# A Scalable Processing-in-Memory Accelerator for Parallel Graph Processing

Junwhan Ahn, Sungpack Hong*, Sungjoo Yoo, Onur Mutlu+, Kiyoung Choi

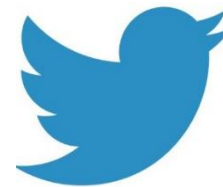Seoul National University    *Oracle Labs    +Carnegie Mellon University

# Graphs

- Abstract representation of object relationships
  - Vertex: object (e.g., person, article, …)
  - Edge: relationship (e.g., friendships, hyperlinks, …)

- Recent trend: explosive increase in graph size



| 36 Million Wikipedia Pages | 1.4 Billion Facebook Users | 300 Million Twitter Users | 30 Billion Instagram Photos |

# Large-Scale Graph Processing

- Example: Google's PageRank

$$R[i] = \alpha + \sum_{j \in \text{Succ}(i)} w_{ji} R[j]$$

```
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;
}
```
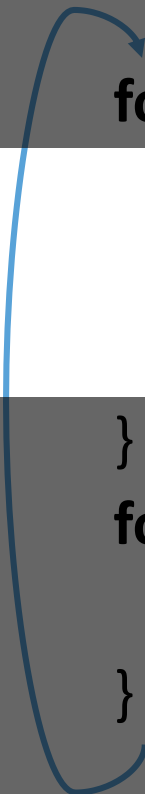
# Large-Scale Graph Processing

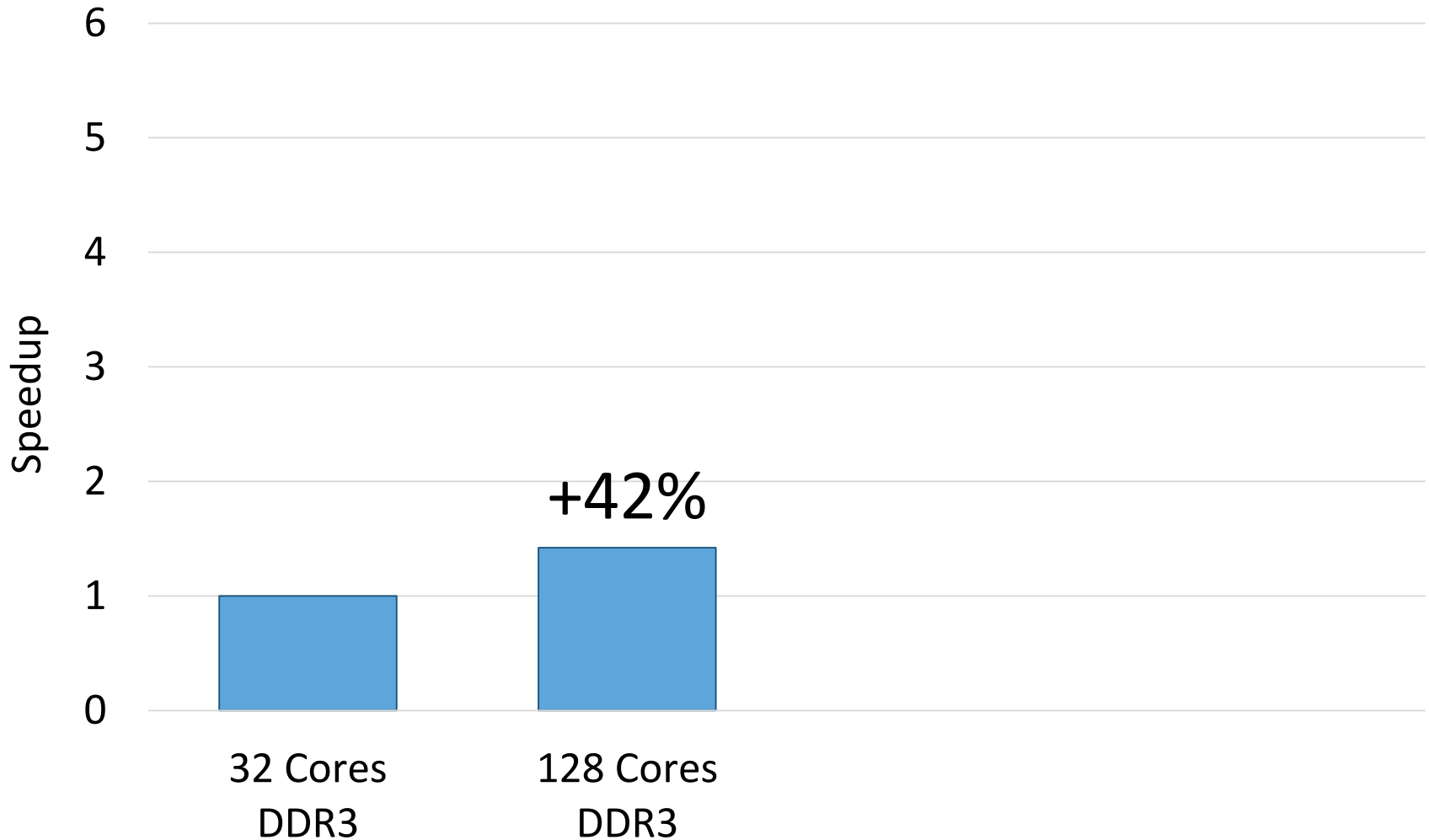- Example: Google's PageRank

Independent to Each Vertex
**Vertex-Parallel Abstraction**

```
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;
}
```
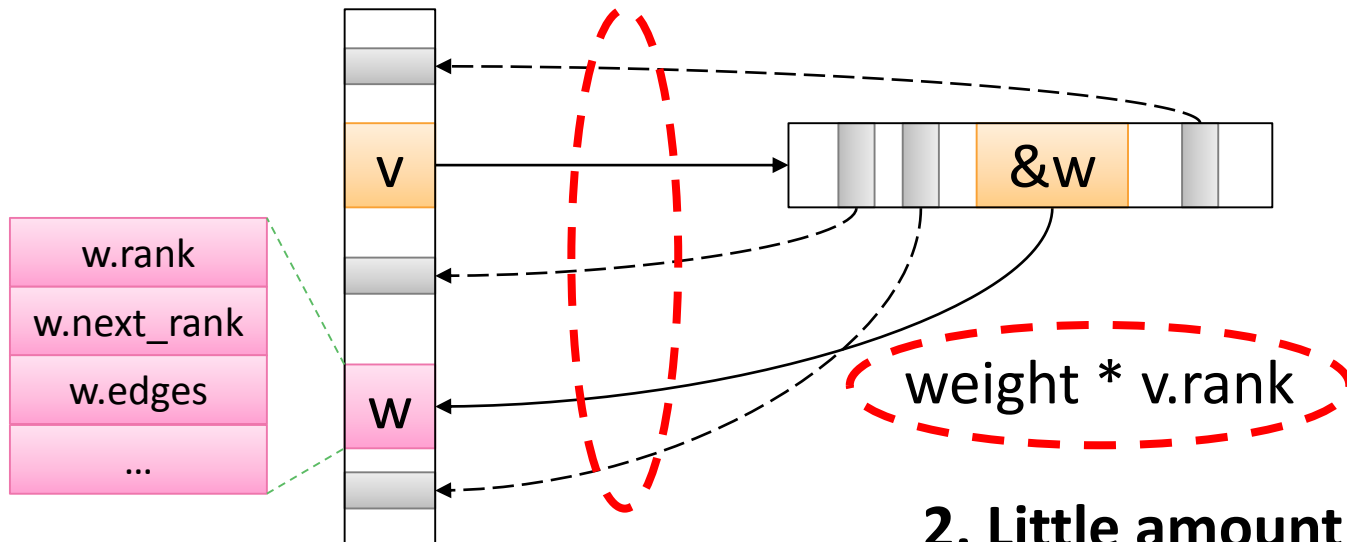
# PageRank Performance

