# Demystifying the Characteristics of 3D-Stacked Memories: A Case Study for the Hybrid Memory Cube (HMC)

Ramyad Hadidi, Bahar Asgari , Burhan Ahmad Mudassar, Saibal Mukhopadhyay, Sudhakar Yalamanchili, and Hyesoon Kim
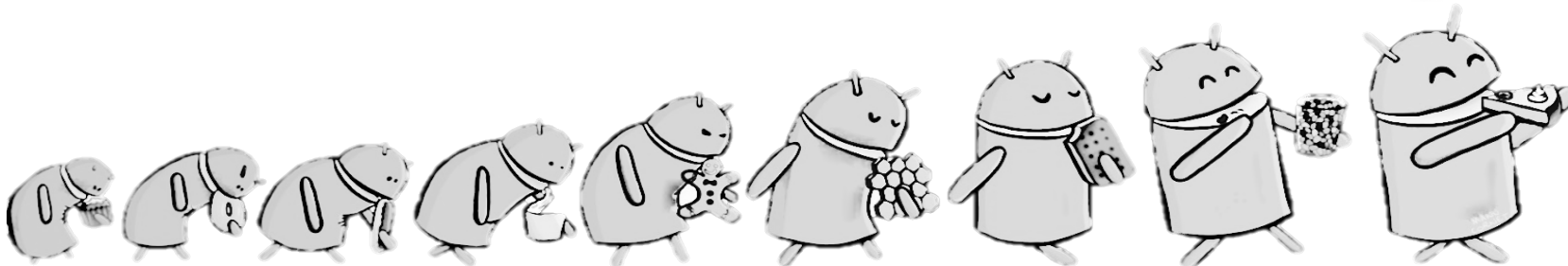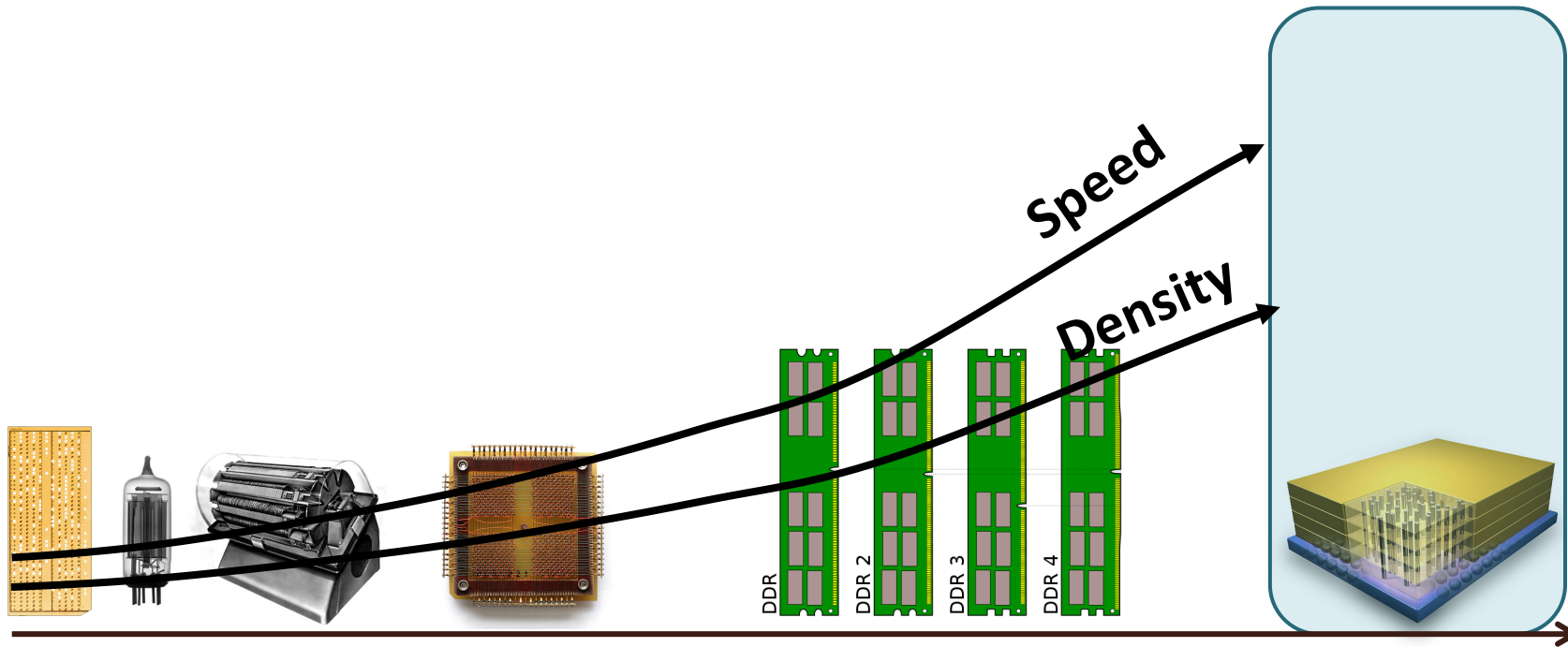
**IISWC'17 Talk**

Georgia Tech · comparch

CASL
Computer Architecture
and Systems Laboratory

# Memory Evolution

Speed

Density

DDR

DDR 2

DDR 3

DDR 4

# 3D-Stacking Technology

Provides opportunities & novel features

3D-DRAMs:

