

Mosaic: A GPU Memory Manager

with Application-Transparent Support for Multiple Page Sizes

Rachata Ausavarungrun, Joshua Landgraf, Vance Miller, Saugata Ghose, Jayneel Gandhi, Christopher J. Rossbach, Onur Mutlu

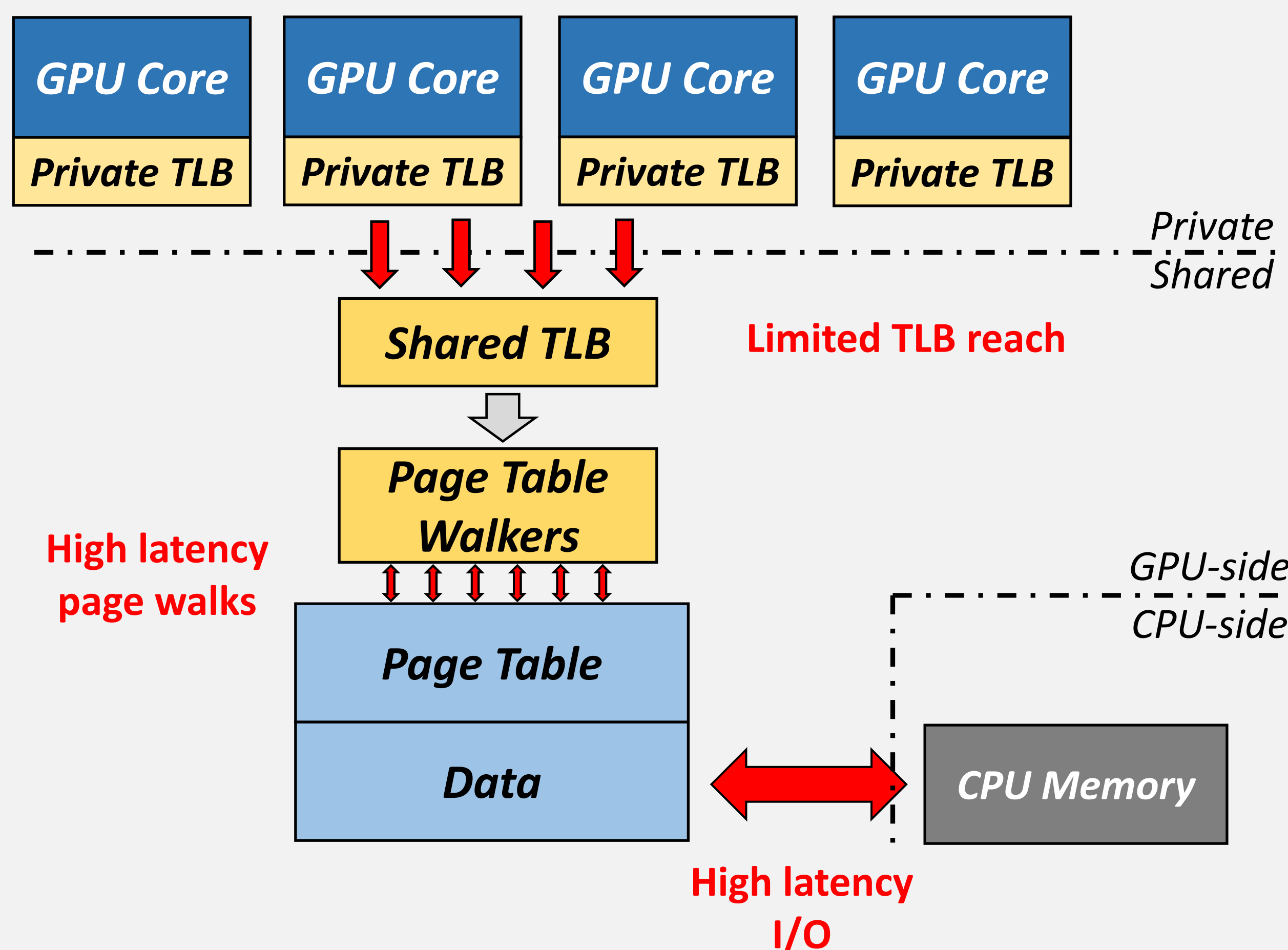


GPU Support for Virtual Memory

