



Memory Harvesting in Multi-GPU Systems with Hierarchical Unified Virtual Memory

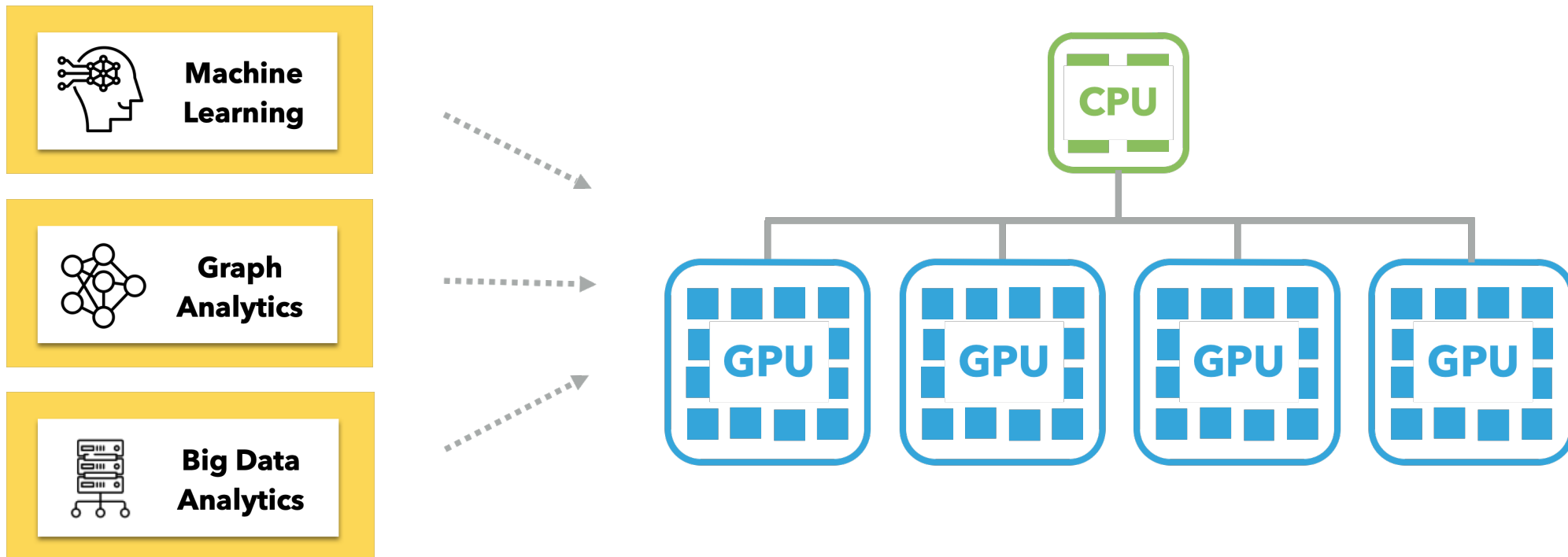
Sangjin Choi*, Taeksoo Kim*, Jinwoo Jeong, Rachata Ausavarungnirun, Myeongjae Jeon, Youngjin Kwon, Jeongseob Ahn