▸ Provide higher bandwidth and density

▸ Enable lower power consumption

▸ Motivate processing-in-memory

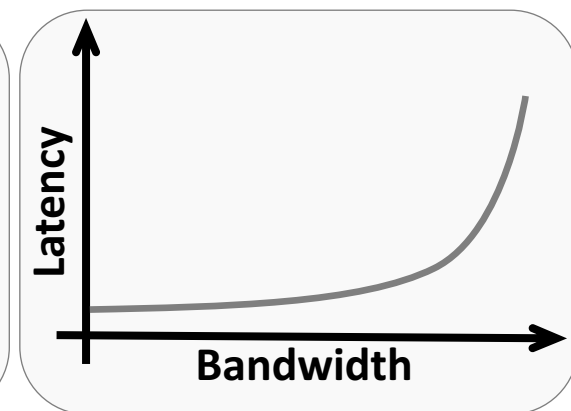HMC is an example of such memories.
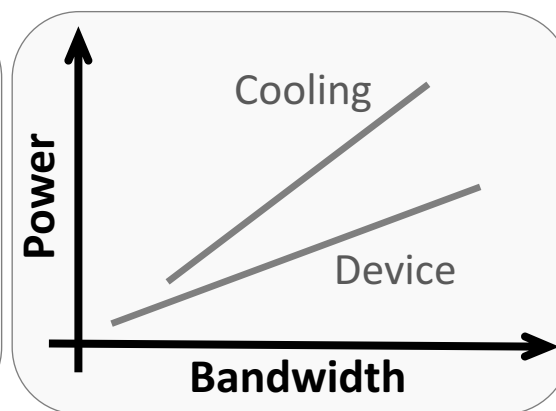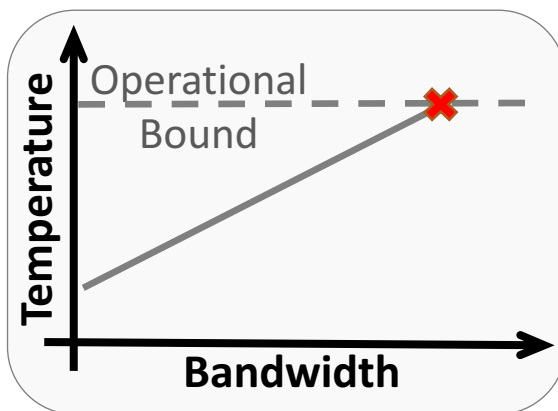
# New Considerations

New **internal organization**

New **thermal** behavior

New **latency** and **bandwidth** hierarchy

New packet-switched **interface**

# Contributions

We evaluate a real system with HMC 1.1 to:

- Study new memory organization
- Present bandwidth, power, and temperature relationships
- Investigate required cooling power
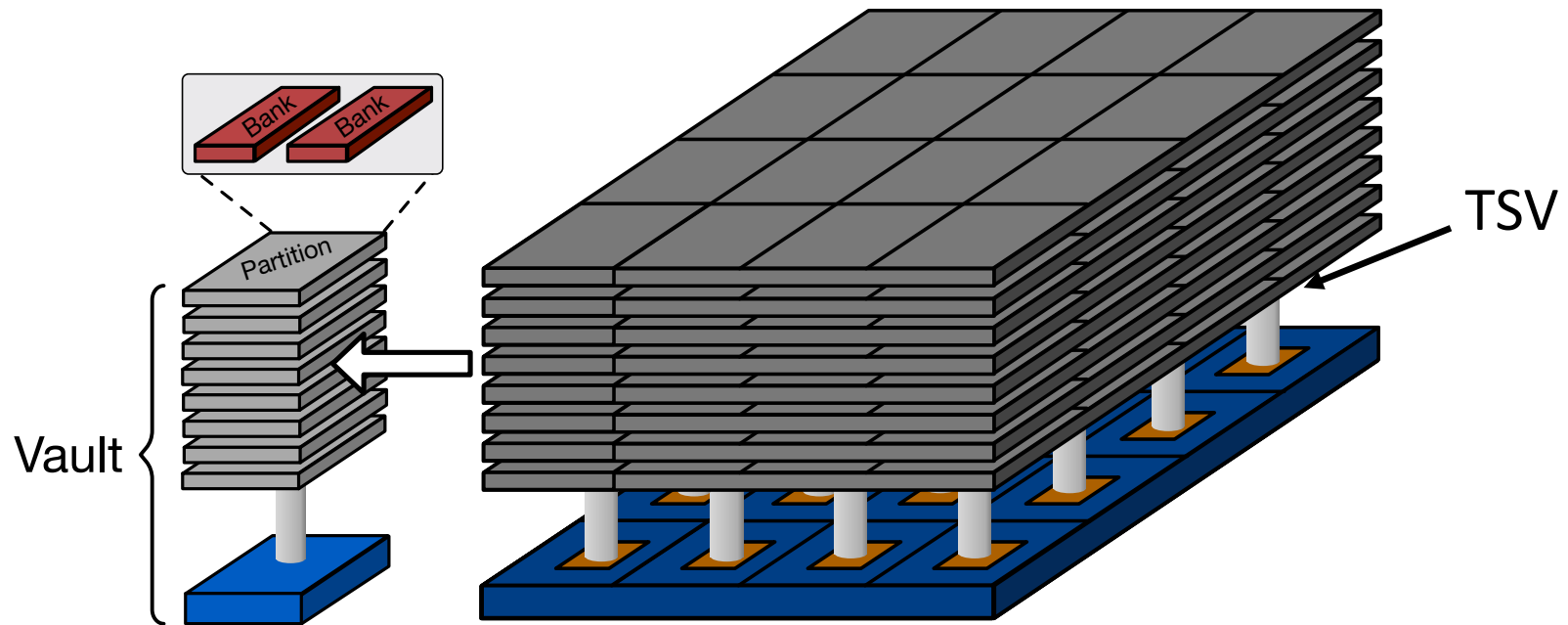- Explore contributing factors to latency

**AC510**

To realize the full-system impact of
3D-stacked memories and HMC in particular.

