# Mosaic: A GPU Memory Manager with Application-Transparent Support for Multiple Page Sizes
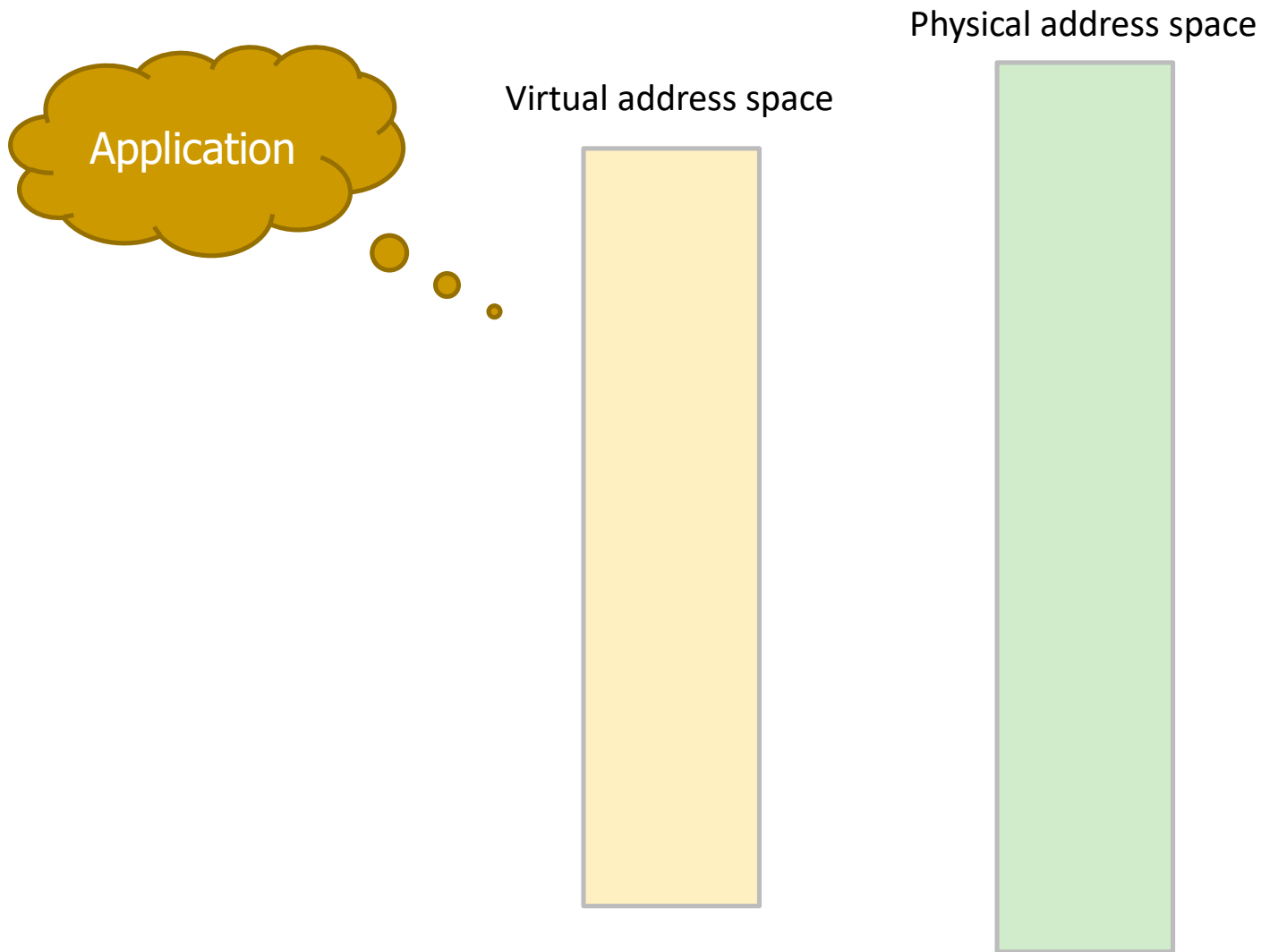
Rachata Ausavarungnirun, Joshua Landgraf, Vance Miller, Saugata Ghose,
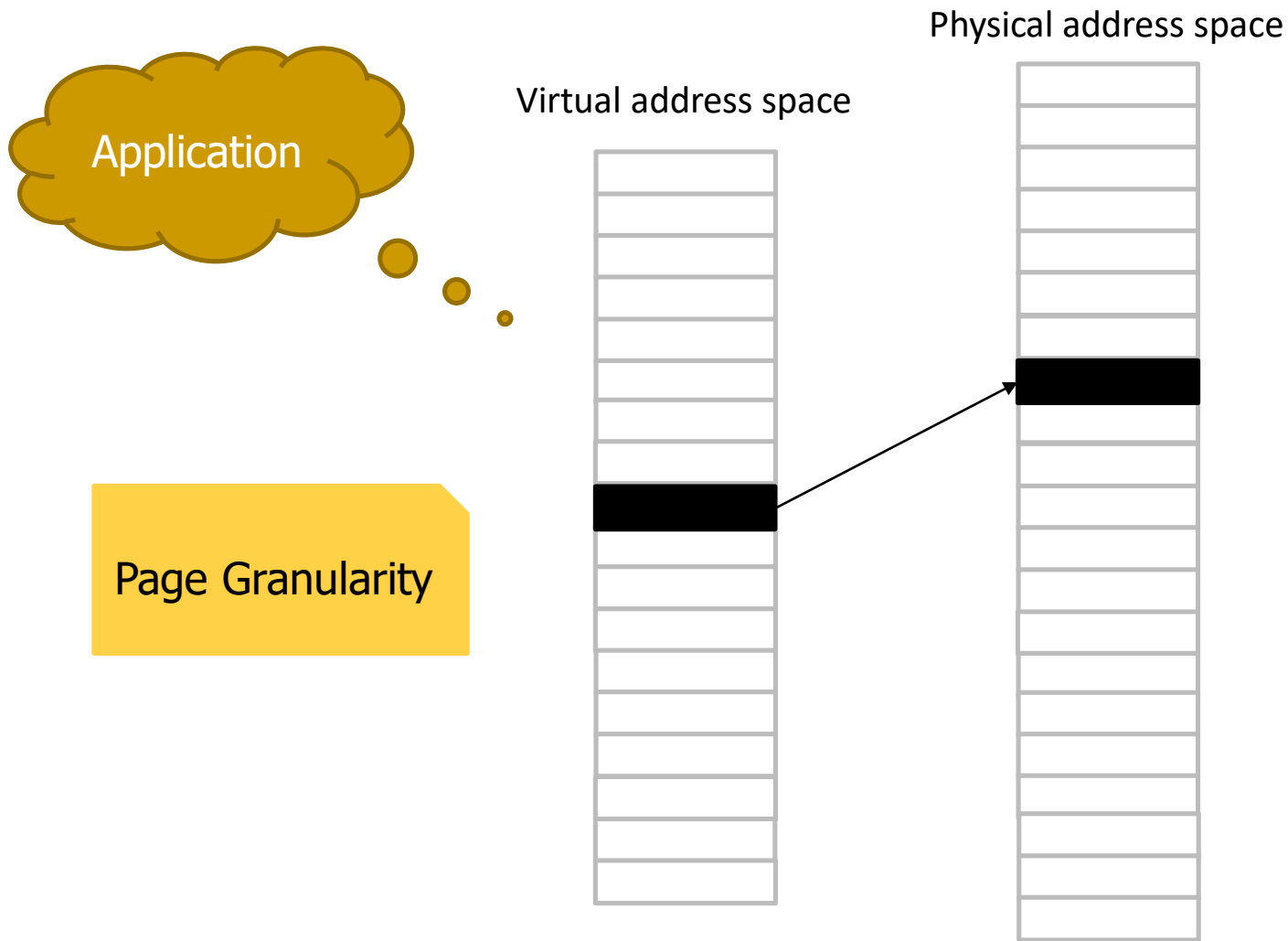Jayneel Gandhi, Christopher J. Rossbach, Onur Mutlu

