# GRIT: Enhancing Multi-GPU Performance with Fine-Grained Dynamic Page Placement

Yueqi Wang*[1], **Bingyao Li*[1]**, Aamer Jaleel[2], Jun Yang[1], Xulong Tang[1]
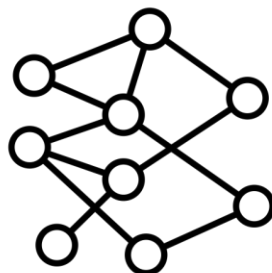
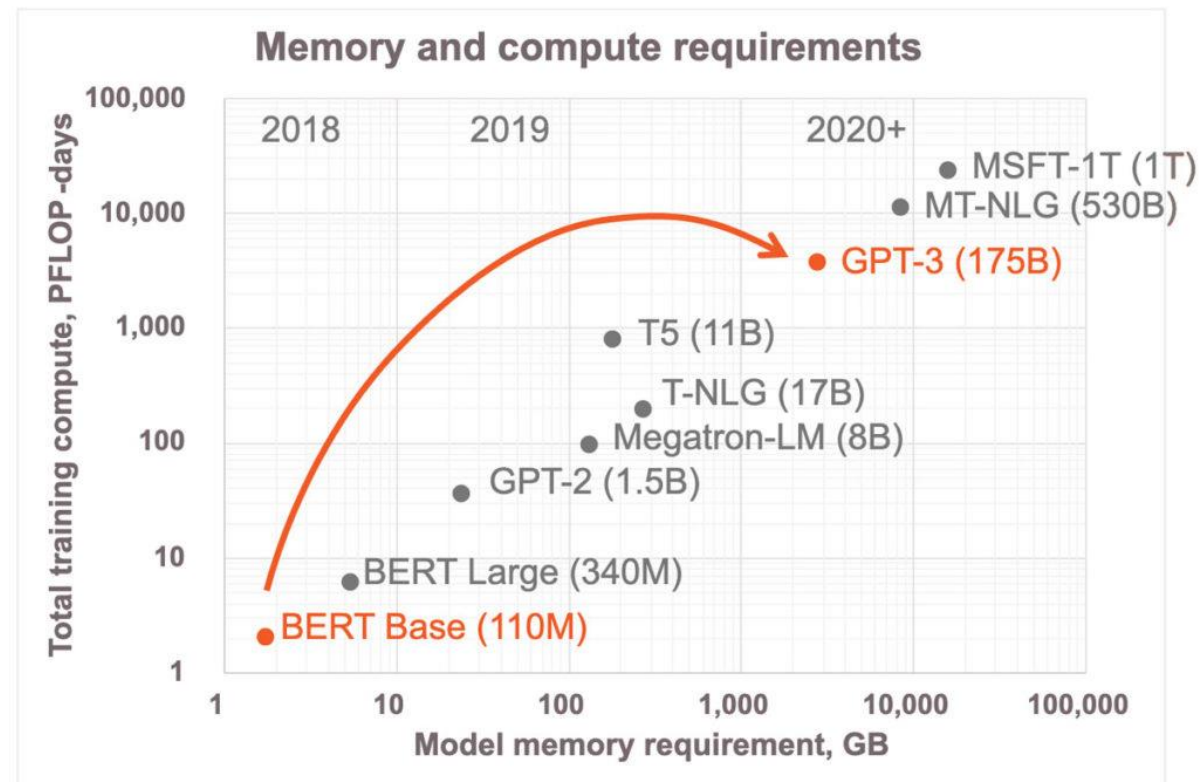[1]University of Pittsburgh, [2]NVIDIA

# Multi-GPU is Popular



Graph Processing

DNN

ChatGPT

Datacenter Workloads

## Memory and compute requirements

2018    2019    2020+

● MSFT-1T (1T)
● MT-NLG (530B)
● GPT-3 (175B)

● T5 (11B)
● T-NLG (17B)
Megatron-LM (8B)
● GPT-2 (1.5B)

BERT Large (340M)
● BERT Base (110M)

Total training compute, PFLOP-days — Model memory requirement, GB

**Ever-growing application complexity and input dataset sizes.**

# Multi-GPU is Popular



**Multi-GPU is here ！ （NVIDIA DGX, Intel Xe)**

**Ever-growing application complexity and input dataset sizes.**

[1] https://www.cerebras.net/blog/harnessing-the-power-of-sparsity-for-large-gpt-ai-models

