# The Locality Descriptor

## A Holistic Cross-Layer Abstraction to Express Data Locality in GPUs

**Nandita Vijaykumar**

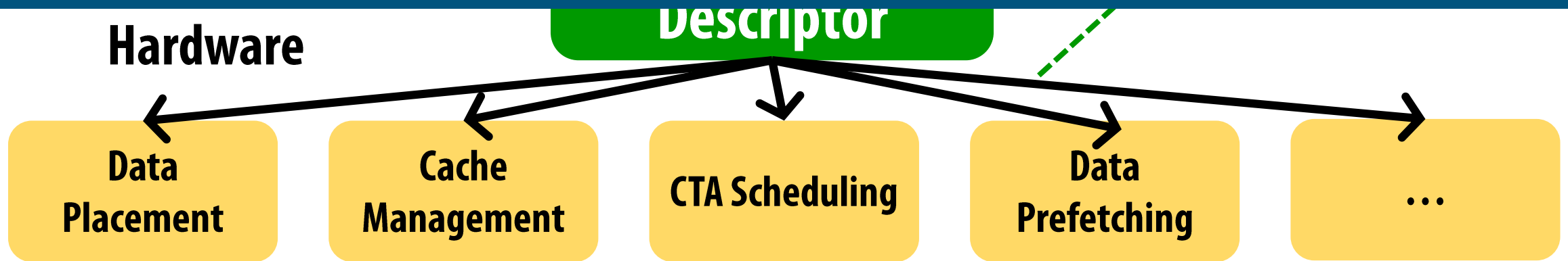**Eiman Ebrahimi, Kevin Hsieh, Phillip B. Gibbons, Onur Mutlu**

Carnegie Mellon University

NVIDIA

ETH Zürich

# Executive Summary

**Exploiting data locality in GPUs is a challenging task**

## Performance Speedups:
### 26.6% (up to 46.6%) from <u>cache locality</u>
### 53.7% (up to 2.8x) from <u>NUMA locality</u>

**Descriptor**

**Hardware**

| Data Placement | Cache Management | CTA Scheduling | Data Prefetching | ... |

2

# Outline

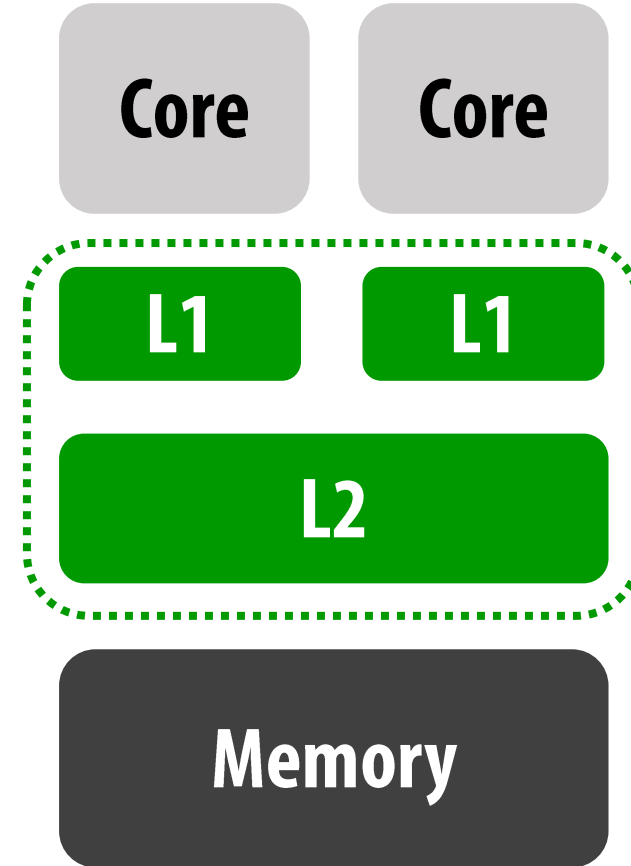Why leveraging data locality is challenging?

