



Reducing Memory Virtualization Overheads in Virtualized Datacenters

Jayneel Gandhi, Arkaprava Basu,
Michael M. Swift, Mark D. Hill

Executive Summary

- Problem: TLB misses in virtual machines
 - Hardware-virtualized MMU has high overheads
 - Up to 280% overhead
- Prior Work: **Direct Segments** – unvirtualized case
- Solution: segmentation to bypass paging
 - **Extend Direct Segments** for virtualization
 - **Three configurations** with different tradeoffs
- Results
 - **Near- or better-than-native** performance

Outline

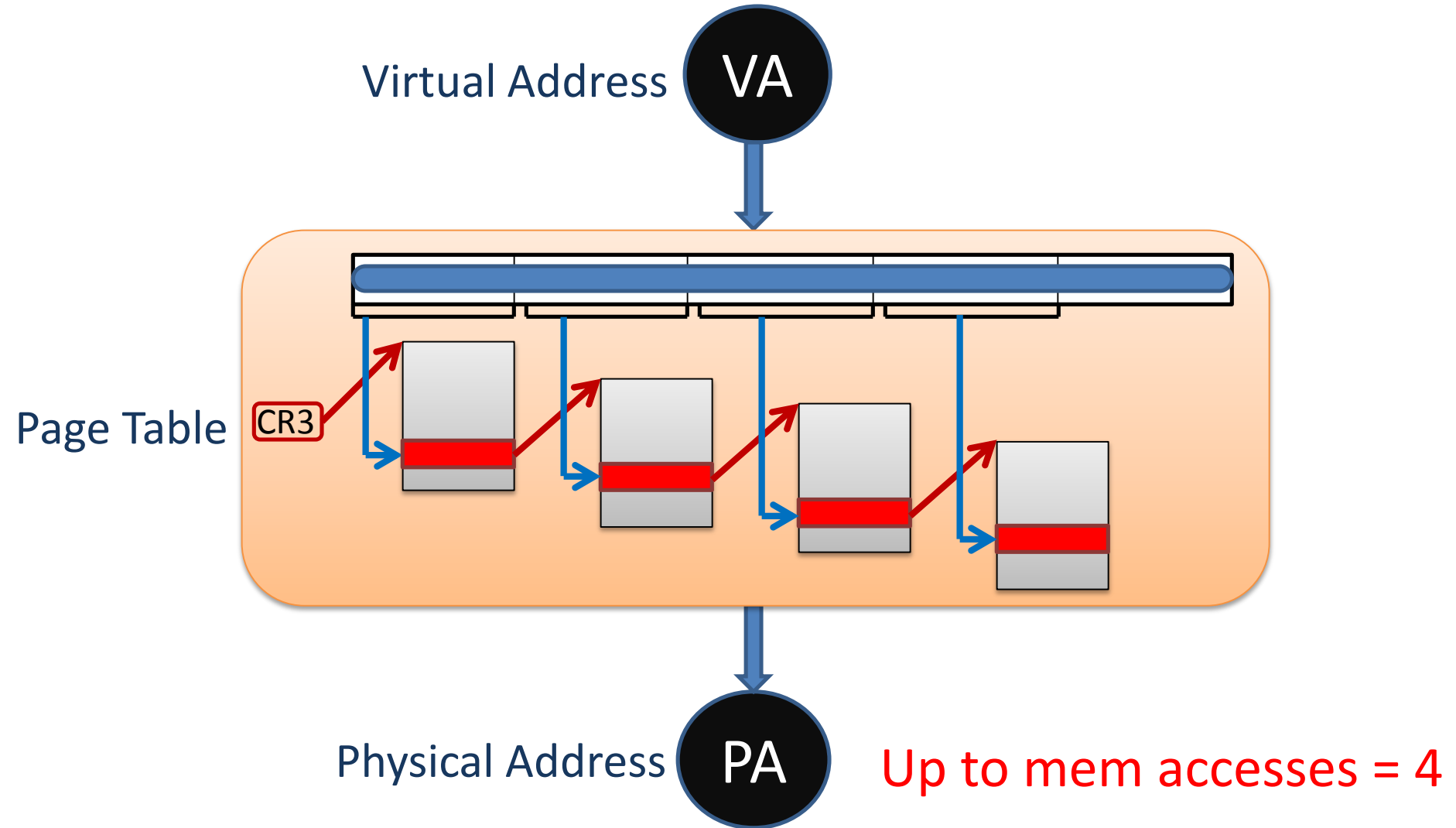
- Motivation ←
- Review: Direct Segments
- Virtualized Direct Segments
- Optimizations
- Evaluation
 - Methodology
 - Results
- Summary

Overheads of Virtualizing Memory

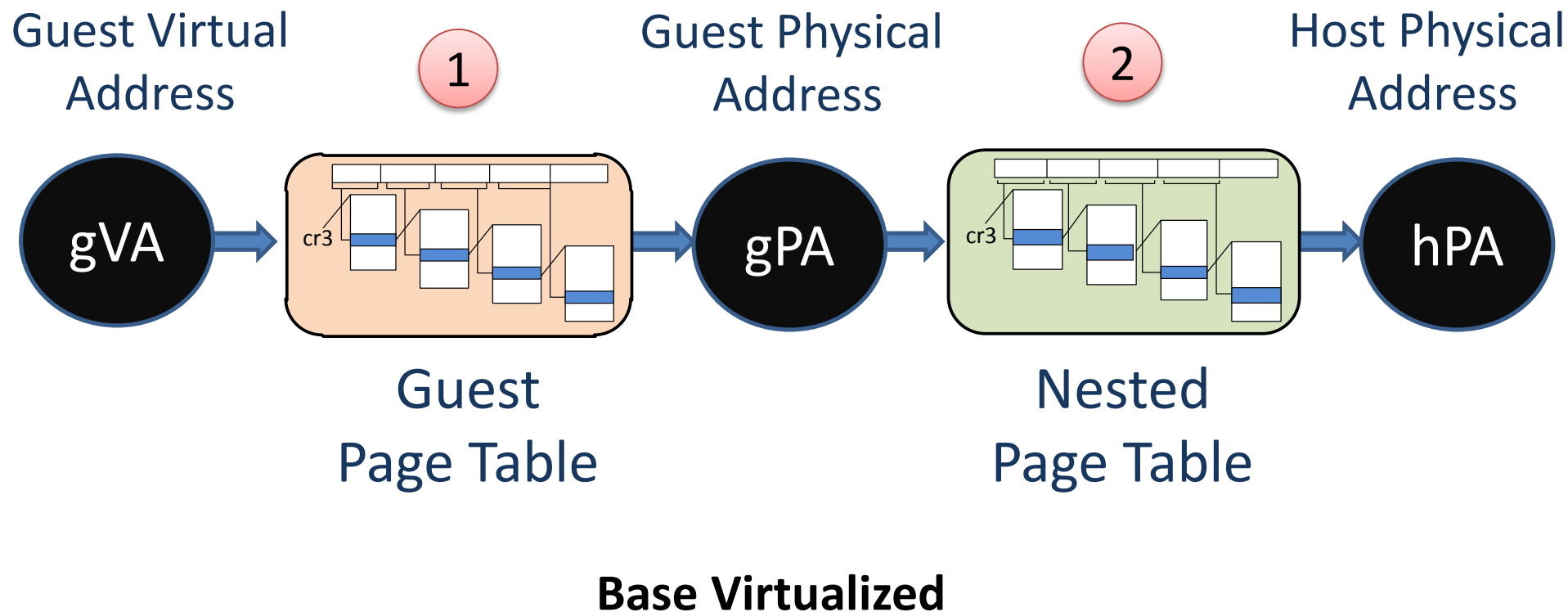
We show that that the increase in translation lookaside buffer (TLB) miss-handling costs due to hard-ware-assisted memory management unit (MMU) is the largest contributor to the performance gap between native and virtual servers.