# Bottleneck of Graph Processing

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

**1. Frequent random memory accesses**
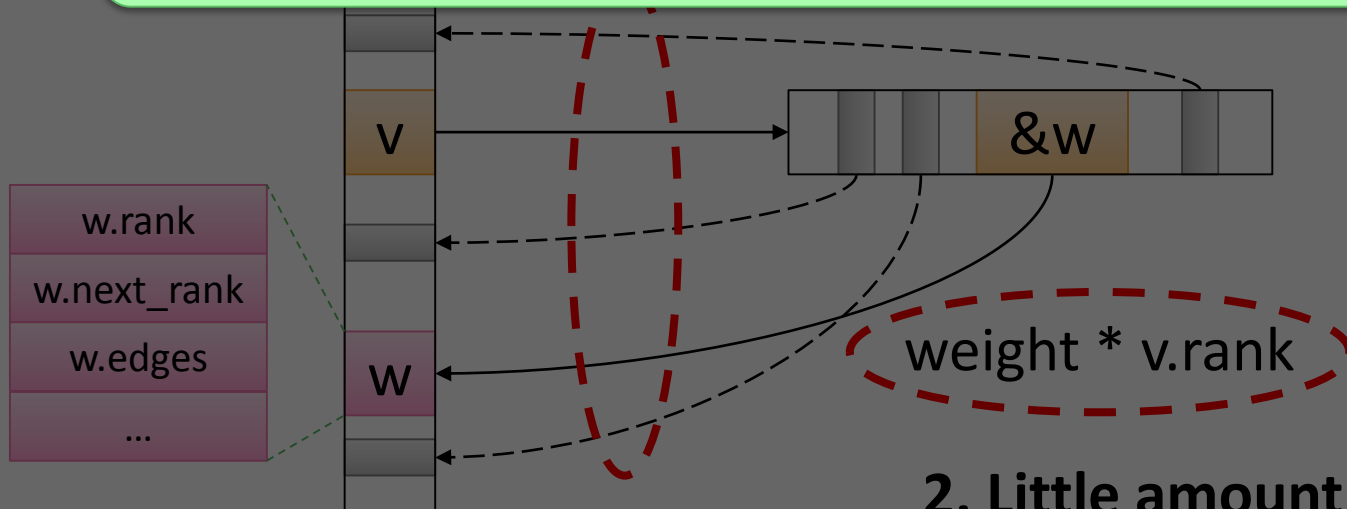
**2. Little amount of computation**

# Bottleneck of Graph Processing

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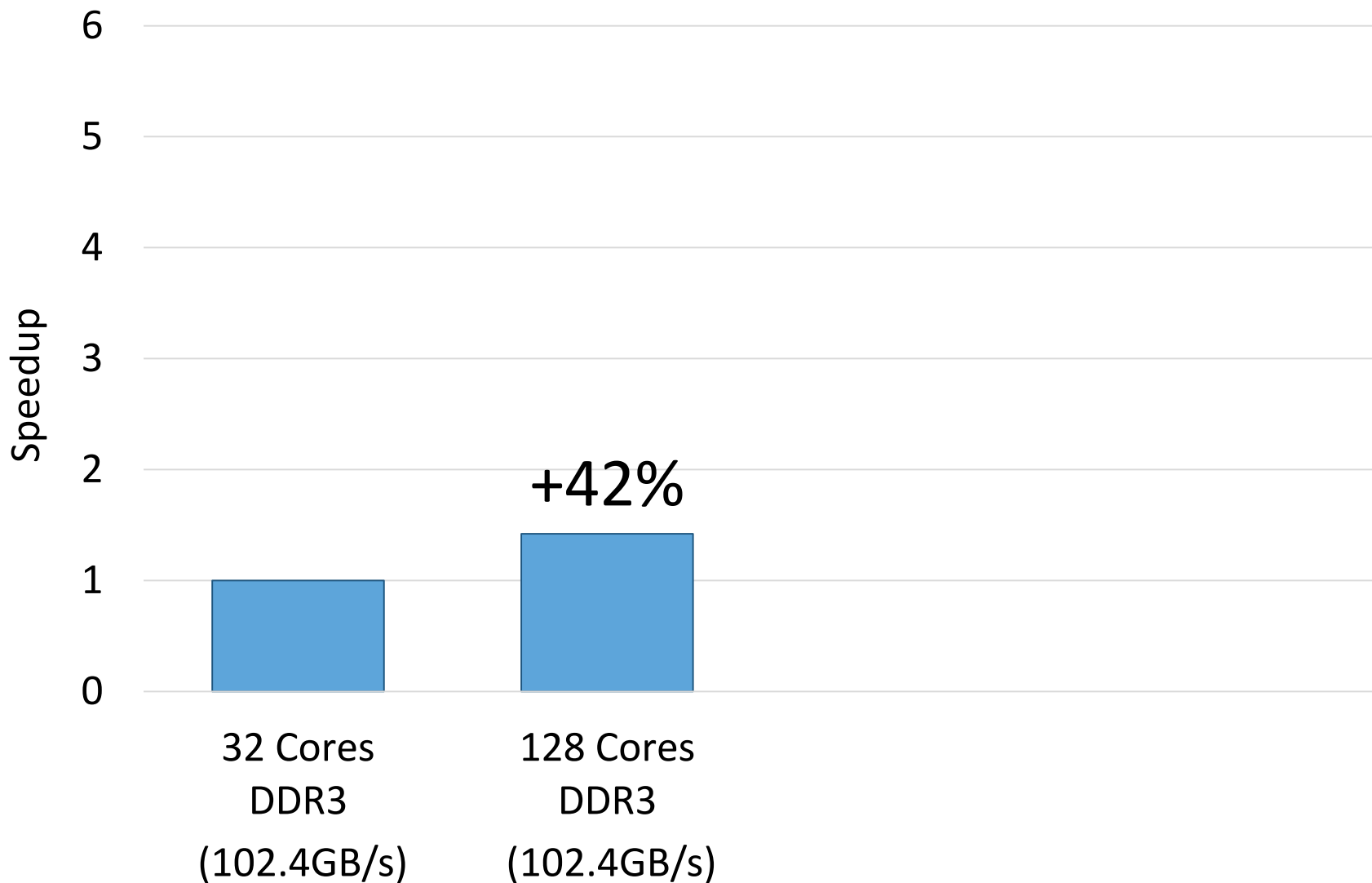
High Memory Bandwidth Demand



w.rank
w.next_rank
w.edges
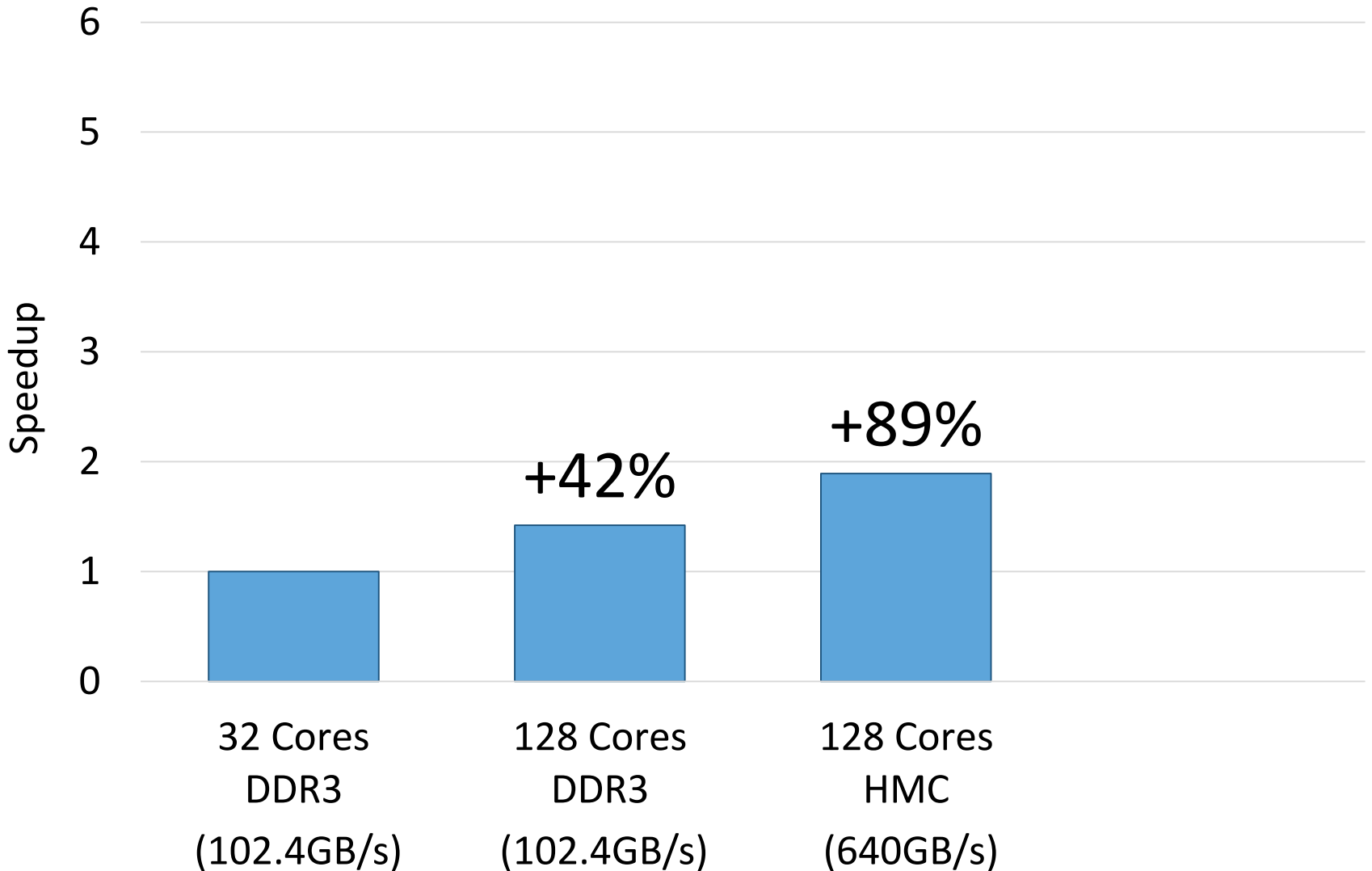…

v

&w

w

weight * v.rank

**2. Little amount of computation**

# PageRank Performance

# PageRank Performance

# PageRank Performance



Speedup

| | 89% |
|---|---|

32 Cores
DDR3
(102.4GB/s)

128 Cores
DDR3
(102.4GB/s)

128 Cores
HMC
(640GB/s)

# PageRank Performance



**Lessons Learned:**

1. *High memory bandwidth* is the key to the scalability of graph processing
2. Conventional systems do not *fully utilize* high memory bandwidth

Speedup

5.3x

+89%

+42%

| | | | |
| 6 | | | |
| 5 | | | |
| 4 | | | |
| 3 | | | |
| 2 | | | |
| 1 | | | |
| 0 | | | |

32 Cores
DDR3
(102.4GB/s)

128 Cores
DDR3
(102.4GB/s)

128 Cores
HMC
(640GB/s)

128 Cores
HMC Internal BW
(8TB/s)

# Challenges in Scalable Graph Processing

- **Challenge 1**: How to provide *high memory bandwidth* to computation units in a practical way?
  - Processing-in-memory based on 3D-stacked DRAM