Designing the Locality Descriptor

Evaluation

# Data locality is critical to GPU performance

**Two forms of data locality:**

**Reuse-based locality (cache locality)**
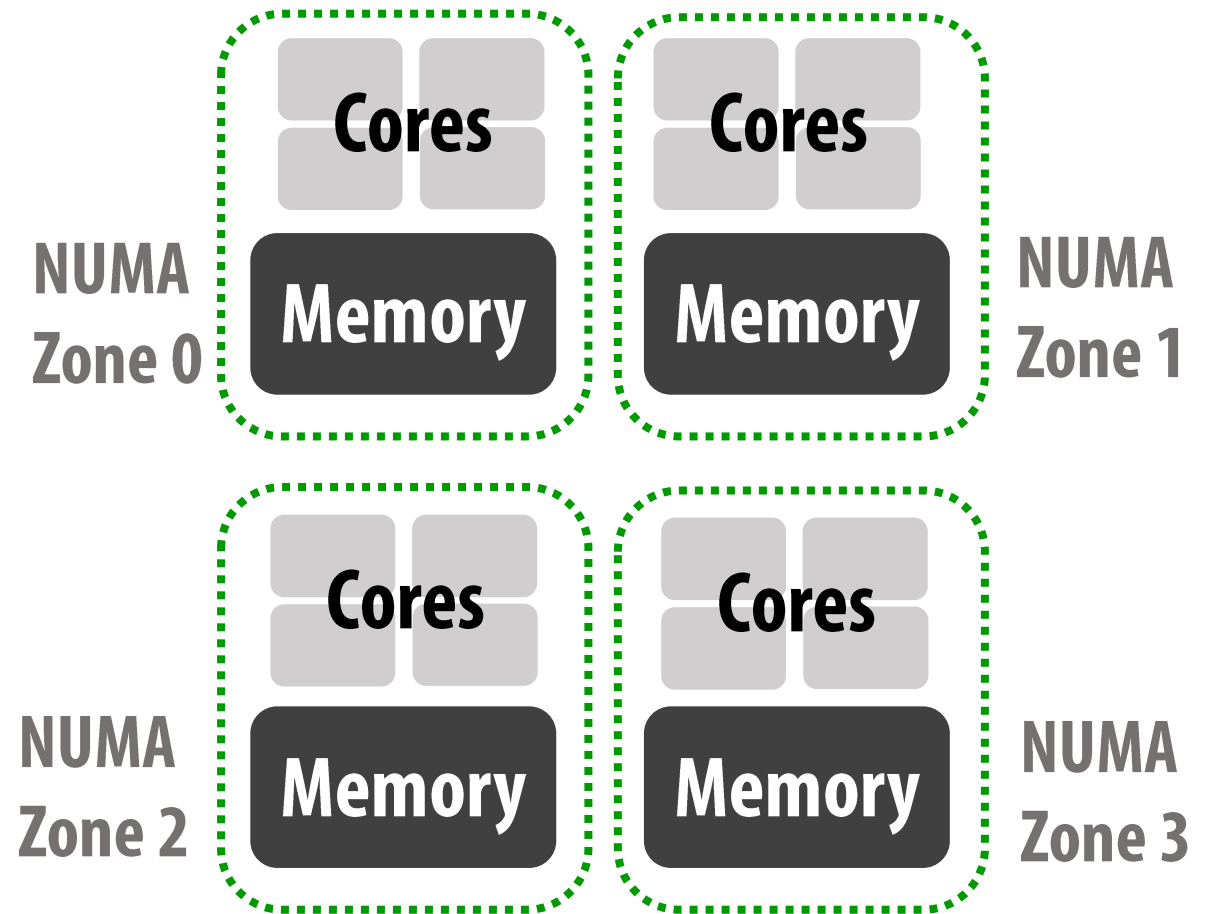
NUMA locality

**Reuse-based (Cache Locality)**

4

# Data locality is critical to GPU performance

**Two forms of data locality:**

Reuse-based locality (cache locality)

**NUMA locality**



NUMA Zone 0 — Cores / Memory
NUMA Zone 1 — Cores / Memory
NUMA Zone 2 — Cores / Memory
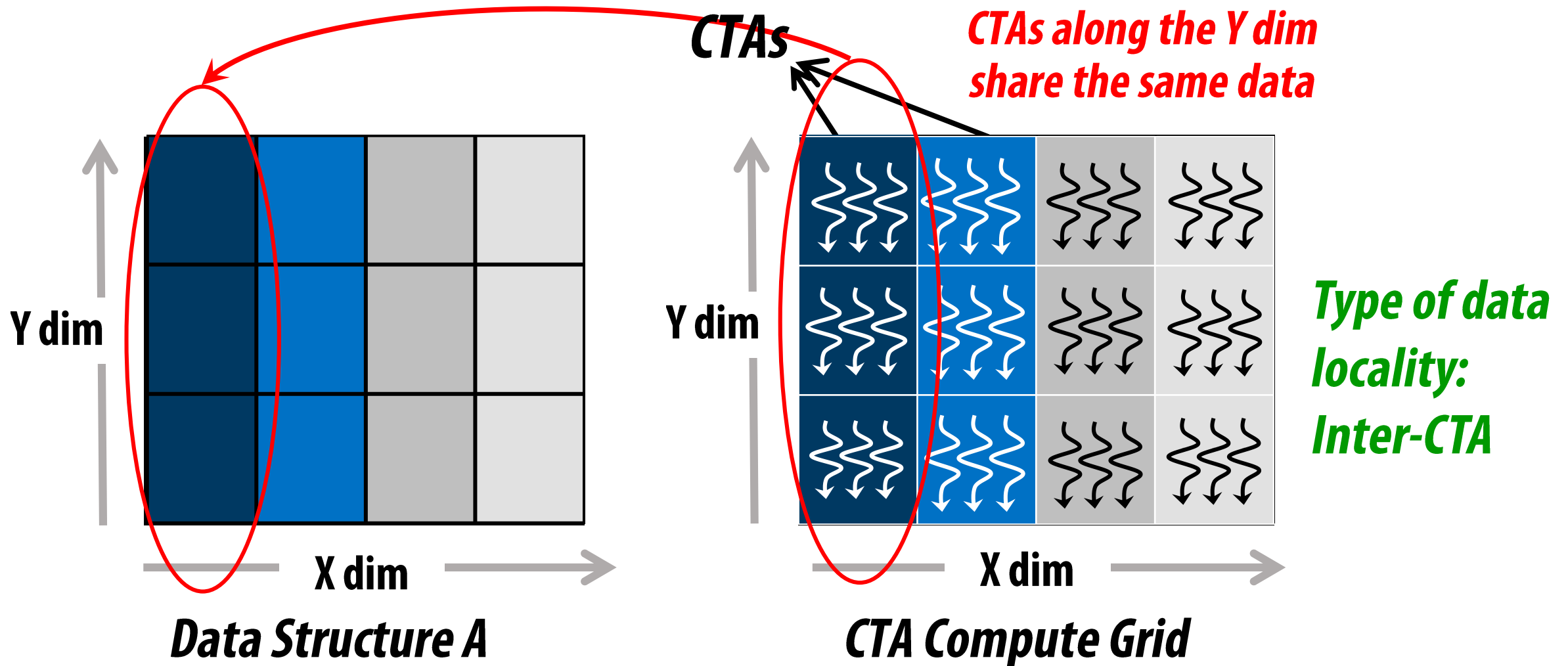NUMA Zone 3 — Cores / Memory

