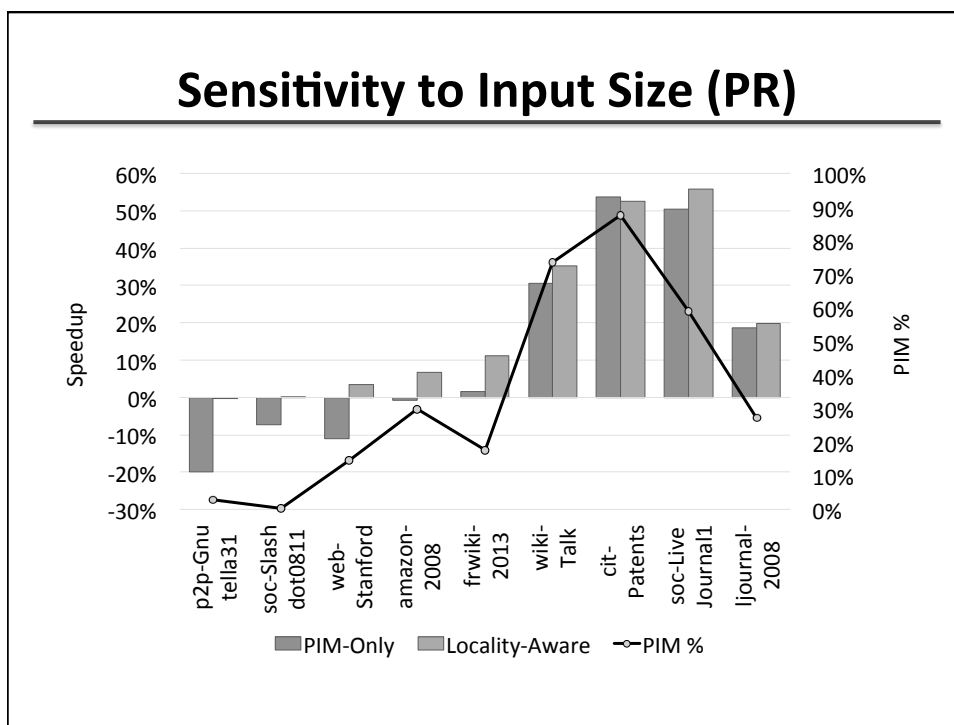
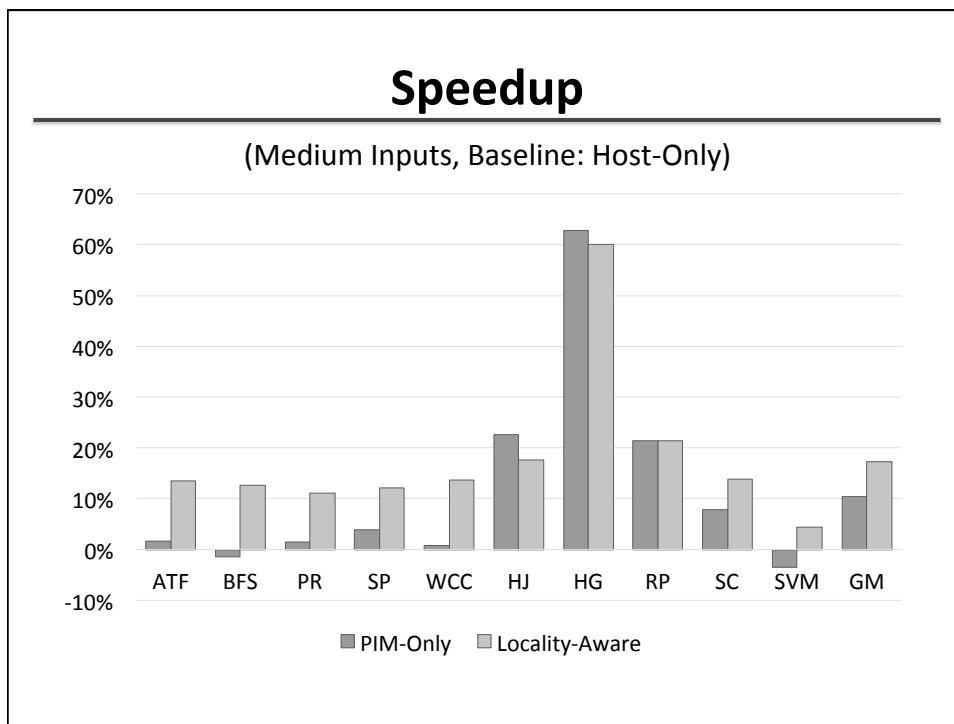
Simulation Configuration