—Buell, et al. VMware Technical Journal 2013

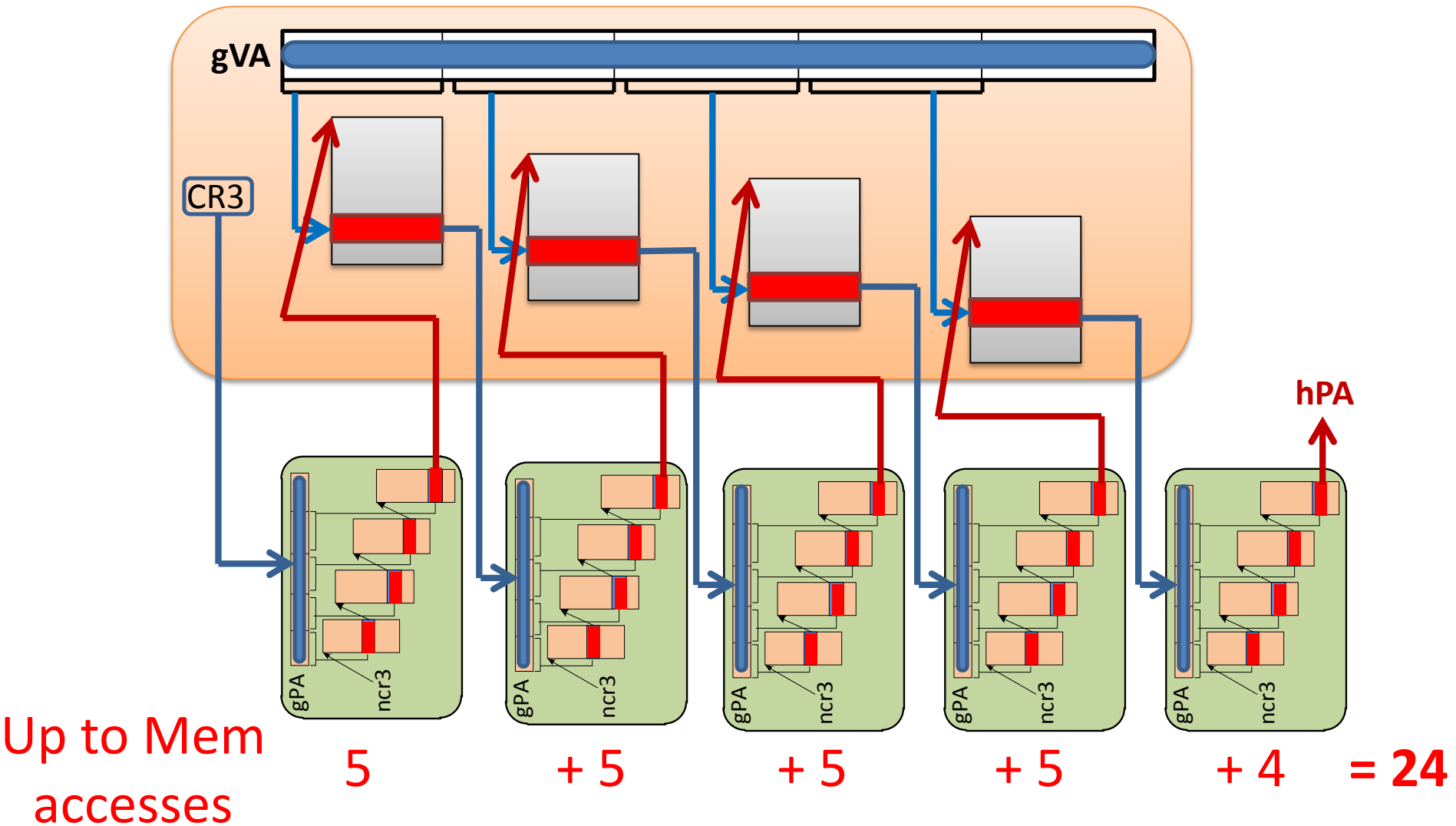
Unvirtualized x86 translation



Two-levels of translation



Support for Virtualizing Memory



Applications



Database



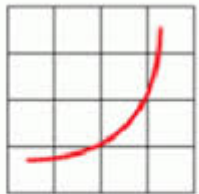
Graph-analytics



Key-value store



HPC Apps



spec

Standard Workloads



Big-memory applications

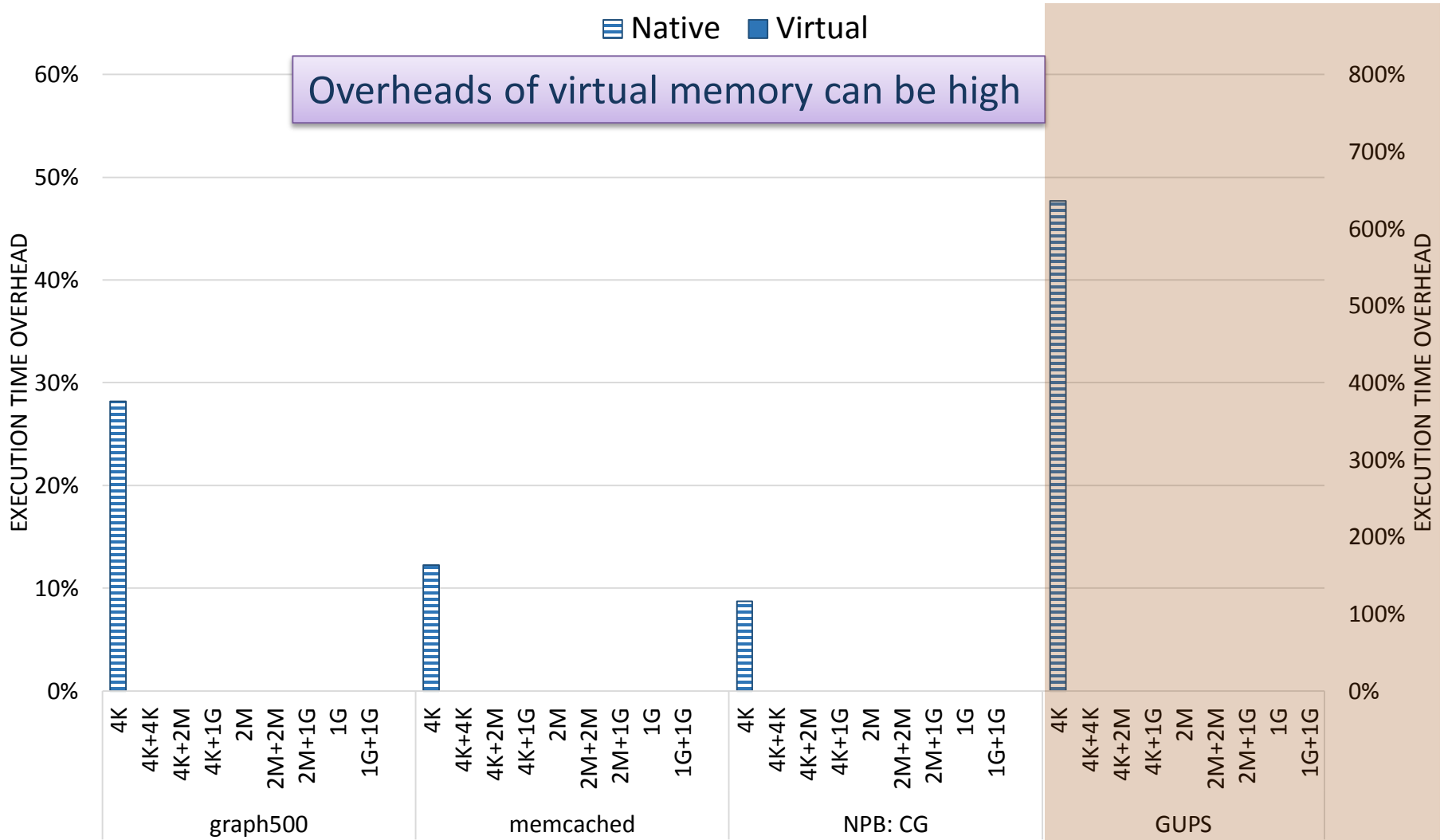
(In-memory apps)