**NUMA Locality**

5

**The GPU execution and programming models are designed to explicitly express <u>parallelism</u>…**
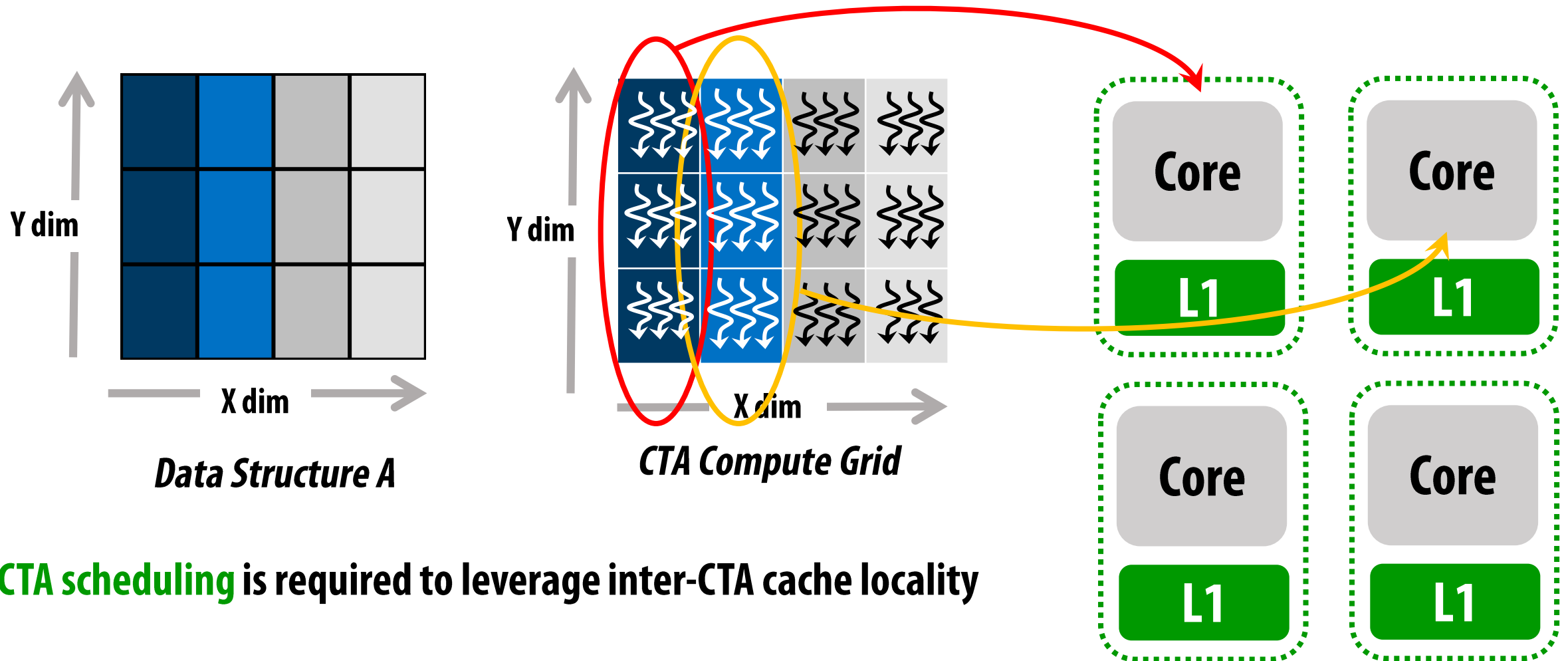
*But there is no explicit way to express <u>data locality</u>*