- **Challenge 2**: How to design computation units that *efficiently exploit large memory bandwidth*?
  - Specialized in-order cores called *Tesseract* cores
    - Latency-tolerant programming model
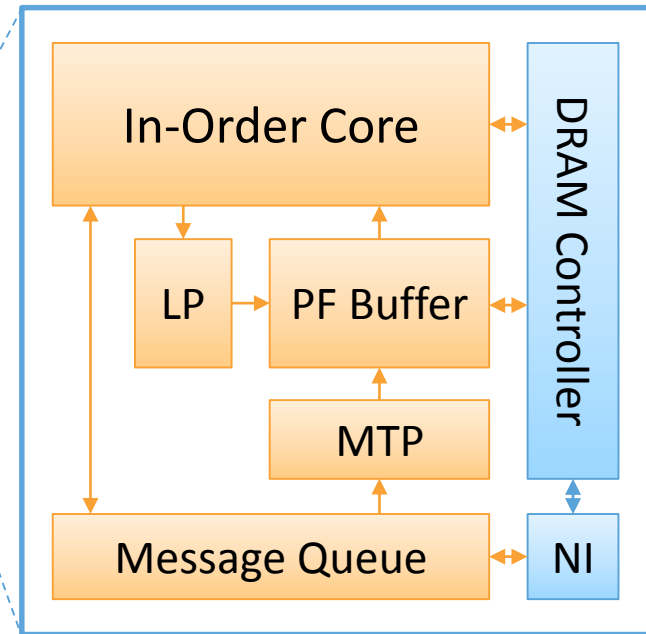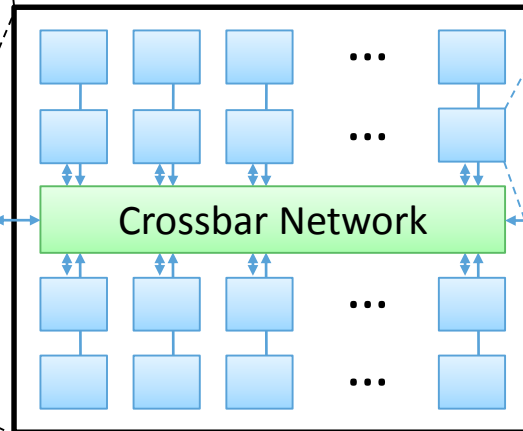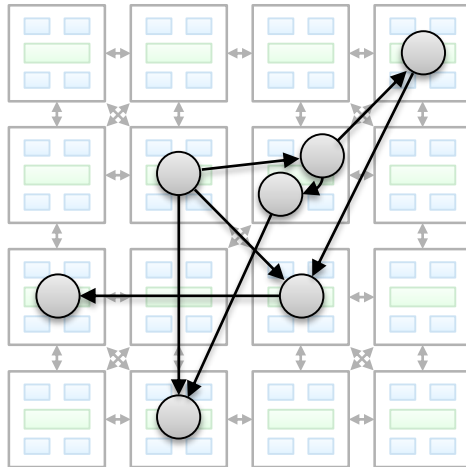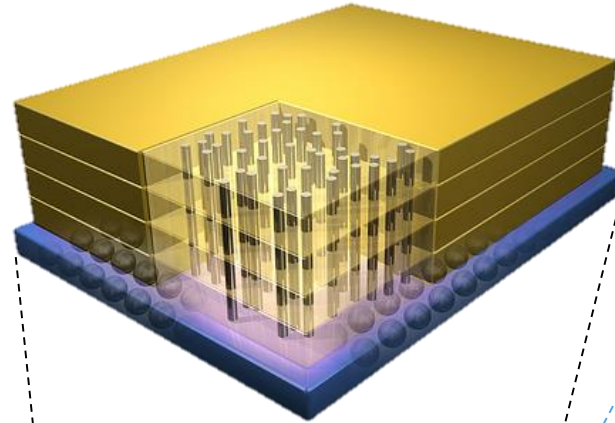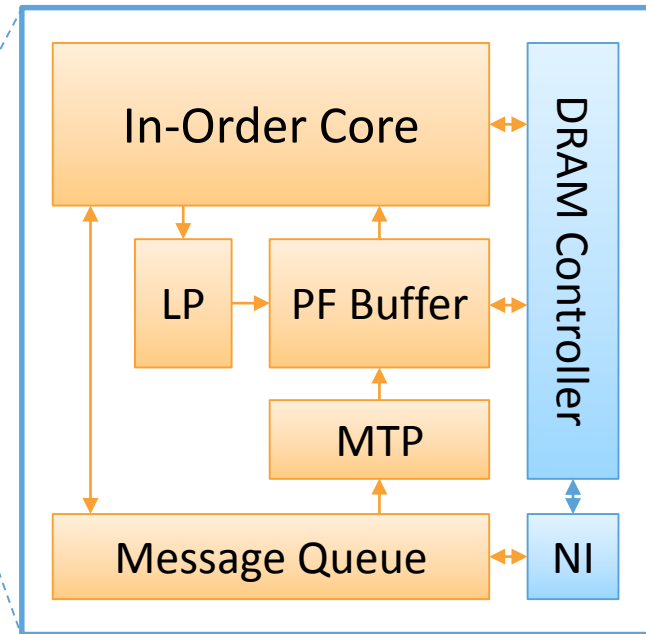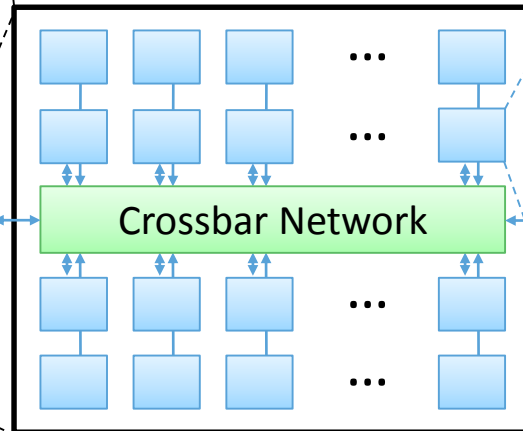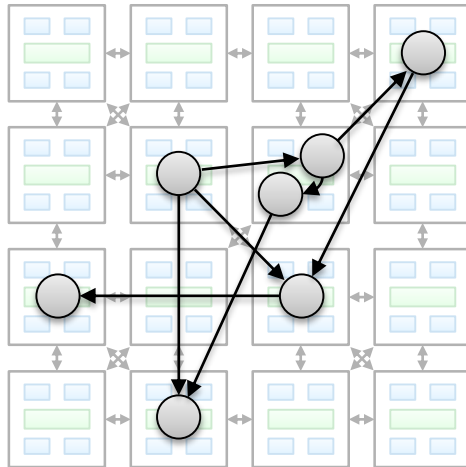    - Graph-processing-specific prefetching schemes
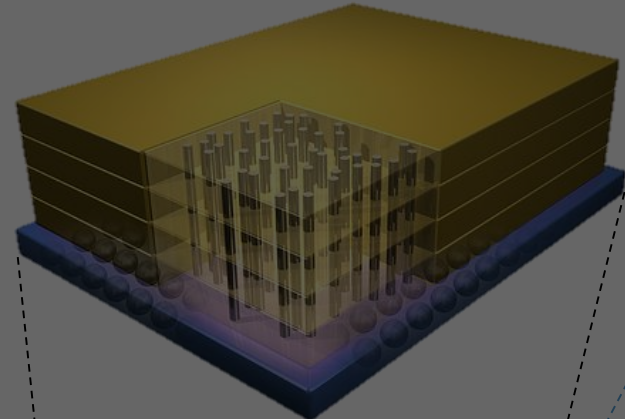
# Tesseract System



Host Processor

Memory-Mapped
Accelerator Interface
(Noncacheable, Physically Addressed)

In-Order Core

DRAM Controller

LP

PF Buffer

MTP

Message Queue
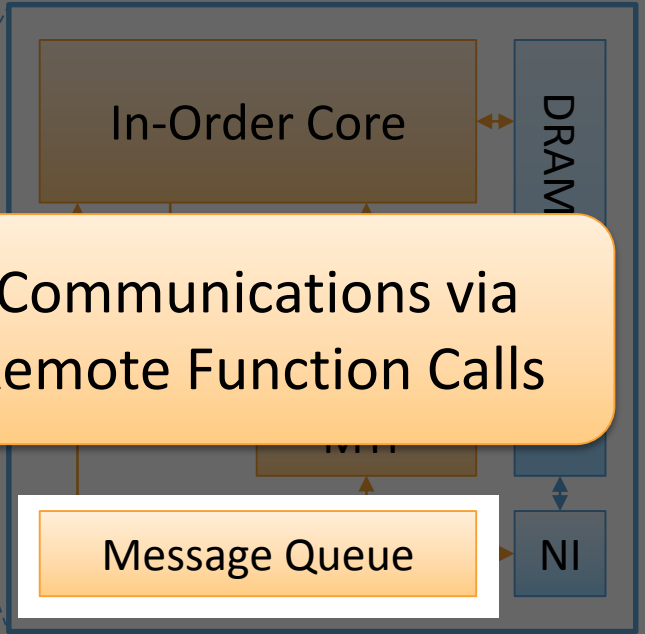
NI

Crossbar Network

# Tesseract System



Host Processor

Memory-Mapped
Accelerator Interface
(Noncacheable, Physically Addressed)

Crossbar Network

In-Order Core

DRAM Controller

LP

PF Buffer

MTP

Message Queue

NI

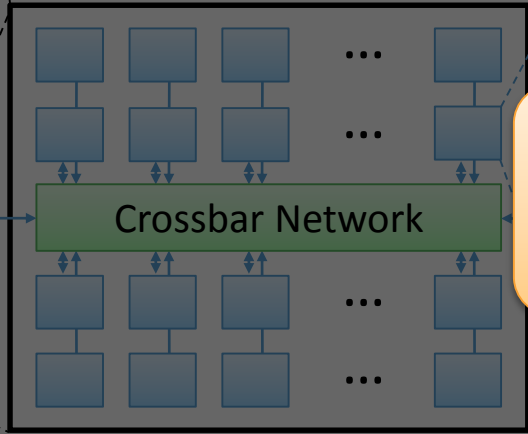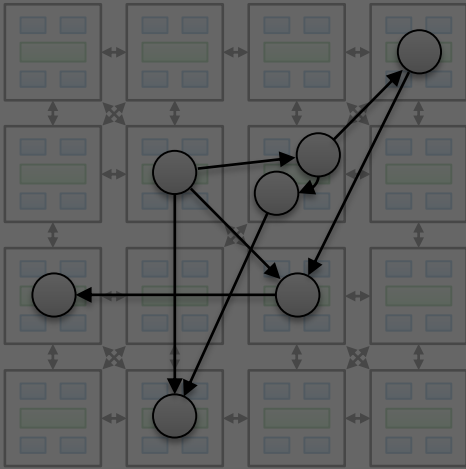# Tesseract System



Host Processor

Memory-Mapped
Accelerator Interface
(Noncacheable, Physically Addressed)