MICRO 2017

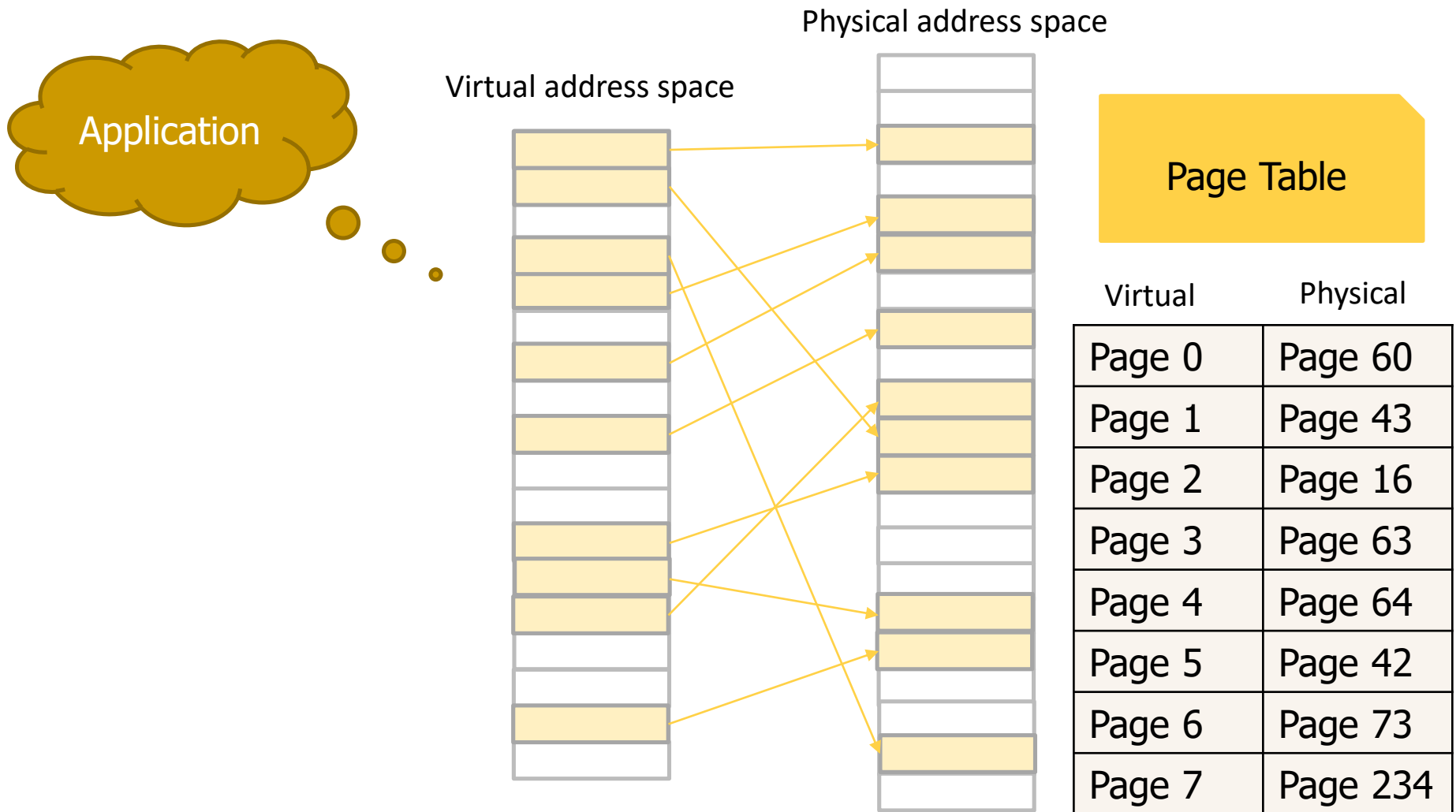Presenter: Christina Giannoula

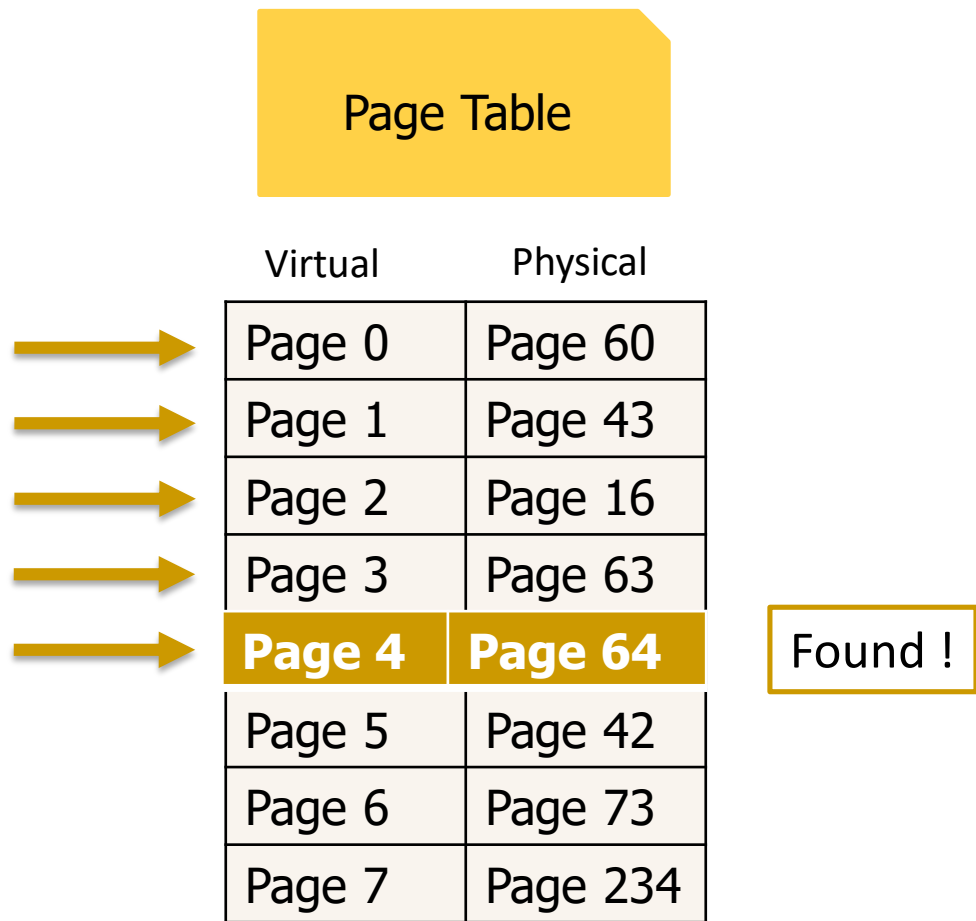# Background and Problem

# Address Translation



Physical address space

Virtual address space

Application

**SAFARI**

# Address Translation

Application

Page Granularity

Virtual address space

Physical address space

SAFARI

# Address Translation

Application

Virtual address space

Physical address space

Page Table

| Virtual | Physical |
|---------|----------|
| Page 0 | Page 60 |
| Page 1 | Page 43 |
| Page 2 | Page 16 |
| Page 3 | Page 63 |
| Page 4 | Page 64 |
| Page 5 | Page 42 |
| Page 6 | Page 73 |
| Page 7 | Page 234 |

# Page Table Walk

- Look up a mapping
- Page Table Walk:
  Ten to hundreds of cycles

Page Table

| Virtual | Physical |
|---|---|
| Page 0 | Page 60 |
| Page 1 | Page 43 |
| Page 2 | Page 16 |
| Page 3 | Page 63 |
| **Page 4** | **Page 64** |
| Page 5 | Page 42 |
| Page 6 | Page 73 |
| Page 7 | Page 234 |

Found !

SAFARI

# Page Table Walk

➤ Look up a mapping

Page Table

| Virtual | Physical |
|---------|----------|
| Page 0 | Page 60 |
| Page 1 | Page 43 |
| Page 2 | Page 16 |
| Page 3 | Page 63 |
| **Page 4** | **Page 64** |
| Page 5 | Page 42 |
| Page 6 | Page 73 |
| Page 7 | Page 234 |

Found !

Page Table Walks:
High Latency

SAFARI

# Translation Lookaside Buffers (TLB)

➢ **TLBs:**
  ➢ Store **recently** used address translations
  ➢ Address translation **cache**

Level 1
Level 2
...

Limited Size

### Page Table

| Virtual | Physical |
|---------|----------|
| Page 0 | Page 60 |
| Page 1 | Page 43 |
| Page 2 | Page 16 |
| Page 3 | Page 63 |
| Page 4 | Page 64 |
| Page 5 | Page 42 |
| Page 6 | Page 73 |
| Page 7 | Page 234 |

### TLB

| Virtual | Physical |
|---------|----------|
| Page 1 | Page 60 |
| Page 3 | Page 63 |
| Page 4 | Page 64 |
| Page 5 | Page 42 |

# State-of-the-art Virtual Memory on GPUs

| GPU core | GPU core | GPU core | GPU core |
|----------|----------|----------|----------|

**SAFARI**

# State-of-the-art Virtual Memory on GPUs

| GPU core | GPU core | GPU core | GPU core | |
|:---:|:---:|:---:|:---:|:---|
| Private TLB | Private TLB | Private TLB | Private TLB | Private |

TLB miss!

**SAFARI**

# State-of-the-art Virtual Memory on GPUs

| GPU core | GPU core | GPU core | GPU core |
|----------|----------|----------|----------|
| Private TLB | Private TLB | Private TLB | Private TLB |

Private

- - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - -

Shared

Shared TLB

TLB miss!

**SAFARI**

# State-of-the-art Virtual Memory on GPUs

| GPU core | GPU core | GPU core | GPU core |
|---|---|---|---|
| Private TLB | Private TLB | Private TLB | Private TLB |

Private

- - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - -

Shared

**Shared TLB**

High
Latency

**Page Table Walkers**

**Page Table (GPU Memory)**

**Data (GPU Memory)**

# State-of-the-art Virtual Memory on GPUs

SAFARI

# State-of-the-art Virtual Memory on GPUs

| GPU core | GPU core | GPU core | GPU core |
|----------|----------|----------|----------|
| Private TLB | Private TLB | Private TLB | Private TLB |

Private

Shared

Shared TLB

GPU Threads

Page Table Walkers

GPU-side memory

CPU-side memory

Page Table
(GPU Memory)

I/O bus

Data
(GPU Memory)

CPU Memory

**SAFARI**

# Address Translation Challenge

Small Pages

Physical address space

Virtual address space

SAFARI

# Address Translation Challenge

Small Pages
(4KB)

Large Pages
(2MB)

Physical address space

Virtual address space

Physical address space

Virtual address space

**SAFARI**

# Address Translation Challenge



Small Pages

Large Pages

Physical address space

Virtual address space

TLB

**Fixed Size**

**3** entries

Physical address space

Virtual address space

# Address Translation Challenge

Small Pages

Large Pages

Physical address space

Virtual address space

TLB

Fixed Size

3 entries

Physical address space

Virtual address space

SAFARI

# Address Translation Challenge

Small Pages

Large Pages

Physical address space

Virtual address space

TLB

Limited TLB reach

Fixed Size

3 entries

Physical address space

Virtual address space

# Address Translation Challenge

Small Pages

Large Pages

Physical address space

Virtual address space

Physical address space

Virtual address space

TLB

Fixed Size

3 entries

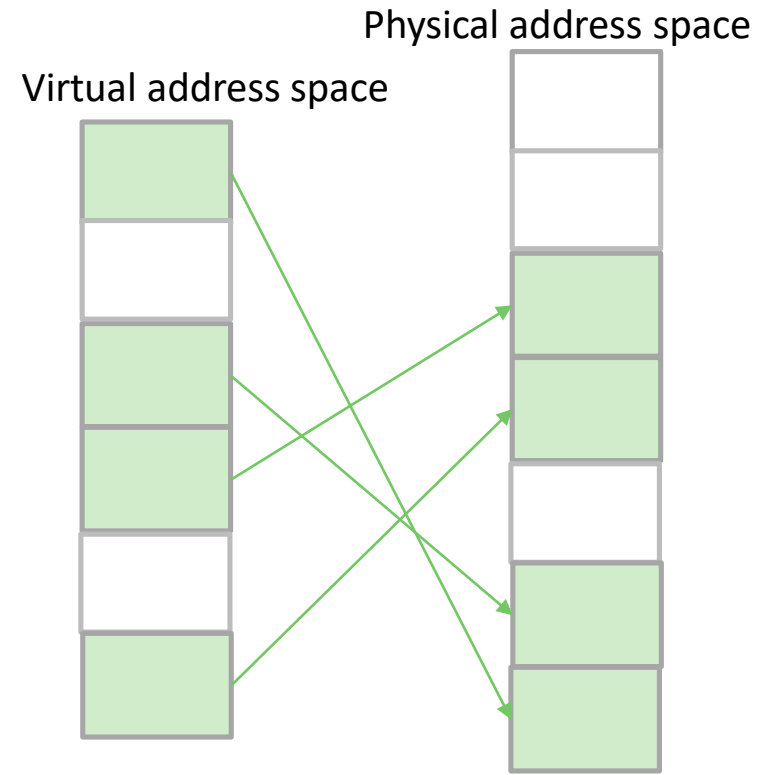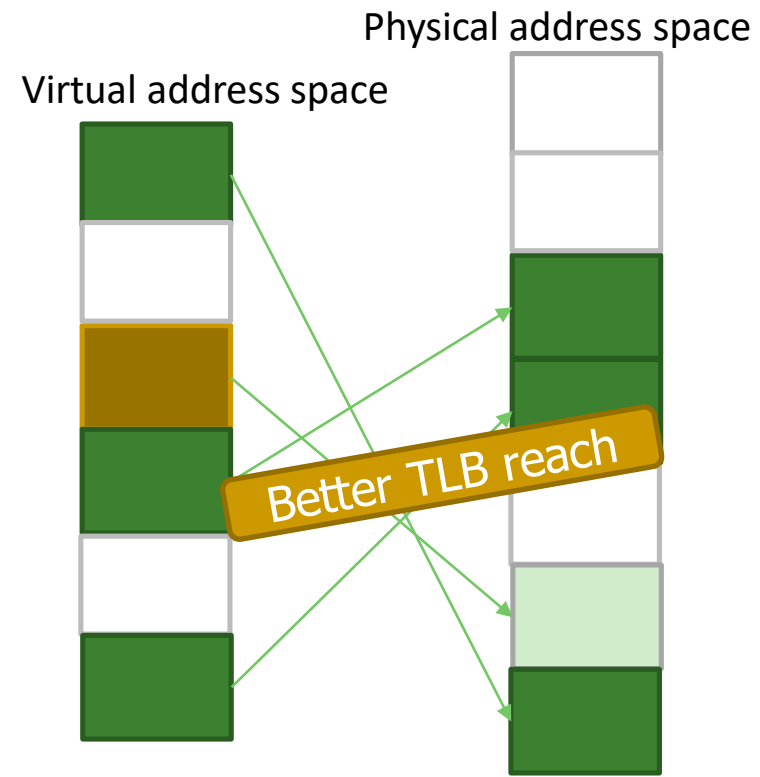SAFARI