**Exploiting data locality in GPUs is a challenging and elusive feat**

# A case study in leveraging data locality: `Histo`



**CTAs along the Y dim share the same data**

**Type of data locality: Inter-CTA**

Y dim

X dim

**Data Structure A**

CTAs

Y dim

X dim

**CTA Compute Grid**

7

# Leveraging cache locality
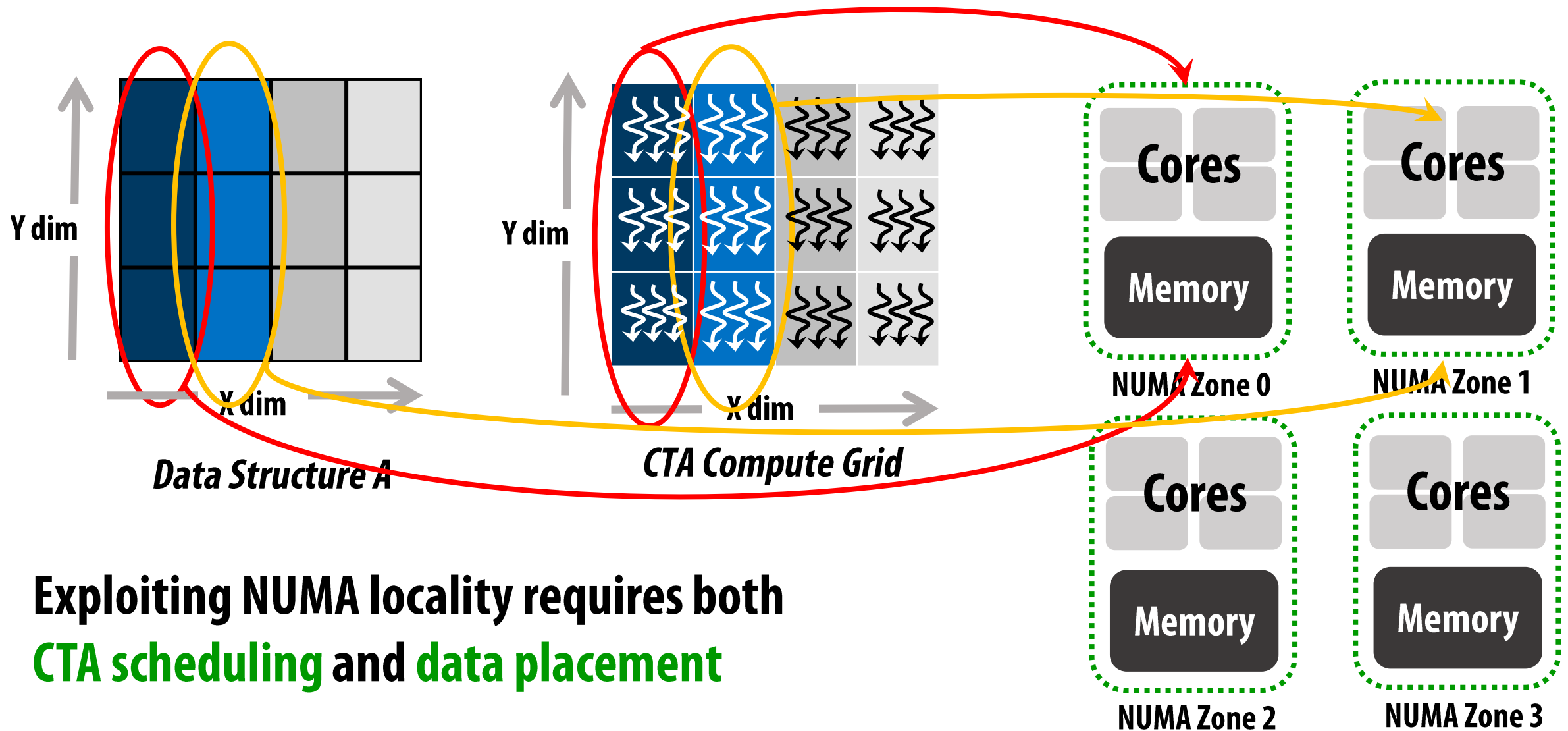


**Data Structure A**

**CTA Compute Grid**

**CTA scheduling** is required to leverage inter-CTA cache locality

CTA scheduling is **insufficient**: we also need other techniques

8

# Leveraging NUMA locality

Y dim

X dim

**Data Structure A**

Y dim

X dim

**CTA Compute Grid**

Cores

Memory

**NUMA Zone 0**

Cores

Memory

**NUMA Zone 1**

Cores

Memory

**NUMA Zone 2**

Cores

Memory

**NUMA Zone 3**

**Exploiting NUMA locality requires both**
**CTA scheduling and data placement**

# Today, leveraging data locality is challenging

**As a programmer:**
- **No easy access** to architectural techniques – CTA scheduling, cache management, data placement, etc.
- Even when using work-arounds, optimization is **tedious** and **not portable**

**As the architect:**
- **Key program semantics** are not available to the hardware

*Where to place data?*
*Which CTAs to schedule together?*

# To make things worse:

There are many different <u>locality types</u>: **Inter-CTA, inter-warp, intra-thread, …**
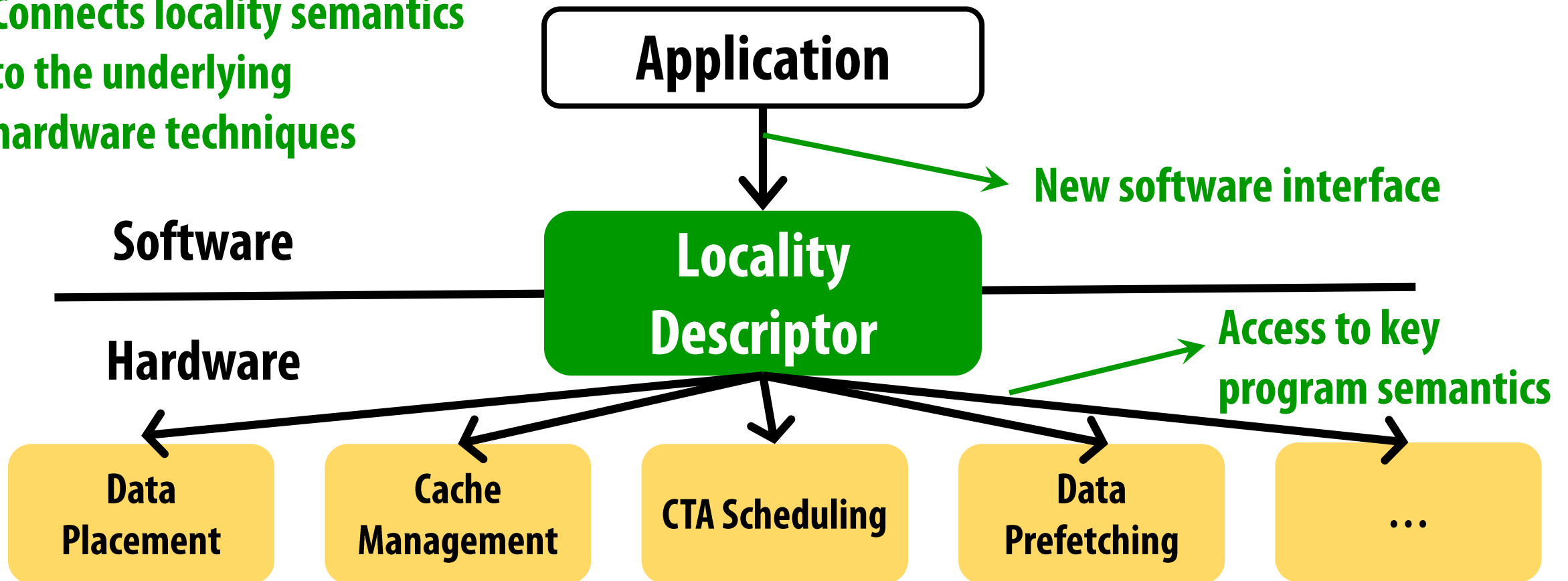
Each type requires a <u>different set</u> of architectural techniques:
- **Inter-CTA locality** requires CTA scheduling + prefetching
- **Intra-thread** locality requires cache management
- …