*Co-first author

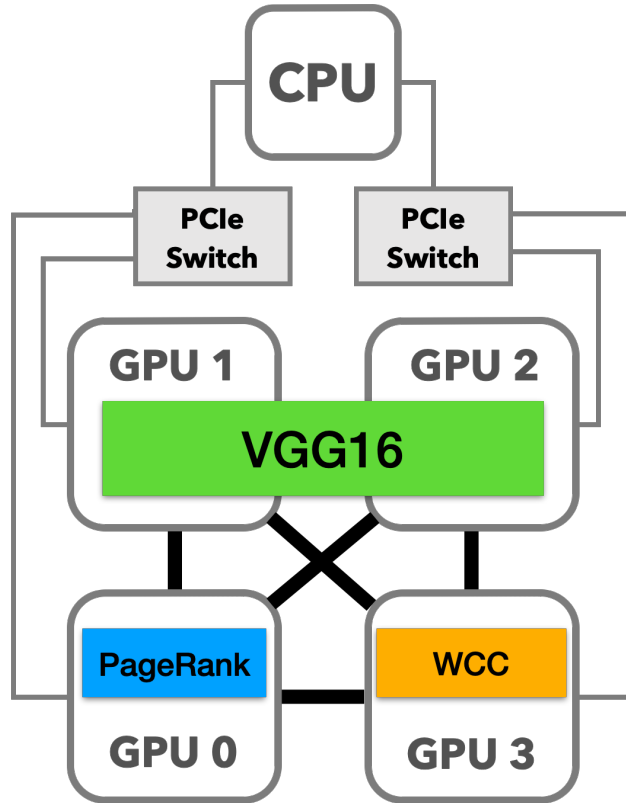
Shared multi-GPU environment

- Shared multi-GPU servers to reduce the cost of infrastructure

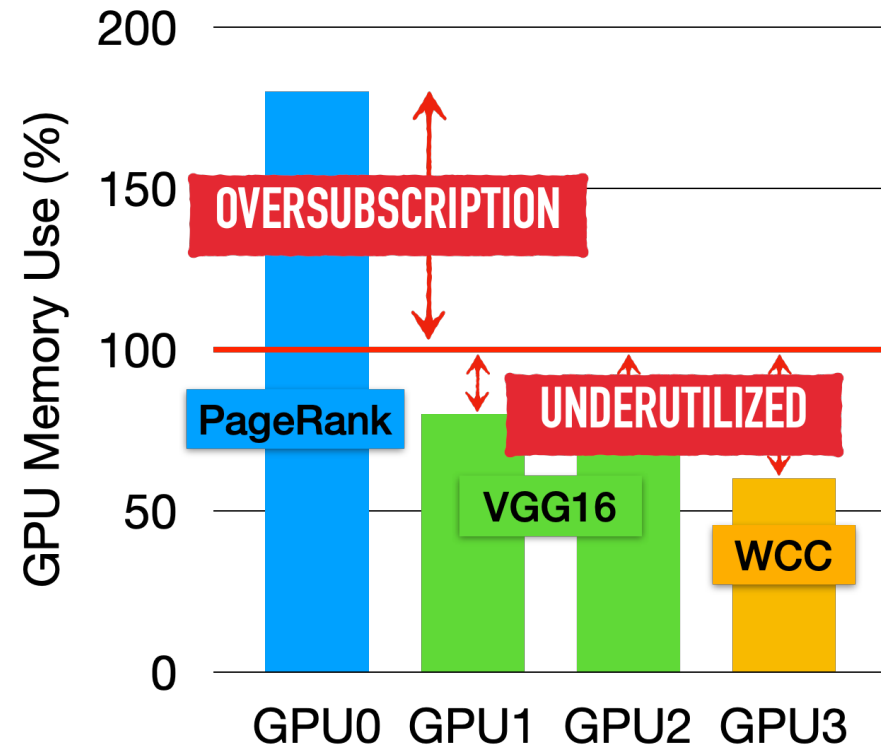


Workloads are consolidated to improve resource utilization

Memory space across GPUs is not fully utilized