Improves **programmability** with a unified address space

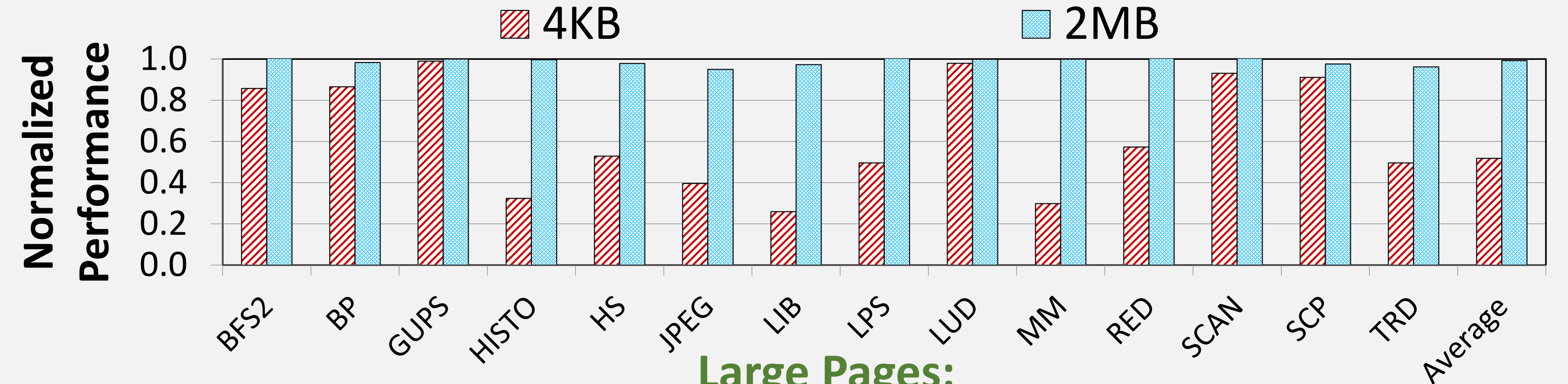
Enables **large data sets** to be processed in the GPU

Allows **multiple applications** to run on a GPU



Page Size Trade-Off

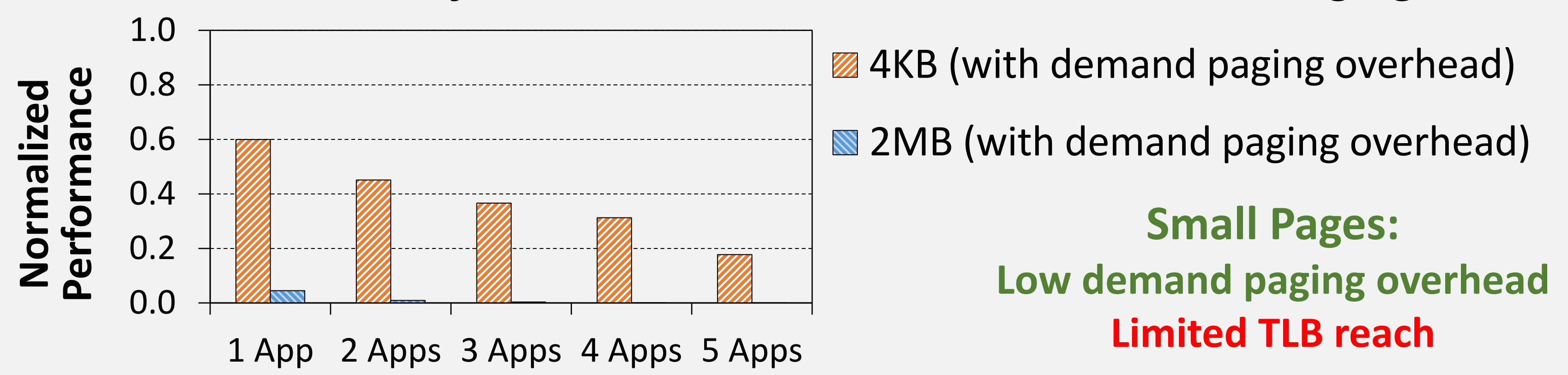
Overhead of Address Translation Without Demand Paging



