# An Instruction Roofline Model for GPUs

***Nan Ding***, Samuel Williams

Computational Research Division
Lawrence Berkeley National Lab
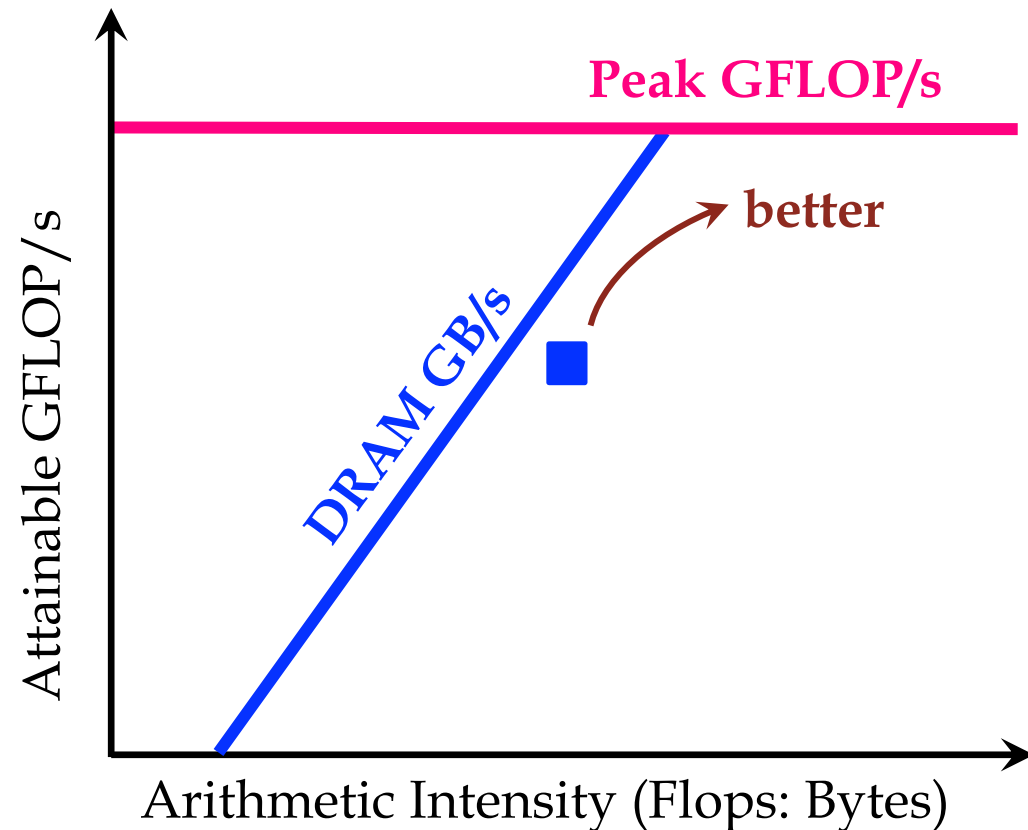{nanding, swwilliams}@lbl.gov
Nov.18th ,2019

BERKELEY LAB
LAWRENCE BERKELEY NATIONAL LABORATORY

U.S. DEPARTMENT OF ENERGY

# History of Roofline Models

- Sustainable performance is bound by

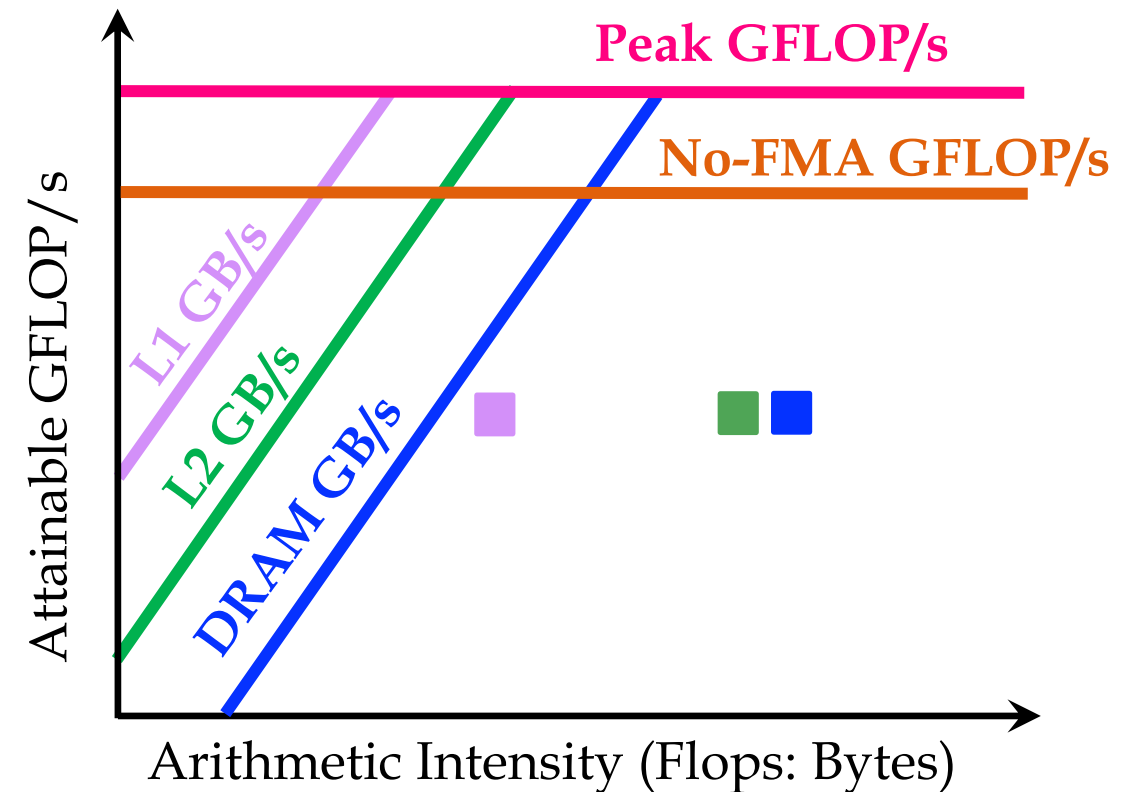$$\text{GFLOP/s} = \min \begin{cases} \text{Peak GFLOP/s} \\ \text{AI} * \text{GB/s} \end{cases}$$