In-Order Core
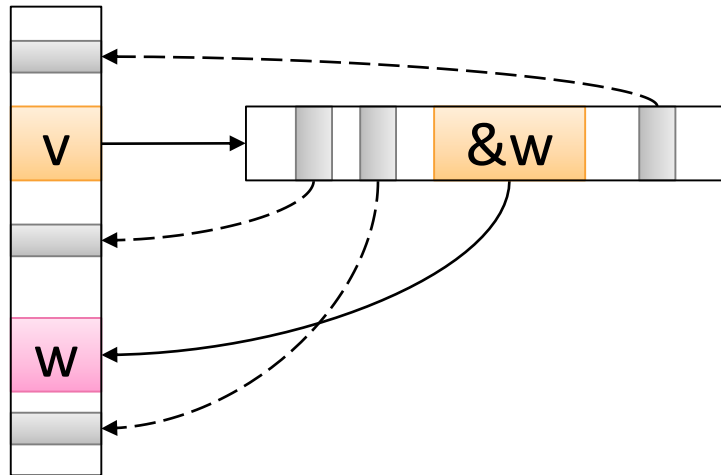
DRAM

Crossbar Network

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

Communications via
Remote Function Calls

Message Queue
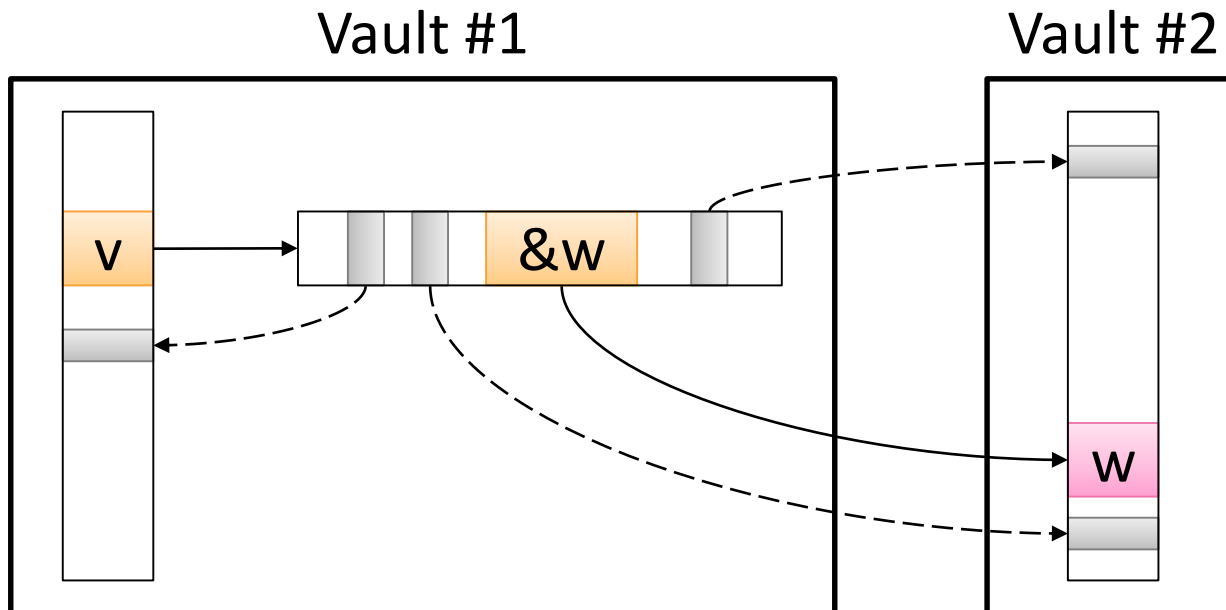
NI

# Communications in Tesseract

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

# Communications in Tesseract

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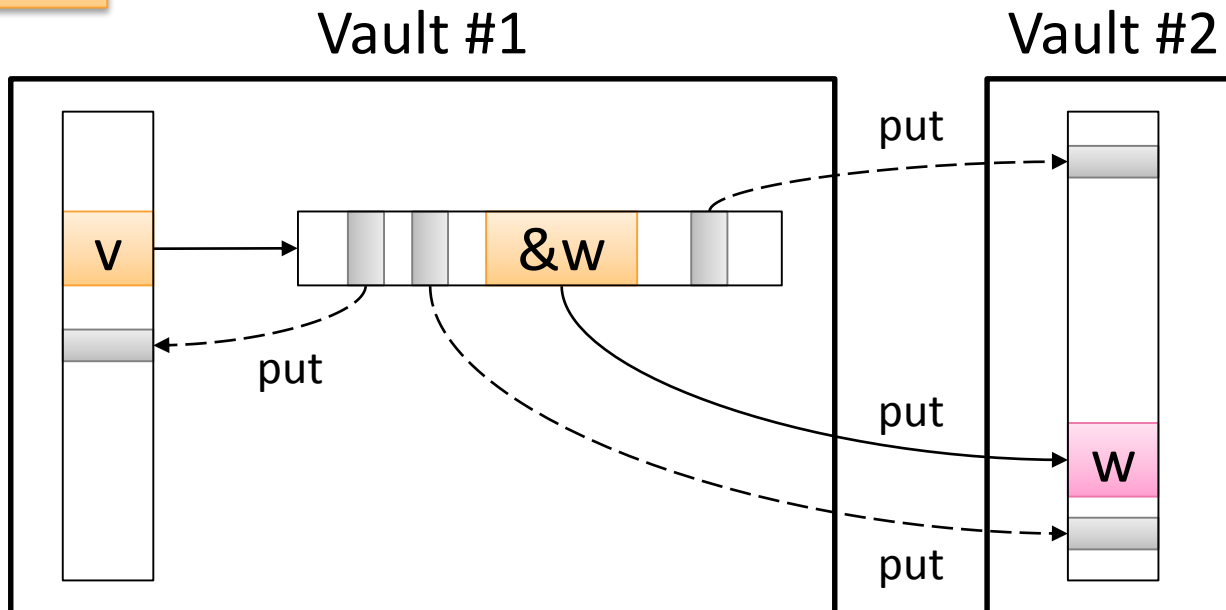


Vault #1    Vault #2

# Communications in Tesseract

```
for (v: graph.vertices) {
    for (w: v.successors) {
        put(w.id, function() { w.next_rank += weight * v.rank; });
    }
}
barrier();
```
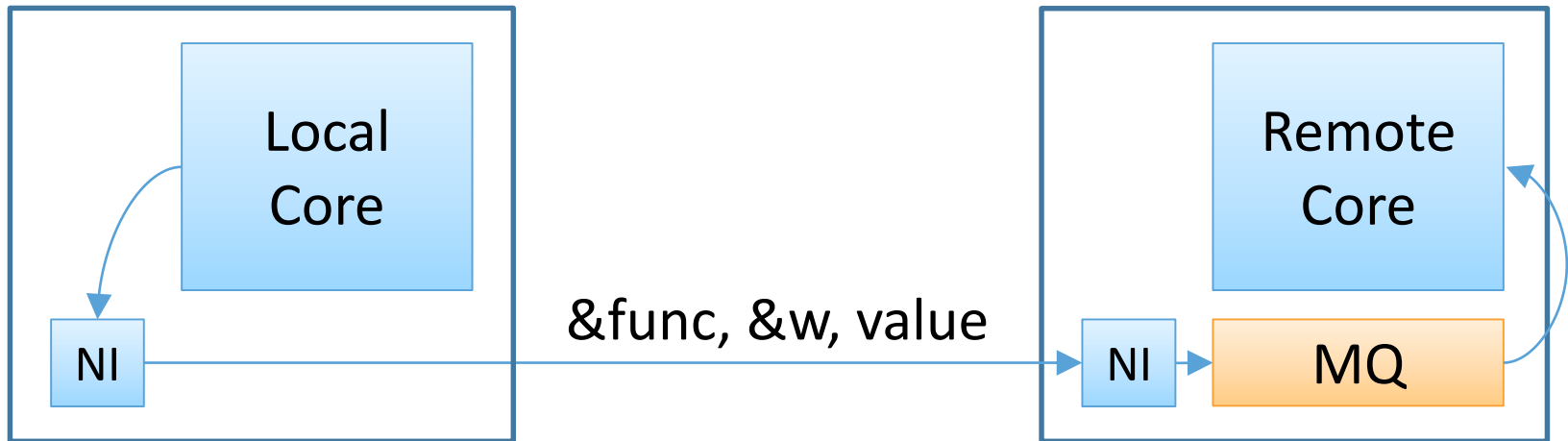
**Non-blocking Remote Function Call**

Can be **delayed**
until the nearest barrier

# Non-blocking Remote Function Call

1. Send function address & args to the remote core
2. Store the incoming message to the message queue
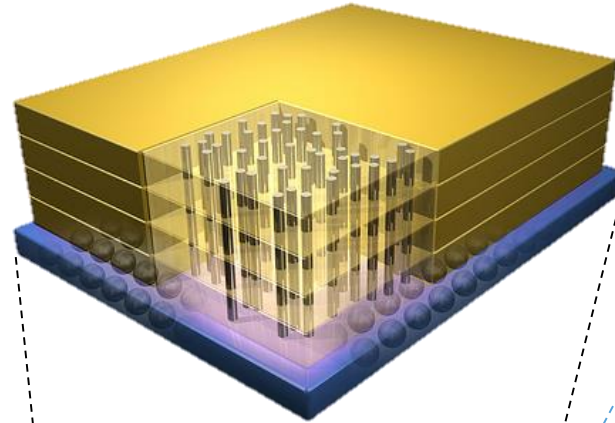3. Flush the message queue when it is full or a synchronization barrier is reached



put(w.id, function() { w.next_rank += value; })

# Benefits of Non-blocking Remote Function Call

- Latency hiding through fire-and-forget
  - Local cores are not blocked by remote function calls

- Localized memory traffic
  - No off-chip traffic during remote function call execution

- No need for mutexes
  - Non-blocking remote function calls are atomic

- Prefetching
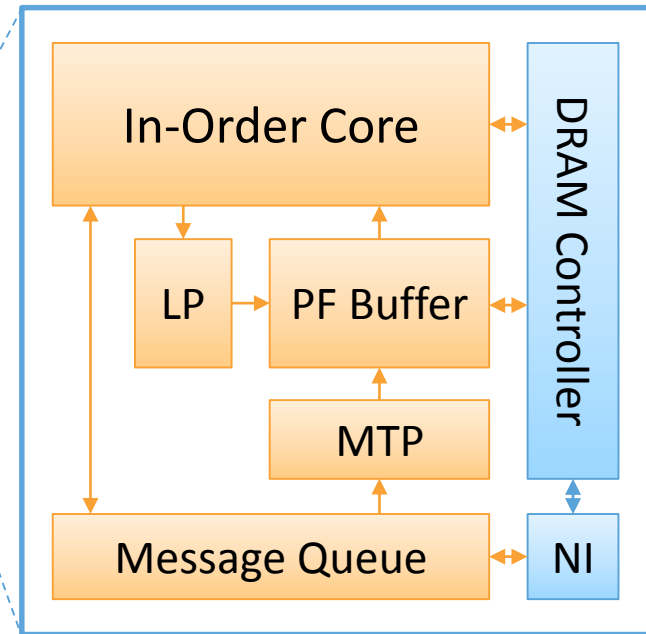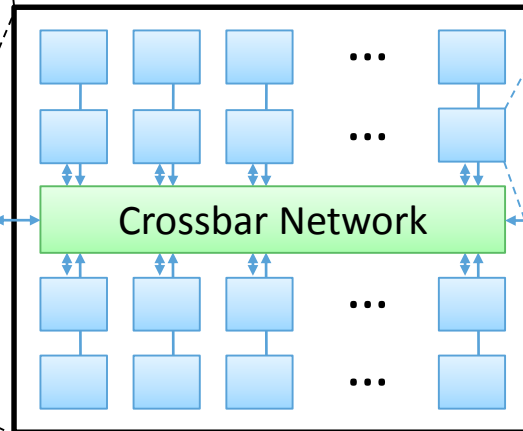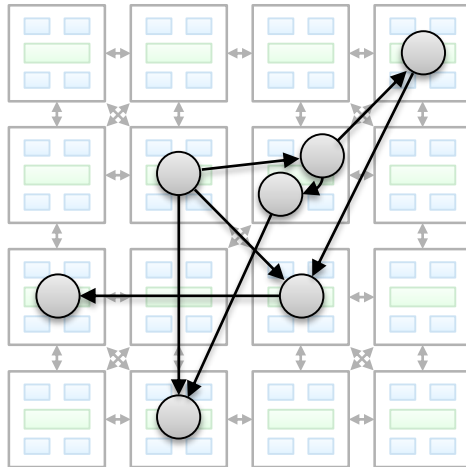  - Will be covered shortly