# Hybrid Memory Cube (HMC)
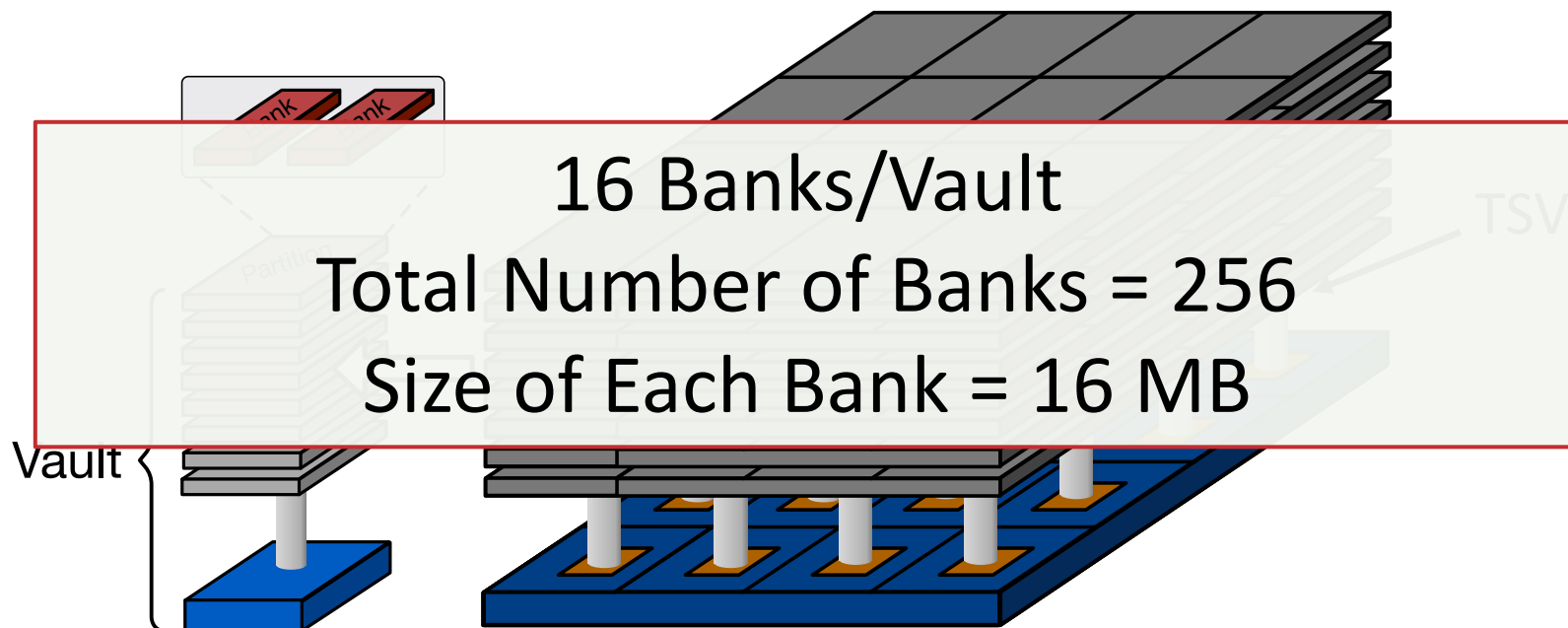
HMC 1.1 (Gen2): 4GB size



TSV

Vault

Partition

Bank  Bank

**Logic Layer**  **Vault Controller**  **DRAM Layer**

# Hybrid Memory Cube (HMC)

HMC 1.1 (Gen2): 4GB size

16 Banks/Vault
Total Number of Banks = 256
Size of Each Bank = 16 MB

Vault

TSV

■ Logic Layer ■ Vault Controller ■ DRAM Layer

CASL
Computer Architecture
and Systems Laboratory

Georgia
Tech

comparch

# HMC Communication I

Follows a serialized **packet-switched** protocol

Partitioned into 16-byte *flit*

Each transfer incurs 1 flit of overhead

| Type | Read | | Write | |
|---|---|---|---|---|
| | Request | Response | Request | Response |
| Data Size | Empty | 1~8 Flits | 1~8 Flits | Empty |
| Overhead | 1 Flit | 1 Flit | 1 Flit | 1 Flit |
| Total Size | 1 Flit | 2~9 Flits | 2~9 Flits | 1 Flit |

# HMC Communication I

Follows a serialized **packet-switched** protocol

Partitioned into 16-byte *flit*

Each transfer incurs 1 flit of overhead

| Type | Read | | Write | |
|------|---------|----------|---------|----------|
| | Request | Response | Request | Response |
| Data Size Overhead | Empty 1 Flit | 1~8 Flits 1 Flit | 1~8 Flits 1 Flit | Empty 1 Flit |
| Total Size | 1 Flit | 2~9 Flits | 2~9 Flits | 1 Flit |

# HMC Communication II

Two/Four full duplex external links:

- Width of 8 or 16 lanes
- Configurable speeds of 10, 12.5, and 15 Gbps
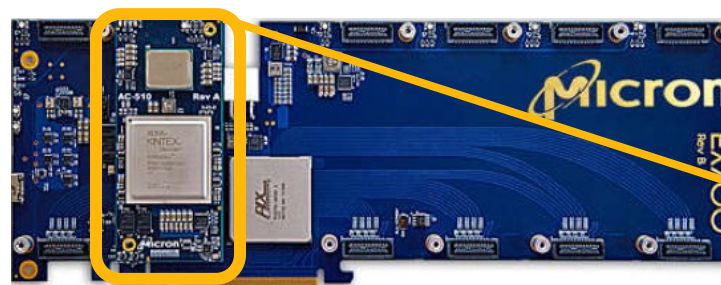


*8 lanes*
*15 Gbps*

Our evaluated system

2 external links – 8 lanes each

# Experimental Setup I
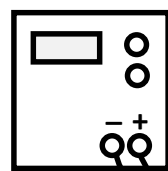
Pico SC6 Mini
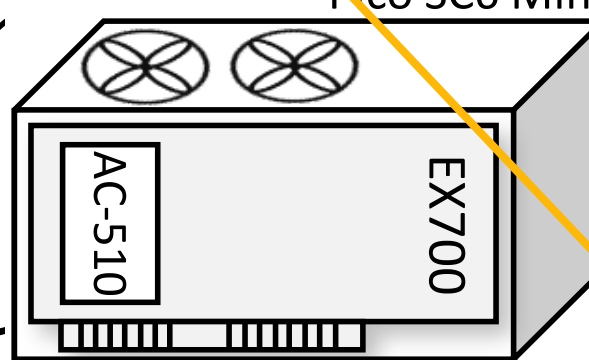
EX700 Backplane

AC510 Module

4GB HMC 1.1



EX700

DC Power Supply:
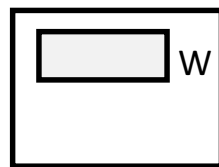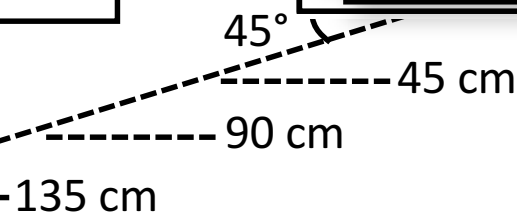Fan Speed Control