# Address Translation Challenge



Small Pages

Large Pages

Physical address space

Virtual address space

TLB

Fixed Size

3 entries

Virtual address space

Physical address space

Better TLB reach

**SAFARI**

# Address Translation Challenge

Small Pages

Large Pages

Physical address space

Virtual address space

TLB

Fixed Size

3 entries

Virtual address space

Physical address space

**Large** page size is better!

Better TLB reach
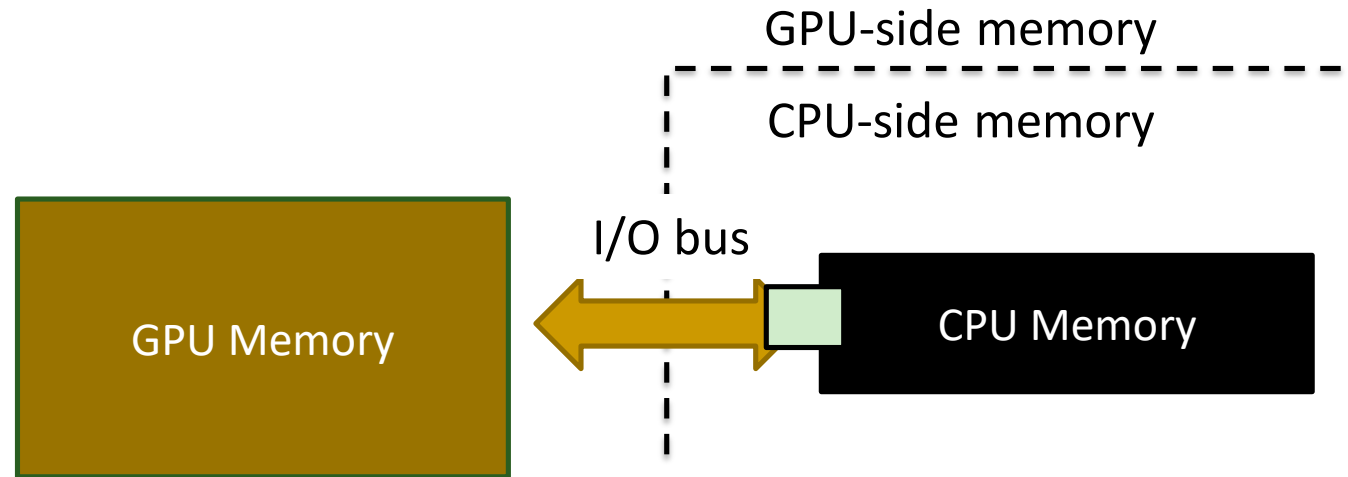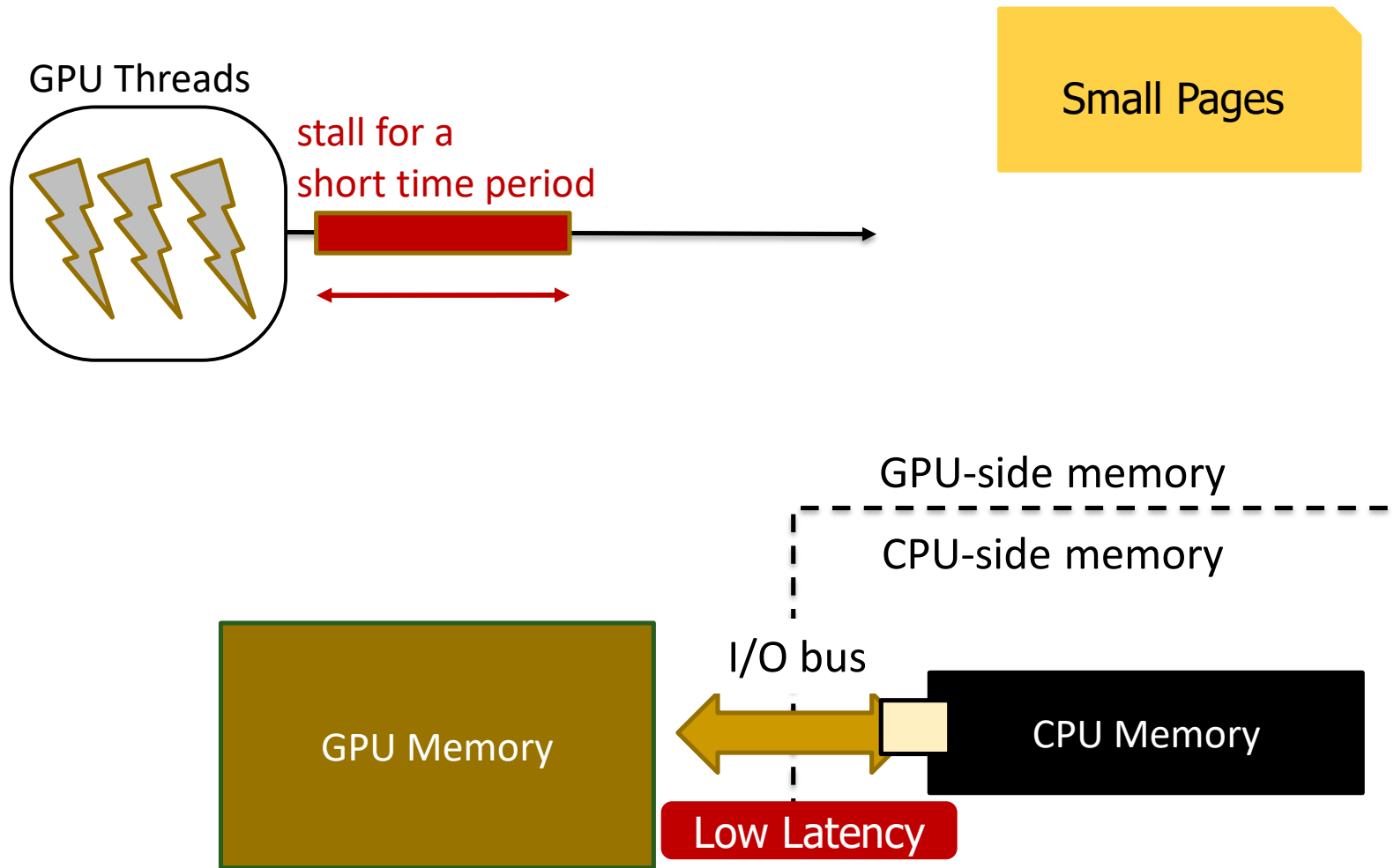
SAFARI

# Demand Paging Challenge

➢ An application requests data that **is not** currently resident in GPU memory

➢ A **Page Fault** is triggered

➢ Transfer data in **page-granularity**

GPU-side memory

CPU-side memory

I/O bus

GPU Memory

CPU Memory

# Demand Paging Challenge

➢ An application requests data that **is not** currently resident in GPU memory

GPU Threads

Small Pages

stall for a
short time period

GPU-side memory

CPU-side memory

I/O bus

GPU Memory

CPU Memory

Low Latency

# Demand Paging Challenge

➢ An application requests data that **is not** currently resident in GPU memory

GPU Threads

stall for a
long time period

Large Pages

GPU-side memory

CPU-side memory

I/O bus

GPU Memory

large emory

High Latency

# Demand Paging Challenge

➤ An application requests data that **is not** currently resident in GPU memory

GPU Threads

Large Pages

stall ...

**Small** page size is better!

GPU-side memory

CPU-side memory

I/O bus

GPU Memory

large emory

High Latency

# Page Size Trade-Off

| Small Pages | vs | Large Pages |
|---|---|---|

Small Pages:
- Low TLB reach
- Low demand paging latency

Large Pages:
- High TLB reach
- High demand paging latency
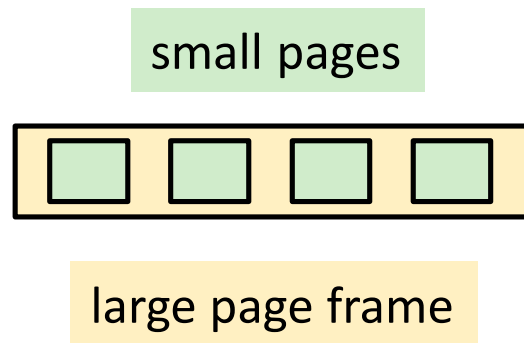
> Can we get the best of both page sizes ?