PARSEC

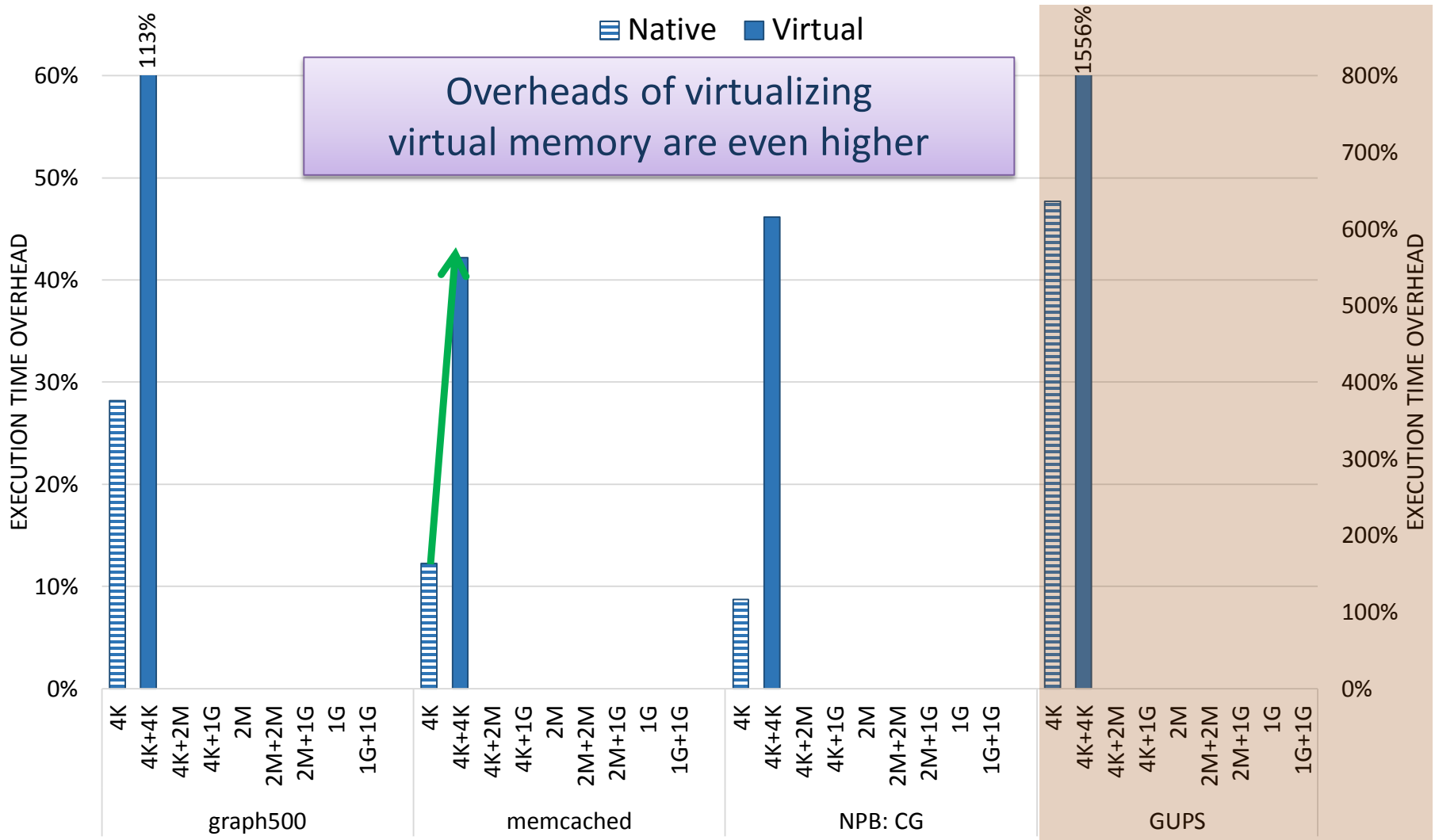


[par-sec] A unit of measure

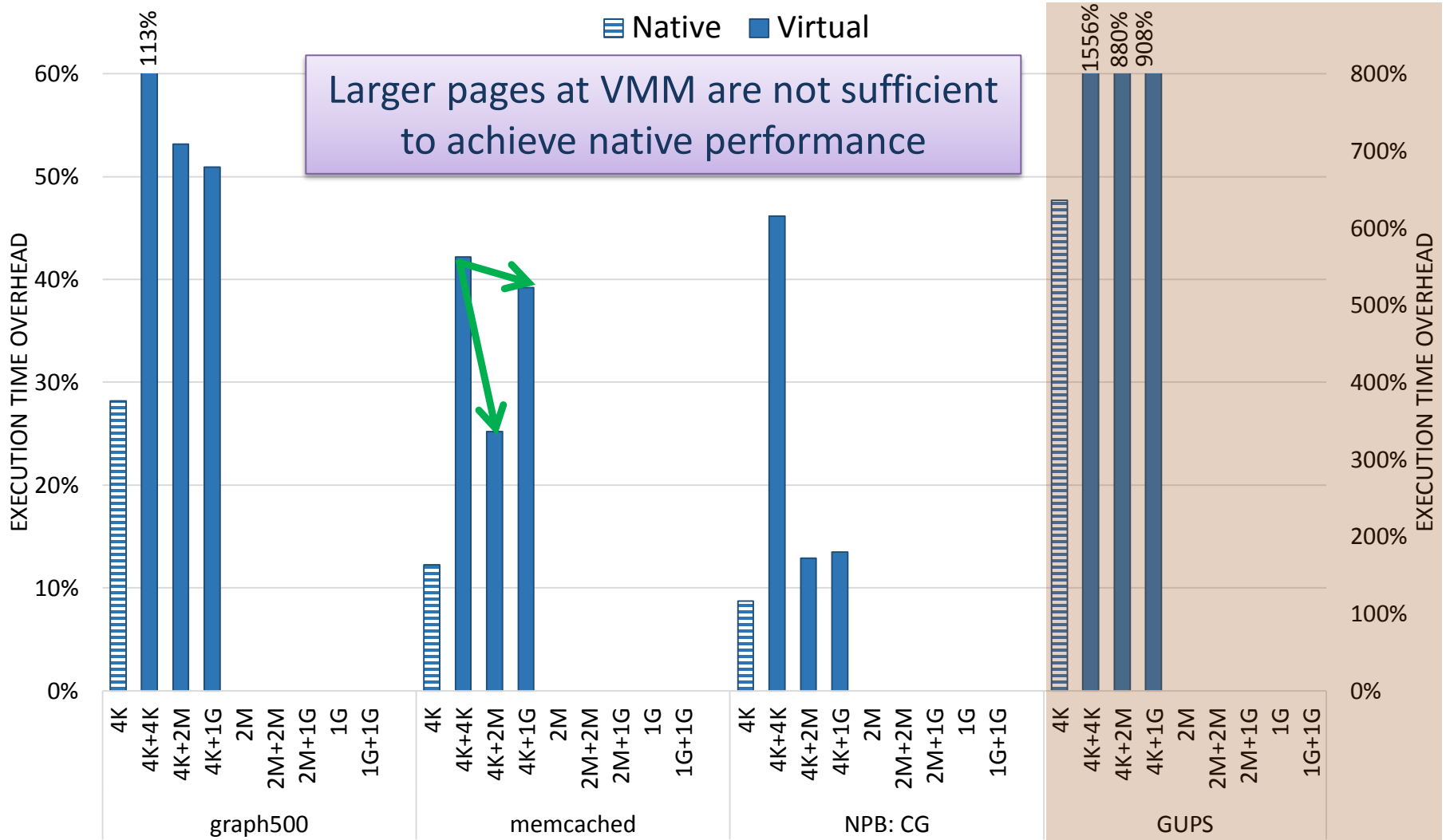
Cost of Virtualization



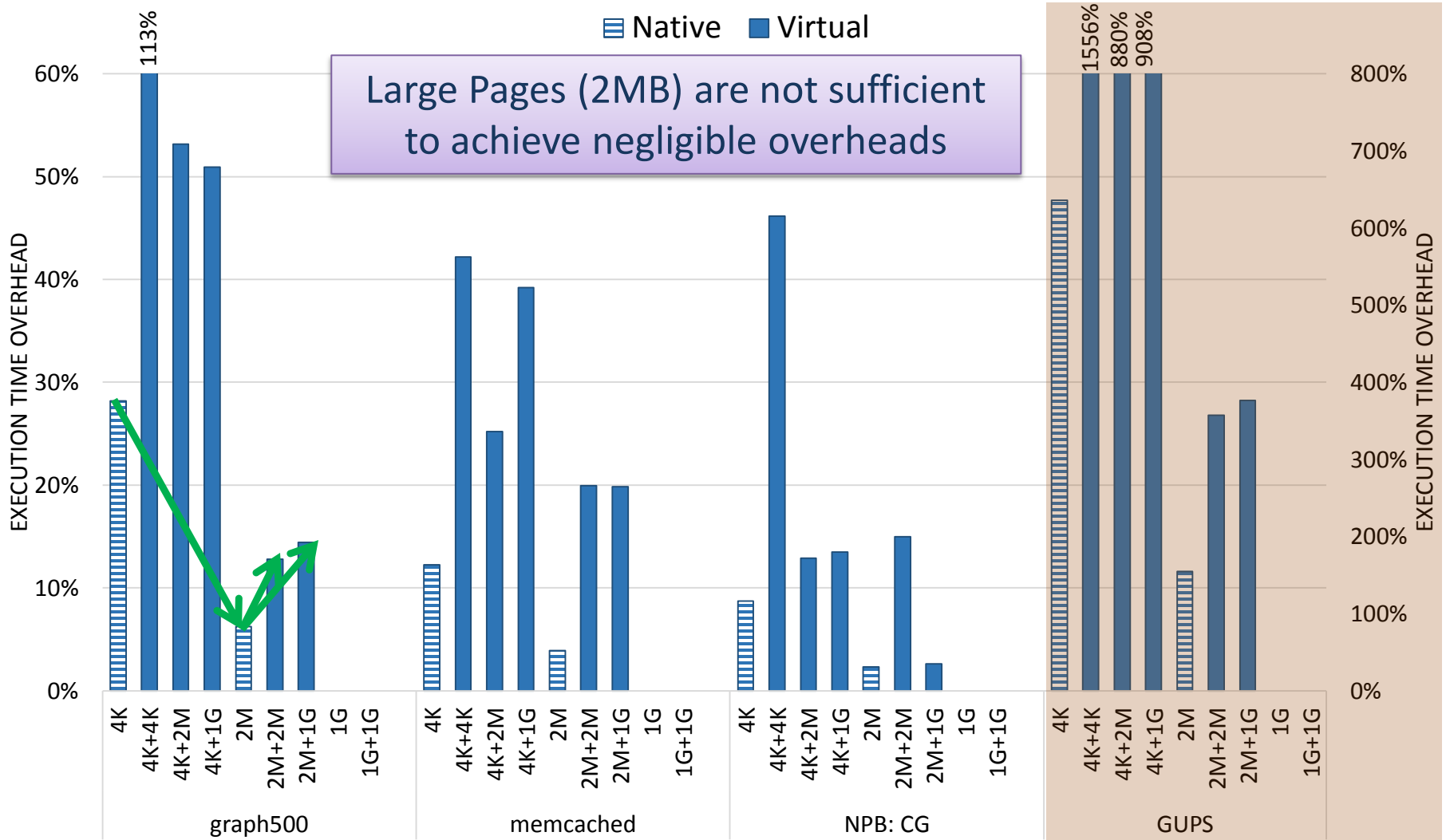
Cost of Virtualization



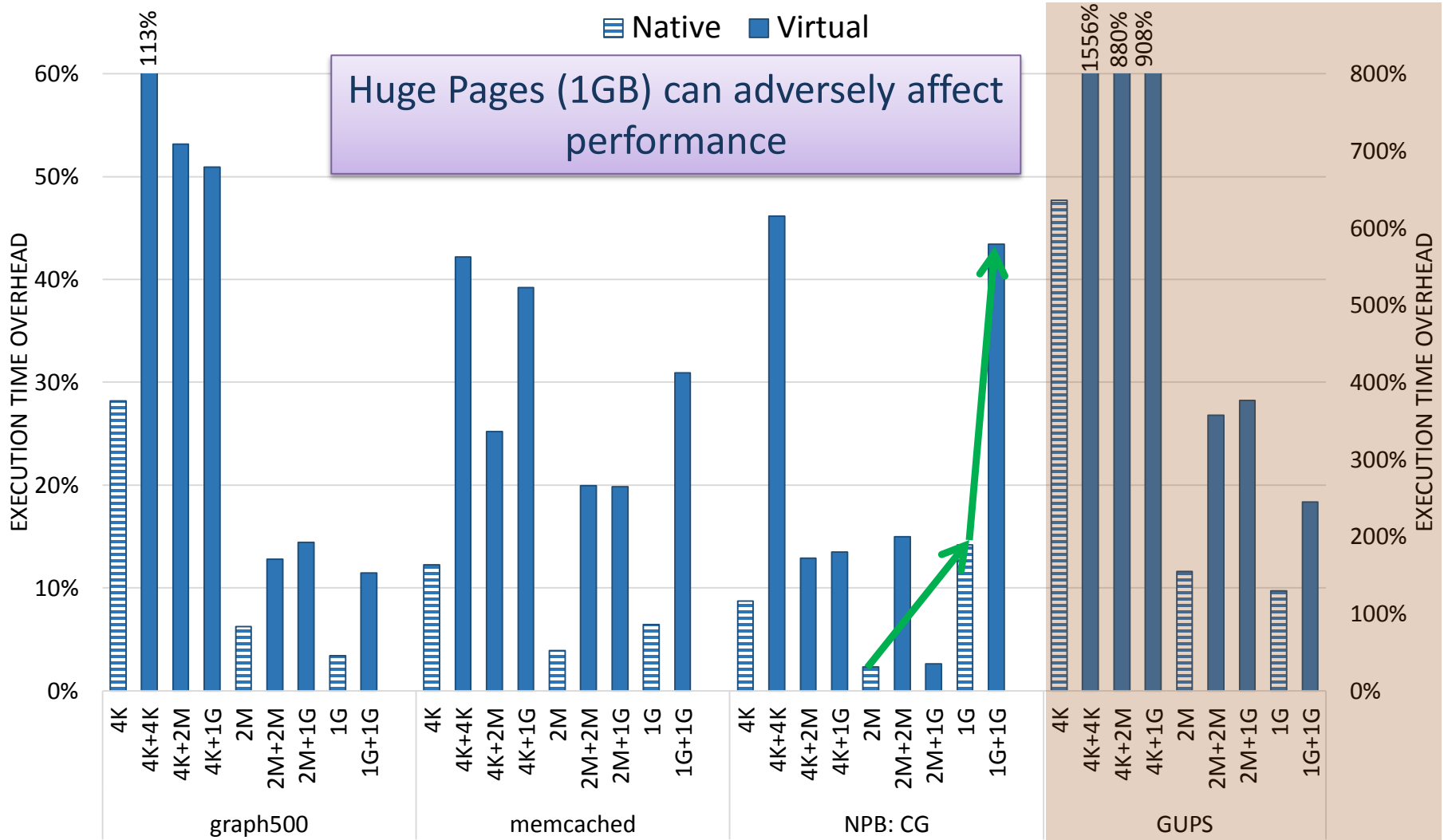
Cost of Virtualization