Large Pages:
Better TLB reach

High demand paging overhead

Overhead of Both Address Translation and Demand Paging

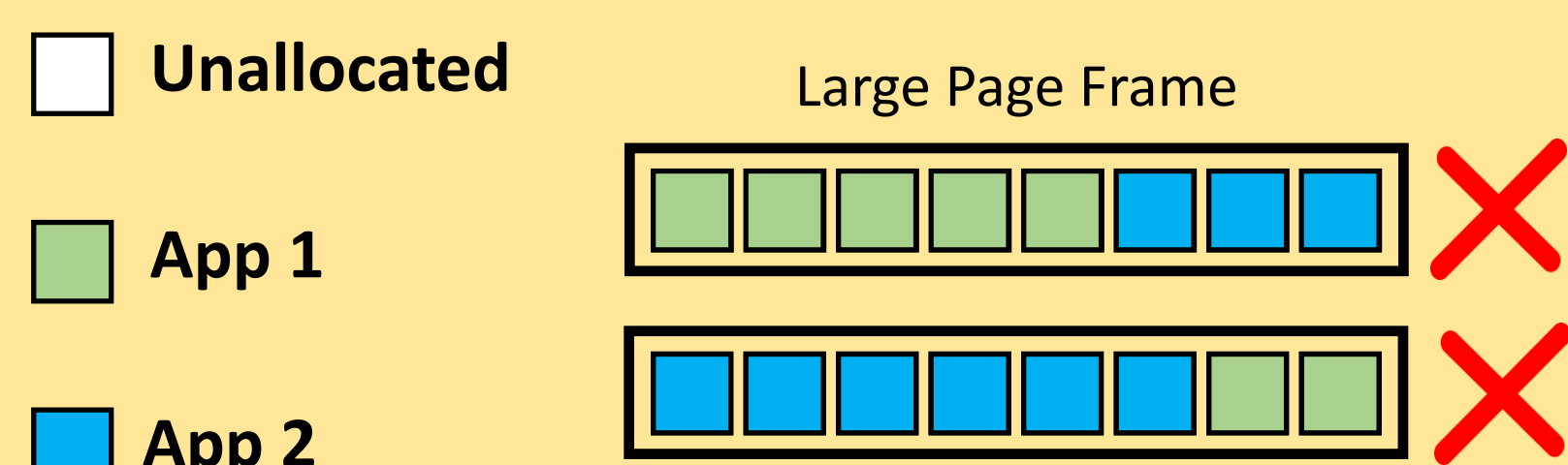


Small Pages:
Low demand paging overhead
Limited TLB reach

How to achieve the best of both page sizes?