# UVM for Multi-GPU

- **Growing trend of multi-GPUs leveraging Unified Virtual Memory (UVM)**

# Multi-GPU Scalability



**Why performance gap:**
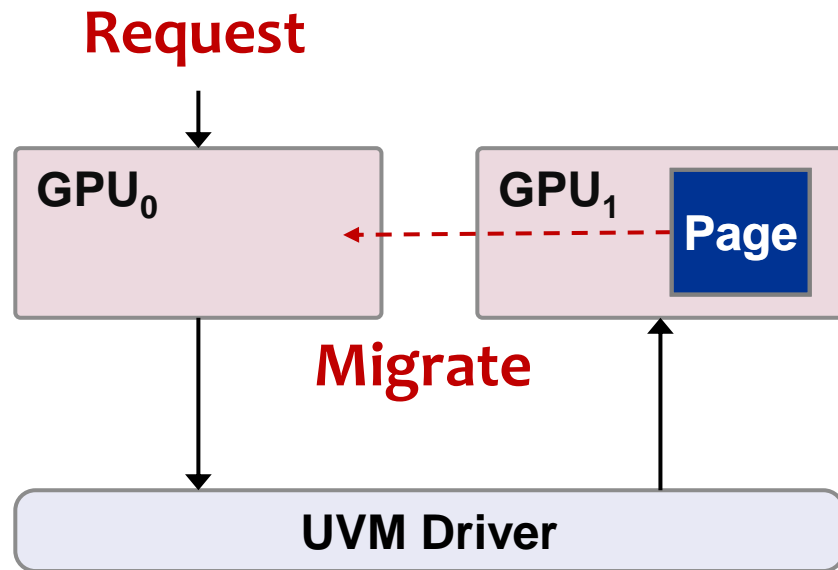
- NUMA data access

- Data transfer
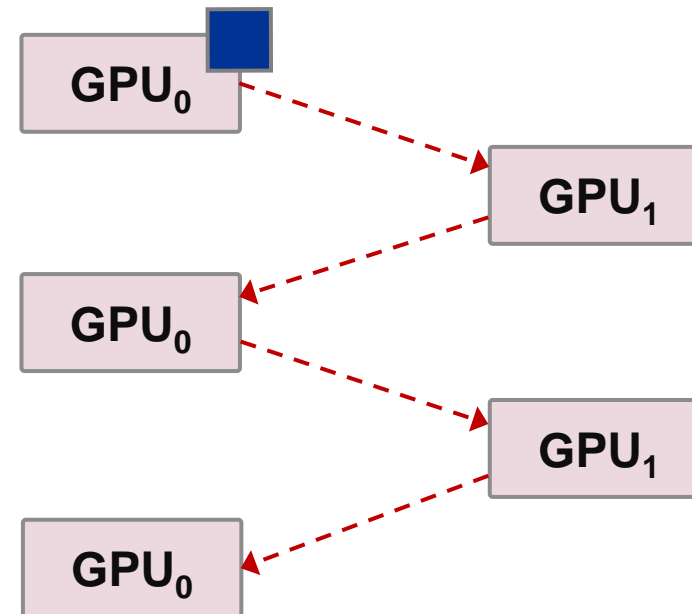
- Address translation

……

[2] https://developer.nvidia.com/blog/easy-multi-gpu-deep-learning-digits-2/

# Multi-GPU Page Placement Schemes

## 1. On-Touch Migration: Request, Migrate.
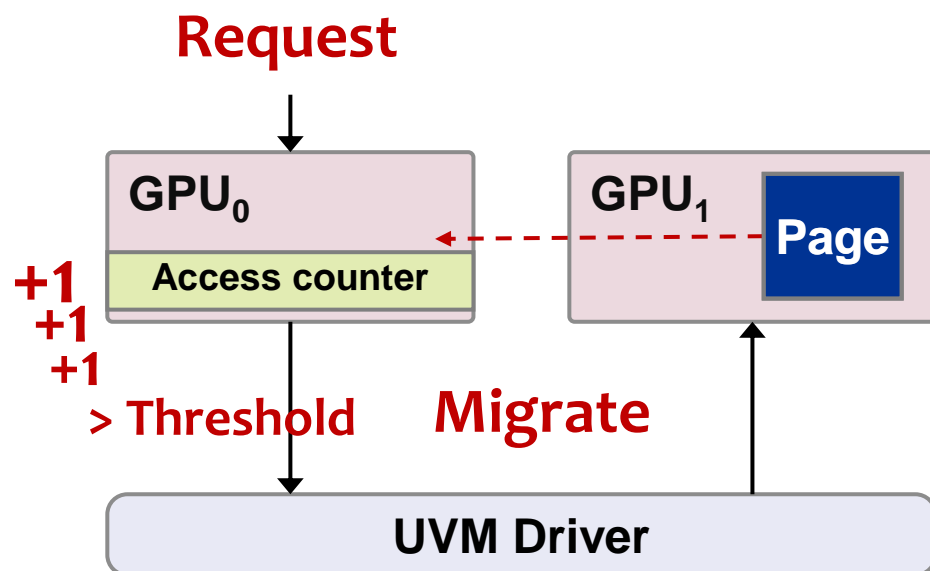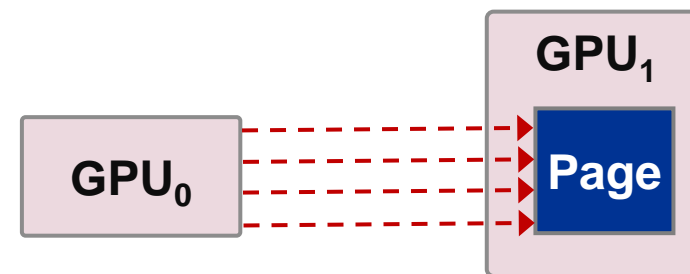


😇 **Every memory access is local**

😣 **Ping-Pong effect**

# Multi-GPU Page Placement Schemes

## 2. Access Counter-based Migration: Reaching threshold, Migrate.



😇 Reduce Ping-Pong migration

😟 Substantial remote access

# Multi-GPU Page Placement Schemes

## 3. Page Duplication: Read, Duplicate.



😇 **Reduce migration & remote access**      😠 **Significant collapsing overhead**

# Performance of Page Placement Schemes



No "one-size-fits-all" page placement scheme

Chart axes: Normalized Performance (y-axis: 0.0, 0.5, 2.0, 2.5); benchmarks (x-axis): BFS, BS, C2D, FIR, GEMM, MM, SC, ST

Legend: ☐ On-touch ▨ Access counter ▨ Duplication

# Page Access Characteristics

➤ **The page-sharing / read-write patterns vary within the same application**



GPU 1    GPU2    GPU3    GPU4

Write    Read

100%

0%    0    5    10    15    20    25    30
**Time**

0    5    10    15    20    25    30
**Time**

**Private / Share**

**Read / Write**

A dynamic page placement scheme that can accommodate variations in page access characteristics

# Problem Summary

*Problem:*

Delivered performance is constrained by **NUMA overhead**

No "**one-size-fits-all**" page placement scheme

*Goal:*

Effectively reduce **NUMA overhead in multi-GPU** by **determining** page placement scheme **at runtime**

University of Pittsburgh

# GRIT (Fine-**GR**ained dynam**I**c page placemen**T**)

Scheme change