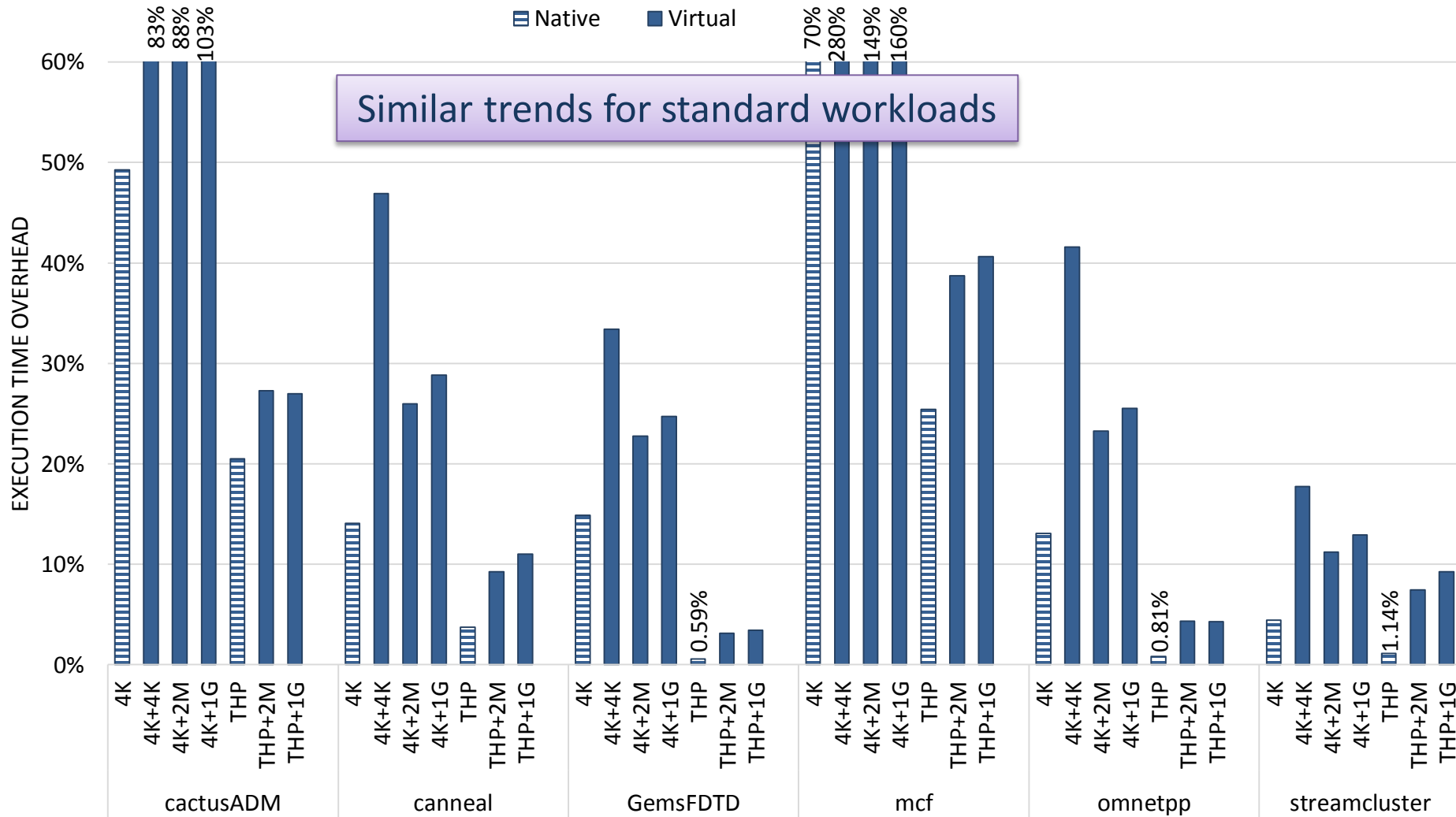
Cost of Virtualization



Cost of Virtualization



Cost of Virtualization



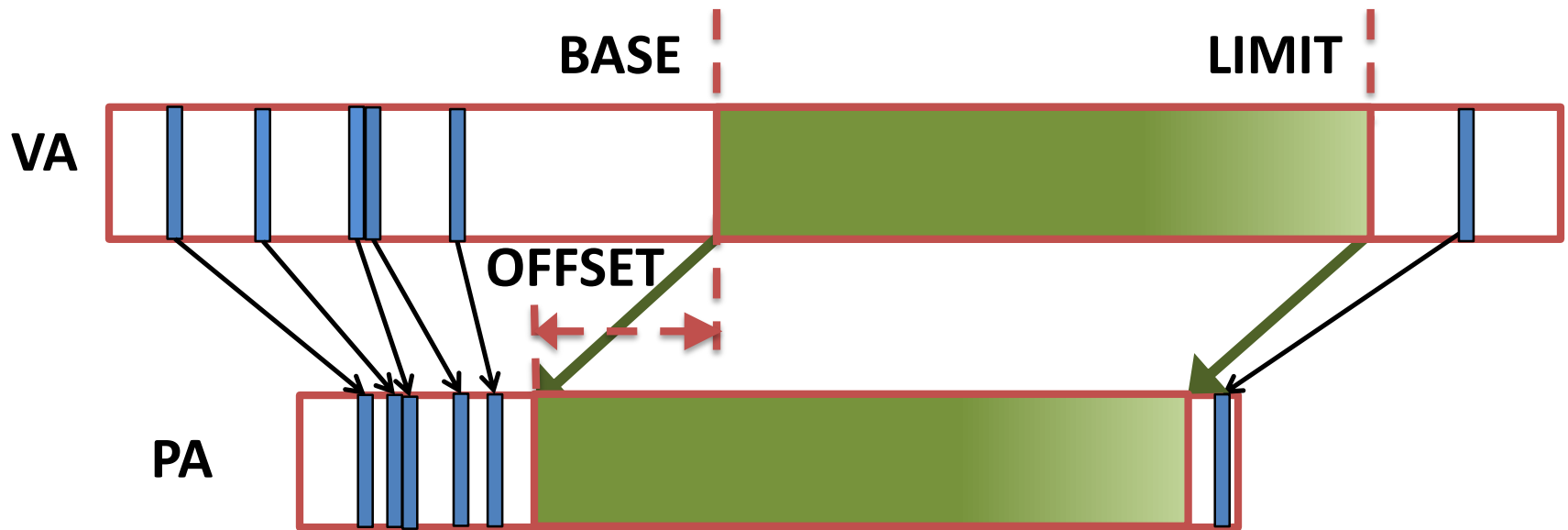
Outline

- Motivation
- Review: Direct Segments ←
- Virtualized Direct Segments
- Optimizations
- Evaluation
 - Methodology
 - Results
- Summary