SAFARI

# Executive Summary

# Executive Summary

- Problem
  - No single best page size for GPU virtual memory (large vs small pages)

- Goal
  - Transparently and efficiently enable <u>both</u> page sizes

# Executive Summary

- **Problem**
  - No single best page size for GPU virtual memory (large vs small pages)

- **Goal**
  - Transparently and efficiently enable both page sizes

- **Key Observation**
  - Can easily coalesce an application's **contiguously-allocated** small pages into a large page
  - GPGPU applications typically allocate large chunks of memory *at once*

small pages

large page frame

**SAFARI**

# Executive Summary

- **Problem**
  - No <span style="color:red">single best page size</span> for GPU virtual memory (large vs small pages)

- **Goal**
  - Transparently and efficiently enable <u>both</u> page sizes

- **Key Observation**
  - Can easily coalesce an application's **contiguously-allocated** small pages into a large page
  - GPGPU applications typically allocate large chunks of memory *at once*
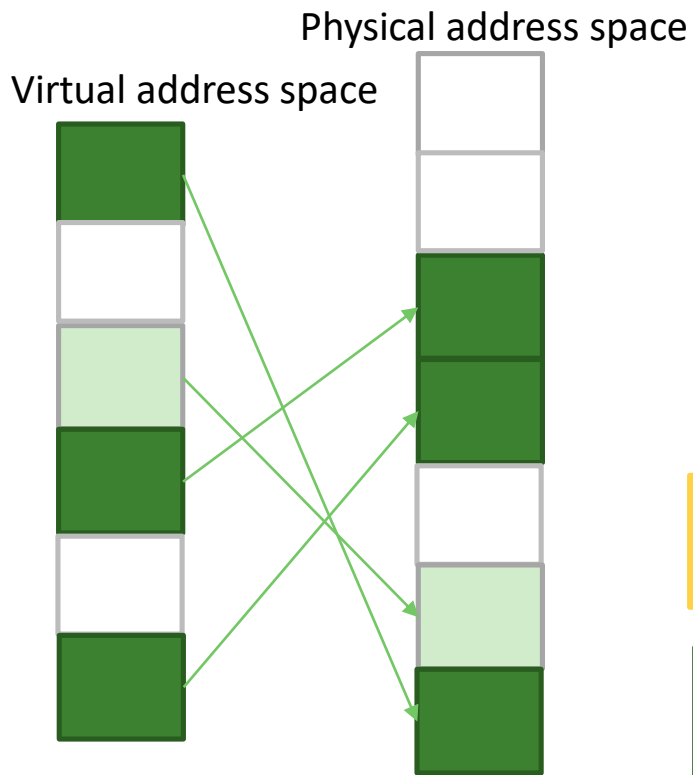
- **Key Idea**
  - *Preserve the virtual address contiguity* of small pages when allocating physical memory to simplify coalescing

- **Mosaic:**
  - A hardware/software cooperative framework
  - Enables the benefits of both small and large pages

- **Key Result: 55% on average performance improvement over state-of-the-art GPU memory management mechanism**

*SAFARI*

# Key Ideas and Challenges

# Key Ideas
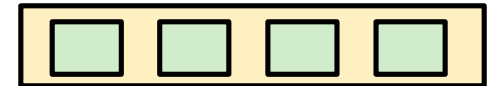
Large Pages

Physical address space

Virtual address space



- Translate using **large page size**
  High TLB reach

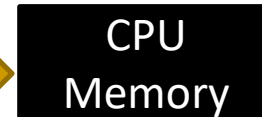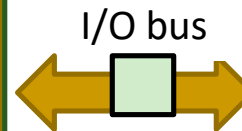- Transfer using **small page size**
  Low demand paging latency

Contiguity

small pages



large page frame

Small Pages

GPU Memory ⟷ I/O bus ⟷ CPU Memory

# Challenges with Multiple Page Sizes

Time

State-of-the-art

App1 Allocation

App2 Allocation

App1 Allocation

App2 Allocation

Coalesce App2 Pages

Large page frame 1 ✓
Large page frame 2
Large page frame 3
Large page frame 4
Large page frame 5 ✓

**Need to search which pages to coalesce**

App1
App2
Unallocated

# Challenges with Multiple Page Sizes

Time

## State-of-the-art

App1
Allocation

App2
Allocation

App1
Allocation

App2
Allocation

Coalesce
App2
Pages

Large page frame 1

Large page frame 2

Large page frame 3

Large page frame 4

Large page frame 5

App1

App2

Unallocated

SAFARI

# Challenges with Multiple Page Sizes

Time

State-of-the-art

App1
Allocation

App2
Allocation

App1
Allocation

App2
Allocation

Coalesce
App2
Pages

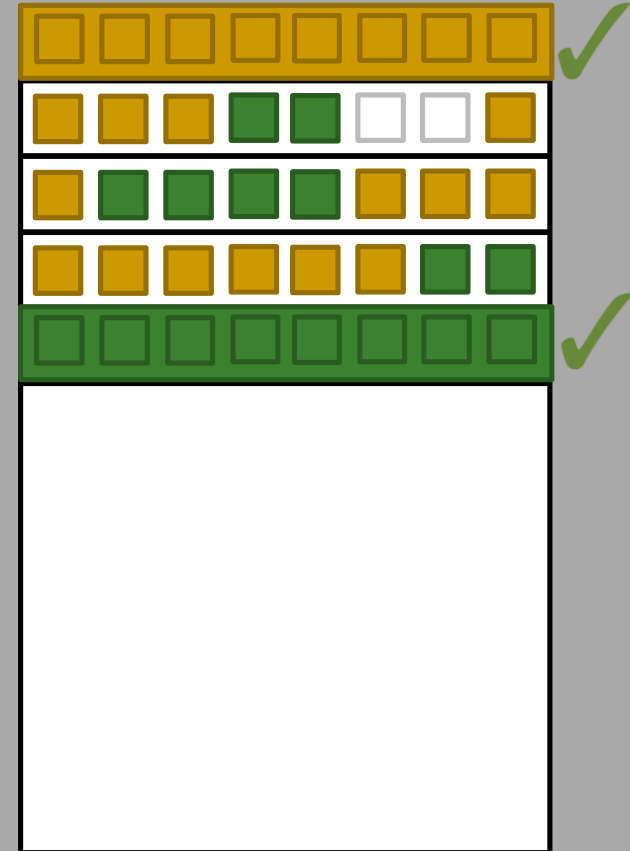Large page frame 1 ✓

Large page frame 2

Large page frame 3 ✓

Large page frame 4 ✓

Large page frame 5 ✓

**Cannot coalesce without migrating multiple pages**

App1

App2

Unallocated

# Why page migration is bad



GPU core — Private TLB (×4)

Private / Shared

Shared TLB

Change Translation

Page Table Walkers

Page Table (GPU Memory)

Data (GPU Memory)

SAFARI

# Why page migration is bad

GPU core

GPU core

GPU core

GPU core

Private TLB

Private TLB

Private TLB

Private TLB

Private

Shared

Shared TLB

TLB flush !

Page Table
Walkers

Page Table
(GPU Memory)

Data
(GPU Memory)

SAFARI

# Why page migration is bad

GPU core
GPU core
GPU core
GPU core

Priva **✗** TLB
Priva **✗** TLB
Priva **✗** TLB
Priva **✗** TLB

Private

- - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - -

Shared

Sha **✗** TLB

TLB flush !

GPU Threads

Stall …

Page Table Walkers

Page Table
(GPU Memory)

Data
(GPU Memory)

SAFARI

# Desirable Allocation

Time

**Desirable Allocation**

App1 Allocation

App2 Allocation

App1 Allocation

App2 Allocation

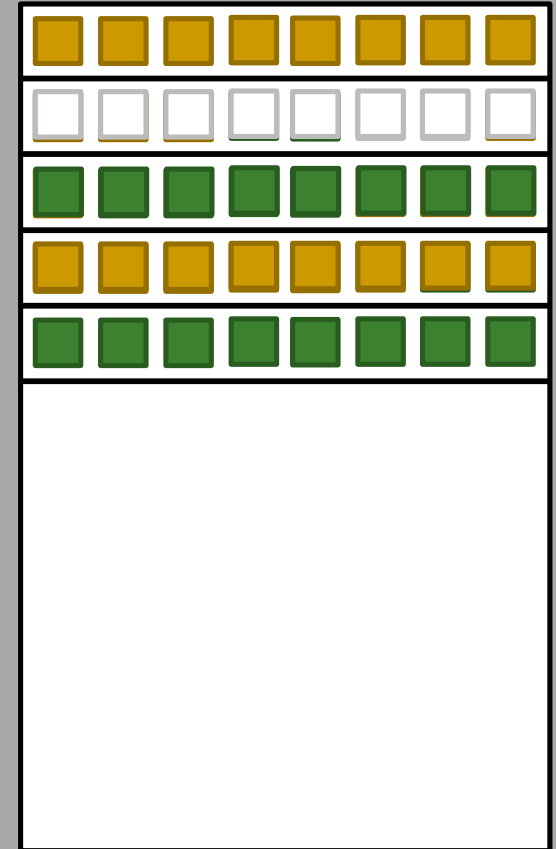Coalesce App2 Pages

Large page frame 1 ✓

Large page frame 2

Large page frame 3 ✓
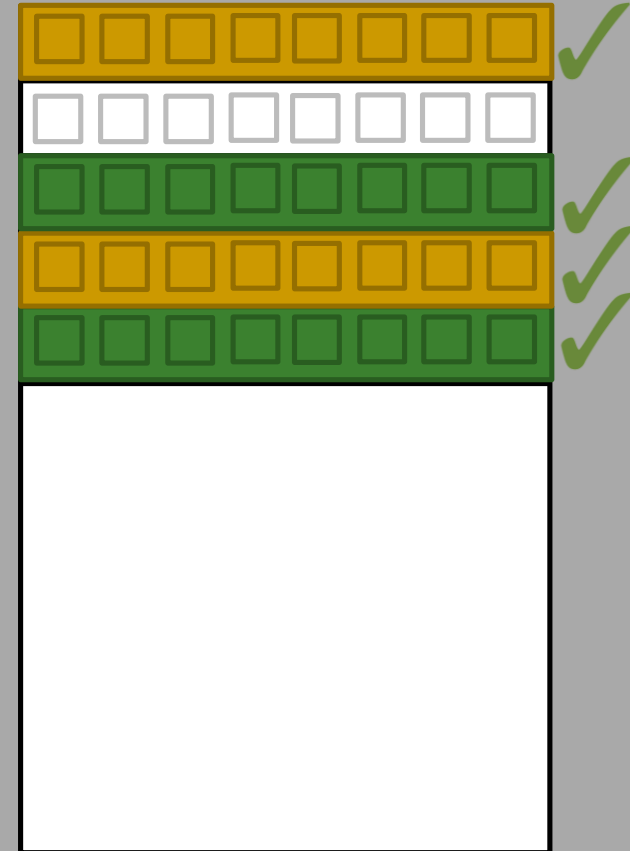
Large page frame 4 ✓