## Roofline Model

- Arithmetic Intensity (AI) : FLOPs/Byte



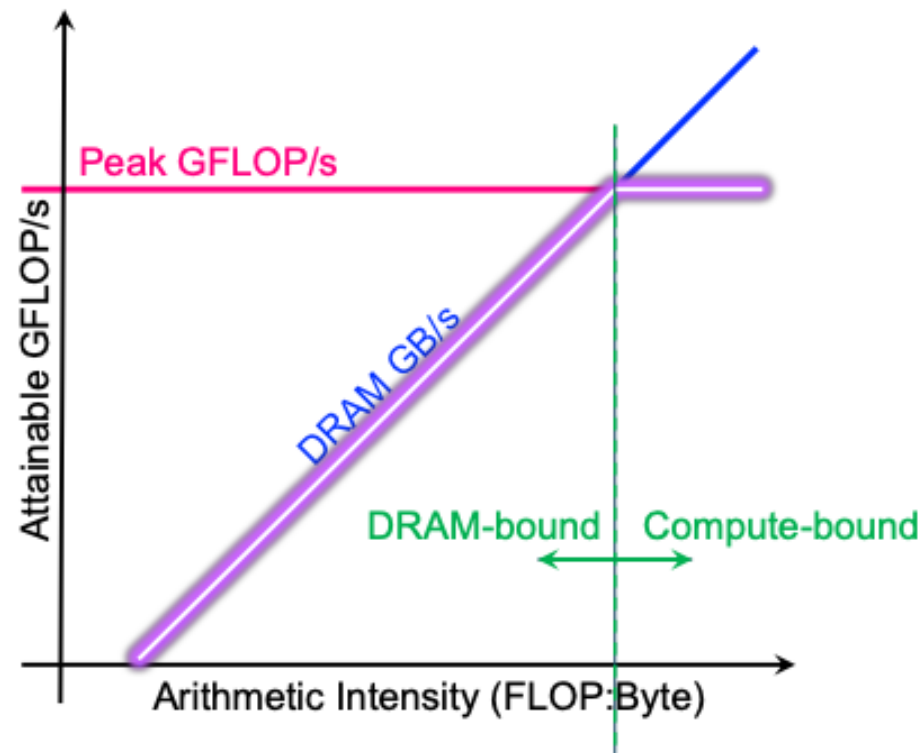## Hierarchical Roofline Model

- Arithmetic Intensity (AI) : FLOPs/Byte(L1/L2/DRAM)
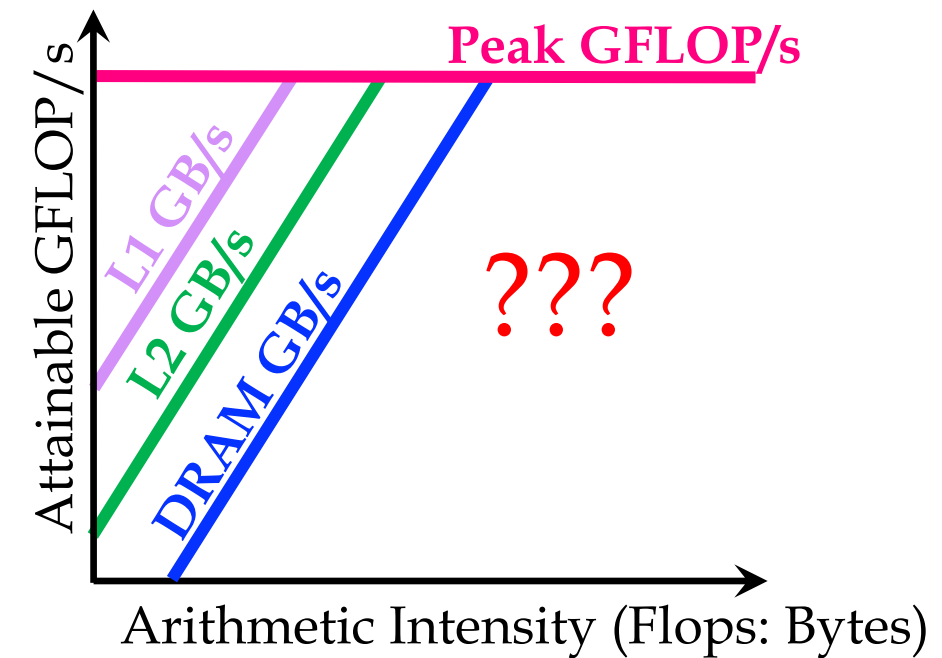- Additional compute ceilings: No-FMA peak

# Roofline is Useful

## Driving performance optimization

3

# However…

- Even with sufficient data locality, one cannot guarantee high performance
    - Pathological memory access patterns?
    - Re-design the data layout?
    - Limited by instruction throughput?

- Many applications perform more integer operations than floating-point/no-flops

# Motivation for the Instruction Roofline Model

## Emerging Domains

- Mixed precision, Integer-heavy

- No floating points operations

**More than Flops**

## Architectural Evolution

- Instruction throughput
  - pipeline utilization

- Warp efficiency
  - Thread predication

- Memory access patterns
  - reduce wasted transactions
  - reduce redundant access

**A New Set of Metrics**

## Practical Use

- What is holding you back?