Unvirtualized Direct Segments

1 Conventional Paging

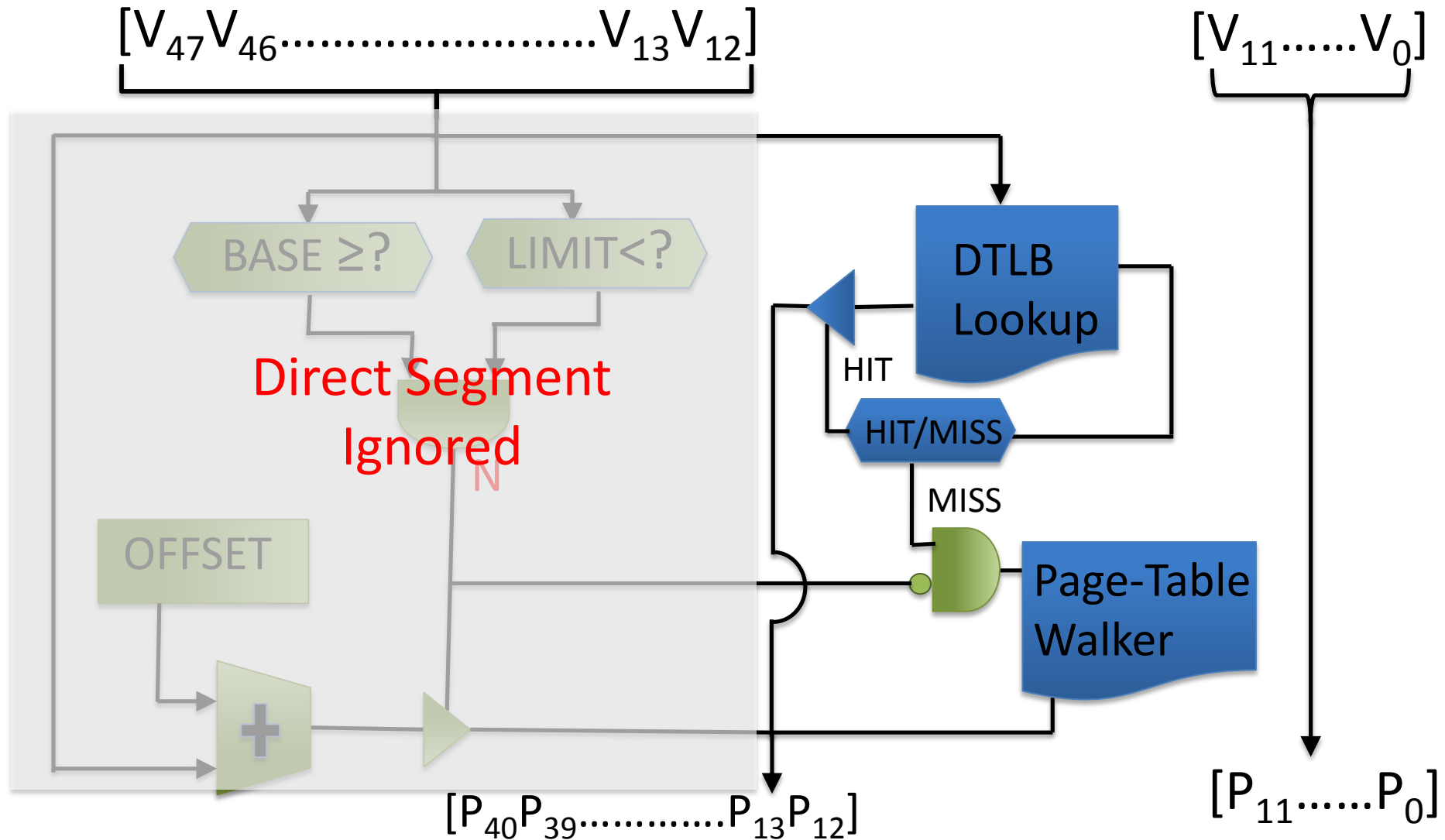
2 Direct Segment



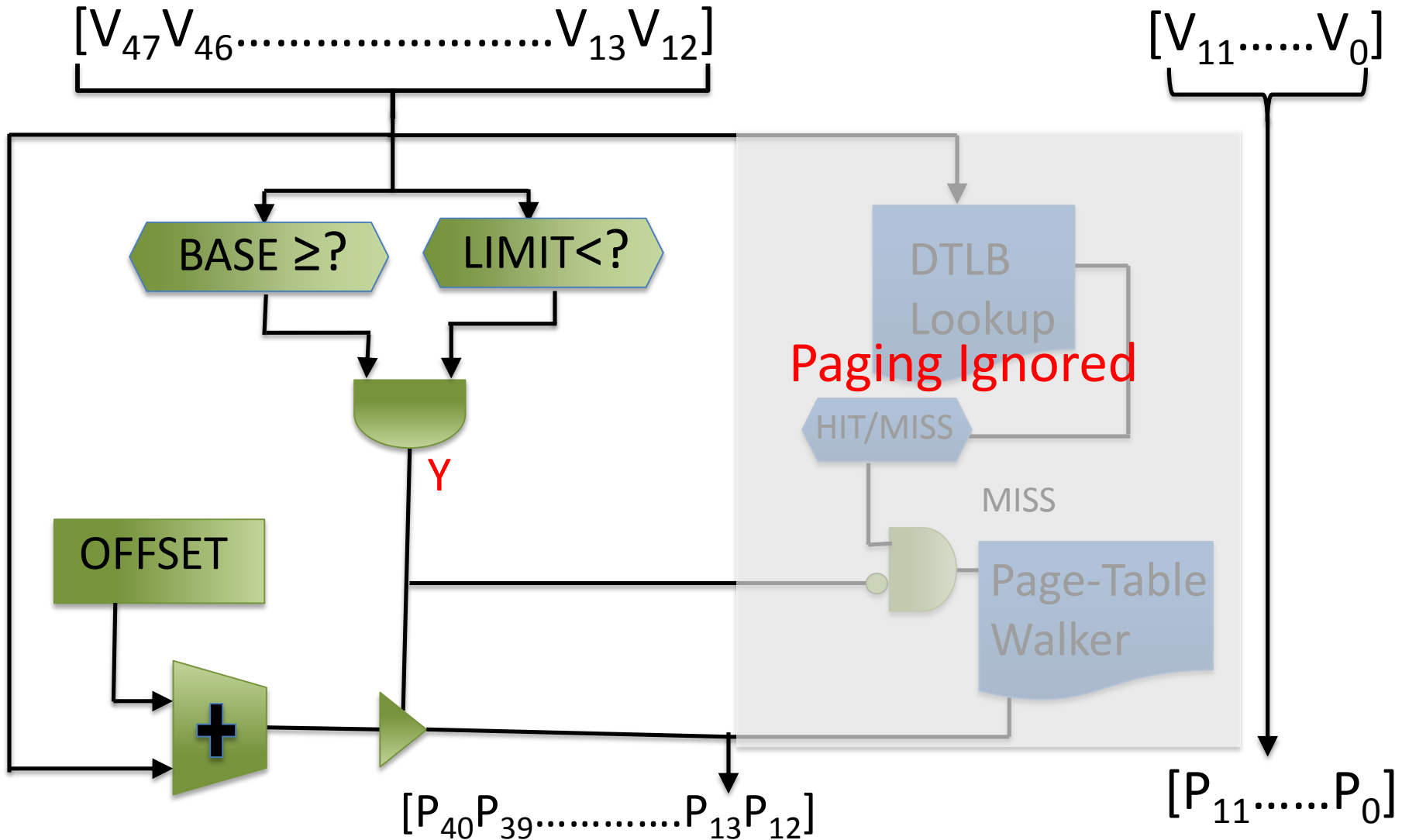
Why Direct Segment?