Power Measurement

15W Fan

Pico SC6 Mini

45°

45 cm

90 cm

135 cm
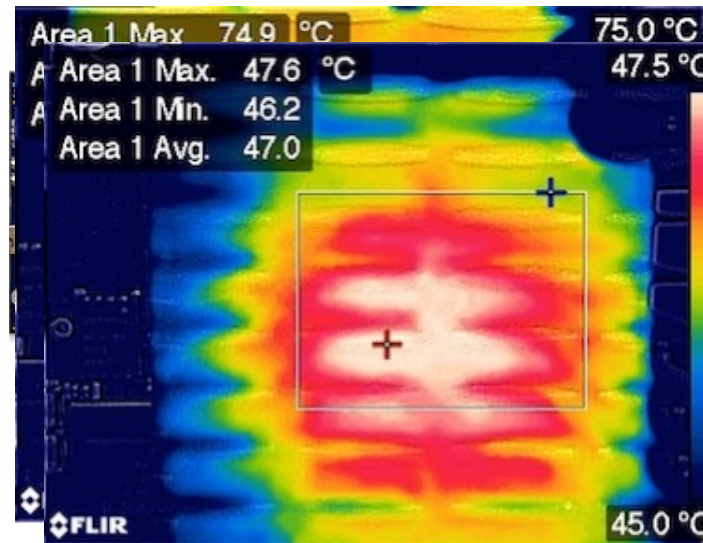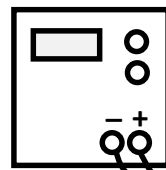
AC510

comparch

# Experimental Setup I

Pico SC6 Mini

EX700 Backplane

AC510 Module

4GB HMC 1.1



DC Power Supply:
Fan Speed Control

Power
Measurement

Pico SC6 Mini

45°

45 cm

90 cm

15W
Fan

135 cm

AC510

comparch

# Experimental Setup II



FPGA frequency: 187.5 MHz

Modified GUPS (giga updates per second) benchmark

Apply different masks to addresses

# Access Patterns

**Accessing Less Banks**

16 vaults | 8 vaults | 4 vaults | 2 vaults | 1 vault (16 banks) | 8 banks | 4 banks | 2 banks | 1 bank

**Access Patterns**

# Access Patterns

**Accessing Less Banks** →

16 vaults | 8 vaults | 4 vaults | 2 vaults | 1 vault (16 banks) | 8 banks | 4 banks | 2 banks | 1 bank

**Access Patterns**

IISWC'17

# Access Patterns

**Accessing Less Banks**

16 vaults  8 vaults  4 vaults  2 vaults  1 vault (16 banks)  8 banks  4 banks  2 banks  1 bank

**Access Patterns**

IISWC'17

# Access Patterns
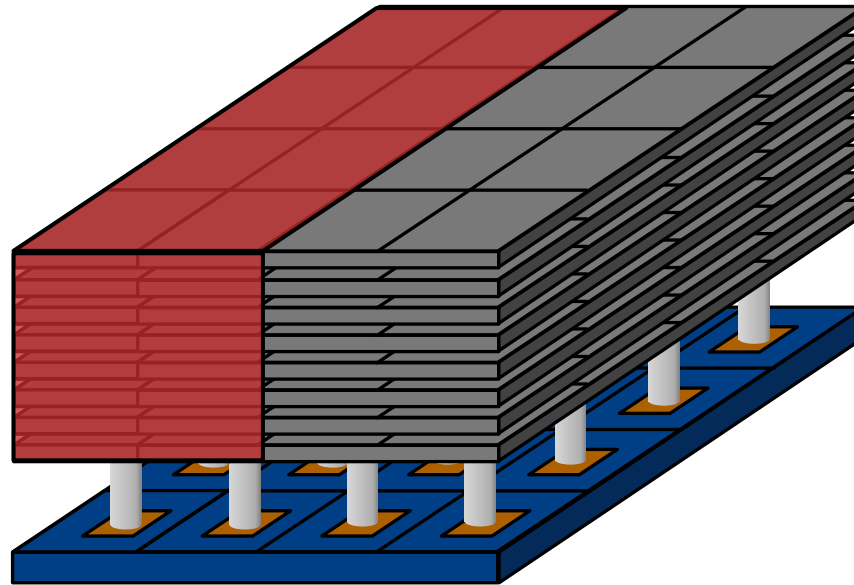
**Accessing Less Banks**

16 vaults     8 vaults     4 vaults     2 vaults     **1 vault (16 banks)**     8 banks     4 banks     2 banks     1 bank
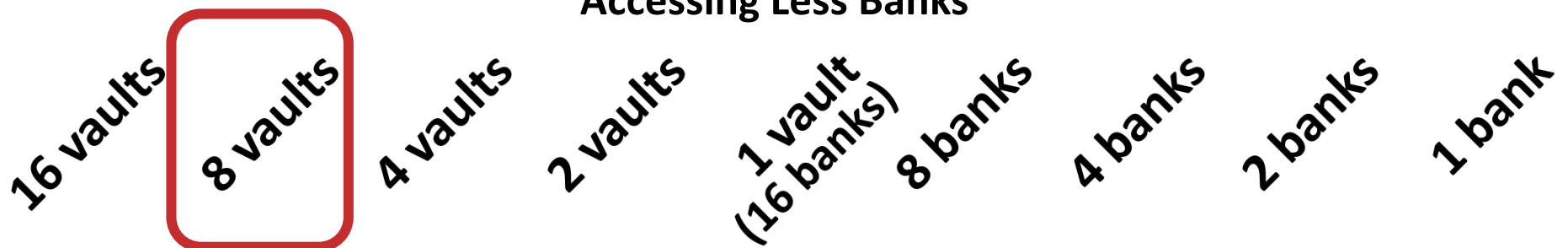
**Access Patterns**

IISWC'17     CASL Computer Architecture and Systems Laboratory     Georgia Tech     comparch

# Bandwidth

# Bandwidth

Accessing 4 banks saturates 1 vault bandwidth. External bandwidth is saturated at 4 vaults.

Type of Accesses: ☐ ro ☐ wo ☐ rw (dependent)



**Bandwidth (GB/s)** vs **Access Pattern** for access patterns: 16 vaults, 8 vaults, 4 vaults, 2 vaults, 1 vault (16 banks), 8 banks, 4 banks, 2 banks, 1 bank.