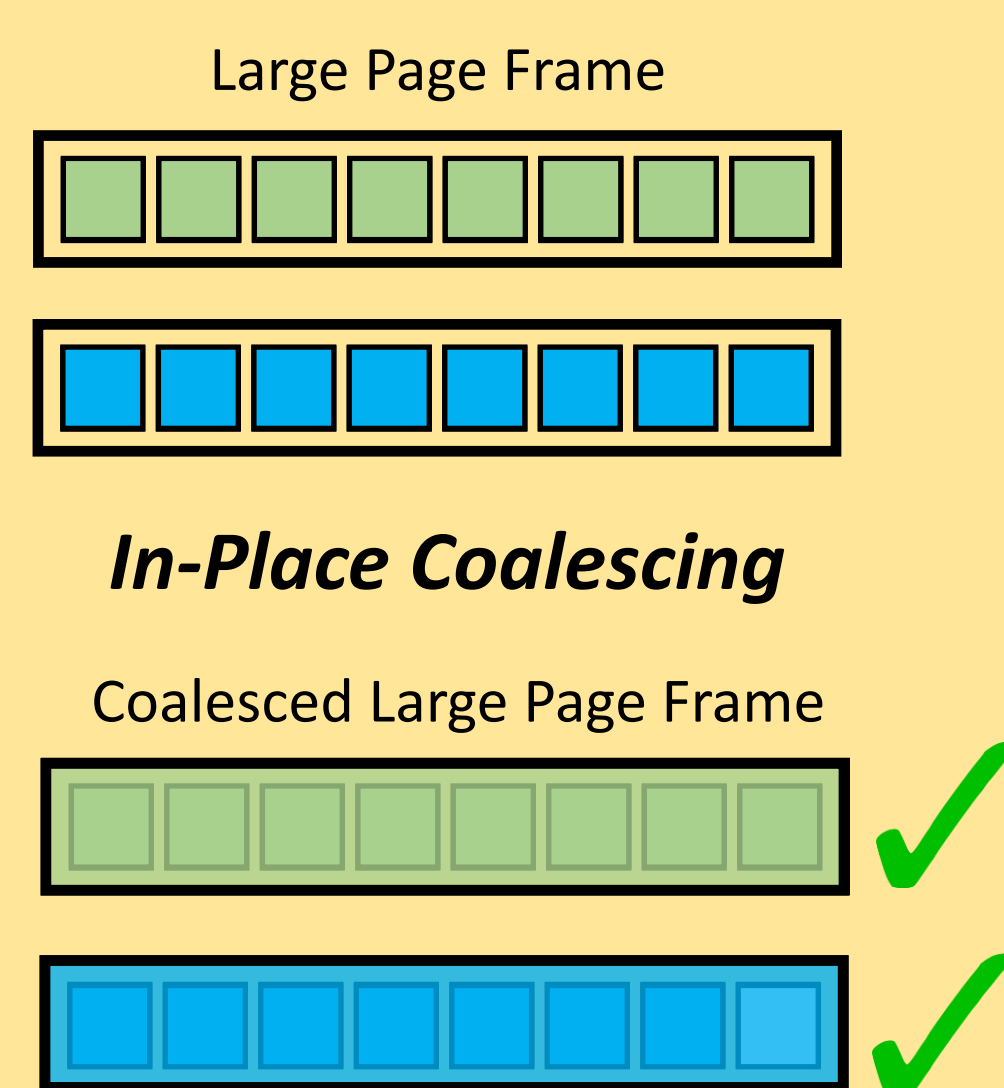
Challenges with Multiple Page Sizes

State-of-the-Art Memory Allocation



Cannot coalesce
(without moving multiple 4K pages)
Need to search which pages to coalesce
High overhead

Mosaic



Design Goals