# Tesseract System
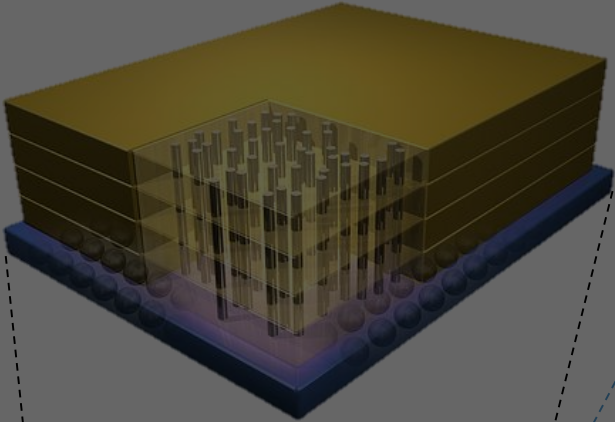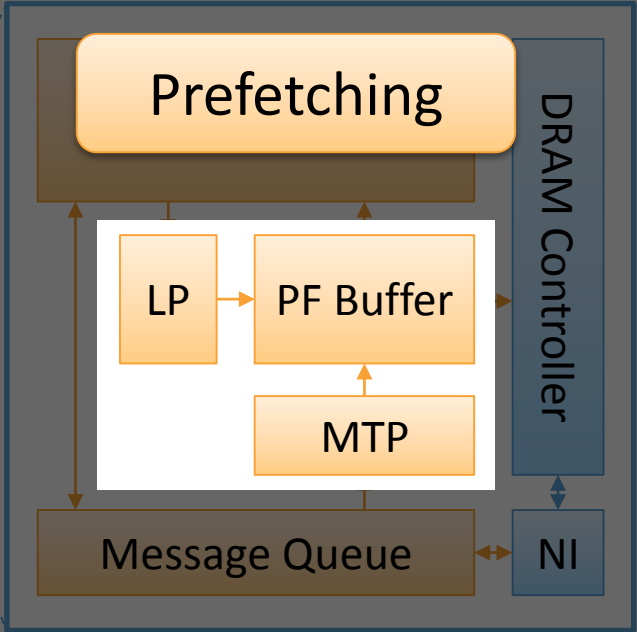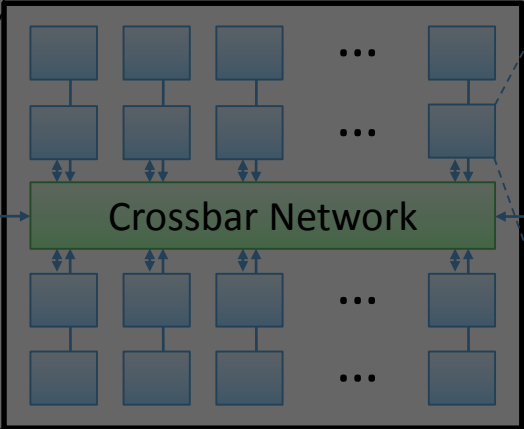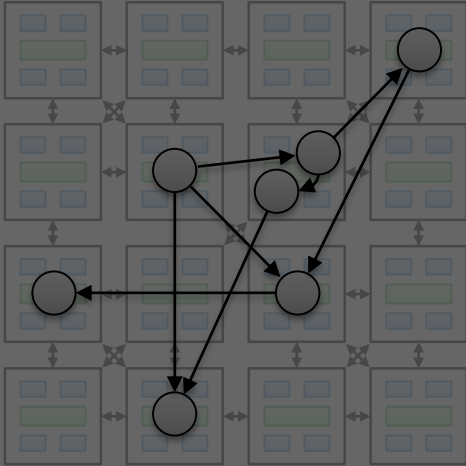
# Tesseract System



Host Processor

Memory-Mapped
Accelerator Interface
(Noncacheable, Physically Addressed)

Crossbar Network

Prefetching

DRAM Controller

LP → PF Buffer

MTP

Message Queue

NI

# Memory Access Patterns in Graph Processing

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

# Message-Triggered Prefetching

- Prefetching random memory accesses is difficult

- Opportunities in Tesseract
  - Domain-specific knowledge of target data address
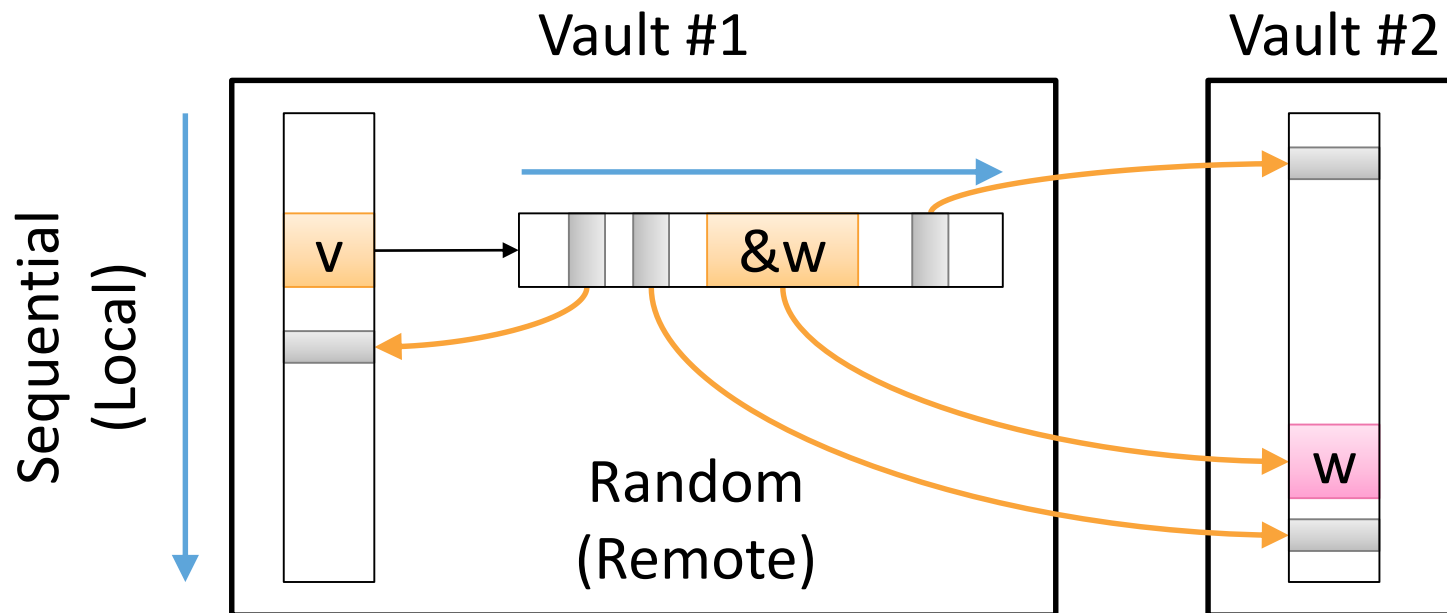  - Time slack of non-blocking remote function calls

```
for (v: graph.vertices) {
    for (w: v.successors) {
        put(w.id, function() { w.next_rank += weight * v.rank; });
    }
}
barrier();
```
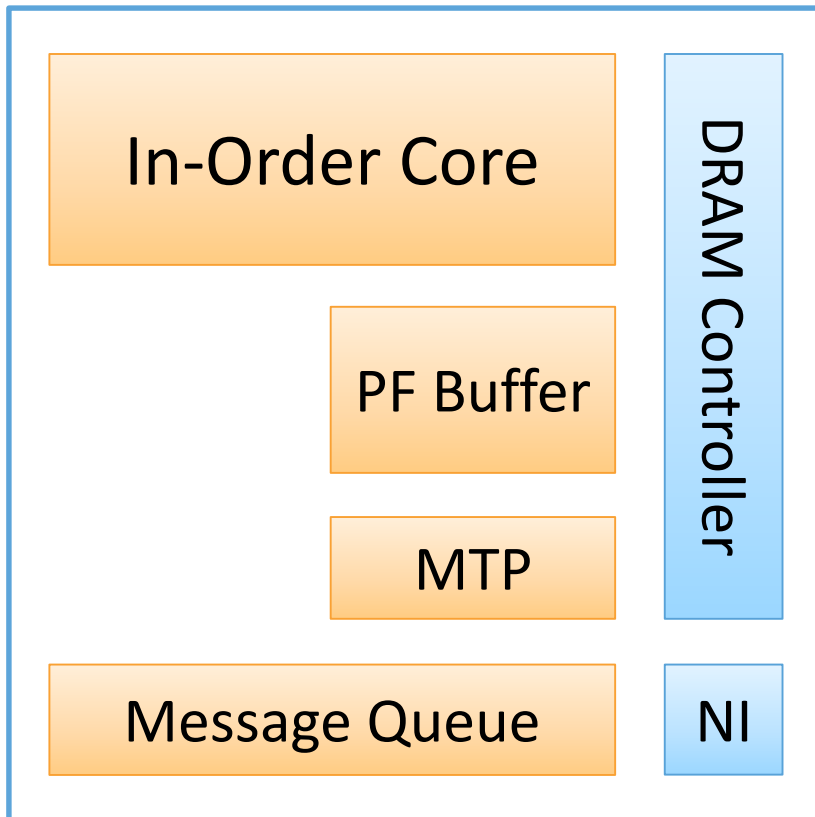
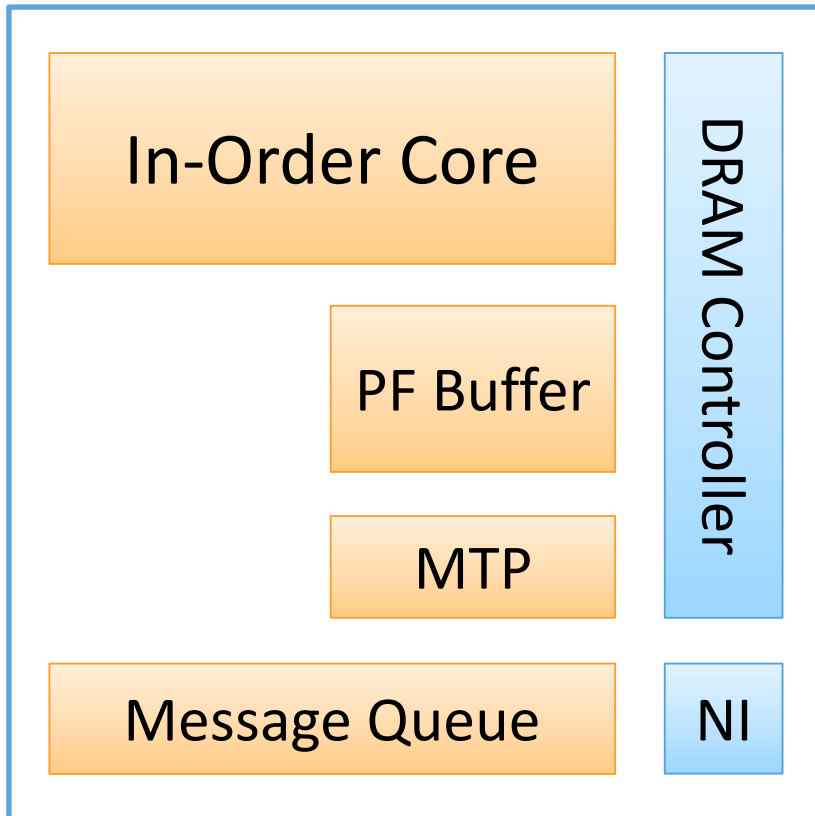Can be **delayed**
until the nearest barrier