IISWC'17    CASL Computer Architecture and Systems Laboratory    Georgia Tech    comparch

# Thermal/Power Experiments

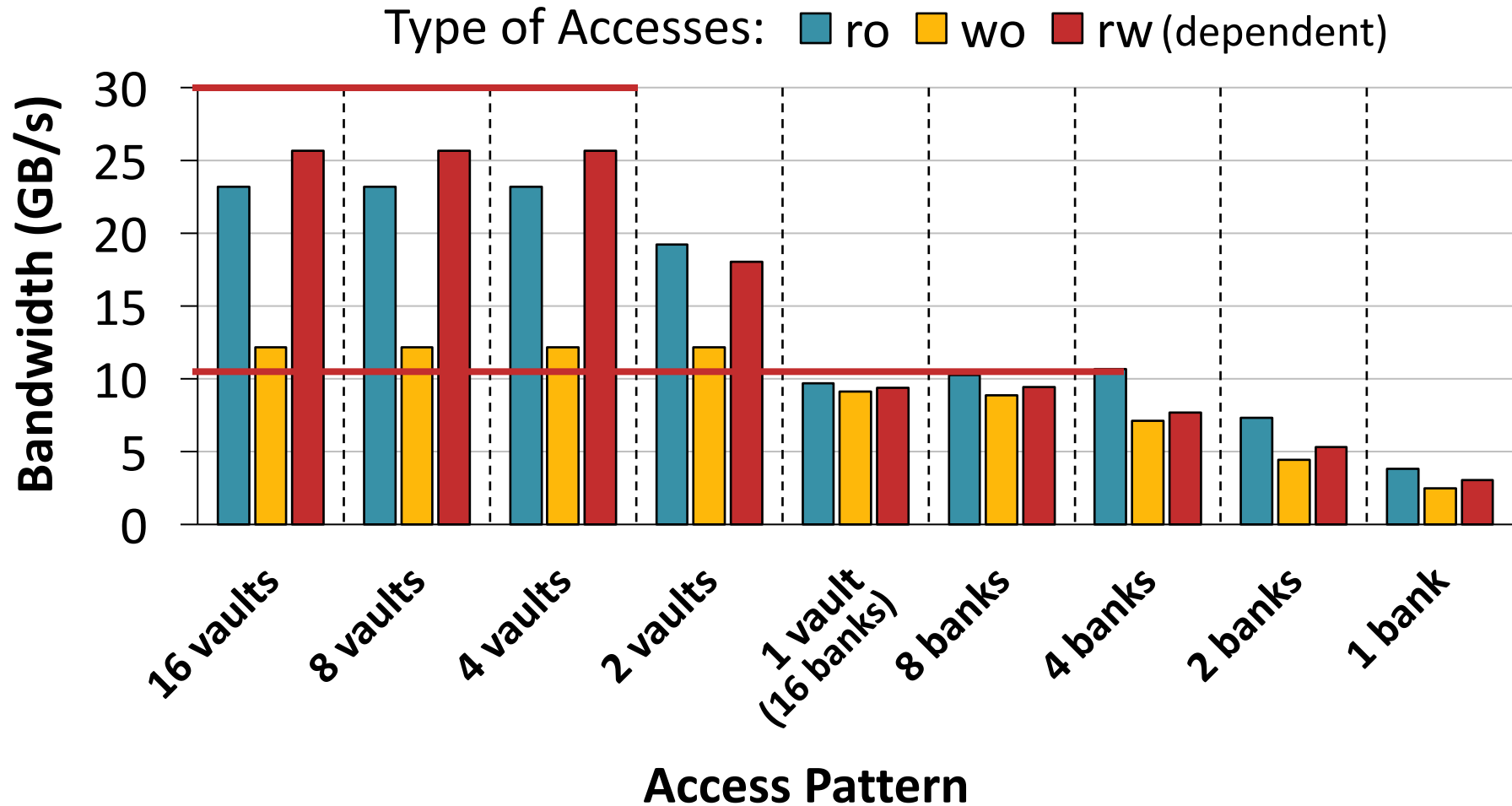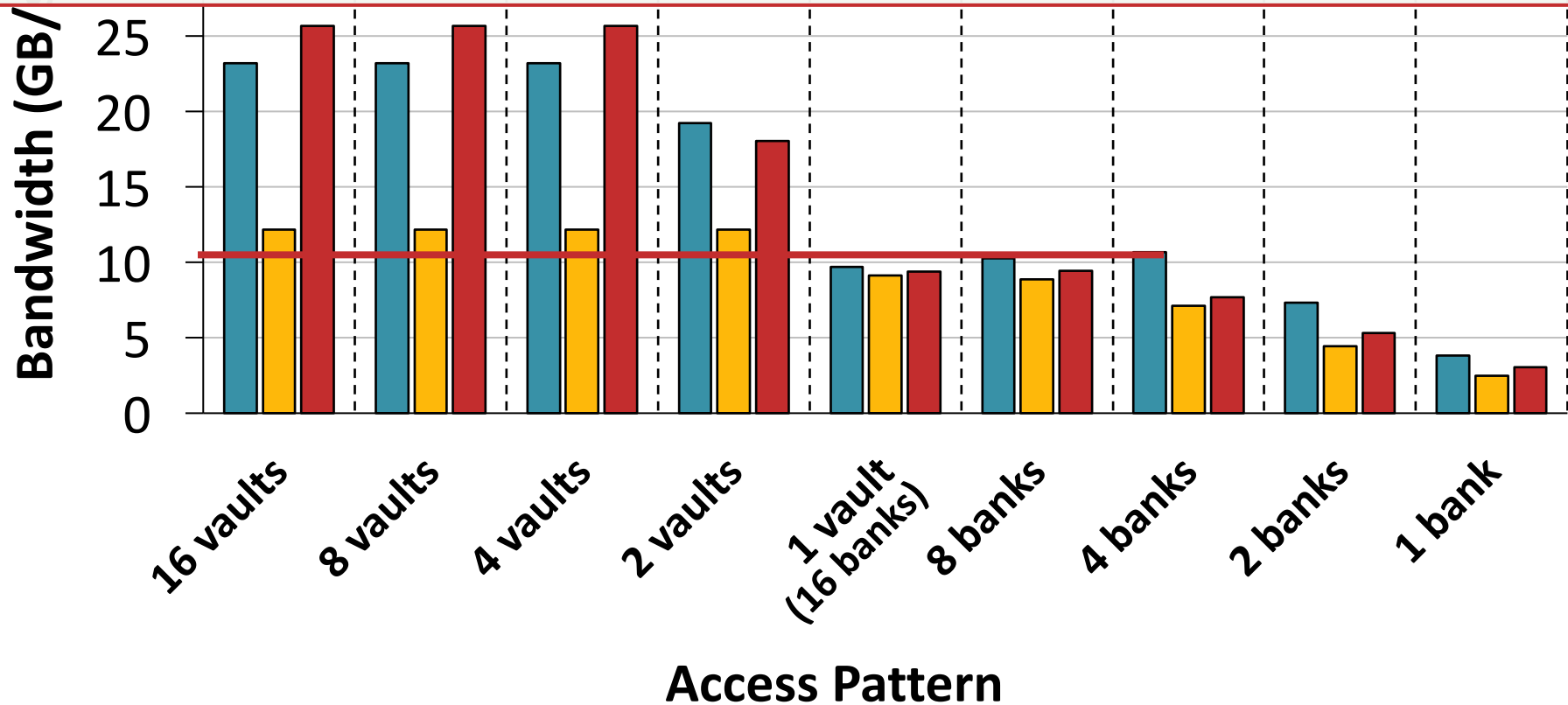| Configuration Name | DC Power Supply: Voltage | Current | 15 W Fan Distance | Average HMC Idle Temperature |
|---|---|---|---|---|
| Cfg1 | 12 V | 0.36 A | 45 cm | 43.1° C |
| Cfg2 | 10 V | 0.29 A | 90 cm | 51.7° C |
| Cfg3 | 6.5 V | 0.14 A | 90 cm | 62.3° C |
| Cfg4 | 6.0 V | 0.13 A | 135 cm | 71.6° C |