- What optimizations should be performed?

- When to stop optimization?

**Drive Code Optimization in a Good Visual Manner**

- Sustainable performance is bound by

$$\textbf{GFLOP/s} = \min \begin{cases} \textbf{Peak GFLOP/s} \\ \text{sAI} * \text{GB/s} \end{cases}$$
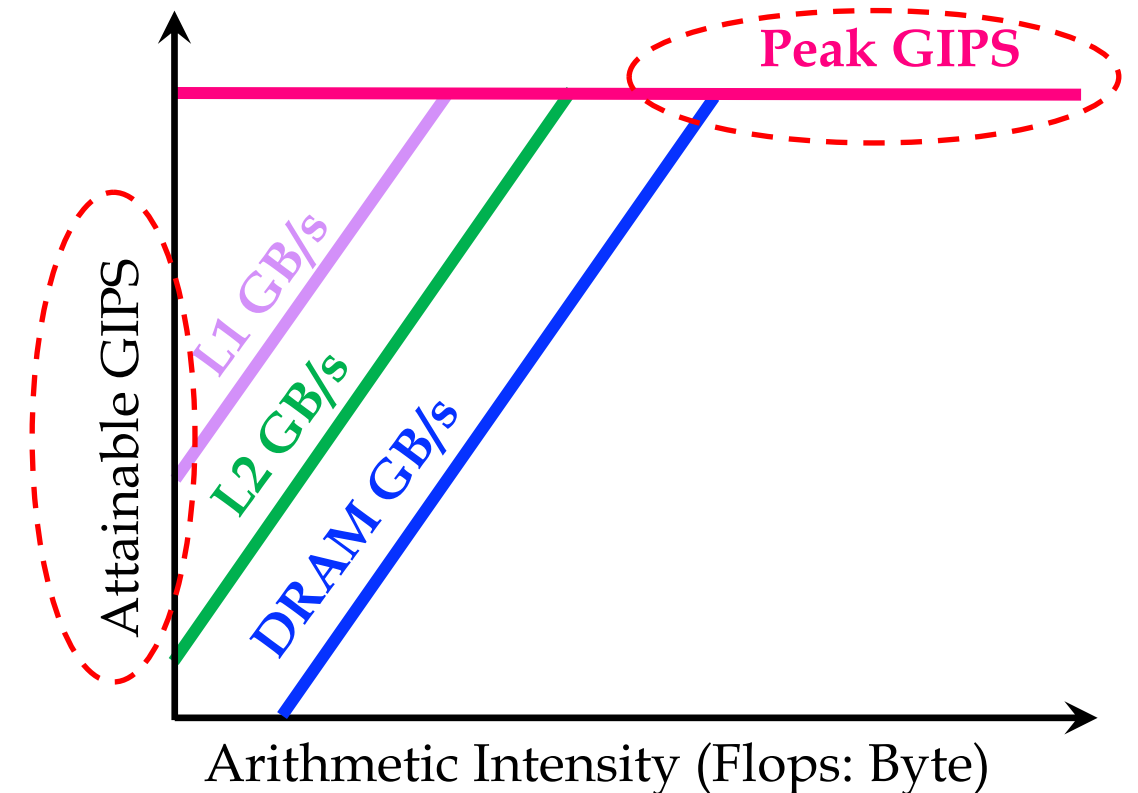
$$\textbf{GIPS} = \min \begin{cases} \textbf{Peak GIPS} \\ \text{II} * \text{GB/s} \end{cases}$$

**Expanding the applicability of roofline to several emerging computational domains**

**Form the basis for several subsequent Instruction Roofline-oriented performance analysis technologies**
- **Identify fetch-decode-issue bottlenecks**
- **Function unit utilization (FPU, tensor, integer, etc...)**



Peak GIPS

Attainable GIPS

L1 GB/s

L2 GB/s

DRAM GB/s

Arithmetic Intensity (Flops: Byte)

# The Second Step to Instruction Roofline Model

- Sustainable performance is bound by

$$\textbf{GFLOP/s} = \min \begin{cases} \textbf{Peak GFLOP/s} \\ \textbf{AI} * \text{GB/s} \end{cases}$$

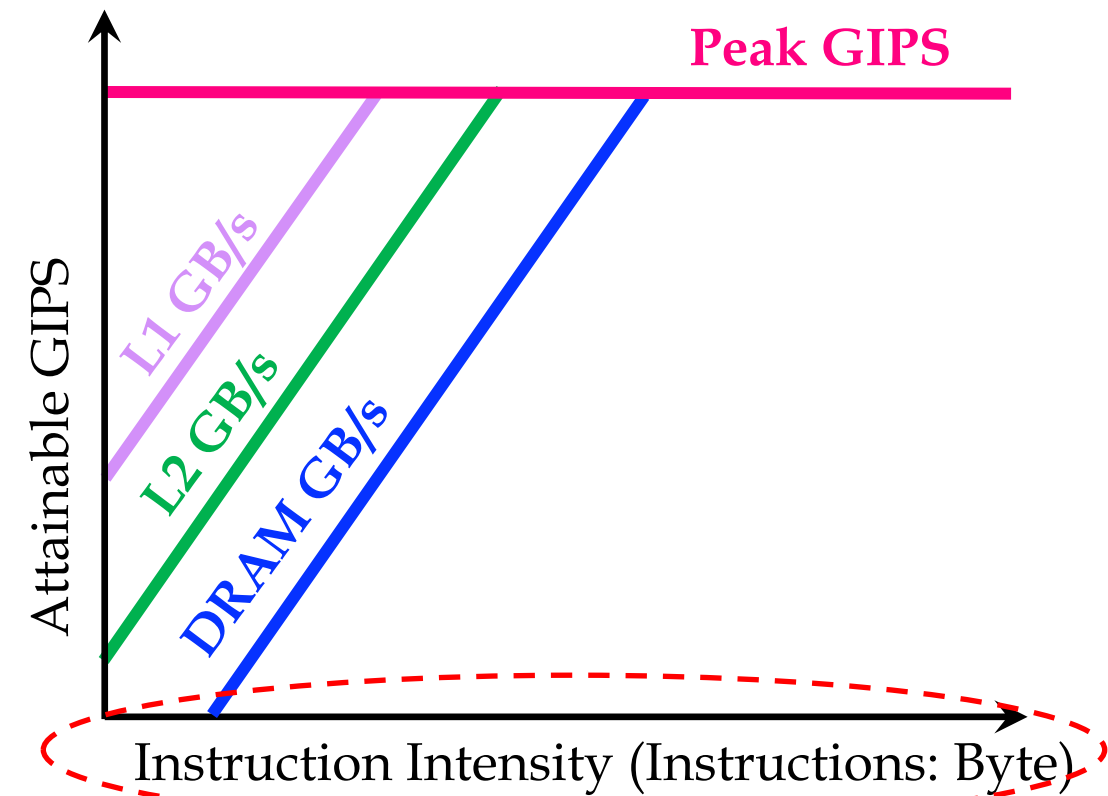$$\textbf{GIPS} = \min \begin{cases} \textbf{Peak GIPS} \\ \textbf{II} * \text{GB/s} \end{cases}$$