# Message-Triggered Prefetching
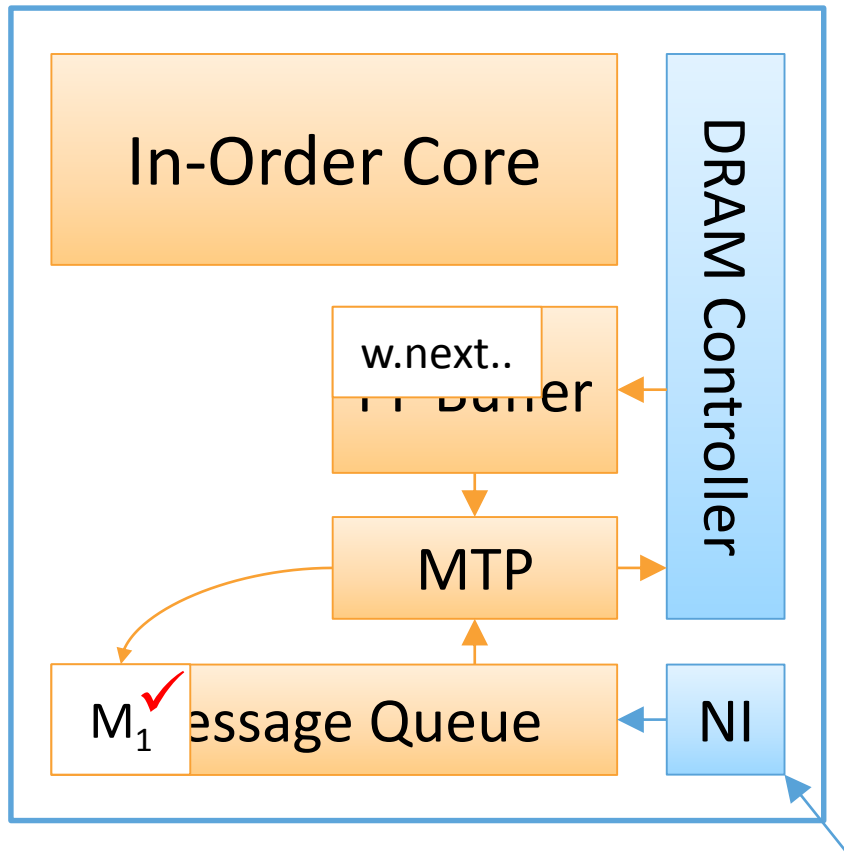


put(w.id, function() { w.next_rank += value; })

# Message-Triggered Prefetching



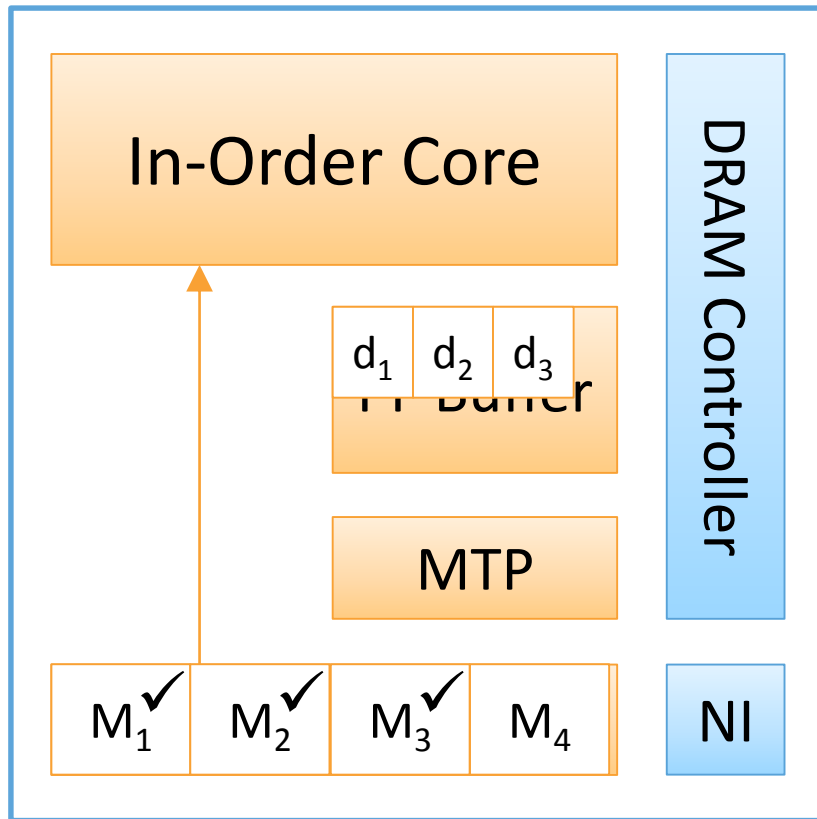put(w.id, function() { w.next_rank += value; }, &w.next_rank)

# Message-Triggered Prefetching



1. Message $M_1$ received
2. Request prefetch
3. Mark $M_1$ as ready when the prefetch is serviced

put(w.id, function() { w.next_rank += value; }, &w.next_rank)

# Message-Triggered Prefetching



1. Message $M_1$ received
2. Request prefetch
3. Mark $M_1$ as ready when the prefetch is serviced
4. Process multiple **ready** messages at once

put(w.id, function() { w.next_rank += value; }, &w.next_rank)
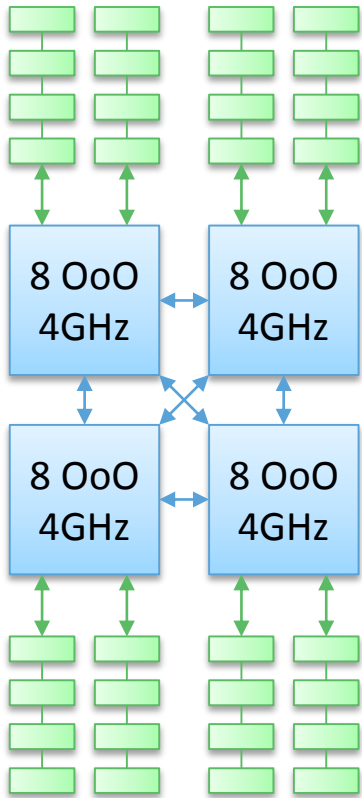
# Other Features of Tesseract

- Blocking remote function calls
- List prefetching
- Prefetch buffer
- Programming APIs
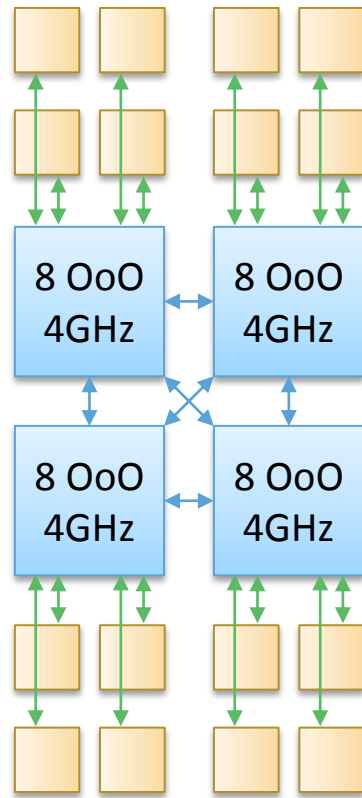- Application mapping

Please see the paper for details

# Evaluated Systems



**DDR3-OoO**
(with FDP)

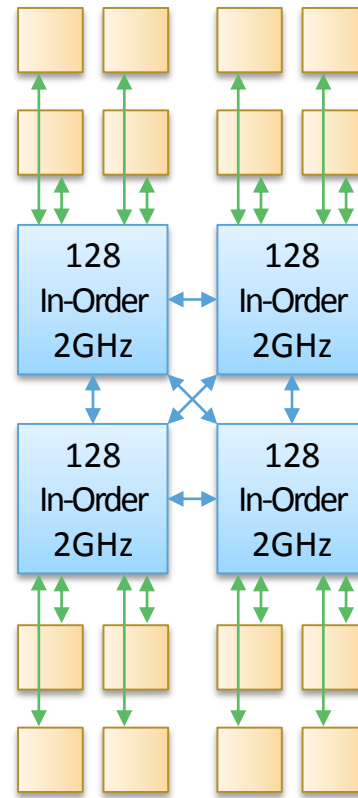**HMC-OoO**
(with FDP)
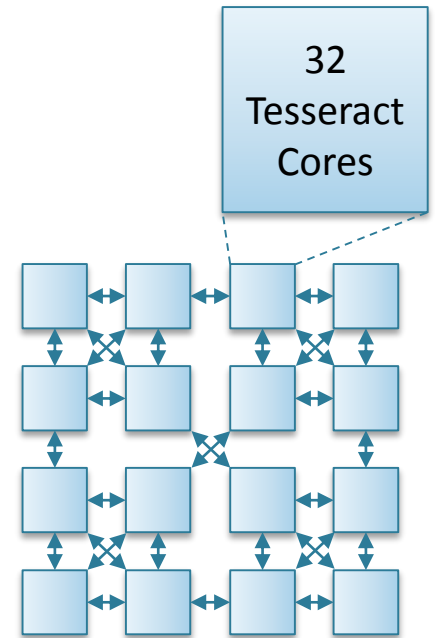
**HMC-MC**

**Tesseract**
(32-entry MQ, 4KB PF Buffer)

| DDR3-OoO | HMC-OoO | HMC-MC | Tesseract |
|---|---|---|---|
| 8 OoO 4GHz | 8 OoO 4GHz | 128 In-Order 2GHz | 32 Tesseract Cores |