**metric**

**Dynamically determine page placement scheme at runtime**

# GRIT – Scheme Change Metric

**Request**    **Write**

**Current scheme is unsuitable**

**GPU$_0$**

~~Page~~

**Page duplication should be avoided**

↑ **local page fault**

**page protection fault** ↑

**UVM Driver**

**Indicator: Number of page faults (local page fault & page protection fault)**

# GRIT: Dynamic Page Placement Scheme

Scheme change metric

How to track information ?

**Dynamically determine page placement scheme at runtime**

**Fault-Aware Initiator**

# GRIT – Page Attribute Table and Cache (PA-Table & Cache)

☹ **Additional memory access**      😇 **Facilitate lookup**

**Write back**

**CPU Memory**

| PA-Table | | |
|---|---|---|
| **VPN** | **Fault Counter** | **Read/Write** |
| 0xA00 | 10 | 0 |
| 0xA01 | 01 | 1 |
| … | … | … |

**PA-Cache**

**Way 0**      …      **Way 3**

| VPT | FC | R/W |
|---|---|---|
| … | … | … |
| … | … | … |
| … | … | … |
| … | … | … |

| VPT | FC | R/W |
|---|---|---|
| … | … | … |
| … | … | … |
| … | … | … |
| … | … | … |

**Write allocate**

University of Pittsburgh

# GRIT: Dynamic Page Placement Scheme

**Scheme change metric**

**How to track information ?**

**Which scheme to change to ?**

## Dynamically determine page placement scheme at runtime

**Fault-Aware Initiator**

**PA-Table & PA-Cache**

# GRIT – Which Scheme

| 63 | 62:54 | 53:52 | 51:12 | 11 | 10:9 |
|----|-------|-------|-------|----|----|
| X D | Unused Bits (UB) | | 4 KB Page Frame Number (PFN) | U B | Scheme Bits |

**On-touch**

**Get access information from PA-Table /PA-Cache**

**FC Reach threshold?**

**True**

**Change Scheme**

**Scheme information is stored in host PTE.**

**All read?**

**True** → **Duplication**

**False** → **Access-counter**

University of Pittsburgh

# GRIT: Dynamic Page Placement Scheme

Scheme change metric

How to track information ?