- **Instruction Intensity**: Instructions per Byte

**Expanding the applicability of roofline to several emerging computational domains**

*Limitation:*

*Hard to motivate more performance analysis techniques, such as memory pattern access*

# A Final Step to Instruction Roofline Model on GPUs

- **Instruction Intensity**
  - Instructions per Transaction

> Expanding the applicability of roofline to more performance analysis technologies GPUs

> Form the basis for several subsequent Instruction Roofline-oriented performance analysis technologies on GPUs:
> - Memory access patterns

- **Memory Transaction**
  - the natural unit to access data on NVIDIA GPUs
  - the natural unit to **analyze memory access**
  - a warp-level load/store ->  1 - 32 transactions

GIPS = min $\begin{cases} \text{Peak GIPS} \\ \text{Instructions/Transaction * GTXN/s} \end{cases}$

# Instruction Roofline Performance Model
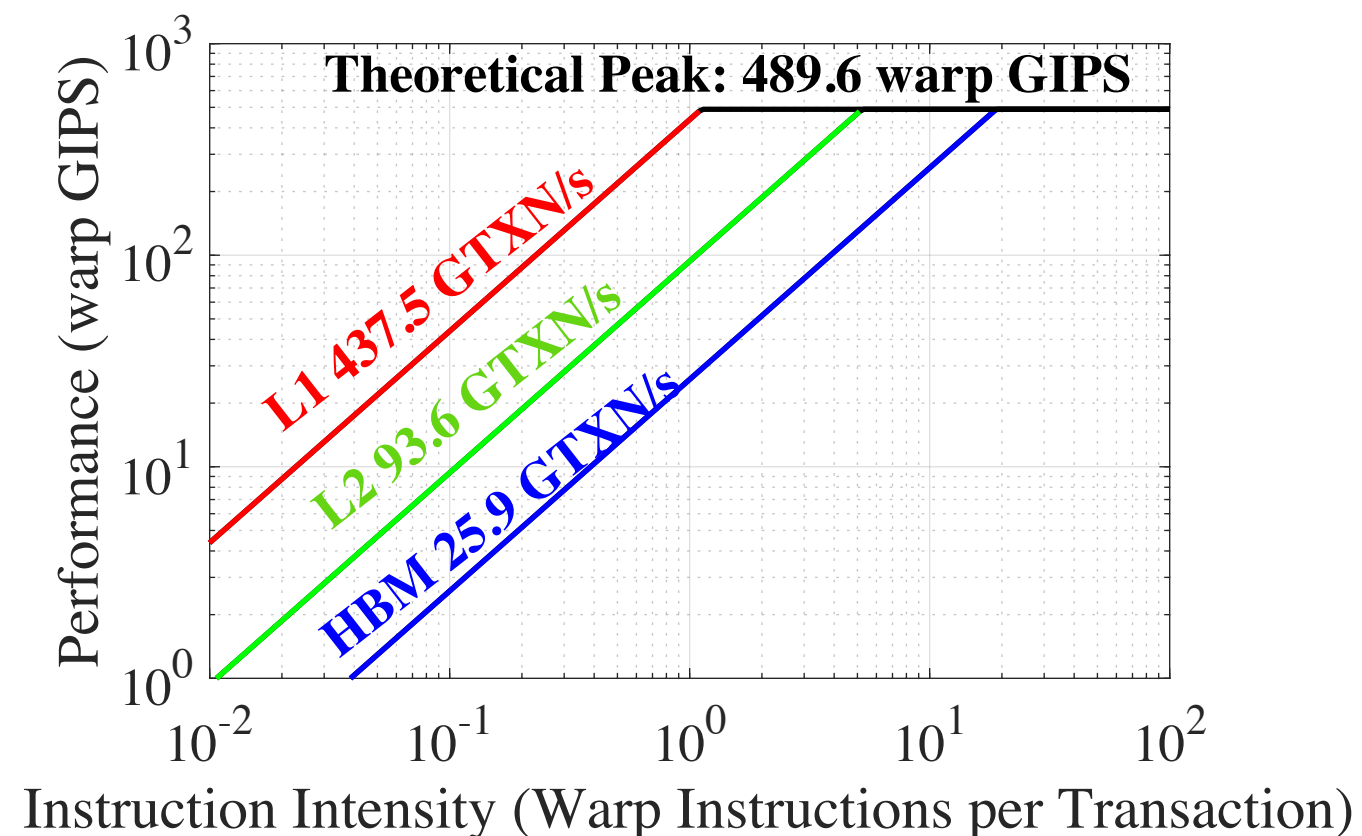
- **Sustainable performance of is bound by**

$$\text{GIPS} = \min \begin{cases} \text{Peak GIPS} \\ \text{Instruction Intensity} * \text{GTransaction/s} \end{cases}$$