- Matches big memory workload needs
- NO TLB lookups => NO TLB Misses

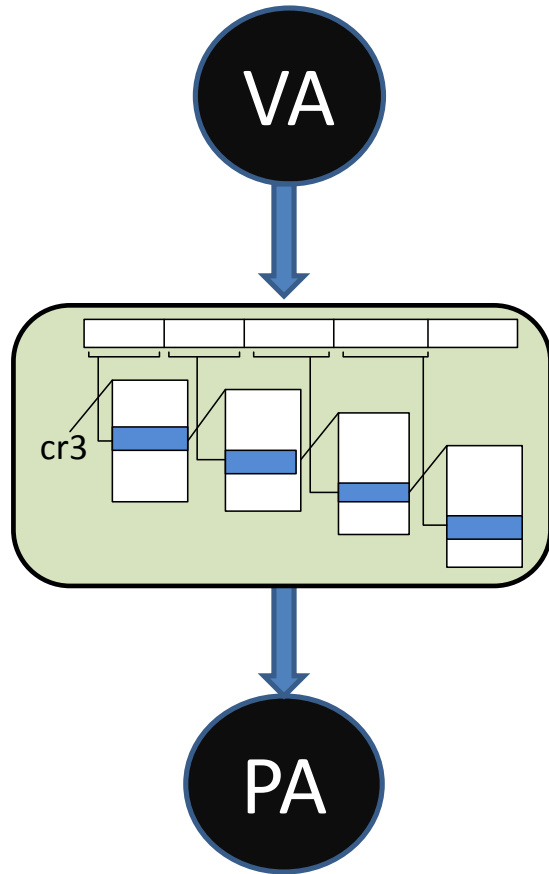
Translation with Direct Segments



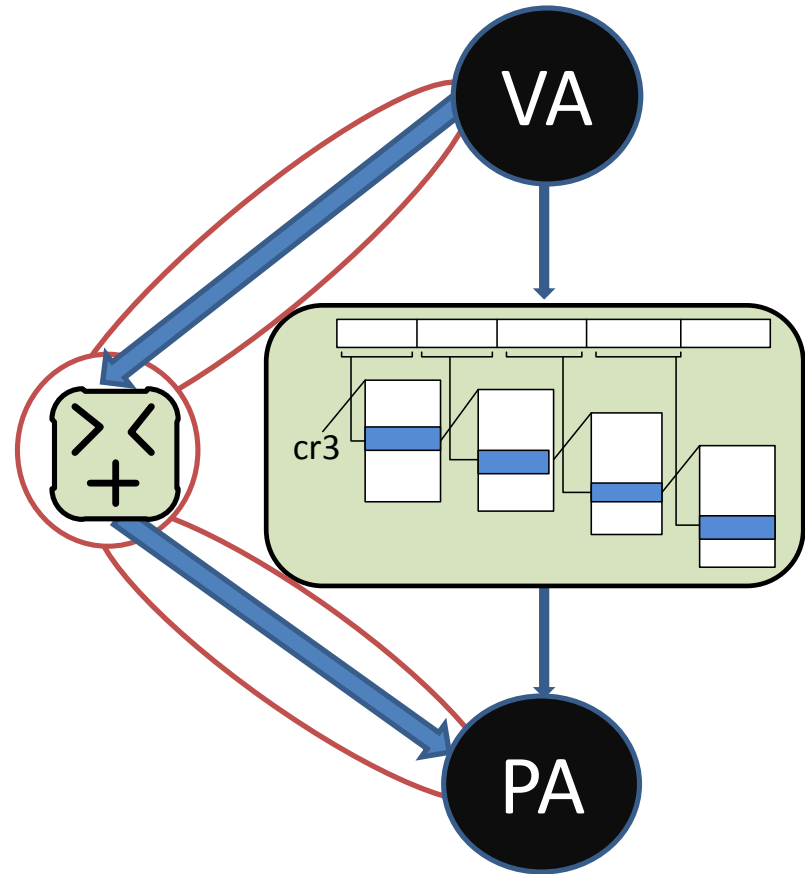
Translation with Direct Segments



Direct Segments



Base Native

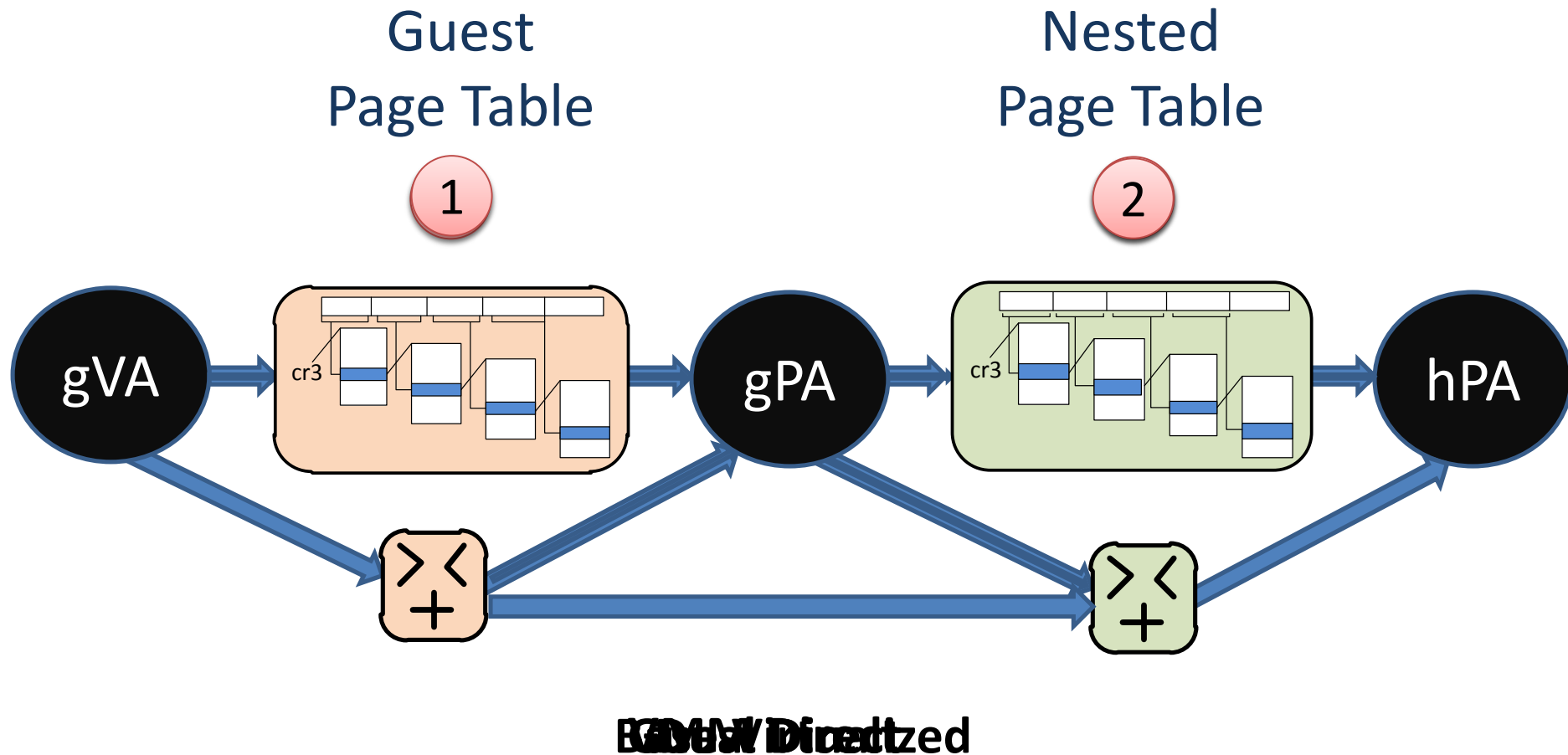


Direct Segments

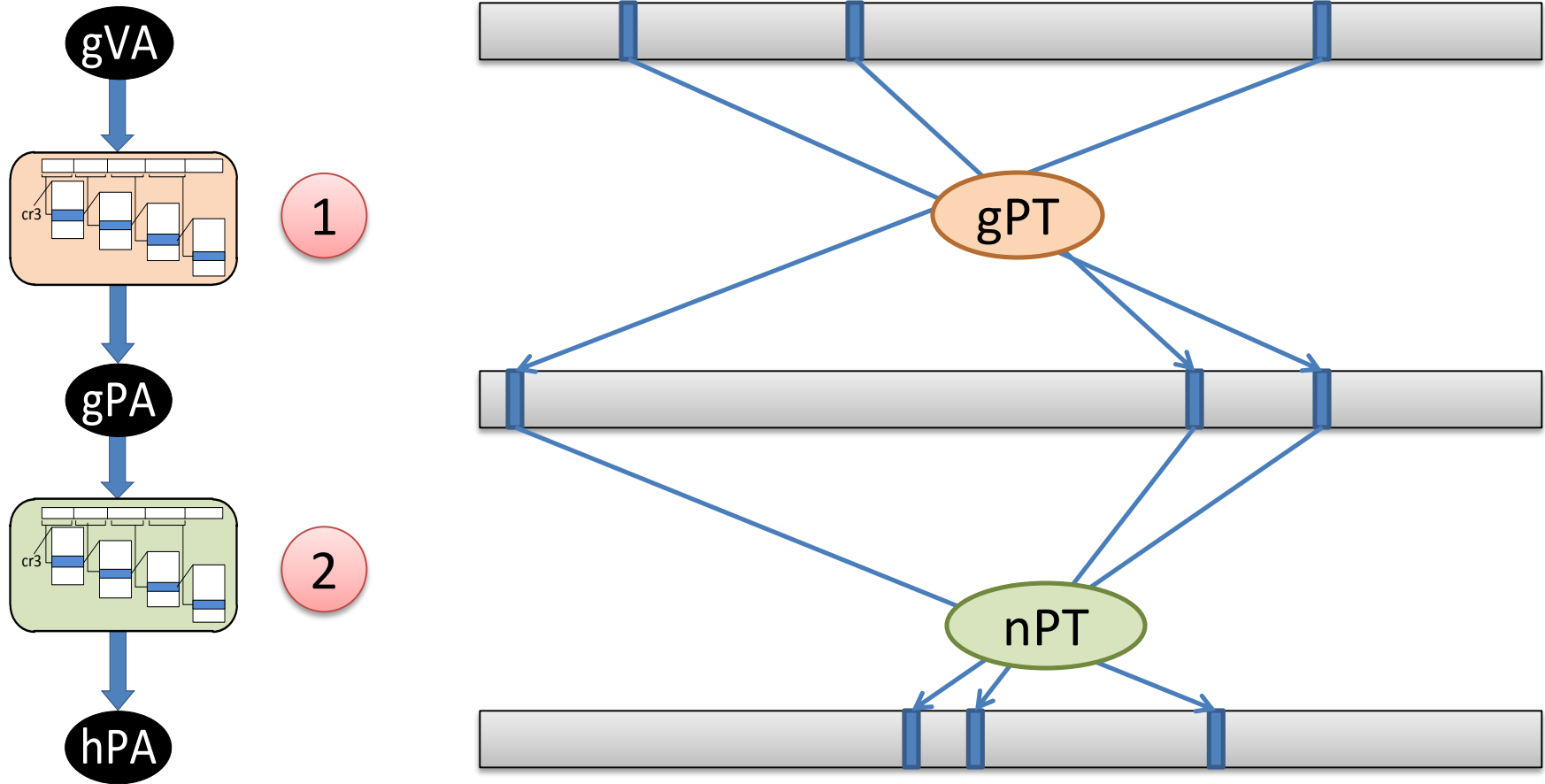
Outline

- Motivation
- Review: Direct Segments
- Virtualized Direct Segments ←
- Evaluation
 - Methodology
 - Results
- Summary