# The Locality Descriptor

## A hardware-software abstraction to express and exploit data locality

**Connects locality semantics to the underlying hardware techniques**
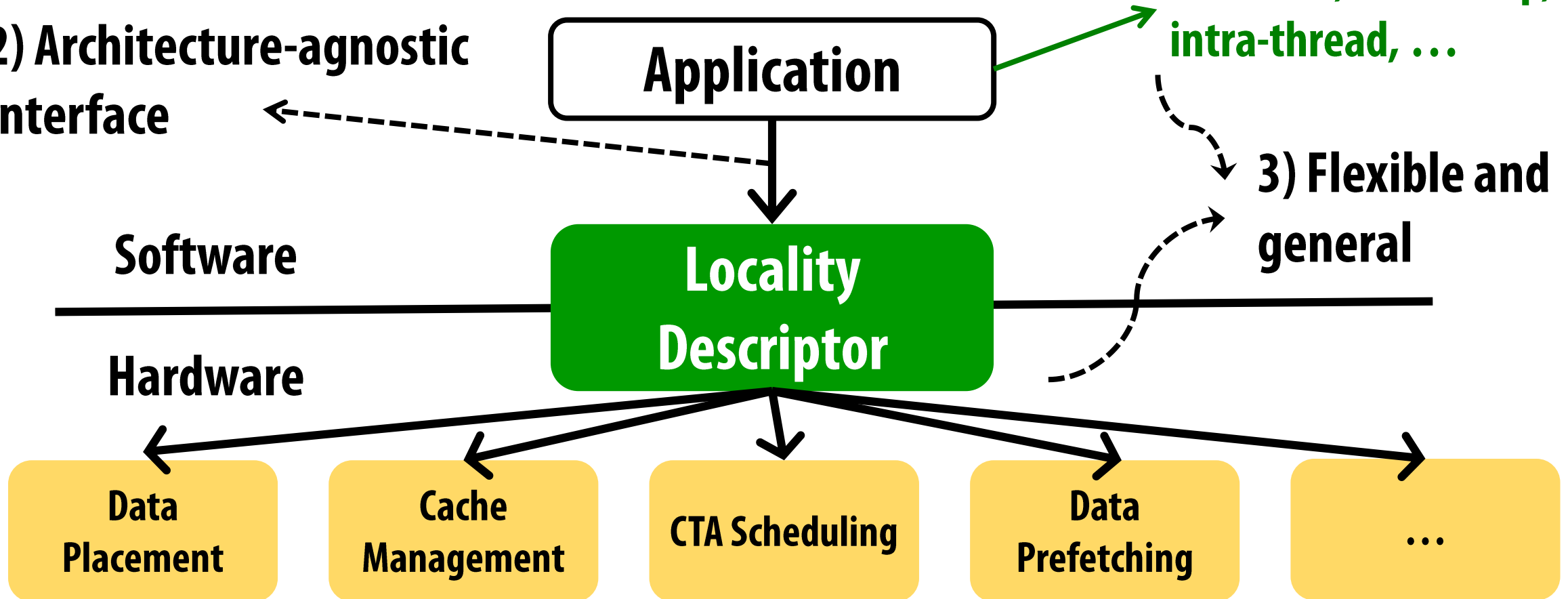
**Application**

→ **New software interface**

**Software**

**Locality Descriptor**

**Hardware**

→ **Access to key program semantics**

| Data Placement | Cache Management | CTA Scheduling | Data Prefetching | … |

# Goals in designing the Locality Descriptor

**1) Supplemental and hint-based**

**2) Architecture-agnostic interface**

**Application**

**Inter-CTA, inter-warp, intra-thread, …**

**3) Flexible and general**

**Software**

**Locality Descriptor**

**Hardware**

| Data Placement | Cache Management | CTA Scheduling | Data Prefetching | … |

# Designing the Locality Descriptor

`LocalityDescriptor ldesc(X; Y, Z);`

Hardware

- Data Placement
- Cache Management
- CTA Scheduling
- Data Prefetching
- ...