- **Theoretical Peak on V100:**

80 SMs x 4 warp scheduler x 1 inst/cyc x 1.53GHz = 489.6 GIPS

- **Memory ceilings on V100:**
  - Based on the GB/s from Empirical Roofline Toolkit[1]
  - Calculate the number of equivalent 32-byte transactions



[1] https://bitbucket.org/berkeleylab/cs-roofline-toolkit

# Instruction Throughput

- **Instruction throughput**

    All instruction,  Transactions of each memory level(L1/L2/HBM), runtime



- **Insights:**
    1. Distance between the ceilings and dots can tell memory-bound or instruction-bound
    2. Distance between the two plots (different memory level) can tell the data reuse.

# Capabilities of Instruction Roofline Model -- Instruction Throughput

- **Instruction throughput**

    All instruction, Transactions of each memory level(L1/L2/HBM), runtime
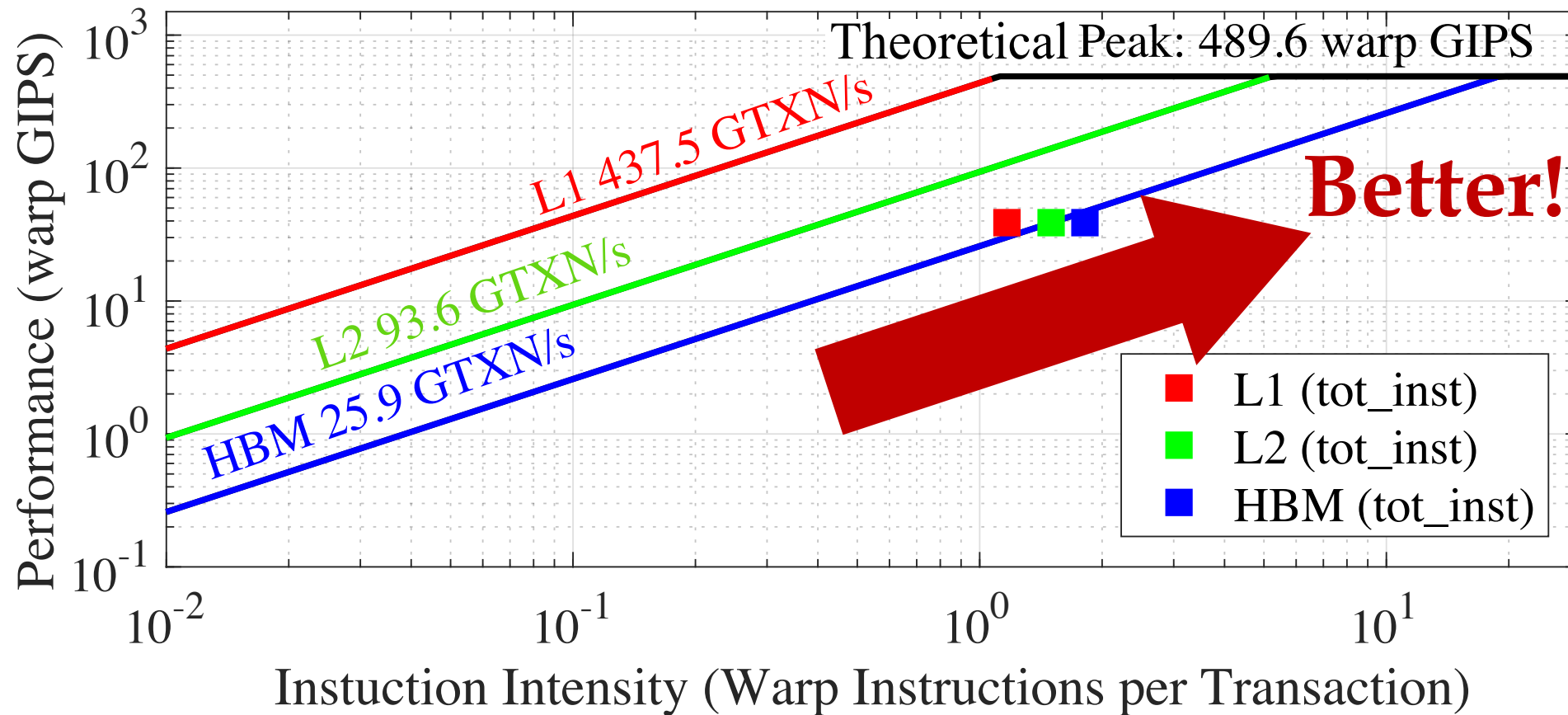


- **Insights:**
    1. Distance between the ceilings and dots can tell memory-bound or instruction-bound
    2. Distance between the two plots (different memory level) can tell the data reuse.