Modes

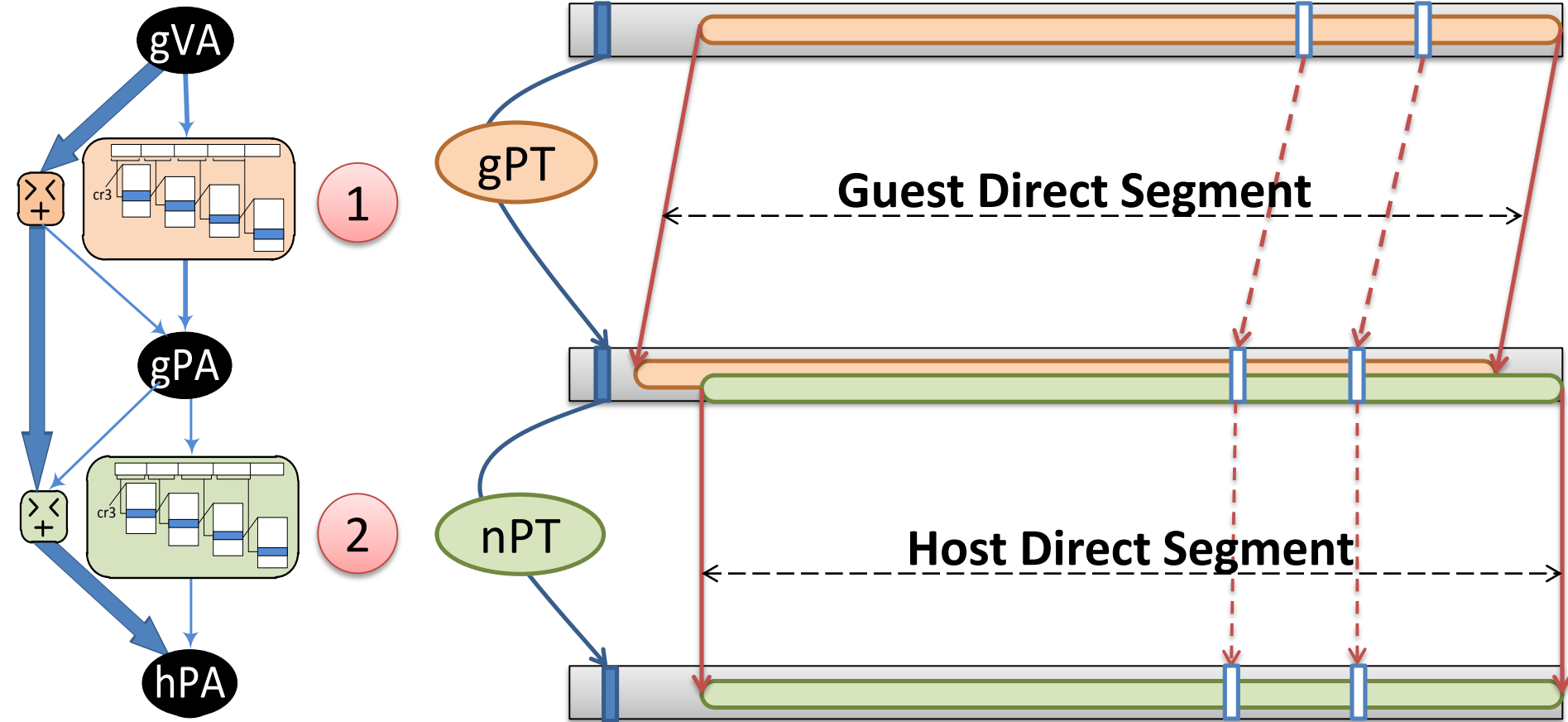


Base Virtualized: Translation



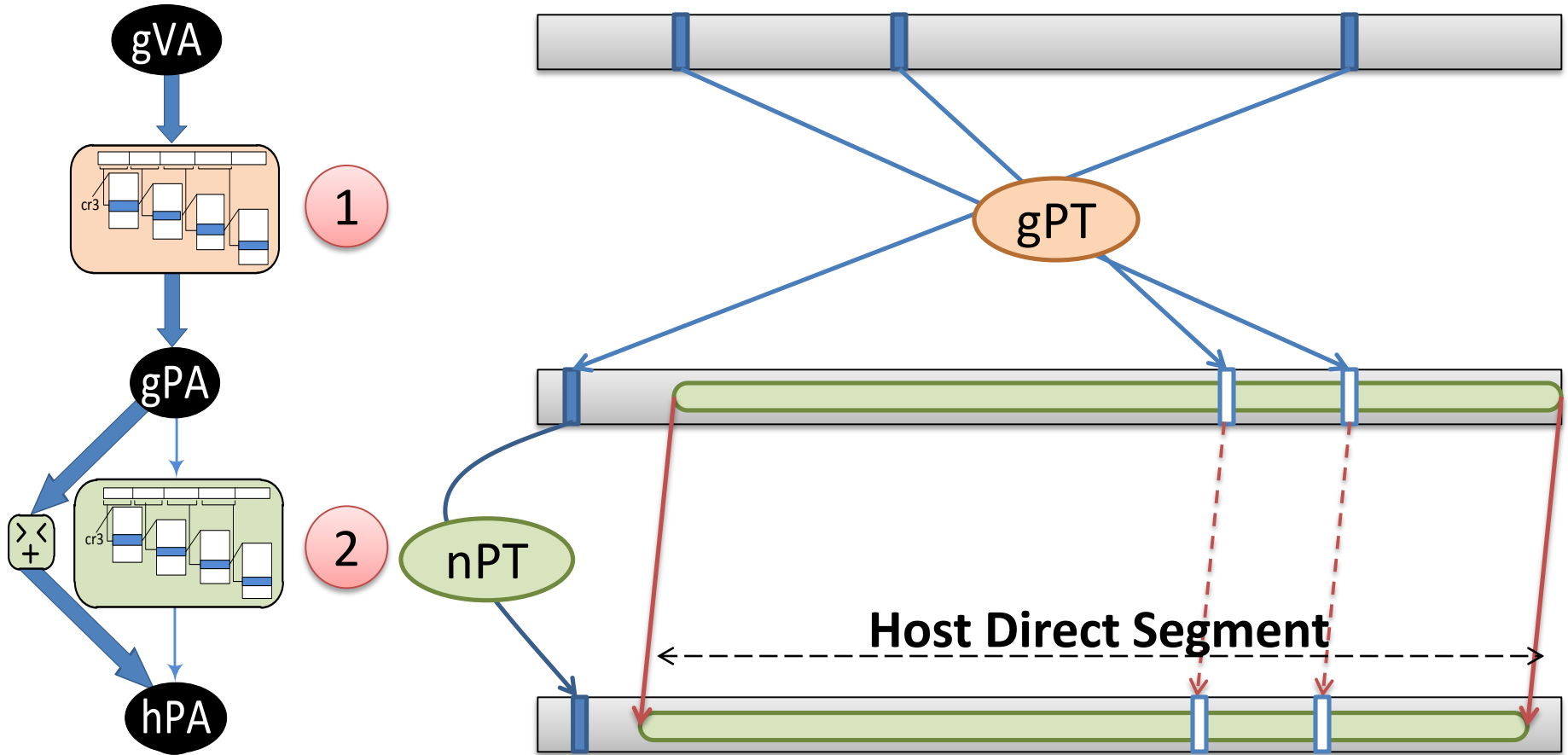
Base Virtualized

Dual Direct: Translation



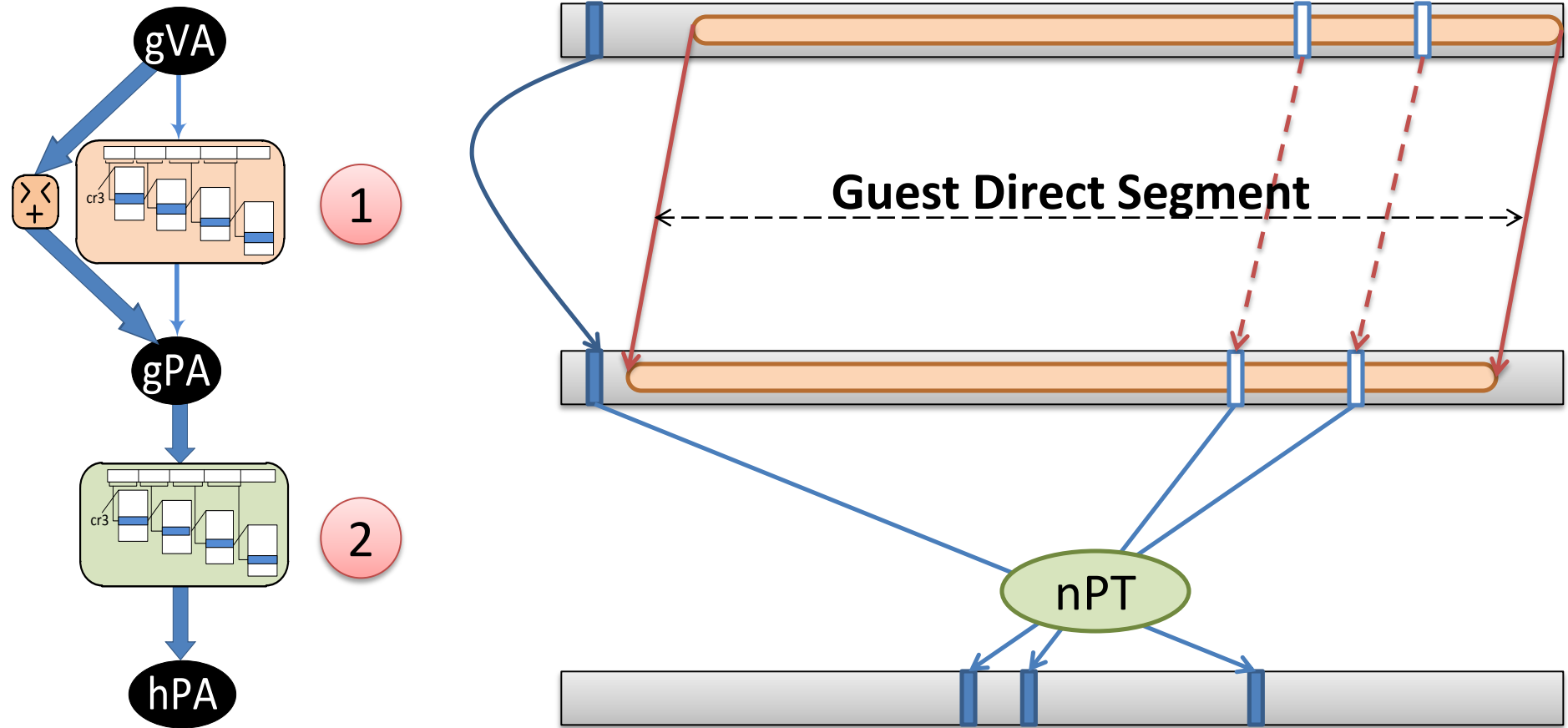
Dual Direct

VMM Direct: Translation



VMM Direct

Guest Direct: Translation



Tradeoffs: Efficiency

Properties	Base Virtualized	Dual Direct	VMM Direct	Guest Direct
Dimension of page walk	2D	0D	1D	1D
# of mem. accesses for most page walks	24	0	4	4
# of base-bound computation for page walks	0	1	5	1

Tradeoffs: Compatibility

Properties	Base Virtualized	Dual Direct	VMM Direct	Guest Direct
Guest OS modifications	✗	✓	✗	✓
VMM modifications	✗	✓	✓	minimal
Application Category	Any	Big-memory	Any	Big-memory

Tradeoffs: Memory Overcommit

Properties	Base Virtualized	Dual Direct	VMM Direct	Guest Direct
Page Sharing	✓	limited	limited	✓
Ballooning	✓	limited	✓	limited
Guest OS Swapping	✓	limited	✓	limited
VMM Swapping	✓	limited	limited	✓

Outline

- Motivation
- Review: Direct Segments
- Virtualized Direct Segments
- Optimizations ←
- Evaluation
 - Methodology
 - Results
- Summary

Optimizations

- **Issue 1:** Guest/Host memory fragmentation
 - Guest physical memory: Self-ballooning
 - Host physical memory: *Compaction*
- **Issue 2:** Permanent “hard” memory faults
 - Escape filter to provide alternate translation
 - Bloom filter stores small number of faulty pages
 - Filter checked in parallel with segment registers

For more details: Come to the poster

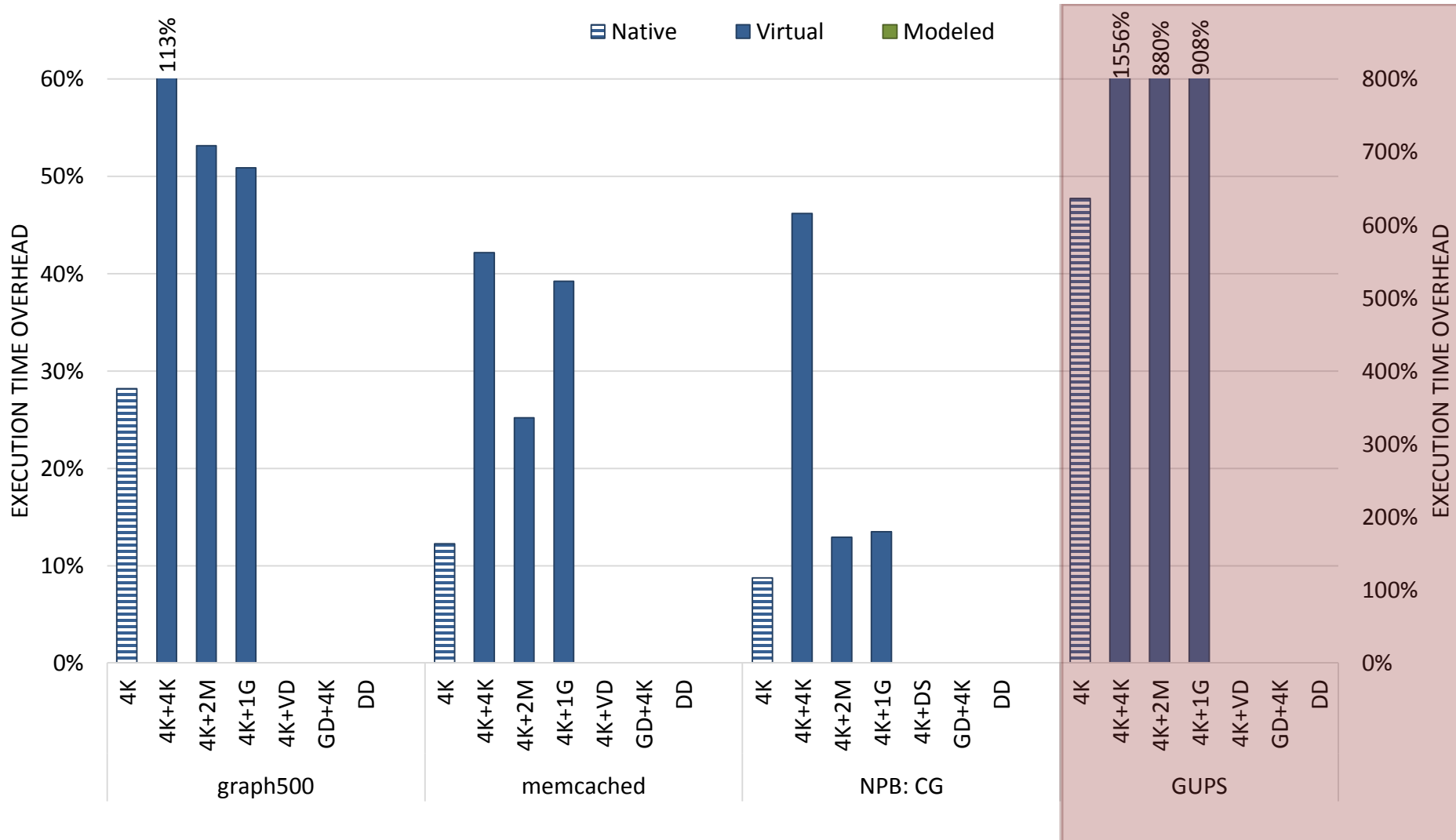
Outline

- Motivation
- Review: Direct Segments
- Virtualized Direct Segments
- Optimizations
- Evaluation ←
 - Methodology
 - Results
- Summary