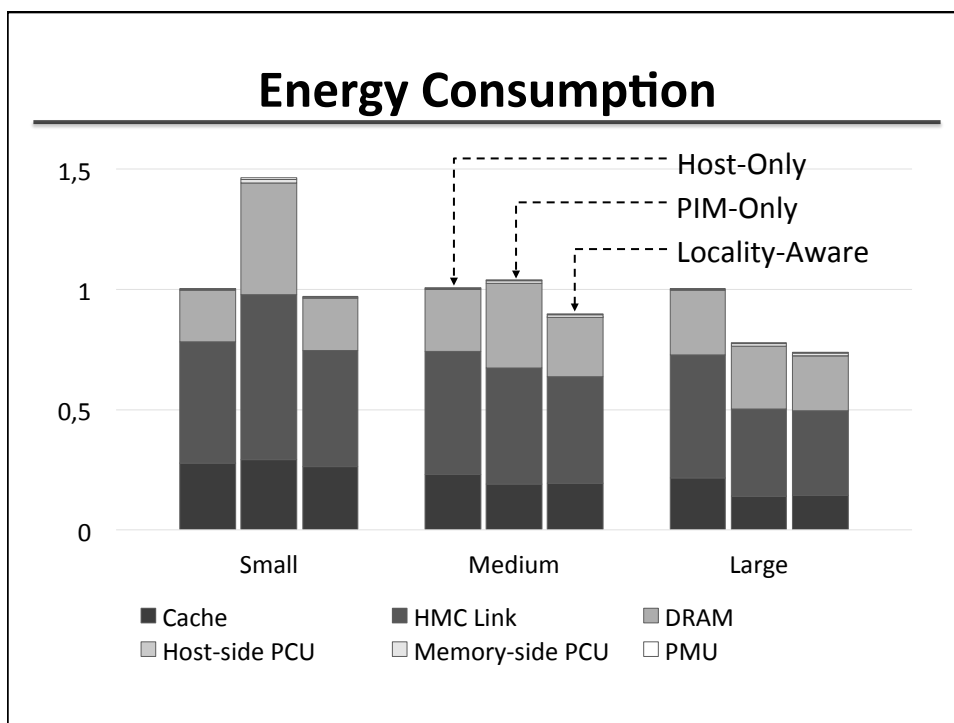
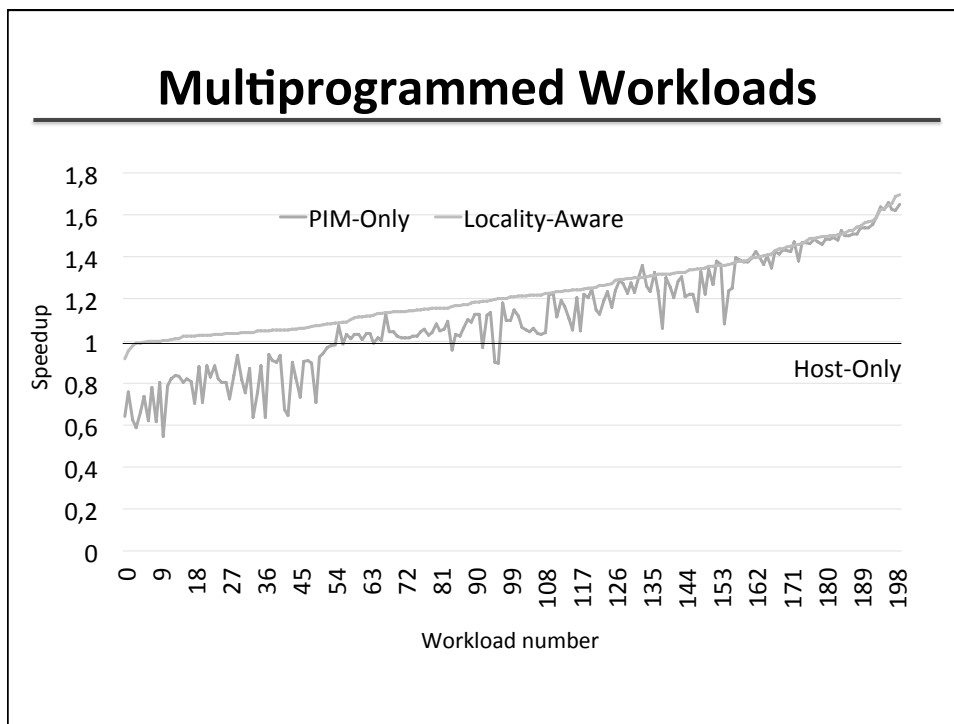
- In-house x86-64 simulator based on Pin
 - 16 out-of-order cores, 4GHz, 4-issue
 - 32KB private L1 I/D-cache, 256KB private L2 cache
 - 16MB shared 16-way L3 cache, 64B blocks
 - 32GB main memory with 8 daisy-chained HMCs (80GB/s)
- PCU
 - 1-issue computation logic, 4-entry operand buffer
 - 16 host-side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







Conclusion

- Challenges of PIM architecture design
 - Costly integration of logic and memory
 - Unconventional programming models
 - Lack of interoperability with caches and virtual memory
- PIM-enabled instruction
 - low-cost PIM implementation
 - 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