IISWC'17   CASL Computer Architecture and Systems Laboratory   Georgia Tech   comparch

# Temperature (read only)

# Temperature (read only)



Thermal Configurations: ●—Cfg4  ○··Cfg3  ▲-·Cfg2  △—Cfg1

**BW**

Access patterns affect temperature.

# Temperature & Bandwidth

Type of Accesses: ● ro ▲ wo ◆ rw



A Bandwidth increment of 15 GB/s

⇩

About 4°C increment in temperature

**Temperature (°C)** (y-axis: 48, 50, 52, 54, 56, 58, 60)

**Bandwidth (GB/s)** (x-axis: 0, 5, 10, 15, 20, 25, 30)

# Temperature & Bandwidth

Type of Accesses: ● ro ▲ wo ◆ rw



Greater slope for writes

⬇

Writes are more sensitivity to temperature

**Temperature (°C)** (y-axis: 48, 50, 52, 54, 56, 58, 60)

**Bandwidth (GB/s)** (x-axis: 0, 5, 10, 15, 20, 25, 30)

IISWC'17

# Device Power Consumption (read only)

# Device Power & Bandwidth



Type of Accesses: ● ro ▲ wo ◆ rw

# Device Power & Bandwidth



Type of Accesses: ● ro ▲ wo ◆ rw

A Bandwidth increment of 15 GB/s

⇩

About 2 W increment in device power

IISWC'17

Cooling Power Consumption (read only)

# Cooling Power Consumption (read only)

Required Cooling Power to
Keep Temperature at (°C):     ◆ 50     ○ 55     ● 60     △ 65     ■ 70

**Cooling Power (W)** — y-axis: 12, 13, 14, 15, 16, 17, 18, 19, 20

**Bandwidth (GB/s)** — x-axis: 5, 10, 15, 20, 25

A Bandwidth increment of 15 GB/s

About 1.5 W increment in cooling power

# Closed-Page Policy

Payload Size:
- ■ 128B  ■ 112B  ■ 96B  ■ 80B
- ■ 64B  ■ 48B  ■ 32B  ■ 16B



**Bandwidth (GB/s)** vs **Access Pattern**

16 vaults: linear, random
1 vault: linear, random

# Closed-Page Policy

Payload Size: ■ 128B  ■ 112B  ■ 96B  ■ 80B
■ 64B  ■ 48B  ■ 32B  ■ 16B

**Bandwidth (GB/s)**

Applications benefit from bank-level parallelism
not from spatial locality

linear          random          linear          random

**16 vaults**                    **1 vault**

**Access Pattern**

IISWC'17  CASL  Computer Architecture and Systems Laboratory  Georgia Tech  comparch

# Achieving High Bandwidth

▶ Promote bank-level parallelism

▶ Remap data to avoid internal organization bottlenecks

▶ Concatenate requests to use bandwidth effectively

# Latency Deconstruction

# Latency Deconstruction

# Latency Deconstruction Summary

TX Path: **287 ns**

**260 ns**

**547 ns**

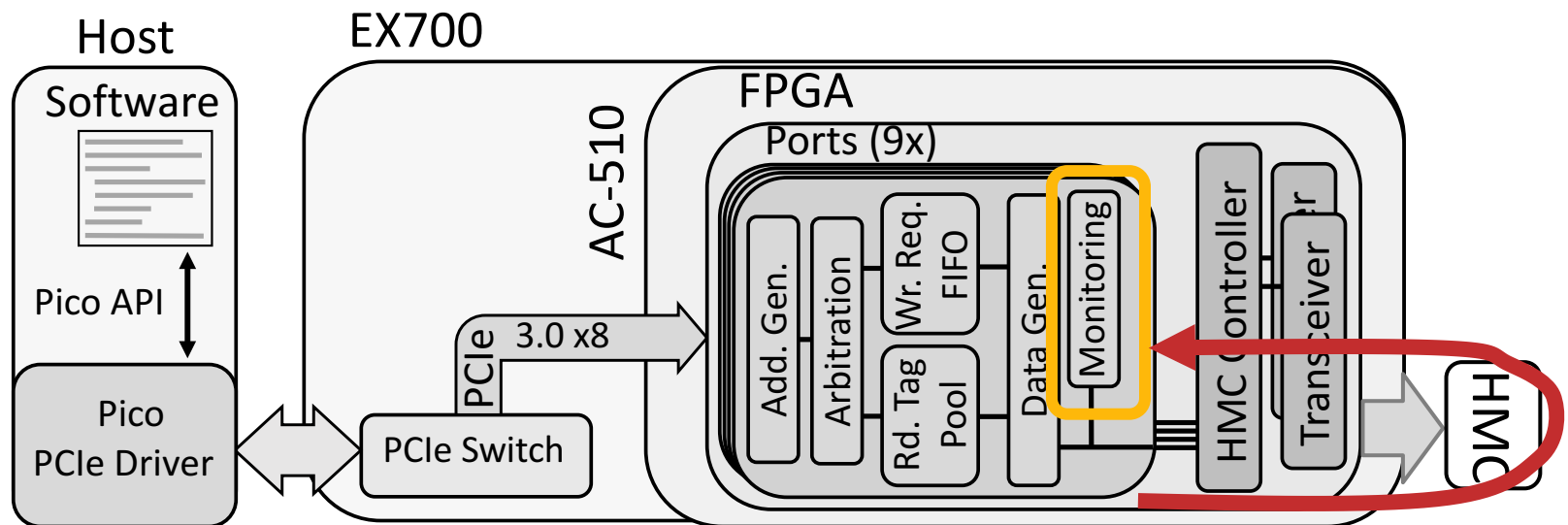| | |
|---|---|
| Conversion to flits & buffering | 10 cycles |
| Round-robin arbitration among ports | 2-9 cycles |
| Add packet fields & flow control | 10 cycles |
| Serialization | 10 cycles |
| Transmission (128B) | 15 cycles |