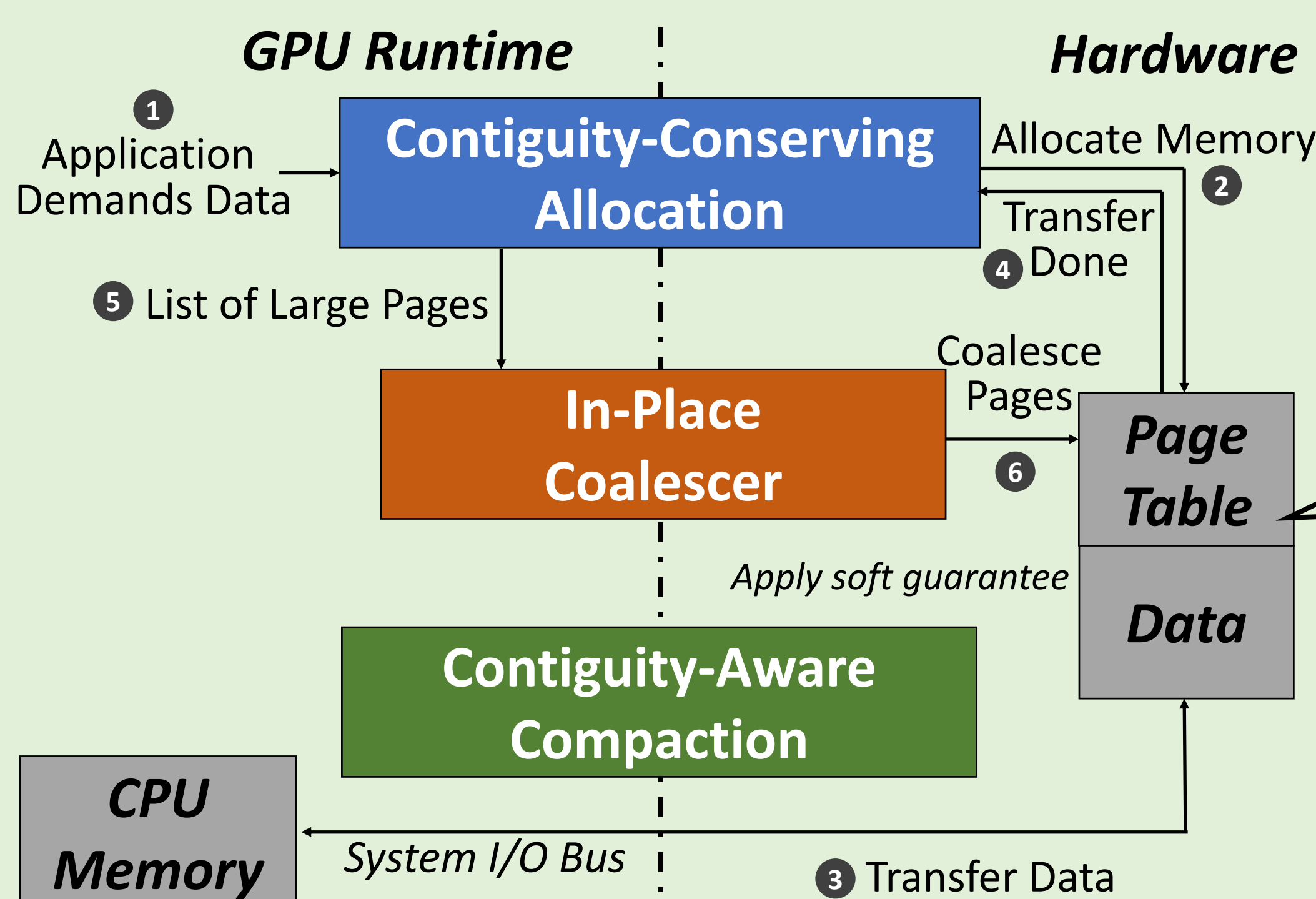
- Exploit benefits of both small and large pages
 - High TLB Reach
 - Low demand paging overhead
- No data movement
- Application transparency
 - Programmers do not need to modify GPGPU applications