Which scheme to change to ?

**Dynamically determine page placement scheme at runtime**

Fault-Aware Initiator

PA-Table & PA-Cache

Scheme Decision Mechanism

University of Pittsburgh

# GRIT: Dynamic Page Placement Scheme



**Performance gap due to improper scheme before trigger scheme change**

Chart axis labels (y-axis): Normalized Performance — 0.0, 2.0, 2.5, 3.0

x-axis categories: BFS, BS, C2D, FIR, GEMM, MM, SC, ST

Legend: ☐ On-touch   ☐ Access counter   ☐ Duplication   ☐ GRIT-Dynamic
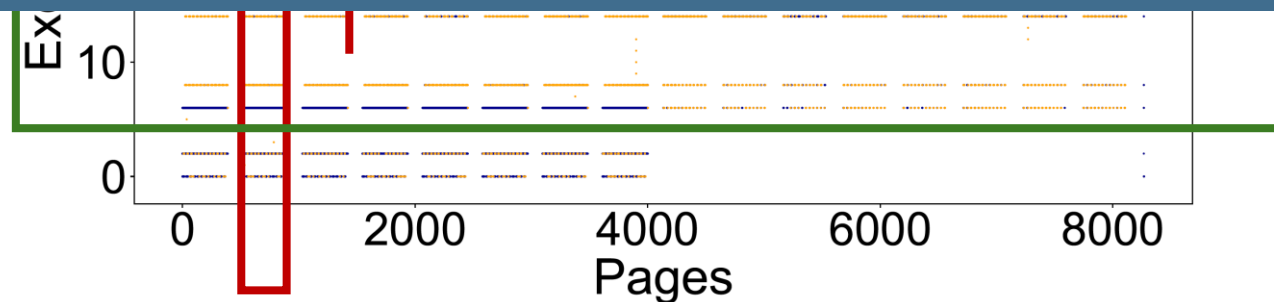
University of Pittsburgh
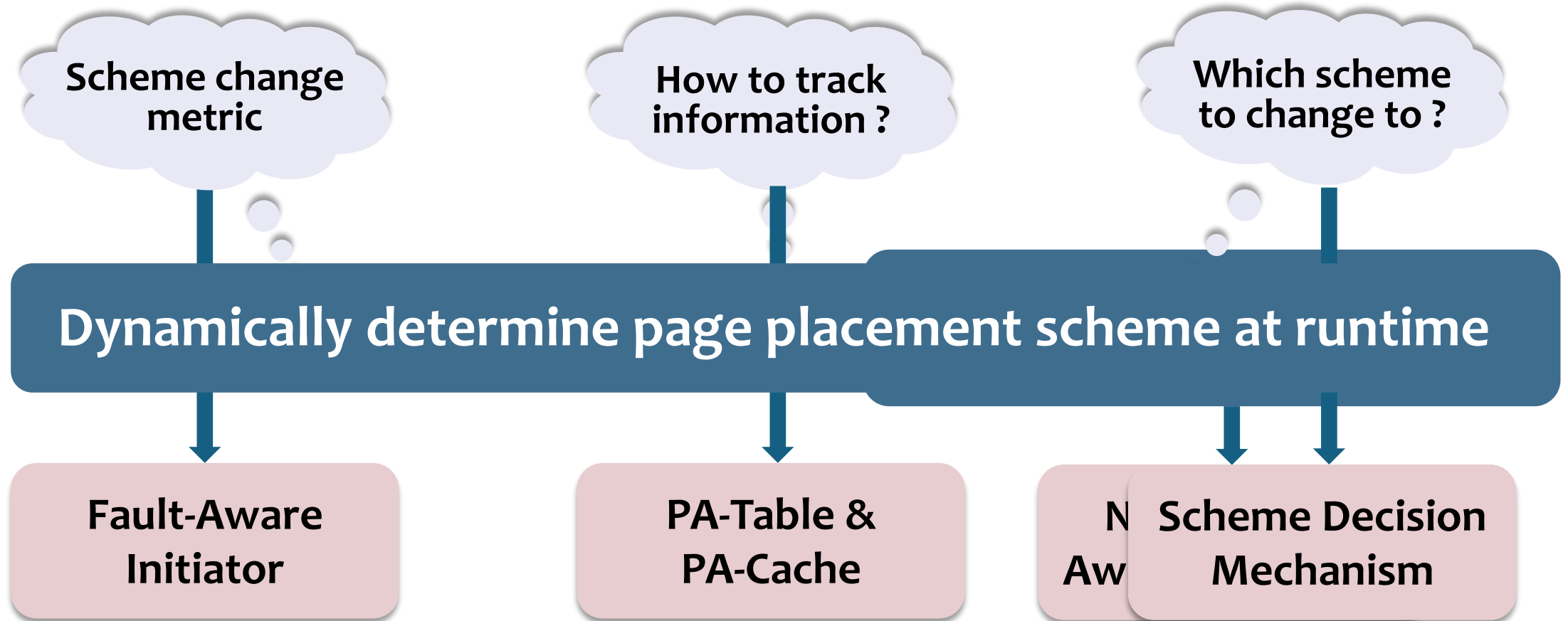
19