Freq.: 187.5 MHz
Cycle: 5.3 ns

IISWC'17   CASL — Computer Architecture and Systems Laboratory   Georgia Tech   comparch

# Low-Load Latency

# Low-Load Latency

**Size 16B**

**Size 32B**

Latency (us)

1.40
1.20
1.00
0.80
0.60
0.40

2 4 6 8 10 12 14 16 18 20 22 24 26 28

**Larger request size**

**Faster latency increment**

Size 64B    **Max**    Size 128B    **Avg.**    **Min**

Latency (us)

2.00
1.80
1.60
1.40
1.20
1.00
0.80
0.60
0.40

2 4 6 8 10 12 14 16 18 20 22 24 26 28

**Number of Read Requests**

**Number of Read Requests**

IISWC'17    CASL Computer Architecture and Systems Laboratory    Georgia Tech    comparch

# Low-Load Latency

**Average latency increases because of maximum latency increments**

Avg.

Min

**Number of Read Requests**

**Number of Read Requests**

# Low-Load Latency

**125 ns is spent in the HMC**

Avg.
Min

Number of Read Requests    Number of Read Requests

# High-Load Latency

# Latency-Bandwidth

# Latency-Bandwidth

**●size 16B**   **△size 32B**   **✕size 64B**   **●size 128B**

**4-banks**

**2-banks**

Latency (us)

Bandwidth (GB/s)

**Lowest Request Rate**

Each layer/bank has a queue.
Limiting factor can be the queue size.

# Conclusions

▸ Mixing read and write requests and using large request sizes lead to effective use of bi-directional bandwidth.

▸ Distributing accesses prevents internal bottlenecks and exploits bank-level parallelism.

▸ Controlling the request rate to avoid high latency.

▸ Employing fault-tolerant mechanisms and using proper cooling solutions enables temperature-sensitive operations to reach a higher bandwidth.

▸ Reducing latency overhead of the infrastructure will greatly benefit latency.

CASL
Computer Architecture
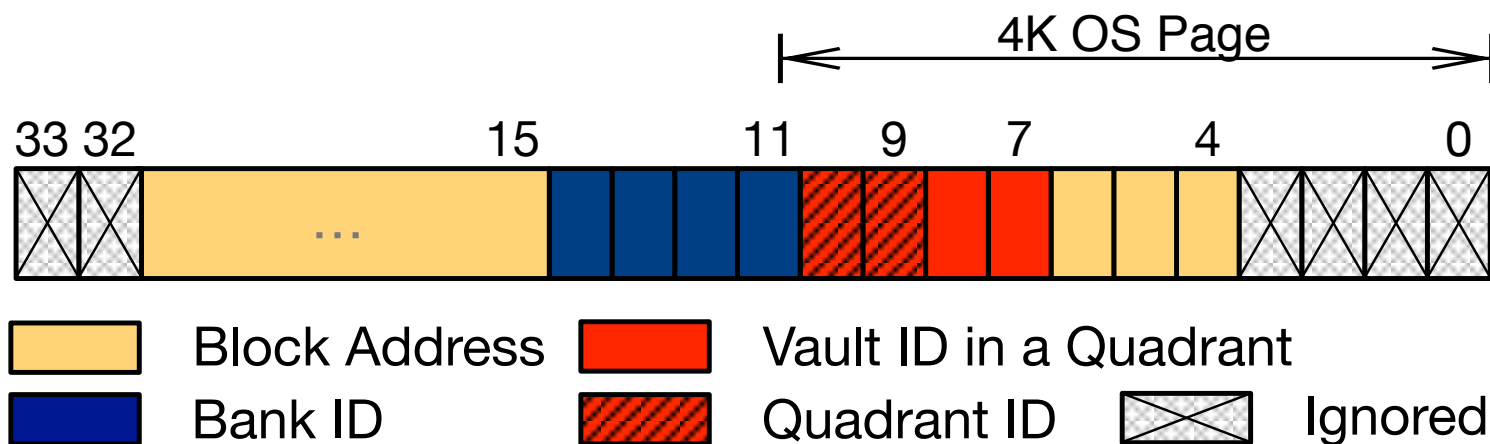and Systems Laboratory

Georgia
Tech

comparch

# Backup Slides

# HMC Memory Addressing

Closed-page policy            Page Size = 256 B

Low-order-interleaving address mapping policy

34-bit address field:



4K OS Page

| 33 32 | | 15 | | 11 | 9 | 7 | | 4 | | 0 |

Block Address — Vault ID in a Quadrant

Bank ID — Quadrant ID — Ignored

# Experimental Setup III



| | Full-scale GUPS | Small-scale GUPS | Stream GUPS |
|---|---|---|---|
| **Addresses** | Random Configurable Mask | Random Configurable Mask | Defined by User |
| **Request Rate** | Maximum | Configurable | Minimum |
| **Experiment** | Bandwidth Power Temperature High-Load Latency | Latency-Bandwidth | Integrity Check Low-Load Latency |

IISWC'17

# Thermal Configurations

| Configuration Name | DC Power Supply: Voltage | Current | 15 W Fan Distance | Average HMC Idle Temperature |
|---|---|---|---|---|
| Cfg1 | 12 V | 0.36 A | 45 cm | 43.1° C |
| Cfg2 | 10 V | 0.29 A | 90 cm | 51.7° C |
| Cfg3 | 6.5 V | 0.14 A | 90 cm | 62.3° C |
| Cfg4 | 6.0 V | 0.13 A | 135 cm | 71.6° C |