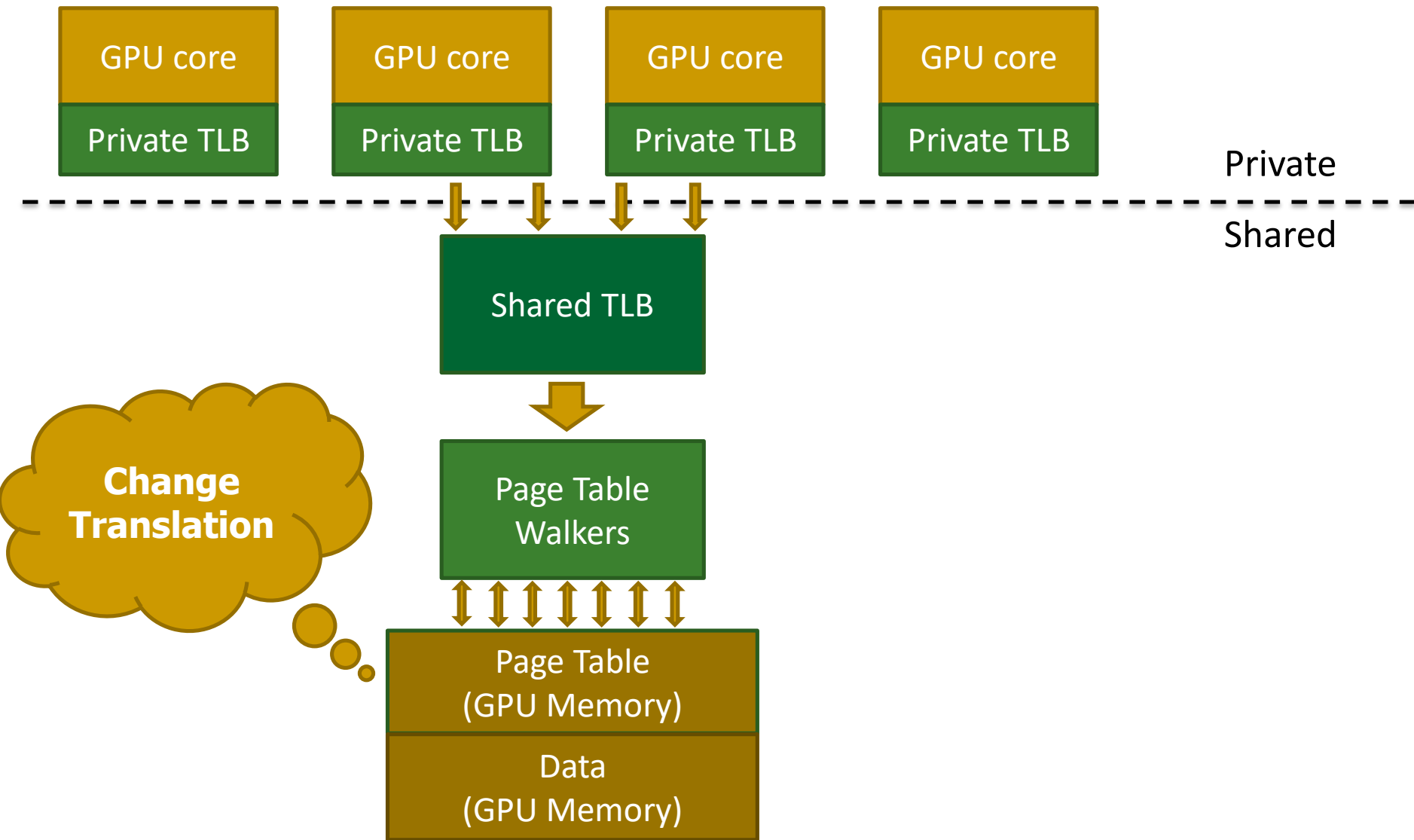Large page frame 5 ✓

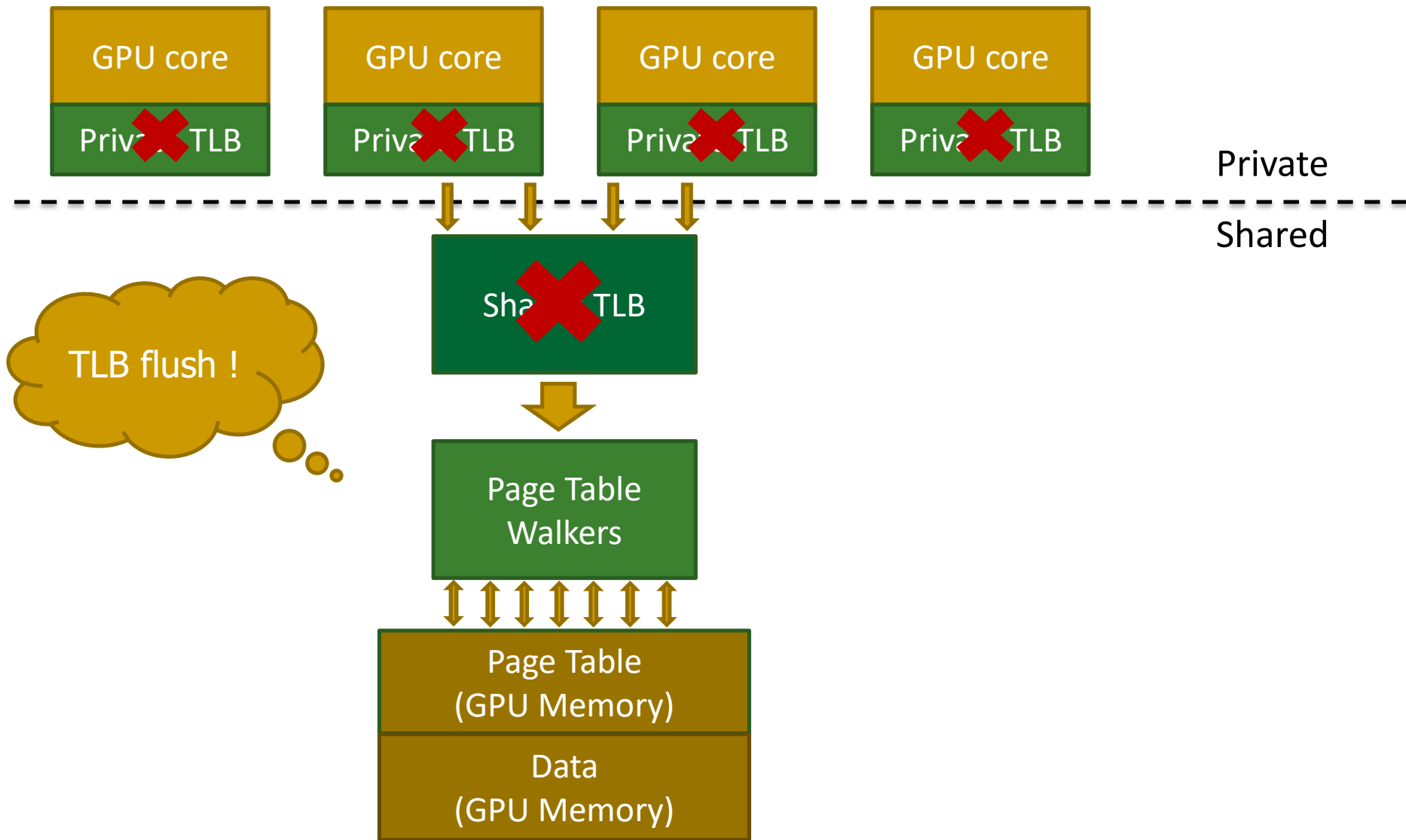**Can coalesce without moving data**

App1

App2

Unallocated

# Key Ideas

- Use **multiple** page sizes

- Translate using **large page size**
    High TLB reach

- Transfer using **small page size**
    Low demand paging latency

- Allocate physical pages in a way that **avoids the need to migrate data**

# Mechanism

# Mosaic

- **3 components**
- **Contiguity-Conserving Allocation**
- **In-Place Coalescer**
- **Contiguity-Aware Compaction**

GPU Runtime

| Contiguity-Conserving Allocation | In-Place Coalescer | Contiguity-Aware Compaction |
|---|---|---|

Hardware

**SAFARI**

# Mosaic

- 3 components

- **Contiguity-Conserving Allocation**

- In-Place Coalescer

- Contiguity-Aware Compaction

GPU Runtime



Hardware

**SAFARI**

# Mosaic: Data Allocation

GPU Runtime

Application demands data

**1**

Contiguity-Conserving Allocation

In-Place Coalescer

Contiguity-Aware Compaction

Hardware

Allocate memory

- A typical GPGPU application allocates a large **number of base pages**

# Mosaic: Data Allocation

GPU Runtime

Application demands data

Contiguity-Conserving Allocation

In-Place Coalescer

Contiguity-Aware Compaction

Hardware

Allocate memory

**2**

Large Page Frame

Page Table

Data

Key observation:
**en masse** memory allocation, i.e., applications allocate a large number of base pages at once

- **Conserve contiguity of base pages** - Virtual memory are contiguous within a large page frame in physical memory

SAFARI

# Mosaic: Data Allocation

GPU Runtime    Application demands data

**Contiguity-Conserving Allocation**

In-Place Coalescer

Contiguity-Aware Compaction

Hardware

Allocate memory

**2**

Large Page Frame

Page Table

Data

Soft guarantee:
A large page frame contains pages from only a single address space

- **Conserve contiguity of base pages** - Virtual memory are contiguous within a large page frame in physical memory

# Mosaic: Data Allocation



- Transfer data/pages at a **smal page** granularity
  - ❑ A page that is transferred is immediately ready to use – low latency

# Mosaic: Data Allocation

Application demands data

**Contiguity-Conserving Allocation**

In-Place Coalescer

Contiguity-Aware Compaction

Hardware

Allocate memory

GPU Threads

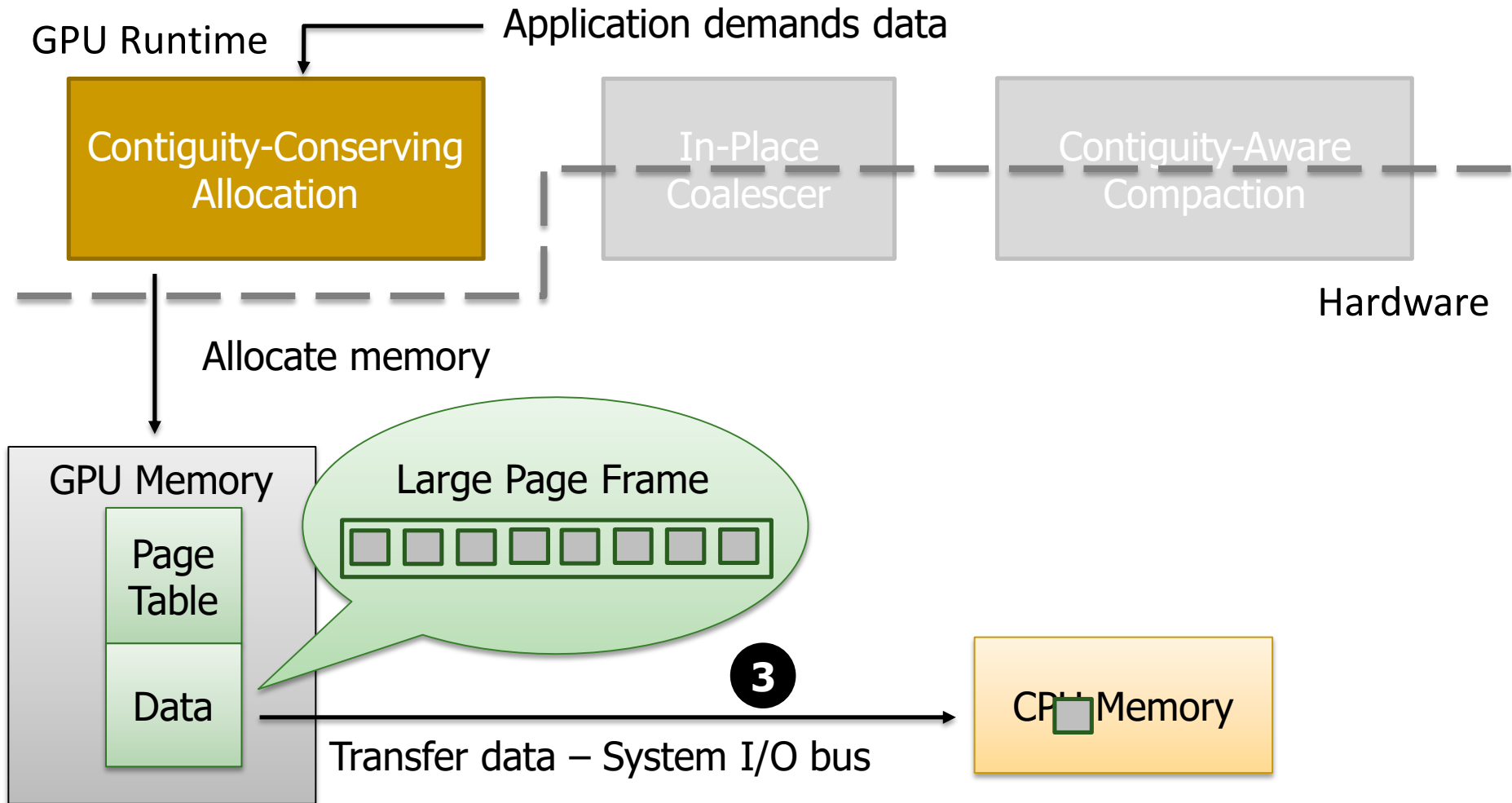GPU Memory

Page Table

Large Page Frame

stall

**3**

Data

CPU Memory

Transfer data – System I/O bus

- Transfer data/pages at a **smal page** granularity
  - A page that is transferred is immediately ready to use – low latency