Mosaic Soft Guarantee

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

High-Level Overview of Mosaic

Data Allocation



Coalescing

Fully-allocated large page frames → Coalescible

Application transparency

Data can be accessed using either page size

No TLB flush

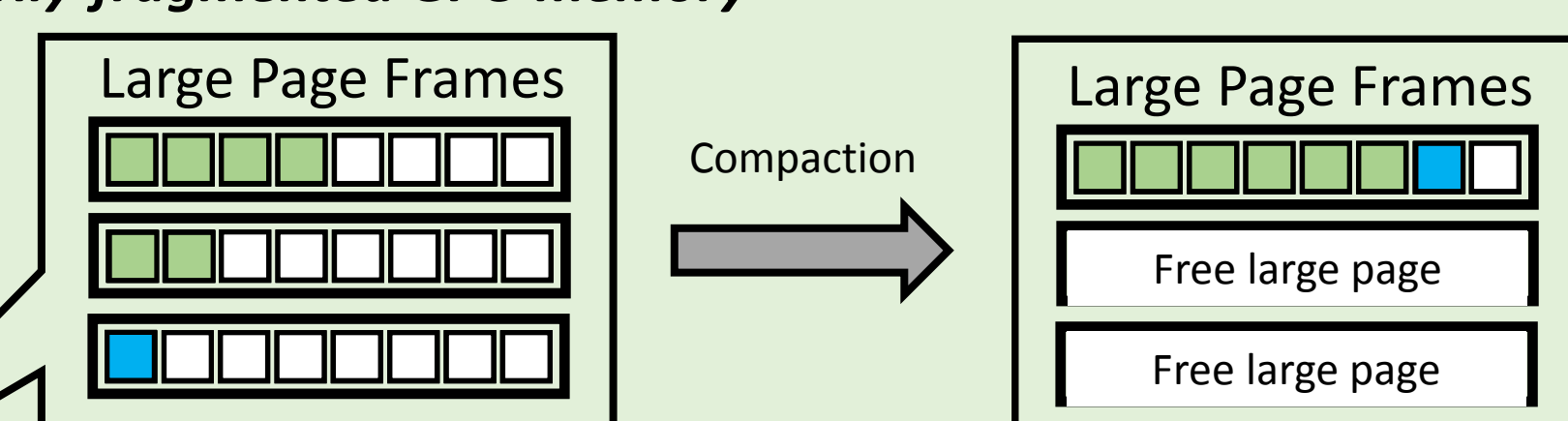
Compaction

Only triggered when memory is heavily fragmented

Goals

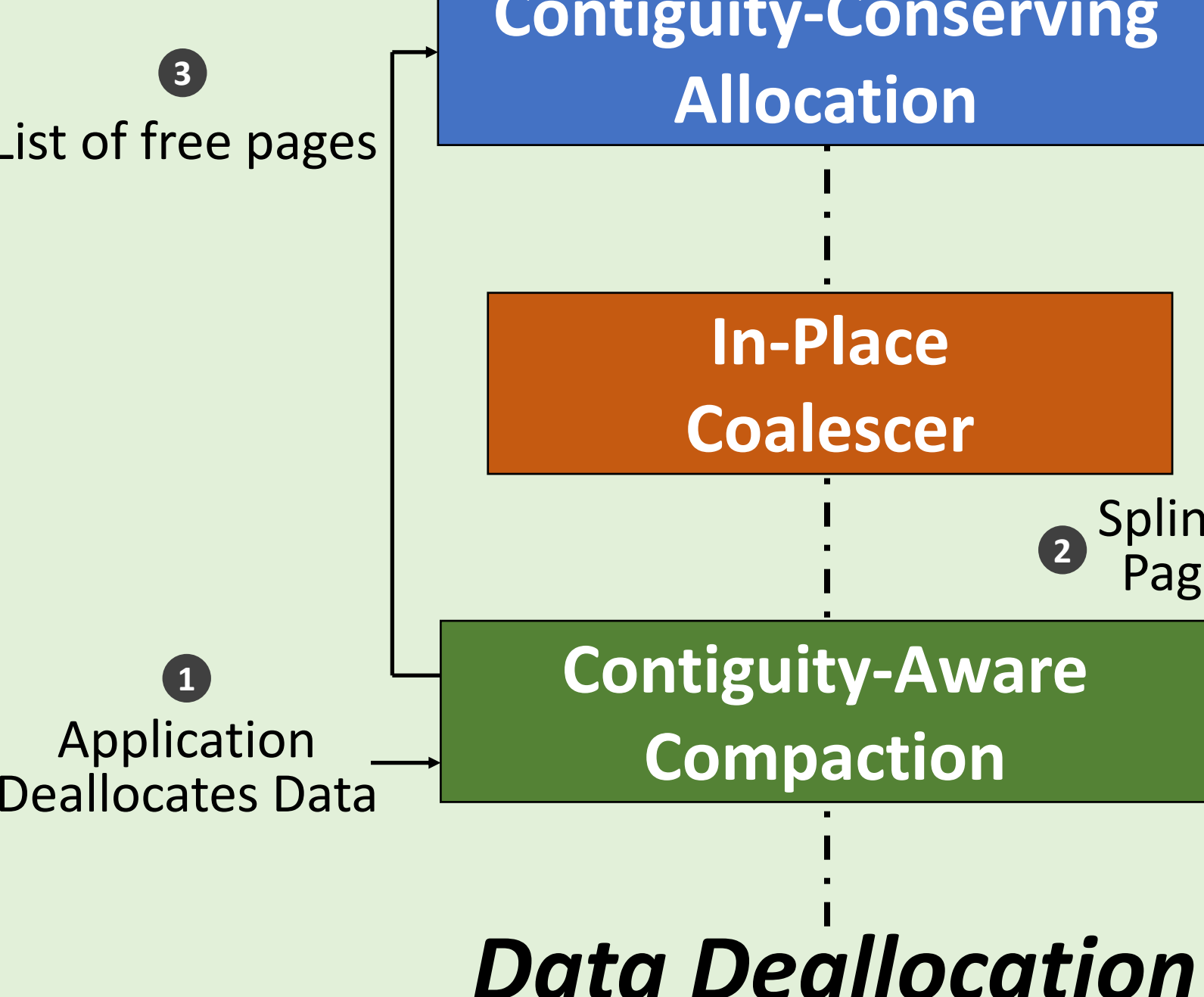
- Reduce memory fragmentation
- Free up large page frames

Heavily-fragmented GPU memory



Compacted pages has no virtual contiguity

No longer coalescible



Methodology

- GPGPU-Sim (MAFIA) configured to a GTX 750 Ti
- Multiple GPGPU applications can execute concurrently
- Model page walks and page tables
- Model virtual-to-physical address mapping
- Available at: <https://github.com/CMU-SAFARI/Mosaic>

Results