# An Overview: The components of the Locality Descriptor

*1) Data Structure*

`LocalityDescriptor` `ldesc`(A, len,

`INTER-THREAD`,

*3) Tile Semantics*

`tile,`

`loc`);

*2) Locality Type*

*4) Locality Semantics*

15

# Outline

**Why leveraging data locality is challenging?**

**Designing the Locality Descriptor**
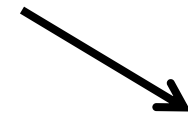
**Evaluation**

# 1. How to choose the basis of the abstraction?

*Key Idea: Use the <u>data structure</u> as the basis to describe data locality*

- Architecture-agnostic
- Each data structure is accessed the same way by all threads

A new instance is required for each important data structure

```
LocalityDescriptor ldesc(A);
```

*Data Structure*

# 2. How to communicate with hardware?

**Locality type drives architecture mechanisms**

Inter-CTA, inter-warp, intra-thread, …

**Architecture-agnostic interface**

Application

**Architecture-specific optimizations**

Software

Hardware

**Locality Descriptor**

Data Placement

Cache Management

CTA Scheduling

Data Prefetching

…

19

# 2. How to communicate with hardware?

*Key Idea:  Use <u>locality type</u> to drive underlying architectural techniques*

Origin of locality (or locality type) causes the challenges in <u>exploiting it</u>

E.g.:
>    Inter-CTA locality requires CTA scheduling as reuse is <u>across threads</u>
>    Intra-thread locality requires cache management to avoid thrashing

Locality type is application-specific and known to the programmer

# 2. How to communicate with hardware?

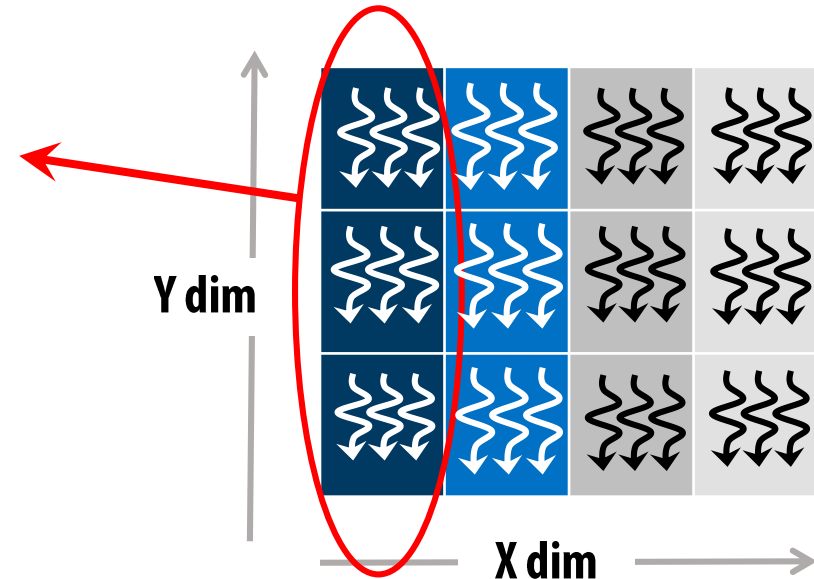*Key Idea:  Use <u>locality type</u> to drive underlying architectural techniques*

**Three fundamental types:**

    **INTER-THREAD**
    **INTRA-THREAD**
    **NO-REUSE**

**INTER-THREAD**

Y dim

X dim

*CTA Compute Grid*

```
LocalityDescriptor ldesc(A);INTER-THREAD);
```

# Driving underlying architectural techniques



**Locality Type?**

INTER_THREAD → Determined based on more information

INTRA_THREAD → CTA Scheduling / Cache Soft Pinning / Memory Placement (if NUMA)

NO_REUSE → Cache Bypassing / CTA Scheduling (if NUMA) / Memory Placement (if NUMA)

# 3. How to describe locality?

Key Idea: Partition the data structure and compute grid into tiles

Basic unit of locality

```
LocalityDescriptor ldesc(A, INTER-THREAD);tile);
```

Data Tile

Compute Tile

```
tile((X_tile, Y_len, 1),
     (1, GridSize.y, 1),
     (1, 0, 0));
```