102.4GB/s | 640GB/s | 640GB/s | 8TB/s
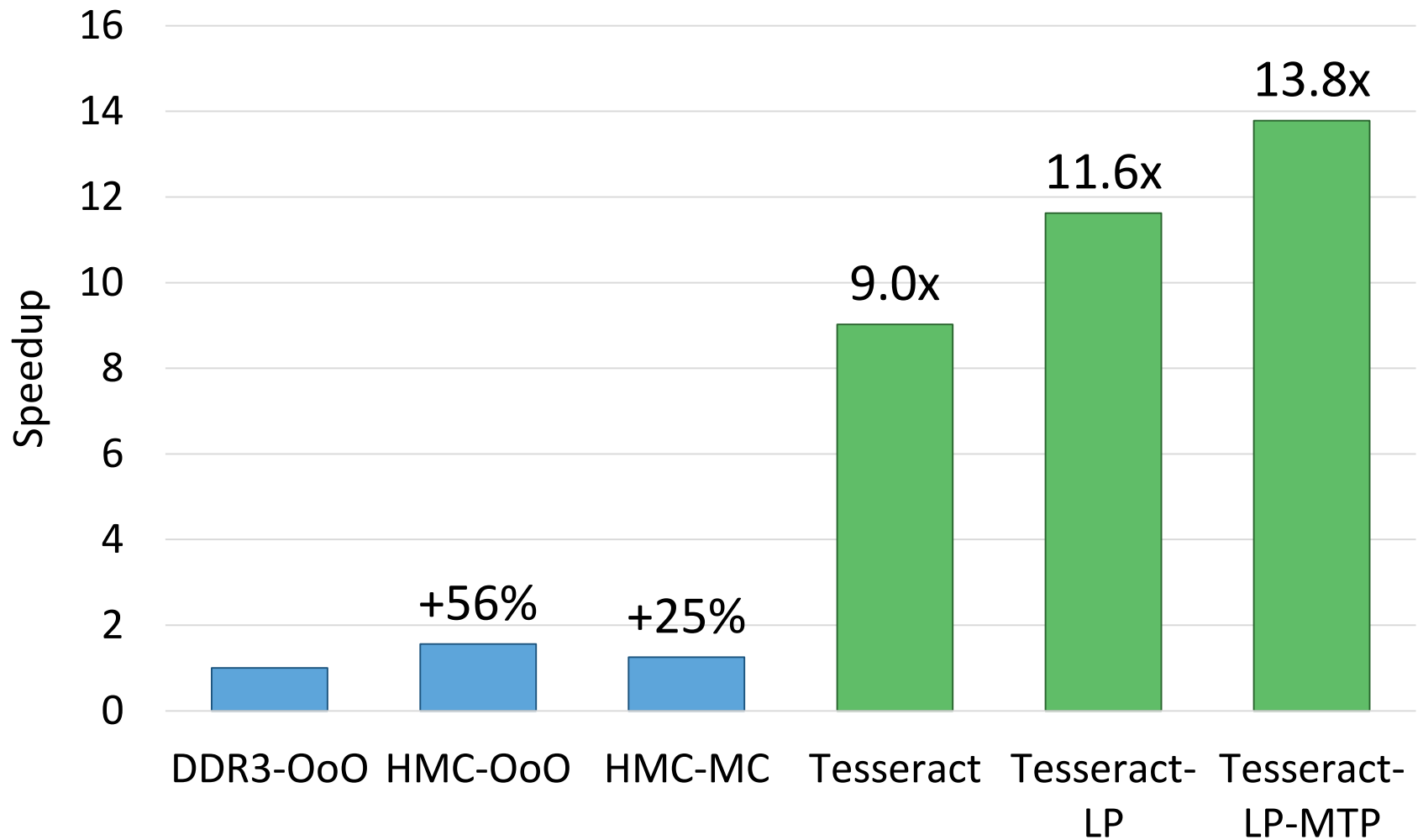
# Workloads

- Five graph processing algorithms
  - Average teenage follower
  - Conductance
  - PageRank
  - Single-source shortest path
  - Vertex cover

- Three real-world large graphs
  - ljournal-2008 (social network)
  - enwiki-2003 (Wikipedia)
  - indochina-0024 (web graph)
  - 4~7M vertices, 79~194M edges

# Performance

# Performance
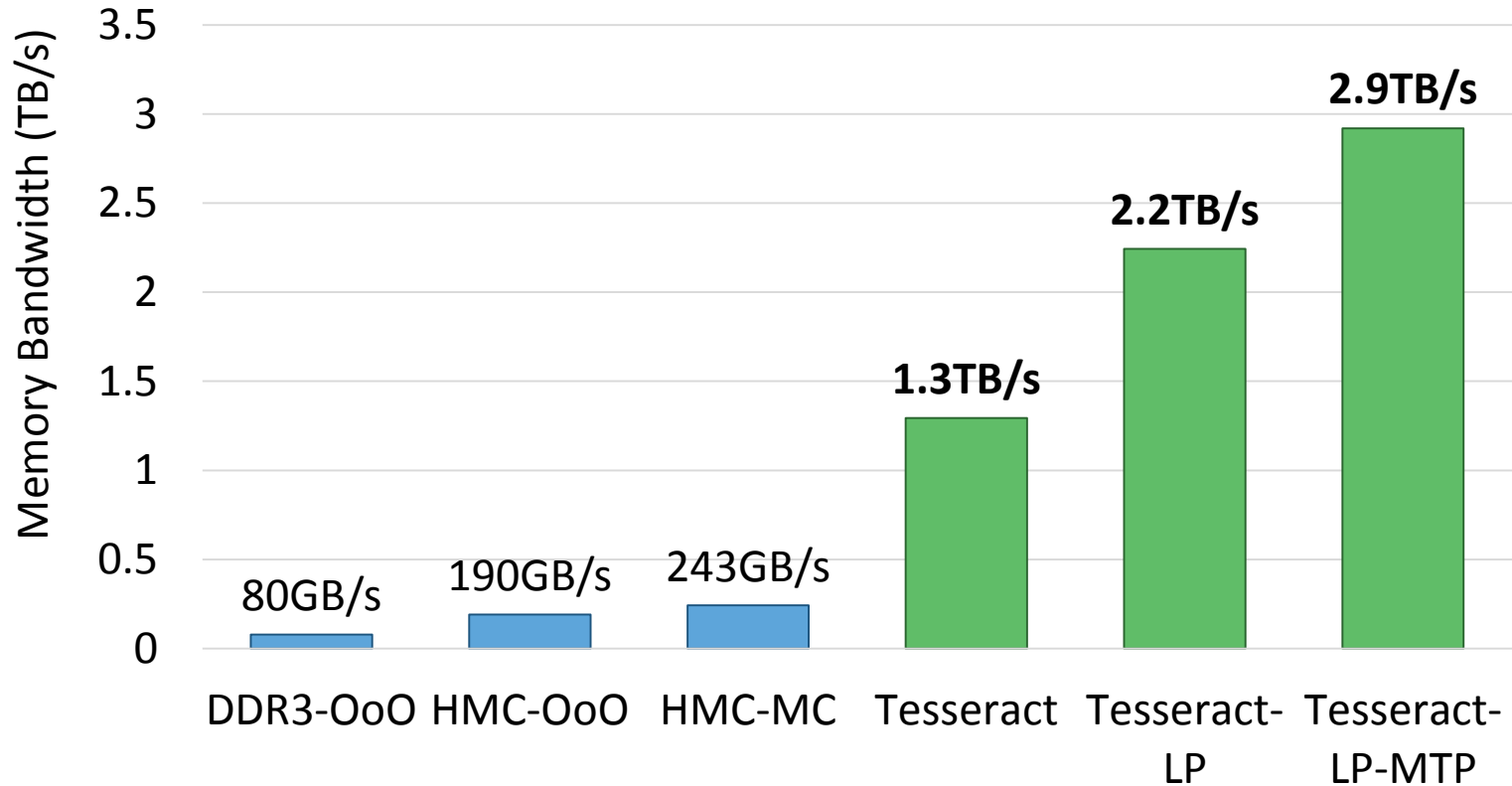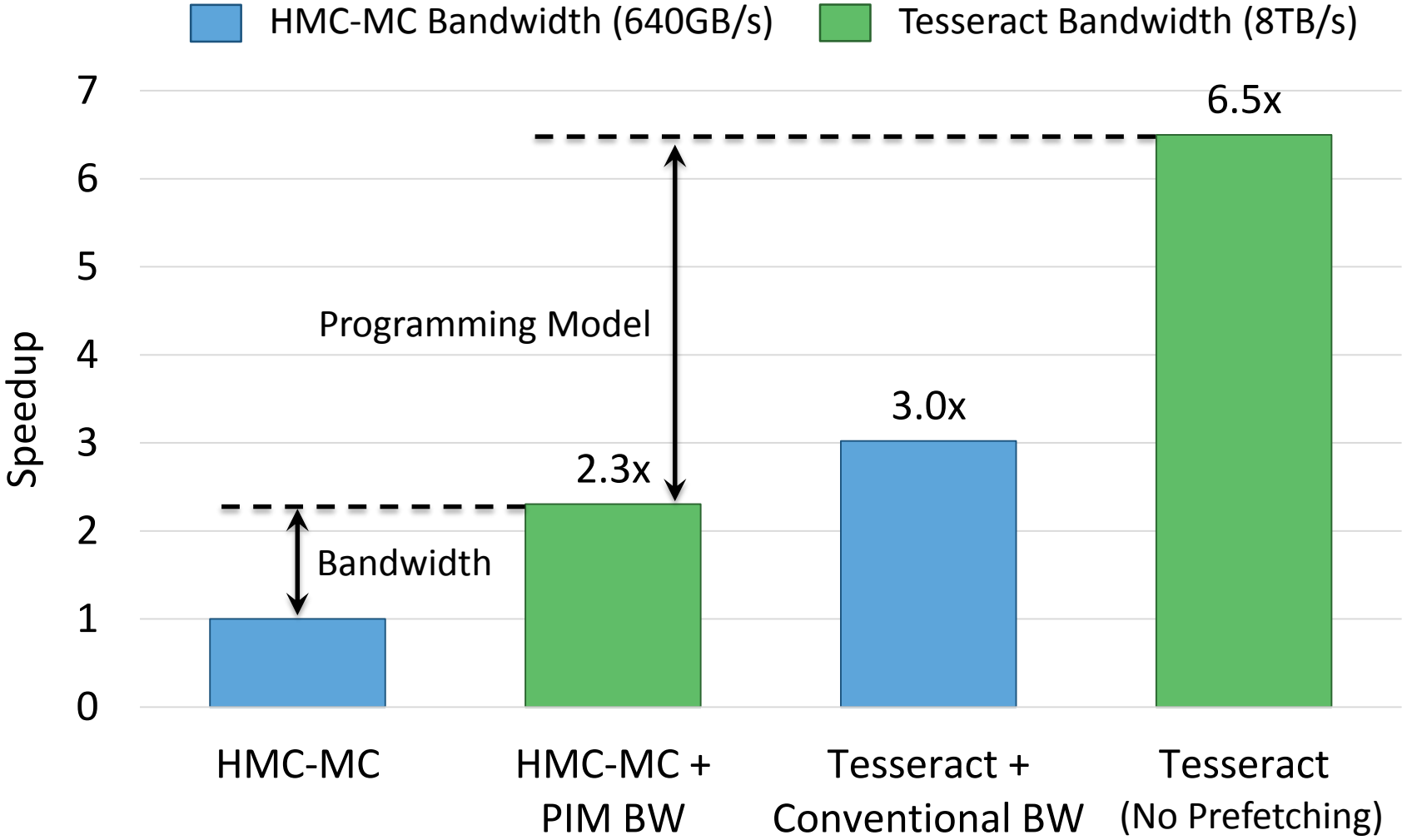


**Memory Bandwidth Consumption**

Memory Bandwidth (TB/s)