# Mosaic: Data Allocation

GPU Runtime

Application demands data

**Contiguity-Conserving Allocation**

In-Place Coalescer

Contiguity-Aware Compaction

Hardware

Allocate memory

GPU Memory

Page Table

Data

Large Page Frame

❸

CPU Memory

Transfer data – System I/O bus

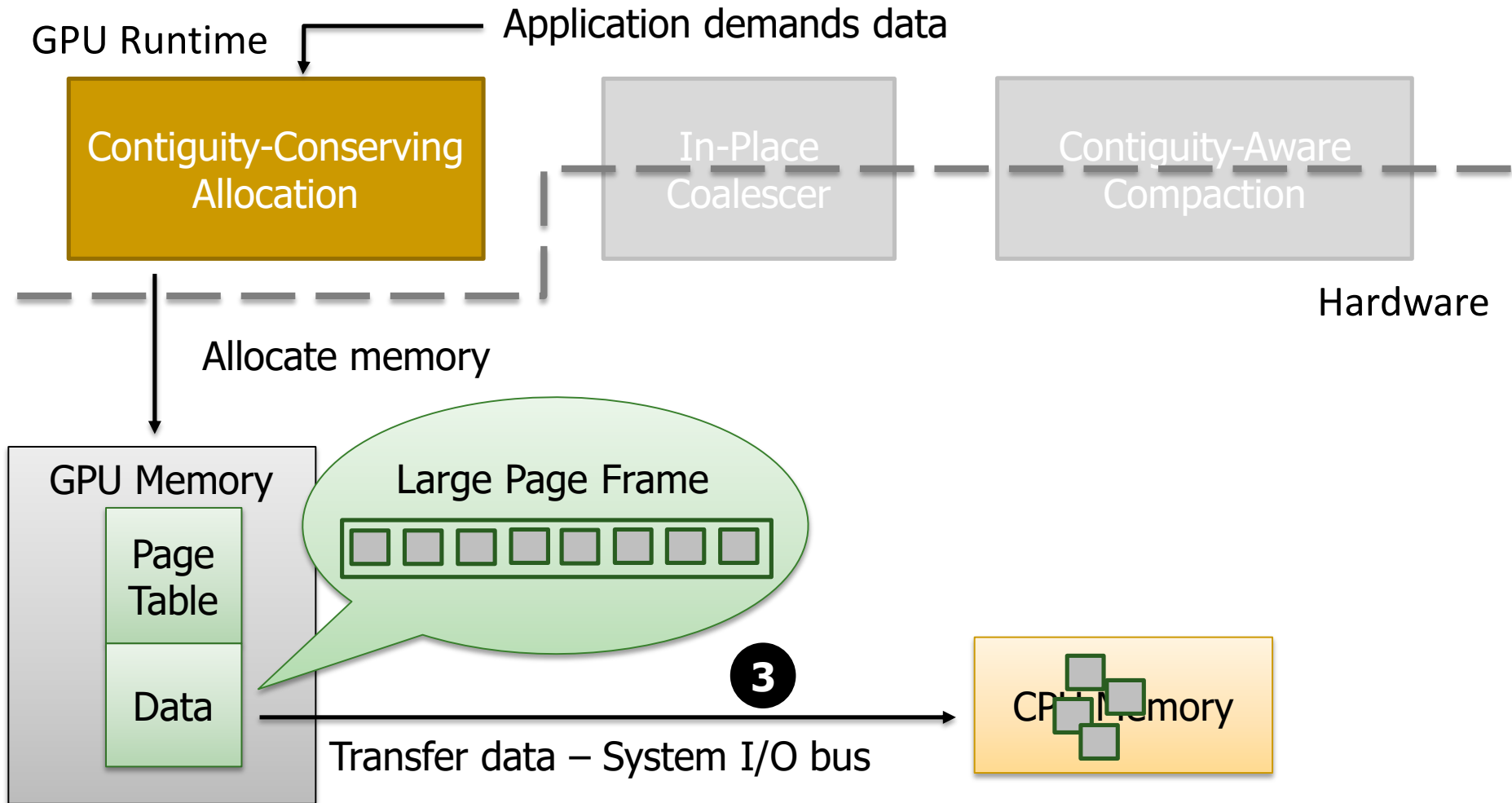- Transfer data/pages at a **smal page** granularity
  - ❑ A page that is transferred is immediately ready to use – low latency

# Mosaic: Data Allocation

GPU Runtime

| Contiguity-Conserving Allocation | In-Place Coalescer | Contiguity-Aware Compaction |
|---|---|---|

Hardware

**Transfer done**

**4**

Page Table

Data

CPU Memory

Transfer data – System I/O bus

- Send to In-Place Coalescer a list of the large page frame addresses that were allocated

# Mosaic

- 3 components
- Contiguity-Conserving Allocation
- **In-Place Coalescer**
- Contiguity-Aware Compaction

GPU Runtime

| Contiguity-Conserving Allocation | | In-Place Coalescer | | Contiguity-Aware Compaction |

Hardware

# Mosaic: Coalescing

GPU Runtime



- Fully-allocated large page frames = coalesceable
  - ❑ **Contiguous** in both virtual and physical memory
  - ❑ All base pages within the large page frame have been allocated and **belong** to the same address space

# Mosaic: Coalescing

GPU Runtime



- **Coalesce without moving data**
  - Simply update the page tables
  - No need for TLB flush
- **With an application transparent way**

# Mosaic: Coalescing

GPU Runtime

Contiguity-Conserving Allocation

**In-Place Coalescer**

Contiguity-Aware Compaction

Hardware

Page Table

Data

List of coalesceable pages

Update page tables

Large Page Table

Small Page Table

1

Coalesced Bit

Large page

| page number | page offset |
|---|---|

Small page

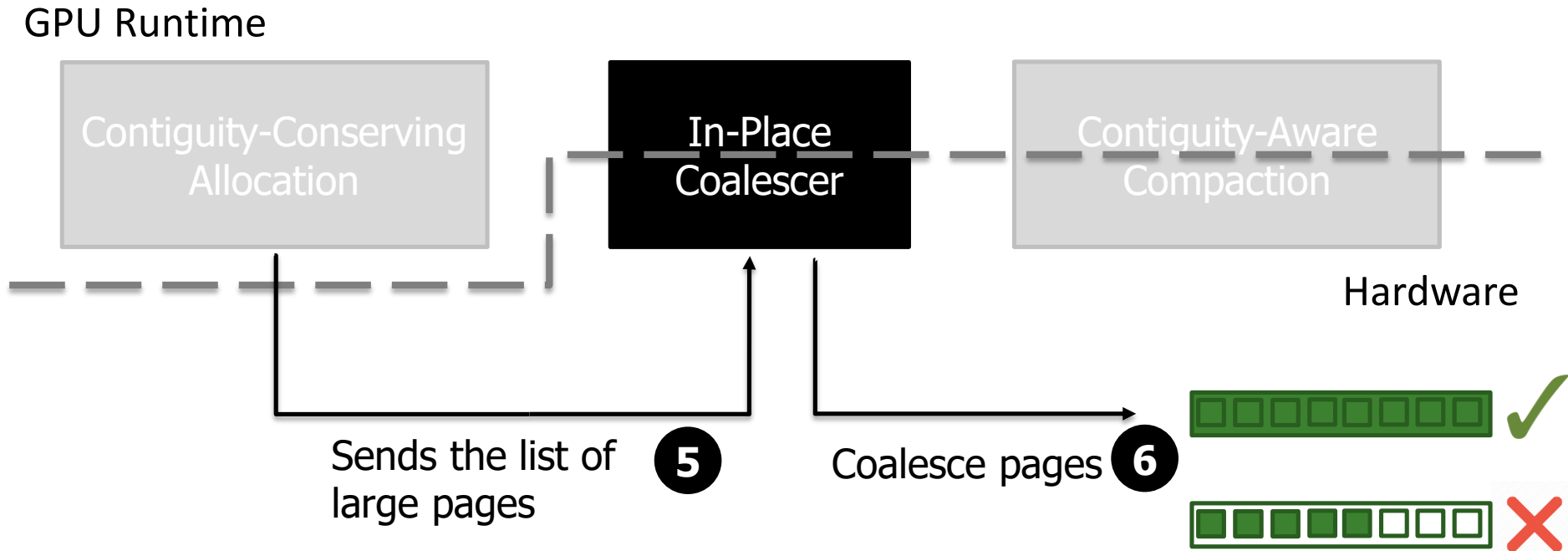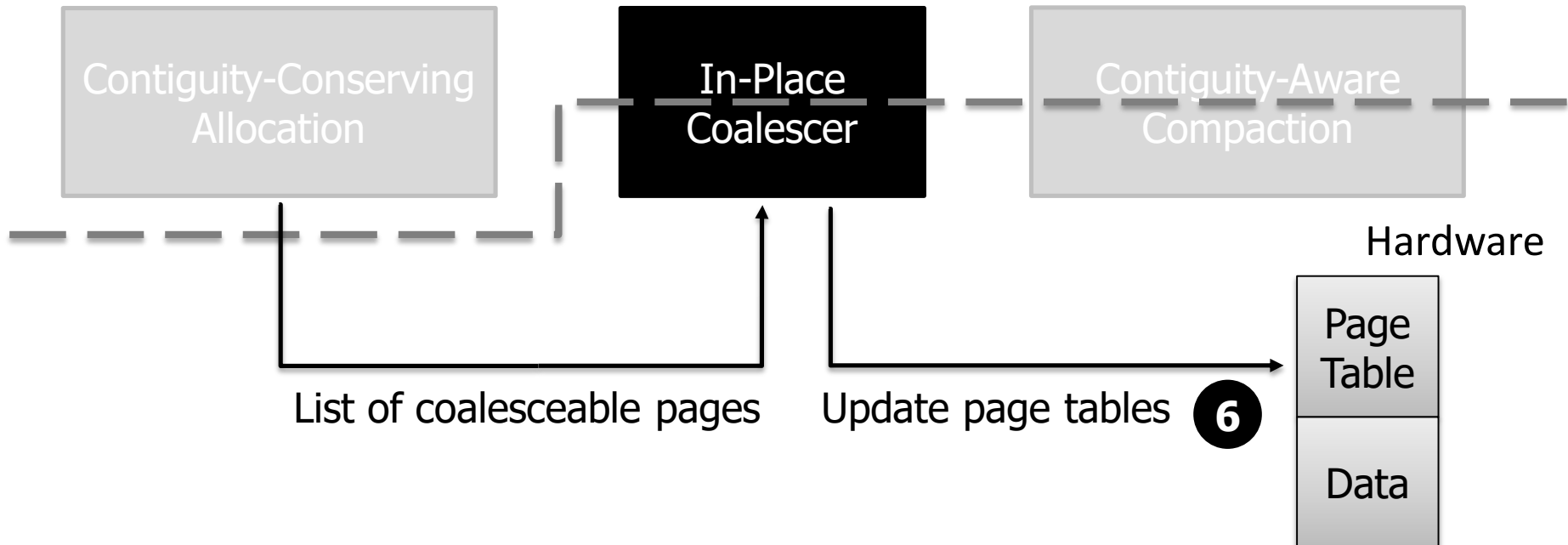| small page number | page offset |
|---|---|