Data Tile Dimensions

Compute Tile Dimensions

Compute-Data Map

# Additional features of the Locality Descriptor

**Locality type insufficient to inform underlying architectural techniques**

      **(INTER-THREAD, INTRA-THREAD, NO-REUSE)**

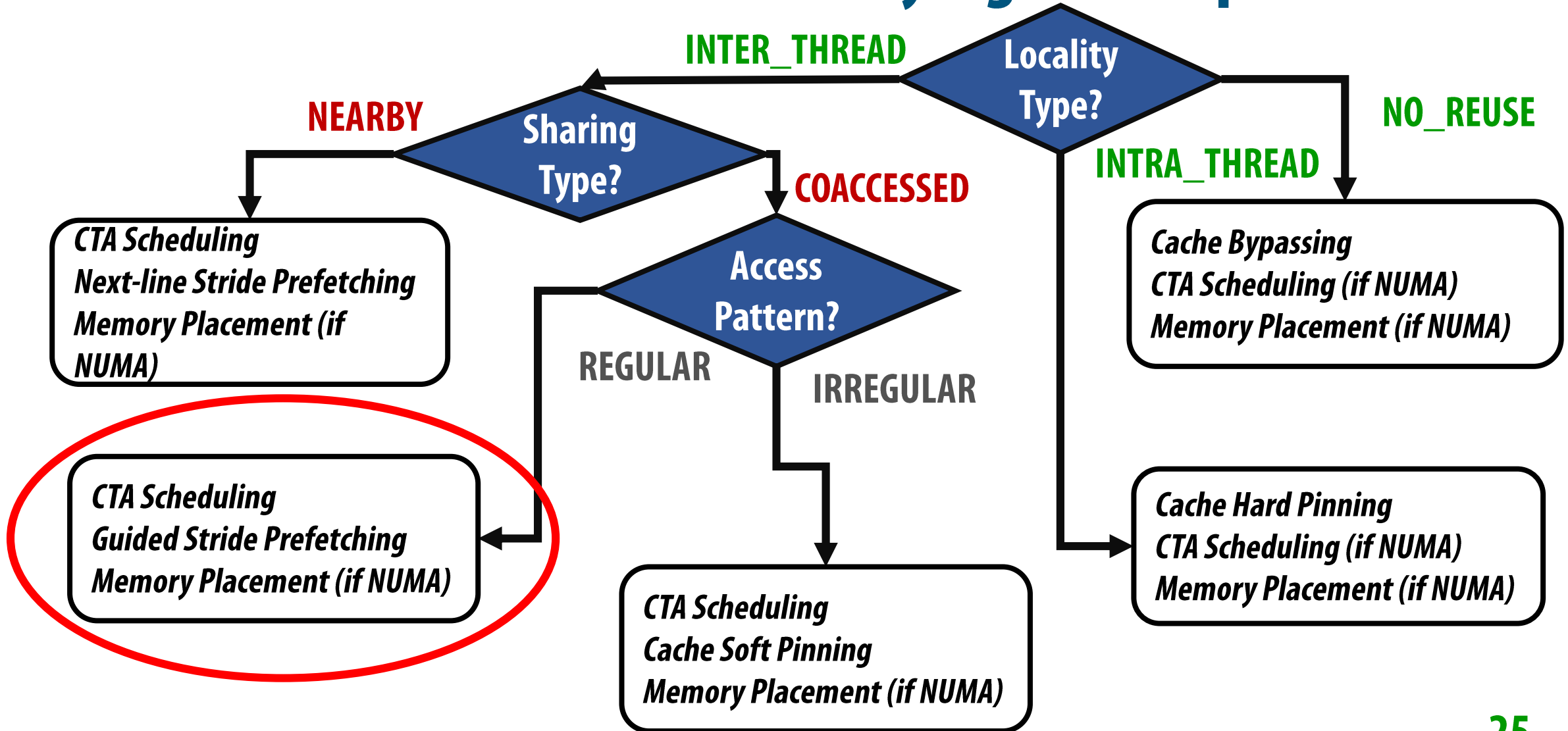**In addition, we also have Locality Semantics to include:**

    **- Sharing Type**

    **- Access Pattern**

```
                    (COACCESSED, REGULAR, X_len)



LocalityDescriptor ldesc(A, INTER-THREAD, tile);loc);
```

# A decision tree to drive underlying techniques

**INTER_THREAD**

**Locality Type?**

**NO_REUSE**

**NEARBY**

**Sharing Type?**

**INTRA_THREAD**

**COACCESSED**

CTA Scheduling
Next-line Stride Prefetching
Memory Placement (if NUMA)

Cache Bypassing
CTA Scheduling (if NUMA)
Memory Placement (if NUMA)

**Access Pattern?**

**REGULAR**

**IRREGULAR**

CTA Scheduling
Guided Stride Prefetching
Memory Placement (if NUMA)

Cache Hard Pinning
CTA Scheduling (if NUMA)
Memory Placement (if NUMA)

CTA Scheduling
Cache Soft Pinning
Memory Placement (if NUMA)

25

# Leveraging the Locality Descriptor

`LocalityDescriptor` **ldesc**`(A, `**`INTER-THREAD`**`, `**`tile, loc`**`);`



**Y dim**

**X dim**

*Data Structure A*

**Y dim**

**X dim**

*CTA Compute Grid*

**Architectural techniques:**
1) **CTA Scheduling**
2) **Prefetching**
3) **Data Placement**

26

# CTA Scheduling



Cluster 0   Cluster 1   Cluster 2   Cluster 3

Compute Tile 0   Compute Tile 1   Compute Tile 2   Compute Tile 3

Y dim

X dim

**Data Structure A**

Y dim

X dim

**CTA Compute Grid**

| Cluster 0 | Cluster 1 | Cluster 2 | Cluster 3 |
|-----------|-----------|-----------|-----------|
| Core 0 | Core 1 | Core 2 | Core 3 |
| L1D | L1D | L1D | L1D |

# Leveraging the Locality Descriptor

`LocalityDescriptor ldesc(A, INTER-THREAD, tile, loc);`



Y dim

X dim

**Data Structure A**

Y dim

X dim

**CTA Compute Grid**

**Architectural techniques:**
1) **CTA Scheduling**
2) **Prefetching**
3) **Data Placement**

# Data Placement

Cluster 0  Cluster 1  Cluster 2  Cluster 3

Cluster Queues

Y dim

X dim

**Data Structure A**

Y dim

X dim

**CTA Compute Grid**

Cluster 0  Cluster 1  Cluster 2  Cluster 3

SM 0  SM 1  SM 2  SM 3

L1D  L1D  L1D  L1D

Memory  Memory  Memory  Memory

*NUMA Zone 0*  *NUMA Zone 1*  *NUMA Zone 2*  *NUMA Zone 3*

# Leveraging the Locality Descriptor

```
LocalityDescriptor ldesc(A, INTER-THREAD, tile, loc);
```



*All threads stall because
they wait on the same data*

**Architectural techniques:**
1) CTA Scheduling
2) Prefetching
3) Data Placement

*Y dim*

*X dim*

**Data Structure A**

*Y dim*

*X dim*

**CTA Compute Grid**

# Outline

**Why leveraging data locality is challenging?**

**The Locality Descriptor**

**Evaluation**

# Methodology

**Evaluation Infrastructure:** GPGPUSim v3.2.2

**Workloads:** Parboil, Rodinia, CUDA SDK, Polybench

**System Parameters:**

**Shader Core:** 1.4 GHz; GTO scheduler [50]; 2 schedulers per SM, Round-robin CTA scheduler

**SM Resources Registers:** 32768; Scratchpad: 48KB, L1: 32KB, 4 ways

**Memory Model:** FR-FCFS scheduling [59, 60], 16 banks/channel

**Single Chip System:** 15 SMs; 6 memory channels; L2: 768KB, 16 ways

**Multi-Chip System:** 4  NUMA zones, 64 SMs (16 per zone); 32 memory channels;
L2: 4MB, 16 ways; Inter-GPM Interconnect: 192 GB/s;
DRAM Bandwidth: 768 GB/s (192 GB/s per module)