# Page Attributes Characterization

> **The neighboring pages tend to exhibit similar access attributes**



**Proactively determine page placement scheme for neighboring pages**

# GRIT: Neighboring-Aware Prediction

Scheme change metric

How to track information ?

Which scheme to change to ?

**Dynamically determine page placement scheme at runtime**

Fault-Aware Initiator

PA-Table & PA-Cache

Scheme Decision Mechanism

# GRIT: Neighboring-Aware Prediction

*Page Table*

| VPN | Scheme |
|-----|--------|
| 0xF000 | 00 |
| 0xF001 | 00 |
| 0xF002 | 00 |
| 0xF003 | 00 |
| 0xF004 | 01 |
| 0xF005 | 01 |
| 0xF006 | 01 |
| 0xF007 | 10 |

Promote

Recursively Promote

32KB page group
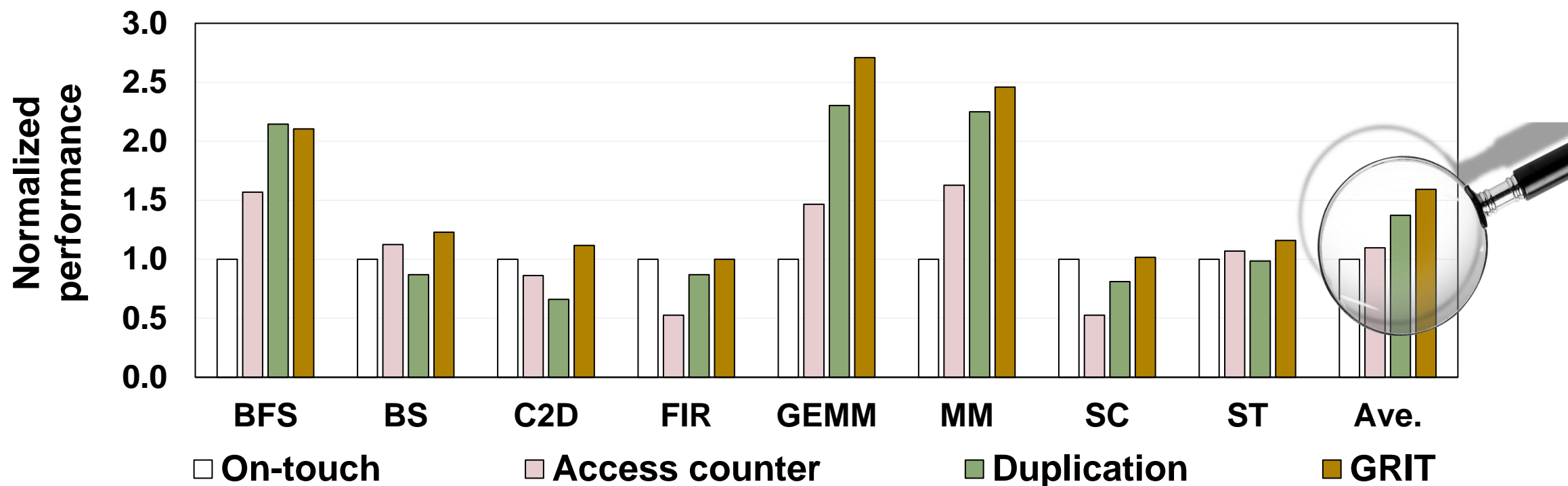
256KB page group

# GRIT – Put All Together

# Methodology

- Simulator:   MGPUSim [ISCA 19']