- Data can be accessed using **either** page size

**SAFARI**

# Mosaic

- 3 components
- Contiguity-Conserving Allocation
- In-Place Coalescer
- **Contiguity-Aware Compaction**

GPU Runtime

| Contiguity-Conserving Allocation | In-Place Coalescer | Contiguity-Aware Compaction |
|---|---|---|

Hardware

# Mosaic: Data Deallocation

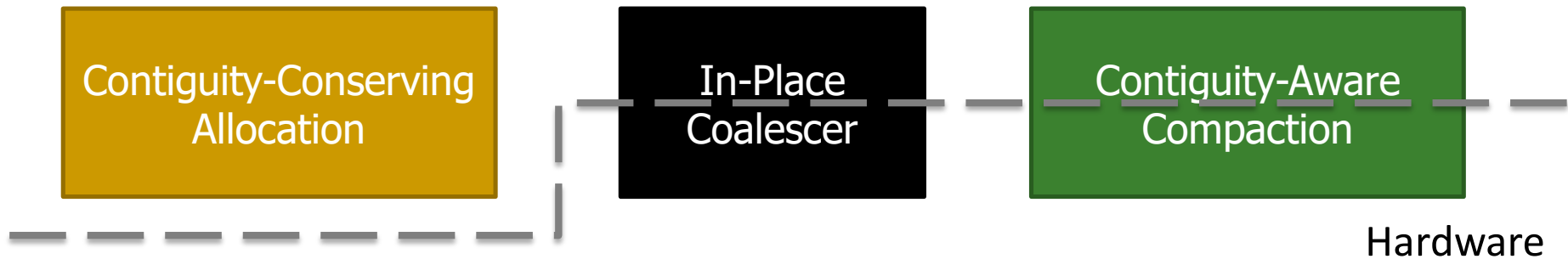GPU Runtime

Application deallocates memory



**Contiguity-Conserving Allocation**

**In-Place Coalescer**

**7**

**Contiguity-Aware Compaction**

Hardware

☐ gpu_free()

- It **sends a deallocation request** to the GPU runtime
- The Runtime invokes Contiguity-Aware Compaction for the corresponding large page

**SAFARI**

# Mosaic: Data Deallocation

GPU Runtime

Application deallocates memory

| Contiguity-Conserving Allocation | | In-Place Coalescer | **7** | Contiguity-Aware Compaction |

Hardware

- Check whether the large page frame has **high degree of internal fragmentation**
- Free-up not fully-used large page frames

**SAFARI**

# Mosaic: Data Deallocation

GPU Runtime

| Contiguity-Conserving Allocation | | In-Place Coalescer | | Contiguity-Aware Compaction |

Hardware

**8** Splinter pages

Page Table

Data

- **Update the page table** to splinter the large page back into its constituent base pages

**SAFARI**

# Mosaic: Data Deallocation

GPU Runtime

| Contiguity-Conserving Allocation | In-Place Coalescer | **Contiguity-Aware Compaction** |
|---|---|---|

Hardware

| Page Table |
|---|
| **Data** |

Compact pages by migrating data ⑨

- Compaction: **Migrating** the remaining base pages to another uncoalesced large page frame that belongs to the same application

**SAFARI**

# Mosaic: Data Deallocation

GPU Runtime



- Compaction: **Migrating** the remaining base pages to another uncoalesced large page frame that belongs to the same application

**SAFARI**

# Mosaic: Data Deallocation

GPU Runtime

| Contiguity-Conserving Allocation | In-Place Coalescer | Contiguity-Aware Compaction |
|---|---|---|

Hardware

**Not coalescable**

Free large page

Page Table
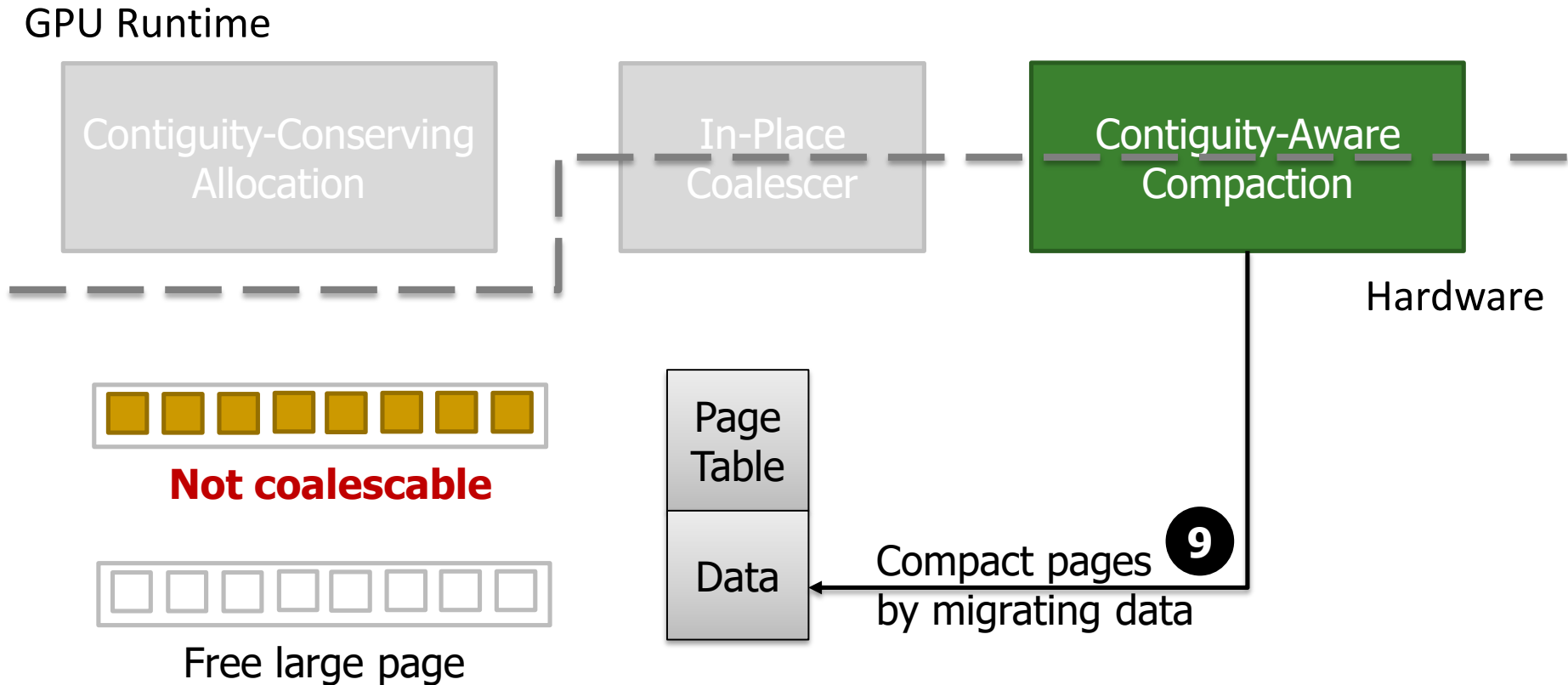
Data

Compact pages by migrating data  **9**

- Compaction: **Migrating** the remaining base pages to another uncoalesced large page frame that belongs to the same application

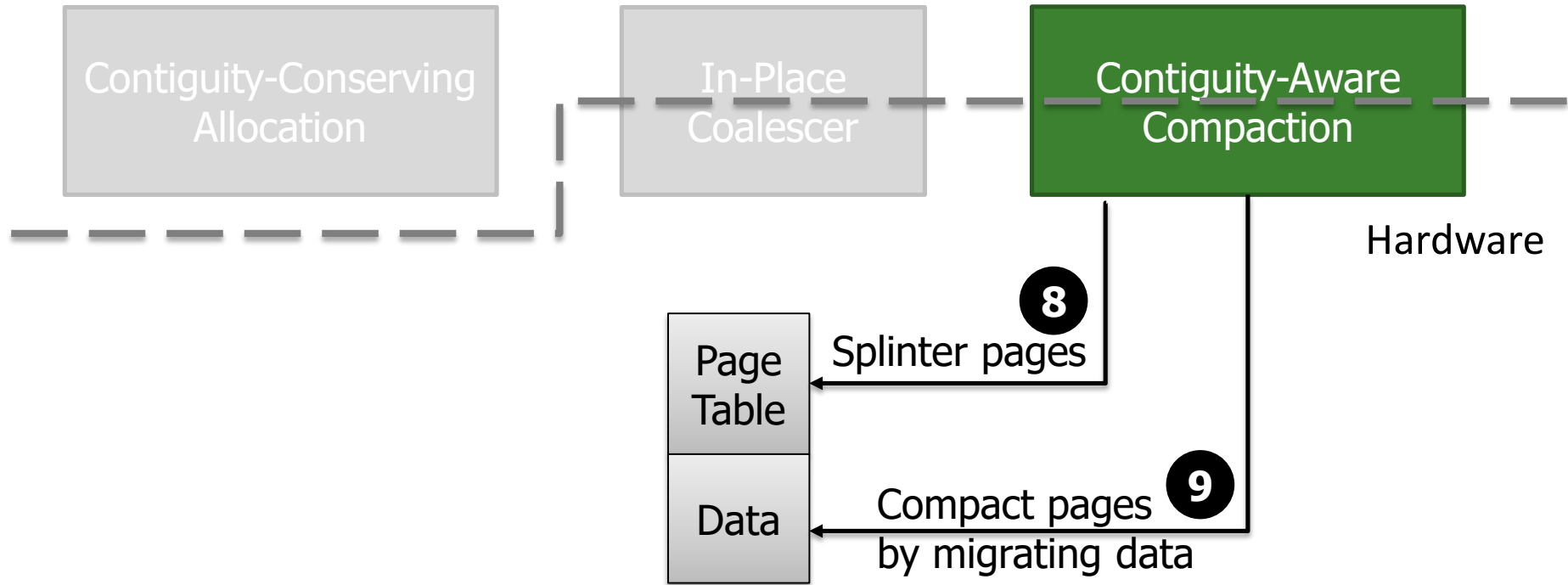# Mosaic: Data Deallocation

GPU Runtime



- Page Migration is required
- **TLB flush** is required

**SAFARI**

# Mosaic: Data Deallocation

**GPU Runtime**

Send list of newly-free pages after compaction  ⑩

| Contiguity-Conserving Allocation | In-Place Coalescer | Contiguity-Aware Compaction |