# Memory Access Patterns

# Memory Access Pattern is Critical to Application Execution Time

Easy to code in an inefficient memory pattern

Low performance

Hidden deep in the code

Time consuming to reason the performance

# Capabilities of Instruction Roofline Model -- Global Memory Patterns

1 warp-level load/store -> 1 to 32 transactions depending on memory patterns

**"Stride-0"**

$$\frac{1\ warp\ Global\ LDST}{1\ Global\ Transaction}$$

**"Stride-8"**

$$\frac{1\ warp\ Global\ LDST}{32\ Global\ Transactions}$$
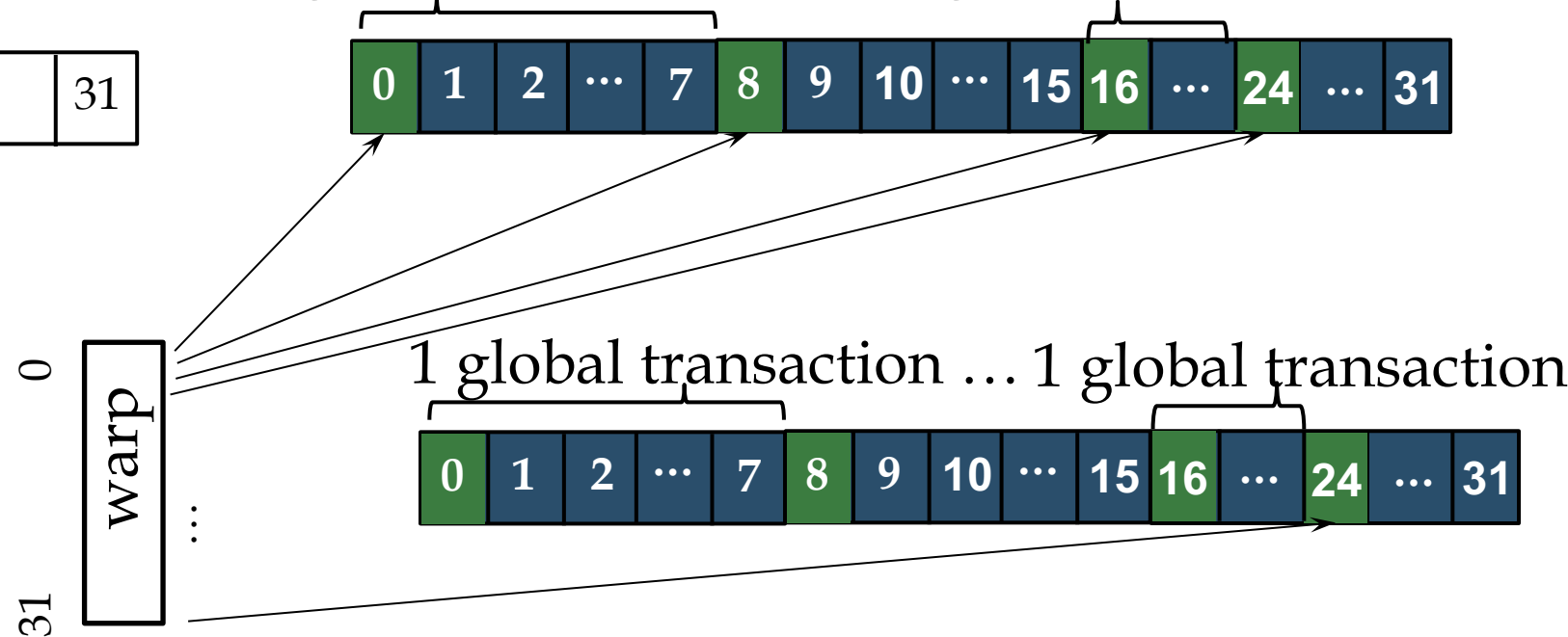
**"Stride-1" (Unit Stride)**