- Workloads:   8 applications from Hetero-Mark, AMDAPPSDK , SHOC, and DNN Mark benchmark suites, including random, adjacent, and scatter-gather access patterns.
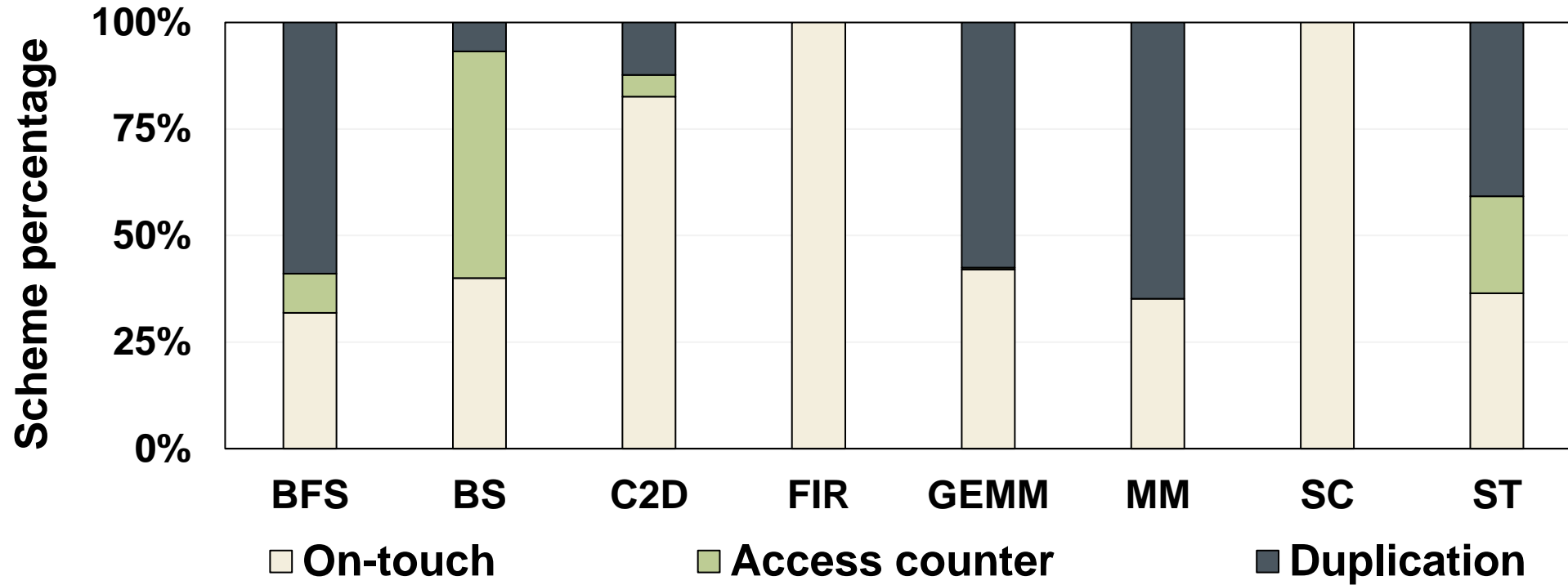
**Detailed page placement scheme modeling in paper**

# Evaluation – Overall Performance



GRIT achieves **60%, 49%, and 29%** performance improvement compared to uniformly employing on-touch, access counter-based, and page duplication scheme.

University of Pittsburgh

# Evaluation – Scheme Breakdown



GRIT is able to distinguish page attributes and consistently select the most suitable scheme accordingly.

# Summary

**Problem:** NUMA overheads in multi-GPU systems

- *No "**one-size-fits-all**" page placement scheme*

**GRIT:**

A. Dynamic page placement scheme **determines schemes in a fine-grained manner**

B. Neighboring-aware prediction **proactively determines adjacent page scheme**

Improves **performance** by **60%** on average.

University of Pittsburgh

# Thanks! Q&A

GRIT: Enhancing Multi-GPU Performance with Fine-Grained Dynamic Page Placement

Yueqi Wang*[1], **Bingyao Li*[1]**, Aamer Jaleel[2], Jun Yang[1], Xulong Tang[1]

[1]University of Pittsburgh, [2]NVIDIA