☐ ☐ Free large page ☐  **Hardware**

- Contiguity-Aware Compaction component notifies Contiguity-Conserving Allocation of the large page frames that are now free after compaction, such that to be used for future memory allocations

# Methodology

# Methodology

- MAFIA framework which uses GPGPU-Sim
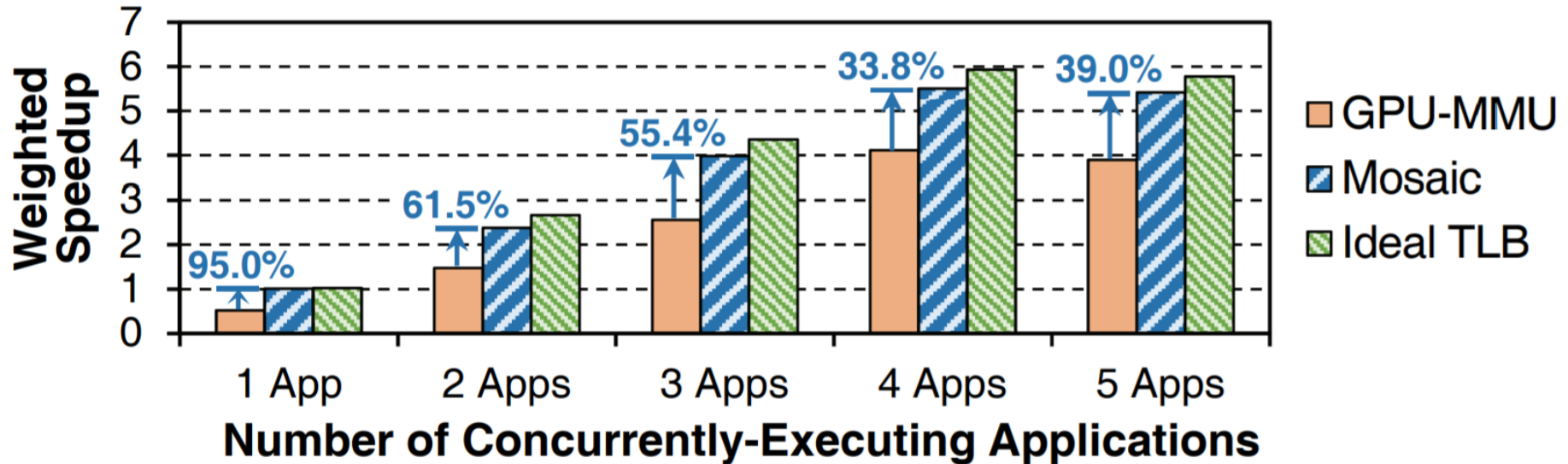  - 30 cores @1020MHz
  - 64KB 4-way L1, 2048KB 16-way L2
  - Private L1 TLB: 128 base pages / 16 large page entries per core
  - Shared L2 TLB: 512 base pages / 256 large page entries
  - DRAM: 3GB GDDR5 @1674 MHz

- Model sequential page walks

- Workloads
  - Homogeneous workloads = multiple copies of the same application
  - Heterogeneous workloads = randomly selected applications
  - Multiple GPGPU applications execute concurrently
  - Test suites: Parboil, SHOC, LULESH, Rodinia, CUDA SDK

- Evaluation metric
  - Weighted Speedup = $\sum_{i=1}^{N} \frac{IPCshared}{IPCalone}$ for each application $i$

SAFARI

# Comparison Points

1. State-of-the-art CPU-GPU memory management
   - GPU-MMU based on [Power et al., HPCA'14]
     - Utilizes parallel page walks, TLB request coalescing and page walk cache to improve performance
     - <u>Limited</u> TLB reach (4KB pages)

2. Ideal TLB: Every TLB access is a L1 TLB hit

SAFARI

# Evaluation

# Homogeneous workloads



Mosaic consistently improves performance:
- **55.5%** averaged over GPU-MMU across 135 workloads
- comes within **6.8%** of the **Ideal** TLB

# Heterogeneous workloads

# Heterogeneous workloads



Mosaic :
- significantly improves performance for TLB-friendly workloads

# Heterogeneous workloads



Mosaic :
- significantly improves performance for TLB-friendly workloads
- provides less benefit in TLB-sensitive workloads

# TLB Hit Rate



Mosaic significantly reduces the TLB miss rate:
- the average miss rate falls **below 1%** in both the L1 and L2 TLBs

**SAFARI**

# More Results in the Paper

- Per-application IPC

- Sensitivity analysis to different TLB sizes

- Memory fragmentation analysis

Mosaic is available at: https://github.com/CMU-SAFARI/Mosaic

# Summary

# Executive Summary

- Problem
  - No single best page size for GPU virtual memory (large vs small pages)

- Goal
  - Transparently and efficiently enable <u>both</u> page sizes

- Key Observation
  - Can easily coalesce an application's **contiguously-allocated** small pages into a large page
  - GPGPU applications typically allocate large chunks of memory *en masse*

- Key Idea
  - <u>*Preserve the virtual address contiguity*</u> of small pages when allocating physical memory to simplify coalescing

- Mosaic:
  - A hardware/software cooperative framework
  - Enables the benefits of both small and large pages

- Key Result: 55% on average performance improvement over state-of-the-art GPU memory management mechanism

# Strengths & Weaknesses

# Strengths

- **Intuitive Idea:**
  - ❑ Exploits the benefits of using both small and large pages
  - ❑ Well-written, insightful paper

# Strengths

- Intuitive Idea:
  - Exploits the benefits of using both small and large pages
  - Well-written, insightful paper
- **Mechanism**:
  - Avoids page migration when memory alocation is requested
  - Application-transparent support for multiple page sizes
  - Data can be accessed using either page size

**SAFARI**

# Strengths

- Intuitive Idea:
  - Exploits the benefits of using both small and large pages
  - Well-written, insightful paper

- Mechanism:
  - Avoids page migration when memory alocation is requested
  - Application-transparent support for multiple page sizes
  - Data can be accessed using either page size

- **Evaluation**:
  - Investigate behavior of multiple GPGPU applications that run concurrently
  - Explore the performance in case of highly fragmented pages
  - High variety of workloads explored

**SAFARI**

# Strengths

- Intuitive Idea:
  - Exploits the benefits of using both small and large pages
  - Well-written, insightful paper

- Mechanism:
  - Avoids page migration when memory alocation is requested
  - Application-transparent support for multiple page sizes
  - Data can be accessed using either page size

- Evaluation:
  - Investigate virtual memory when multiple GPGPU applications run concurrently
  - Explore the performance in case of highly fragmented pages
  - High variety of workloads explored

- **Online available**

SAFARI

# Weaknesses

- **Mechanism**:
  - Provides **soft guarantee** that a large page frame contains pages from only a single address space
  - What is the threshold after which Mosaic splinters a large page frame into small pages?
  - Needs many changes in the system stack
    - Software-hardware cooperative solutions are not always be easy to adopt

SAFARI

# Weaknesses

- **Mechanism:**
  - Provides soft guarantee that a large page frame contains pages from only a single address space
  - What is the threshold after which Mosaic splinters a large page frame into small pages?
  - Needs many changes in the system stack
    - Software-hardware cooperative solutions are not always be easy to adopt
- **Evaluation**:
  - No comparison with an approach that uses large page frames
    - Mosaic mainly benefits from TLB-friendly applications
  - Model sequential page walks in simulation
  - Less benefit in TLB sensitive applications and highly fragmented pages
  - Simulation-based evaluation

**SAFARI**

# Takeaways

# Takeaways

- A novel idea to enable benefits of both small and large pages

- Hardware/software cooperative framework

- Application-transparent support for multiple page sizes
  - No TLB flush when coalescing

- Online available

- Easy to read and understand paper

**SAFARI**

# Open Discussion

# Discussion

- We do not completely avoid data migration !

- Avoid page migration to the critical path :
  - gpu_malloc();
  - … access …

**On allocation: do not move the data**

SAFARI

# Discussion

- We do not completely avoid data migration !

- Avoid page migration to the critical path :
  - gpu_malloc();
  - … access …
  - gpu_free();
  - On deallocation:

**On deallocation: move the data**

SAFARI

# Discussion

- We do not completely avoid data migration !

- Avoid page migration in the critical path :
  - gpu_malloc();
  - ... access ...
  - gpu_free();
  - On deallocation:

**Be Optimistic !**

- Any similar concepts ?

SAFARI

# Discussion

- We do not completely avoid data migration !

- Hardware Transactional Memory (HTM)

  **Be Optimistic !**

- Lazy PIM: CAL 2016

- Other works related to speculation ?

# Discussion

- Mosaic does not significantly improve the performance for **TLB sensitive** workloads

- No comparison with other research works that use <u>large</u> page sizes

- Any ideas to extend this work for TLB-sensitive applications ?
  - **TLB prefetching ?**

**SAFARI**

# Discussion

- TLB is **shared among multiple** concurrently-executing applications. These applications <u>compete</u> for the shared TLB.
    - ❑ Can we improve inter-application interference?

**SAFARI**

# Discussion

- TLB is shared among multiple concurrently-executing applications. These applications compete for the shared TLB.

  - Can we improve inter-application interference?

  - TLB partitioning?

  - Static/Dynamic partitioning?

## MASK: Redesigning the GPU Memory Hierarchy to Support Multi-Application Concurrency

Rachata Ausavarungnirun[1]    Vance Miller[2]    Joshua Landgraf[2]    Saugata Ghose[1]
Jayneel Gandhi[3]    Adwait Jog[4]    Christopher J. Rossbach[2,3]    Onur Mutlu[5,1]

[1]Carnegie Mellon University    [2]University of Texas at Austin    [3]VMware Research
[4]College of William and Mary    [5]ETH Zürich

# Discussion

- Is the 'ideal' page size an <u>application-specific</u> parameter? How can we **predict** the 'correct' page size for each application?

  - ❑ tracks the difference or distance between successive TLB miss virtual pages to identify it

- How to apply such an idea when different applications are executing concurrently?

# Discussion

- Mosaic provides a **soft guarantee** that a large page frame contains pages from only a single address space.

  - ❑ No discussion about the heuristic used
  - ❑  Any good heuristic/idea?

**SAFARI**

printf ( "Thank you" );