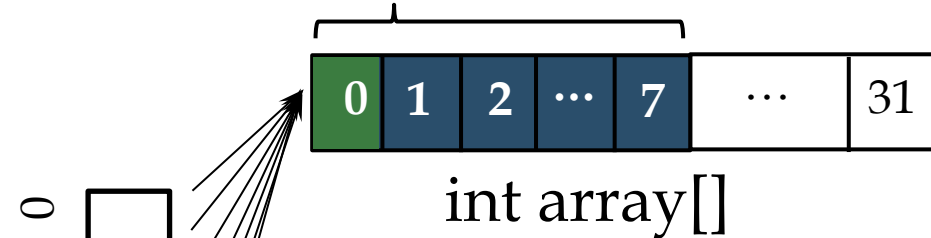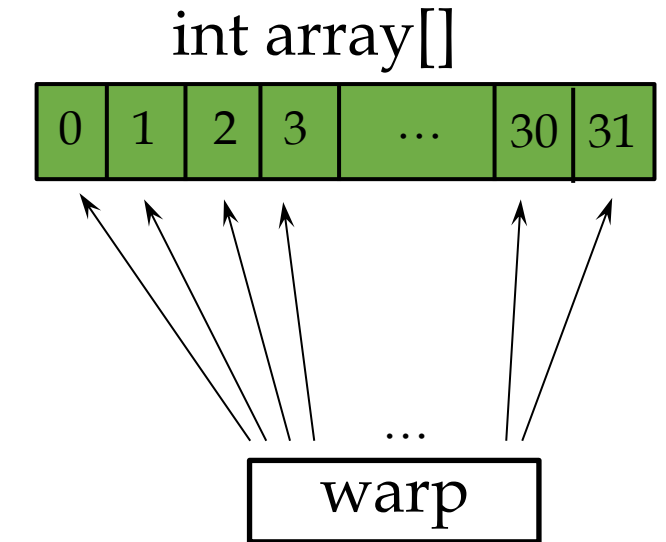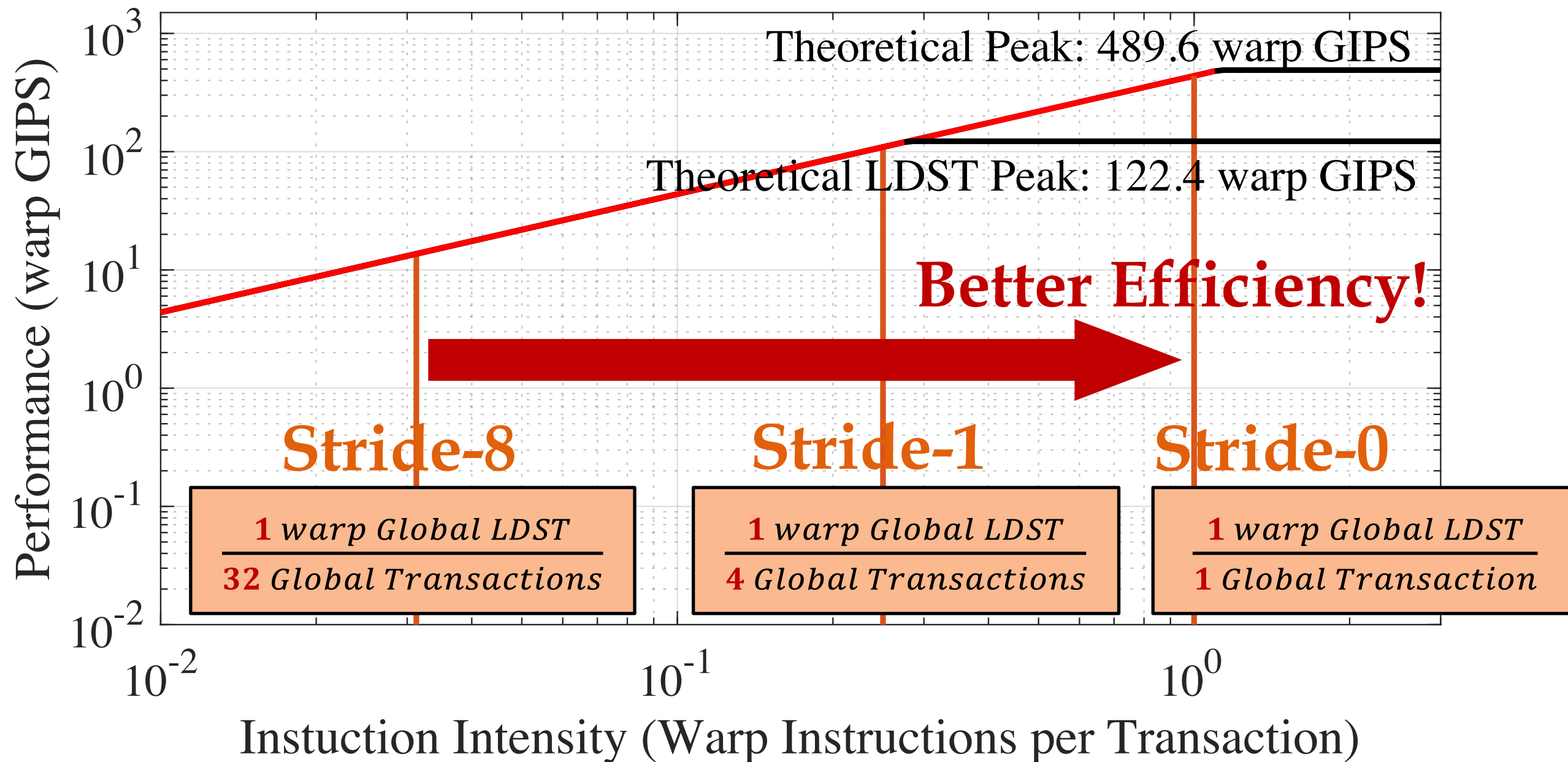$$\frac{1\ warp\ Global\ LDST}{4\ Global\ Transactions}$$

Useful data   Waste data

1 global transaction = 32 Bytes

one cache line: 128 bytes

int array[]

# Capabilities of Instruction Roofline Performance Model
## ---- Three Intensity ``Walls'' for Stride Global Memory Access Patterns

Breakdown the **L1 dot** into **Global Memory Only metrics -> Stride Global Memory Patterns** according to **Global Memory Walls**



X-axis: $\dfrac{Total\ instructions}{Total\ L1\ (global + shared + local)\ transactions}$

Y-axis : $\dfrac{Total\ instructions}{Runtime}$

X-axis: $\dfrac{Global\ LDST\ instructions}{Global\ transactions}$

Y-axis : $\dfrac{Global\ LDST\ instructions}{Runtime}$

**Better Efficiency!**

Performance (warp GIPS)

Global 437.5 GTXN/s

Stride-8    Stride-1    Stride-0

Theoreti...

Thoeretical LD...

Instuction Intensity (Warp Instructions per Transaction)

**Notional**

**Physical**

__shared__ int array[32][32]

32 banks per shared memory row
Each bank is 4 Byte

No bank conflict

__shared__ int array[32][32]

**32-way** bank conflict

**Reduce bank conflicts**
__shared__ int array[32][32+1]

No bank conflict

- "No bank conflict" $= \dfrac{1\ warp\ Shared\ LDST}{1\ Shared\ Transaction}$

  - different 4-byte word, different bank
  - same 4-byte word , same bank

- "32-way bank conflict" $= \dfrac{1\ warp\ Shared\ LDST}{32\ Shared\ Transactions}$

  - different 4-byte words, same bank

Breakdown the **L1 dot** into **Shared Memory Only metrics -> banked Shared Memory Patterns**
according to **Shared Memory Walls**



X-axis: $\dfrac{Total\ instructions}{Total\ L1\ (global + shared + local)\ transactions}$

Y-axis : $\dfrac{Total\ instructions}{Runtime}$

X-axis: $\dfrac{Shared\ LDST\ instructions}{Shard\ transactions}$

Y-axis : $\dfrac{Shared\ LDST\ instructions}{Runtime}$

**Better Efficiency!**

# An example to understand the outputs from Instruction Roofline Model

| Example: Matrix Transpose | | |
|---|---|---|
| **Description** | A-> A$^T$, stored in column major | |
| **Matrix size** | 1024 x 1024 | |
| **Machine** | NVIDIA's latest V100 GPU | |
| **Implementations** | Naive | Coalesced | Coalesced_NoBankConflict |
| | Simple copy | Coaleced global memory access | Based on "Coalesced" Reduce shared memory bank conflicts |
| | Using 32×8 thread blocks operating on 32×32 matrix tiles | | |