# Cooling Power

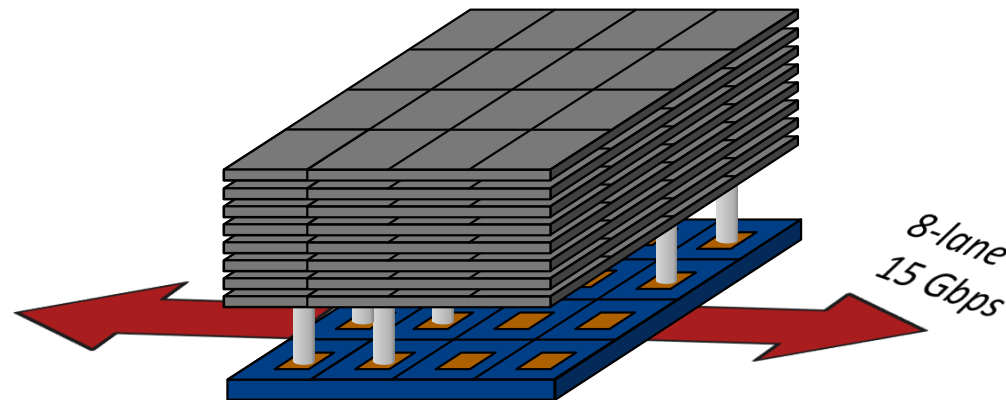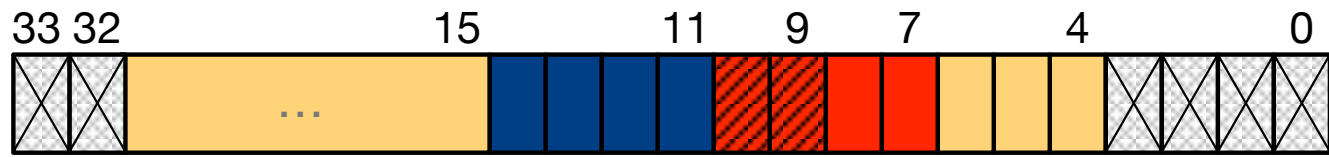| Configuration | Cooling Power |
|---------------|---------------|
| cfg1 | 19.32 W |
| cfg2 | 15.90 W |
| cfg3 | 13.90 W |
| cfg4 | 10.78 W |

# HMC Communication II
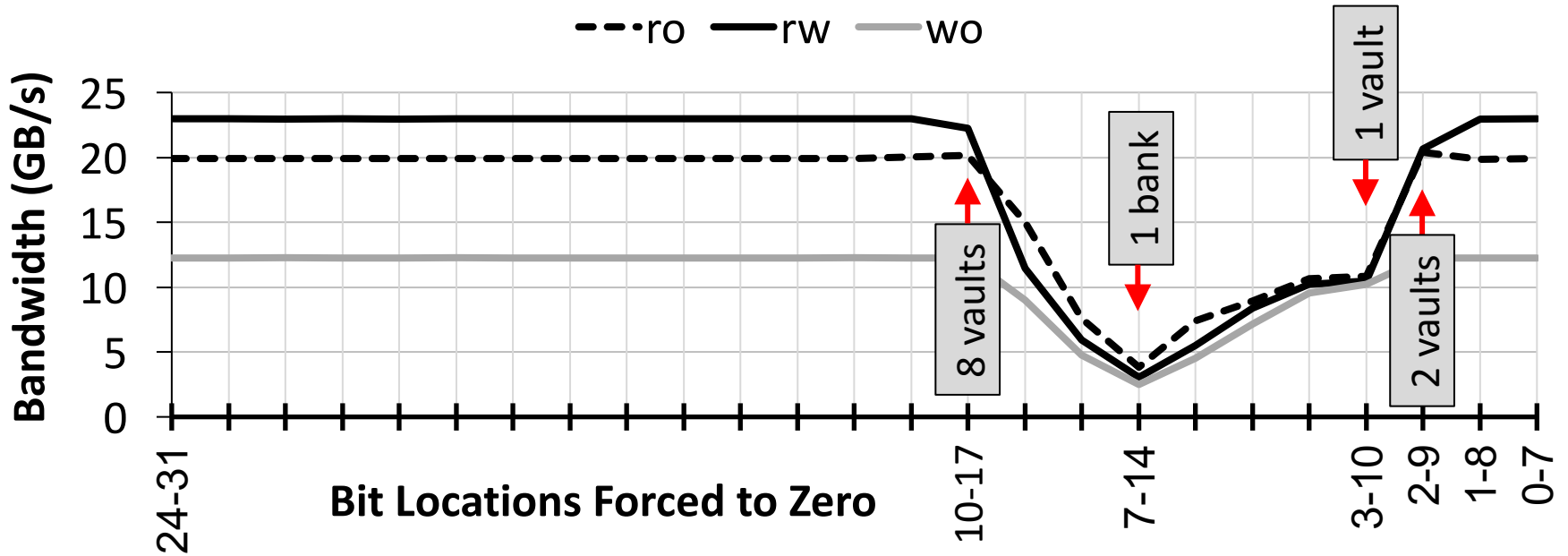
Two/Four full duplex external links:

- Width of 16 or 8 lanes

- Configurable speeds of 10, 12.5, and 15 Gbps



8-lane
15 Gbps

$$BW_{peak} = 2\,\mathrm{link} \times 8\,\mathrm{lanes/link} \times 15\,\mathrm{Gbps} \times 2\,\mathrm{full\ duplex}$$
$$= 480\,\mathrm{Gbps} = 60\,\mathrm{GB/s.}$$

# Address Mapping

# Bandwidth II

Legend: 128B (teal bars), 64B (yellow bars), 32B (red bars), MRPS 128B (dotted line with squares), MRPS 64B (dashed line with triangles), MRPS 32B (solid line with diamonds)

Y-axis (left): Bandwidth (GB/s) — 0 to 25
Y-axis (right): #Req. (M) / Second — 0 to 350
X-axis: Access Pattern — 16 vaults, 8 vaults, 4 vaults, 2 vaults, 1 vault, 8 banks, 4 banks, 2 banks, 1 bank

IISWC'17   CASL Computer Architecture and Systems Laboratory   Georgia Tech   comparch

# Latency-Bandwidth II

# Latency-Bandwidth III



Legend: 1 bank, 2 banks, 4 banks, 8 banks, 1 vault, 2 vaults, 4 vaults, 8 vaults, 16 vaults

Four plots of Latency (us) vs Bandwidth (GB/s): Size 16B, Size 32B, Size 64B, Size 128B.