| | | |
|---|---|---|
| 80GB/s | 190GB/s | 243GB/s |
| DDR3-OoO | HMC-OoO | HMC-MC |

1.3TB/s — Tesseract
2.2TB/s — Tesseract-LP
2.9TB/s — Tesseract-LP-MTP

# Iso-Bandwidth Comparison

# Execution Time Breakdown



Legend: ■ Normal Mode ■ Interrupt Mode ■ Interrupt Switching ■ Network ■ Barrier

Categories: AT.LJ, CT.LJ, PR.LJ, SP.LJ, VC.LJ

# Execution Time Breakdown



**Network Backpressure due to Message Passing**
(Future work: message combiner, ...)

Legend: Normal Mode ■ Interrupt Mode ■ Interrupt Switching ■ Network ■ Barrier

Categories: AT.LJ, CT.LJ, PR.LJ, SP.LJ, VC.LJ

# Execution Time Breakdown



**Workload Imbalance**
(Future work: load balancing)

Legend: Normal Mode, Interrupt Mode, Interrupt Switching, Network, Barrier

X-axis: AT.LJ, CT.LJ, PR.LJ, SP.LJ, VC.LJ

# Prefetch Efficiency

# Scalability

**DDR3-OoO**



**Tesseract with Prefetching**
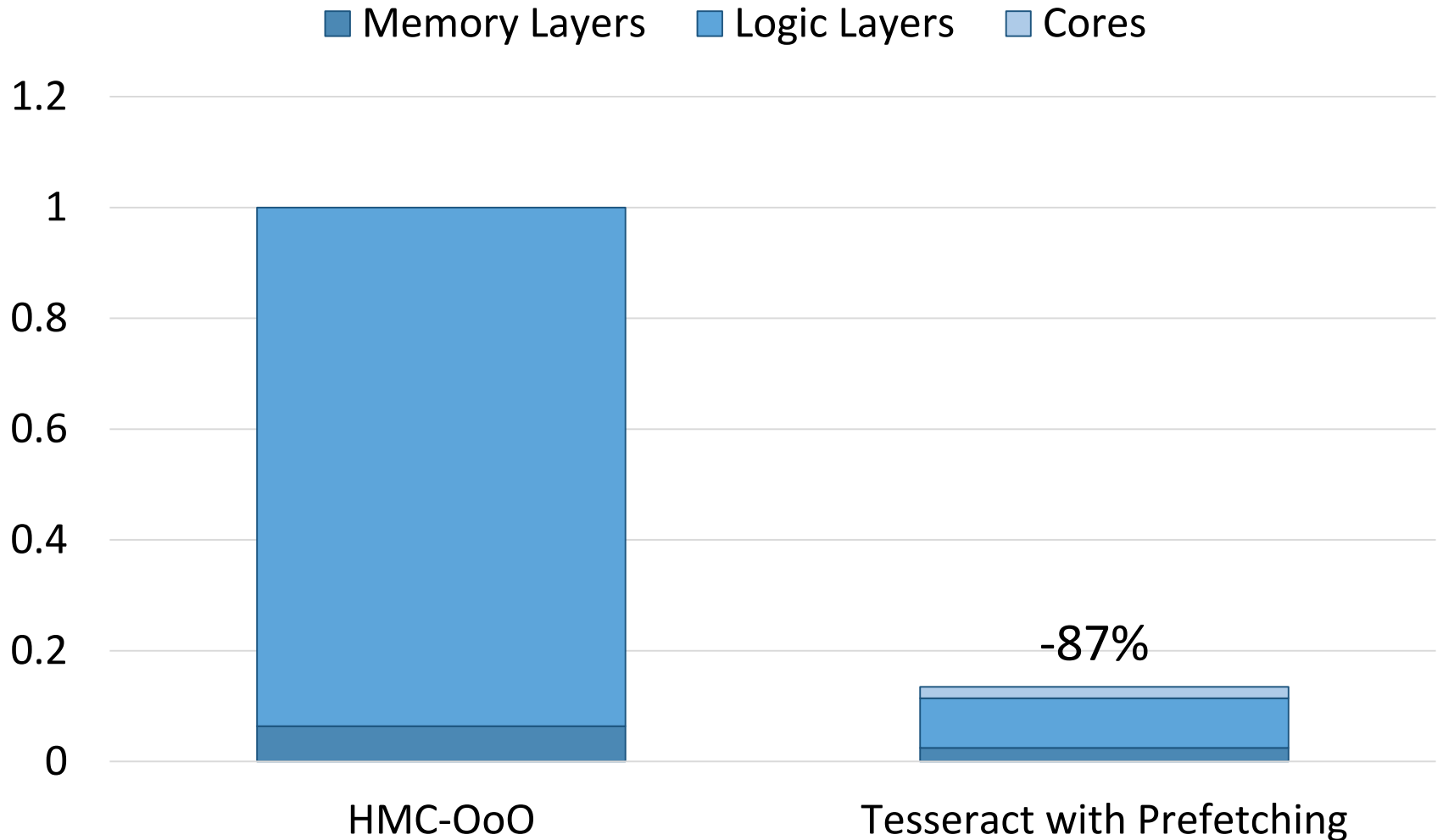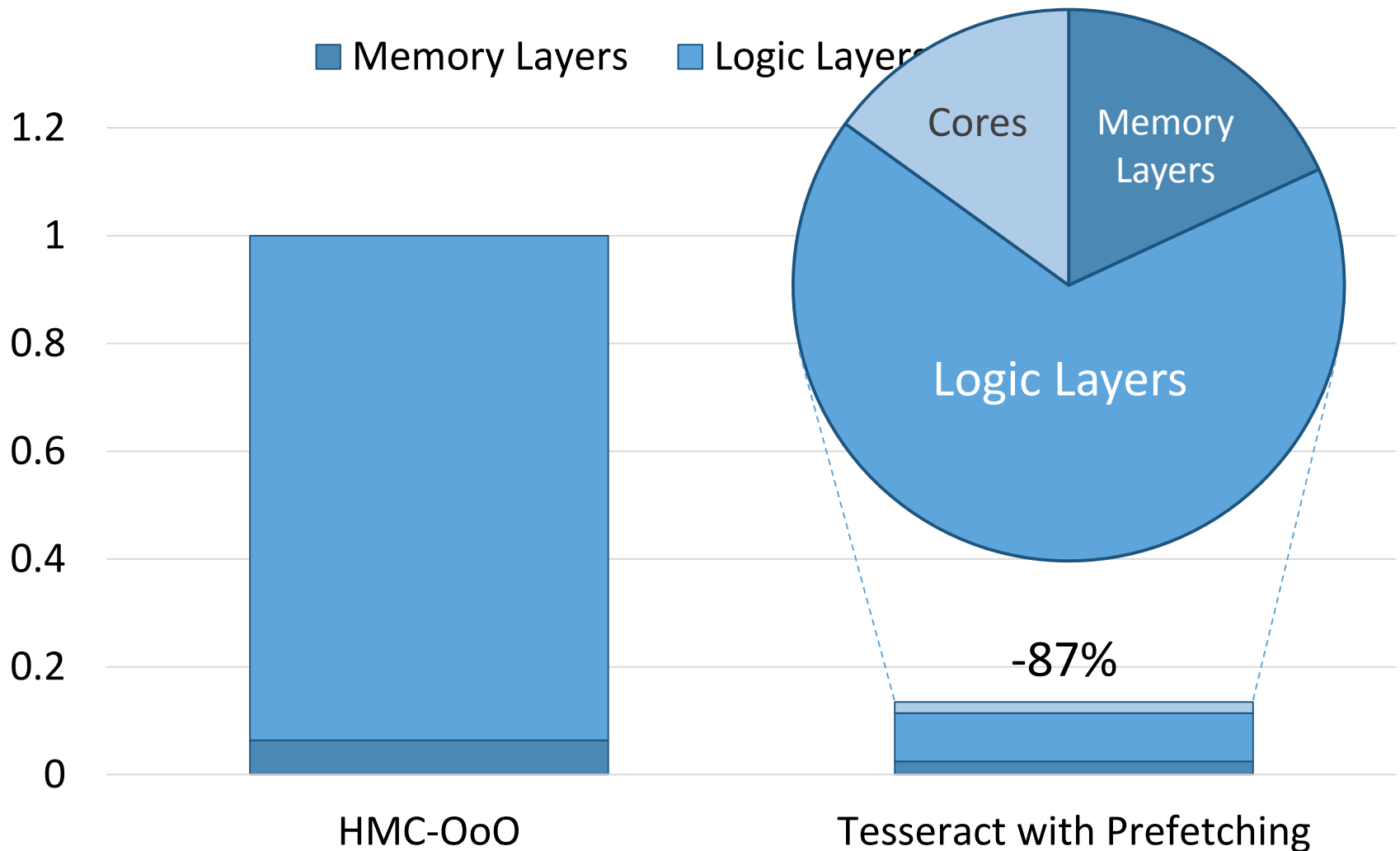
# Memory Energy Consumption

# Memory Energy Consumption

# Conclusion

- Revisiting the PIM concept in a new context
  - Cost-effective 3D integration of logic and memory
  - Graph processing workloads demanding high memory bandwidth

- Tesseract: scalable PIM for graph processing
  - Many in-order cores in a memory chip
  - New message passing mechanism for latency hiding
  - New hardware prefetchers for graph processing
  - Programming interface that exploits our hardware design

- Evaluations demonstrate the benefits of Tesseract
  - 14x performance improvement & 87% energy reduction
  - Scalable: *memory-capacity-proportional* performance

# A Scalable Processing-in-Memory Accelerator for Parallel Graph Processing

Junwhan Ahn, Sungpack Hong*, Sungjoo Yoo,
Onur Mutlu+, Kiyoung Choi

Seoul National University     *Oracle Labs     +Carnegie Mellon University