# *Naive Implementation*

**4096 Bytes (1024 floats)**



**Input Matrix: A (column major)**

**4096 Bytes (1024 floats)**



**Output Matrix: A$^T$ (column major)**

# *Global memory stride access*



# *Instruction Throughput*

# Coalesced Implementation

**Input Matrix: A (column major)**
- 4096 Bytes (1024 floats)
- 4096 Bytes (1024 floats)
- 32 x 32 (floats)
- **Coalesced read**

**Shared memory**
- Recalculating the array index

**Output Matrix: A$^T$ (column major)**
- 4096 Bytes (1024 floats)
- 4096 Bytes (1024 floats)
- **Coalesced write**

# Global memory stride access
- Performance (warp GIPS)
- Theoretical LDST Peak: 122.4 warp GIPS
- Global 437.5 GTXN/s
- **Better Efficiency!**
- Stride-8
- Stride-1
- Stride-0
- Global (ldst_inst)
- Instuction Intensity (Warp Instructions per Transaction)

# Instruction Throughput
- Performance (warp GIPS)
- Theoretical Peak: 489.6 warp GIPS
- L1 437.5 GTXN/s
- L2 93.6 GTXN/s
- HBM 25.9 GTXN/s
- *Coalesced*
- *Naive*
- L1 (tot_inst)
- L2 (tot_inst)
- HBM (tot_inst)
- Instuction Intensity (Warp Instructions per Transaction)

# Coalesced Implementation

**4096 Bytes (1024 floats)**

**4096 Bytes (1024 floats)**

**32 x 32 (floats)**

**Coalesced read**

**Input Matrix: A (column major)**

**Shared memory**

k      k+16

| T0 | T1 |
| T2 | T3 |
| T4 | T5 |
| ... | ... |
| T30 | T31 |

**16-way bank conflict**

**Recalculating the array index**

**4096 Bytes (1024 floats)**

**4096 Bytes (1024 floats)**

**Coalesced write**

**Output Matrix: Aᵀ (column major)**

# Shared memory bank access



Theoretical LDST Peak: 122.4 warp GIPS

Performance (warp GIPS)

32-way bank conflict

1-way bank conflict

**Better Efficiency!**

□ Shared (ldst_inst)

Instruction Intensity (Warp Instructions per Transaction)

# Instruction Throughput



Theoretical Peak: 489.6 warp GIPS

Performance (warp GIPS)

L1 437.5 GTXN/s

L2 93.6 GTXN/s

HBM 25.9 GTXN/s

*Coalesced*

*Naive*

■ L1 (tot_inst)
■ L2 (tot_inst)
■ HBM (tot_inst)

Instruction Intensity (Warp Instructions per Transaction)

# *Coalesced_NoBankConflict Implementation*

**4096 Bytes (1024 floats)**

**4096 Bytes (1024 floats)** (vertical label)

**32 x 32 (floats)**

**Coalesced read**

**Input Matrix: A (column major)**

**Shared memory**

**[32][32+1]**

k          k+16

T0    T1
T2    T3
T4    T5
...   ...
T30   T31

**Recalculating the array index**

**4096 Bytes (1024 floats)**

**4096 Bytes (1024 floats)** (vertical label)

**Coalesced write**

**Output Matrix: A$^T$ (column major)**

# *Shared memory bank access*

Theoretical LDST Peak: 122.4 warp GIPS

32-way bank conflict

No Bank conflict

**Better Efficiency!**

Performance (warp GIPS)

Instuction Intensity (Warp Instructions per Transaction)

□ Shared (ldst_inst)

# *Instruction Throughput*

Theoretical Peak: 489.6 warp GIPS

L1  437.5 GTXN/s

L2 93.6 GTXN/s

HBM 25.9 GTXN/s

*NoBankConflict*

*Coalesced*

*Naive*

Performance (warp GIPS)

Instuction Intensity (Warp Instructions per Transaction)

■ L1 (tot_inst)
■ L2 (tot_inst)
■ HBM (tot_inst)

# Summary

- **Instruction Throughput**
  - Expanding the applicability of roofline to several emerging computational domains

- **Global Memory Access Patterns**
  - Quantify the memory access pattern, e.g. unit-stride vs. gather/scatter

- **Shared Memory Access Patterns**
  - Denote the efficiency of shared memory access.

**There's more in the paper !!**
- **Thread Predication**
- **More Examples**
  - **HPGMG (mixed precision): three implementations.**
  - **BatchSW (integer-only): two implementations**
- **Tensor core**
  - **WMMA**
  - **cuBLAS**

Closing Thoughts

BERKELEY LAB
LAWRENCE BERKELEY NATIONAL LABORATORY

U.S. DEPARTMENT OF
ENERGY

# What the Instruction Roofline Models Tell us...

| Emerging Domains | Architectural Evolution | Practical Use |
|---|---|---|
| • Mixed precision<br>• Integer-heavy<br>• No floating points operations<br><br>**Applicability to several emerging computational domains** | • Instruction throughput<br>   • pipeline utilization<br>• Quantify memory pattern<br>   • Unit-stride, scatter/gather<br>• Efficiency of memory access<br>• Warp efficiency<br>   • Thread predication<br><br>**Applicability of roofline to GPUs with greater insights** | • Unified visualization of bandwidth and access efficiency<br><br>**Rapidly tell how different aspects of modern GPU architectures constrain performance.** |

# Future Work

- Apply our methodology to other accelerated architectures

- Extend the access efficiency concept to networking, I/O, and Lustre file systems.

# Acknowledgement

# Questions?

Backup