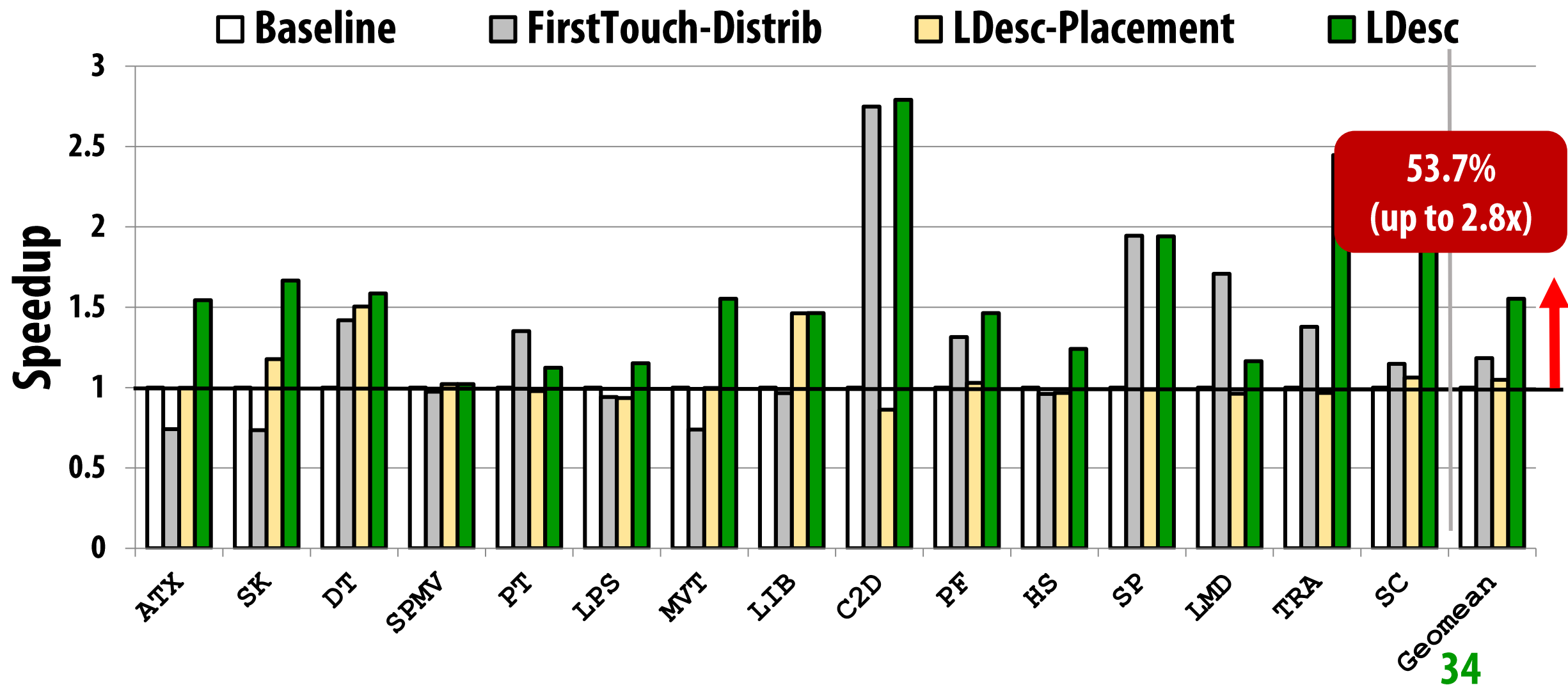# Locality descriptors are an effective means to leverage cache locality

1.5

26.6%

# Different locality types require different optimizations
# A single optimization is often insufficient

| SK | DT | HS | D2D | C2D | SPMV | LIB | LMD | Geomean |
|----|----|----|-----|-----|------|-----|-----|---------|
| INTER-THREAD (COACCESSED) | | | INTER-THREAD (NEARBY) | | | INTRA-THREAD | | |

# Performance Impact: Leveraging NUMA Locality

□ Baseline     ▨ FirstTouch-Distrib     ▨ LDesc-Placement     ▨ LDesc

**Speedup**

53.7%
(up to 2.8x)

ATX   SK   DT   SPMV   PT   LPS   MVT   LIB   C2D   PF   HS   SP   LMD   TRA   SC   Geomean

34

# Conclusion

**Problem:**

GPU programming models have no explicit abstraction to express data locality

Leveraging data locality is a challenging task, as a result

**Our Proposal: The Locality Descriptor**

A HW-SW abstraction to explicitly express data locality

A architecture-agnostic and flexible SW interface to express data locality

Enables HW to leverage key program semantics to optimize locality

**Key Results:**

26.6% (up to 46.6%) performance speed up from leveraging <u>cache locality</u>

53.7% (up to 2.8x) performance speed up from leveraging <u>NUMA locality</u>

# The Locality Descriptor

## A Holistic Cross-Layer Abstraction to Express Data Locality in GPUs

**Nandita Vijaykumar**

**Eiman Ebrahimi, Kevin Hsieh, Phillip B. Gibbons, Onur Mutlu**

Carnegie Mellon University

NVIDIA

ETH Zürich