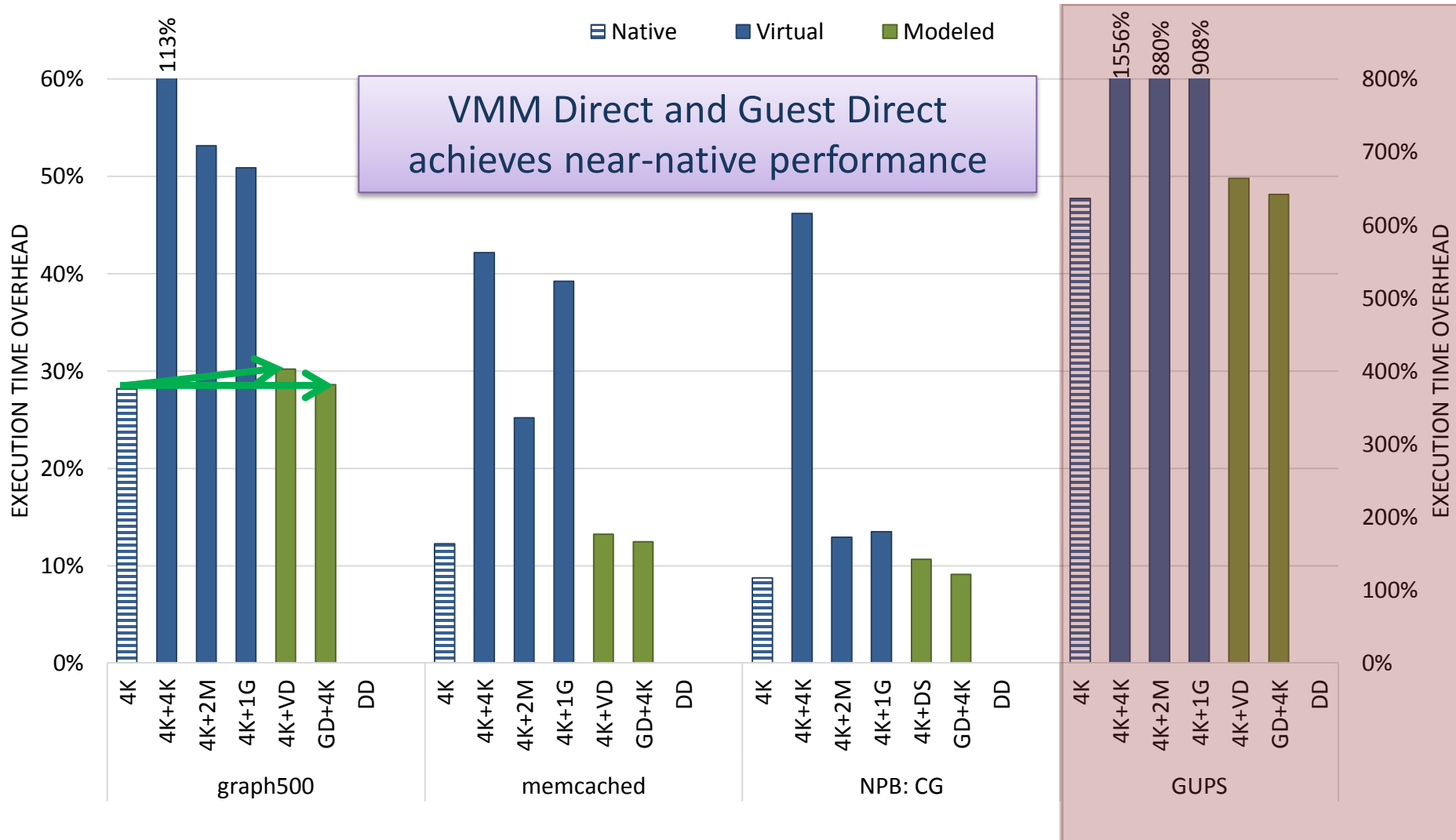
Methodology

- Estimate performance of future hardware
 - Measure fraction of TLB misses to segmented memory
 - Measure TLB miss cost with performance counters
 - Estimate performance gain with linear model
- Prototype
 - Linux v3.12.13 host/guest
 - Qemu-KVM hypervisor
- Intel 12-core Sandy-bridge with 96GB memory

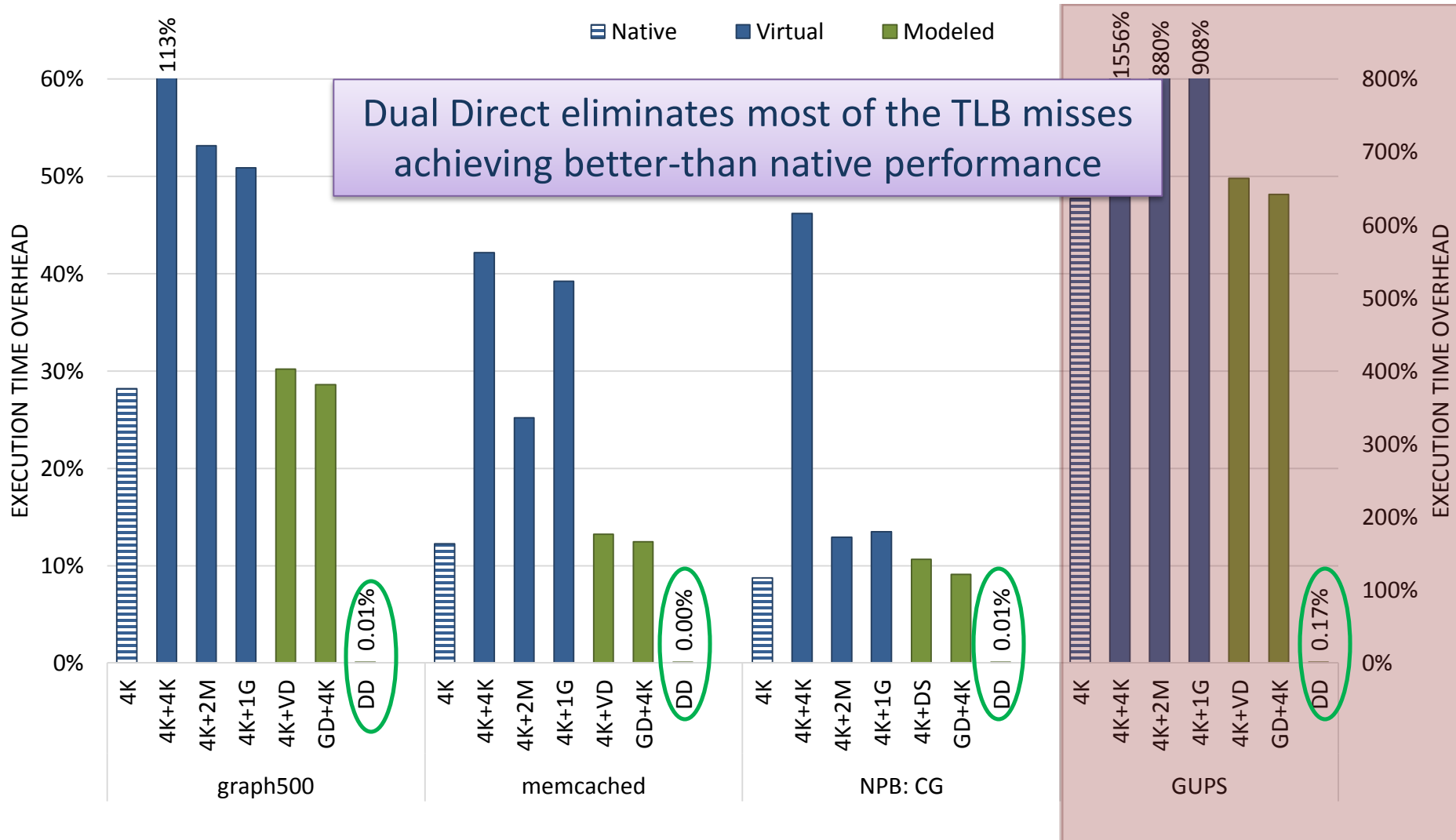
Results



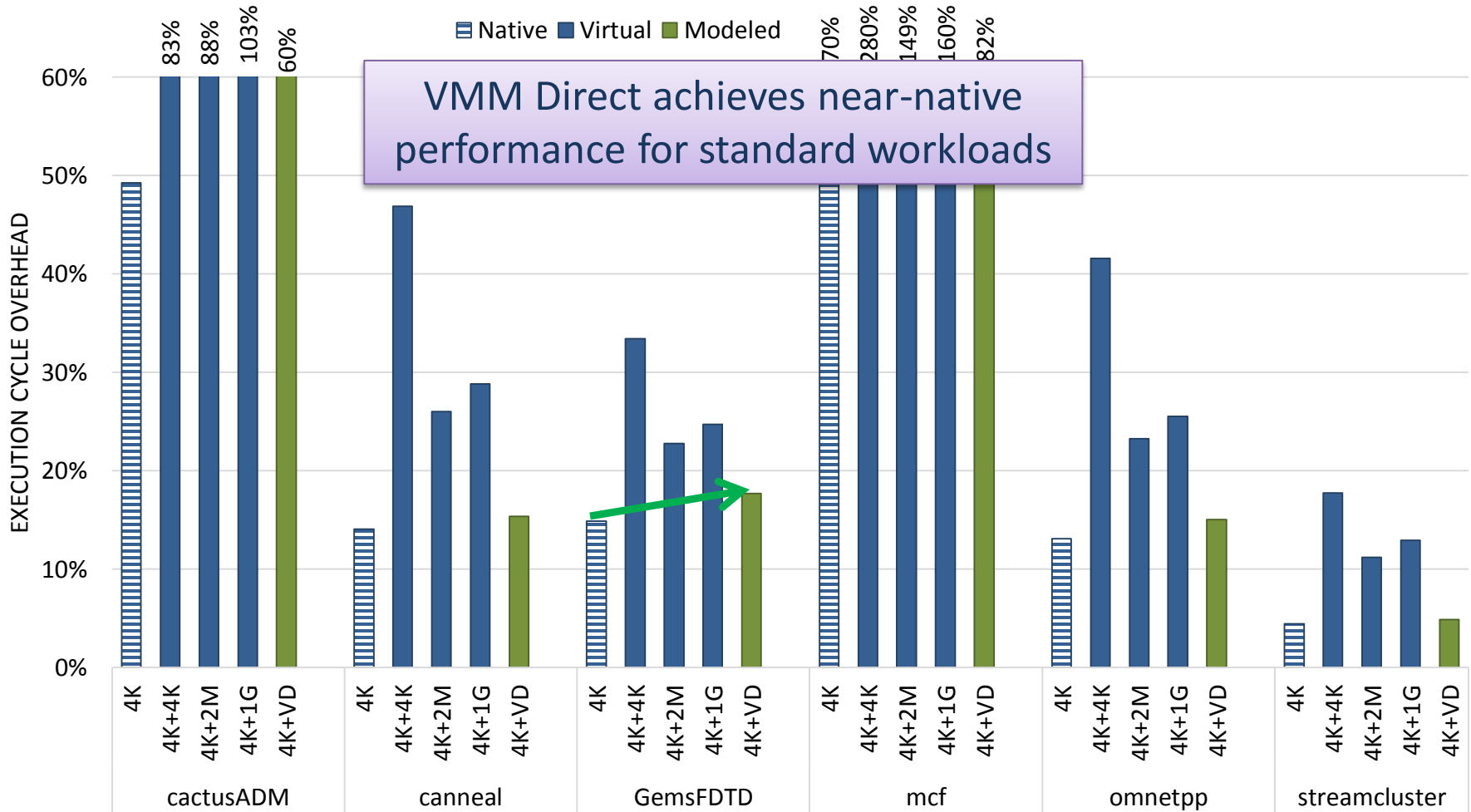
Results



Results



Results



Summary

- Problem: TLB misses in virtual machines
 - Hardware-virtualized MMU has high overheads
 - Up to 280% overhead
- Prior Work: **Direct Segments** – unvirtualized case
- Solution: segmentation to bypass paging
 - **Extend Direct Segments** for virtualization
 - **Three configurations** with different tradeoffs
- Results
 - **Near- or better-than-native** performance

Questions ?



For more details: Come to the poster