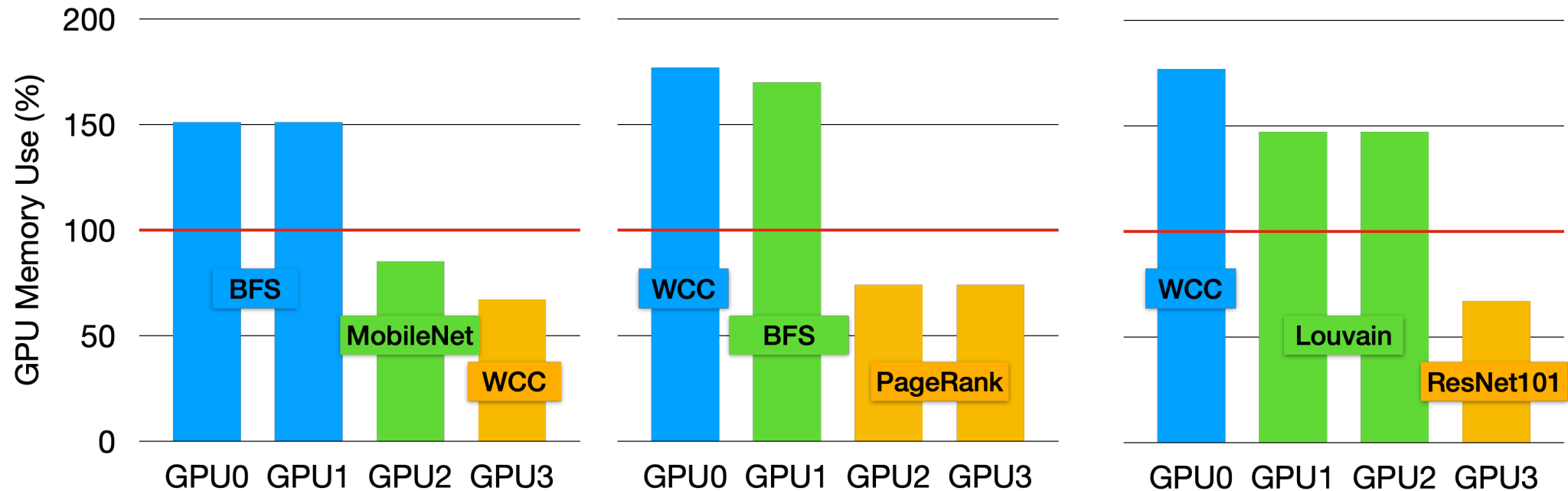


AWS p3.8xlarge instance



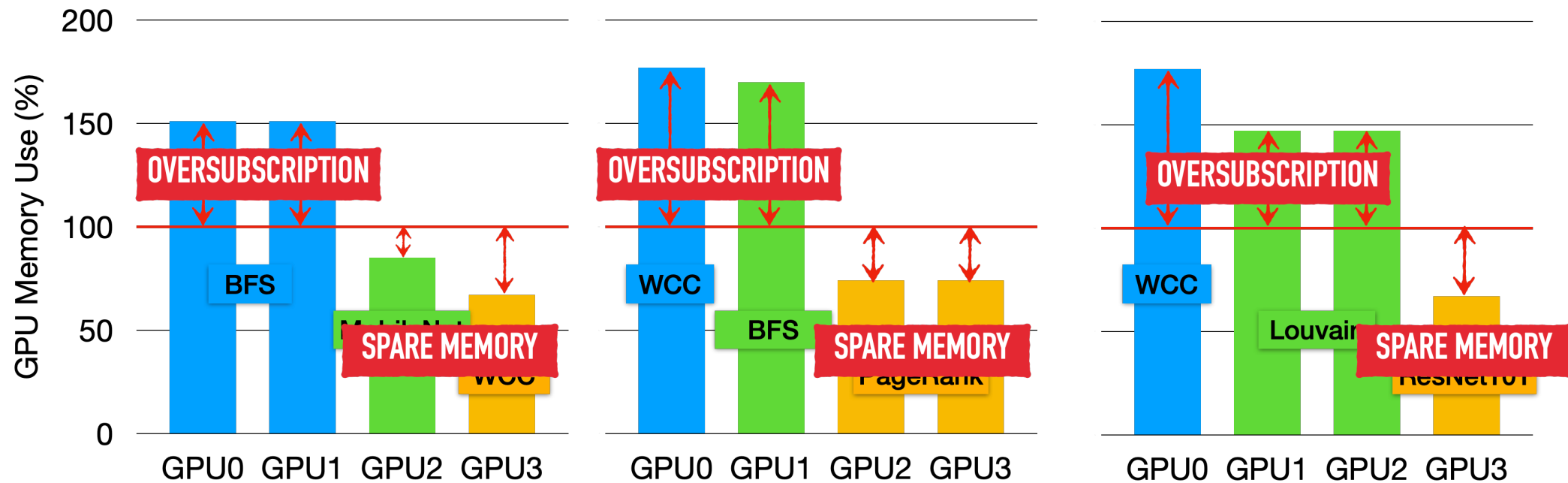
Memory space across GPUs is not fully utilized

- Each workload has highly varying memory demands
- GPU memory is not fully utilized in shared GPU clusters (e.g., Microsoft GPU cluster^[1])



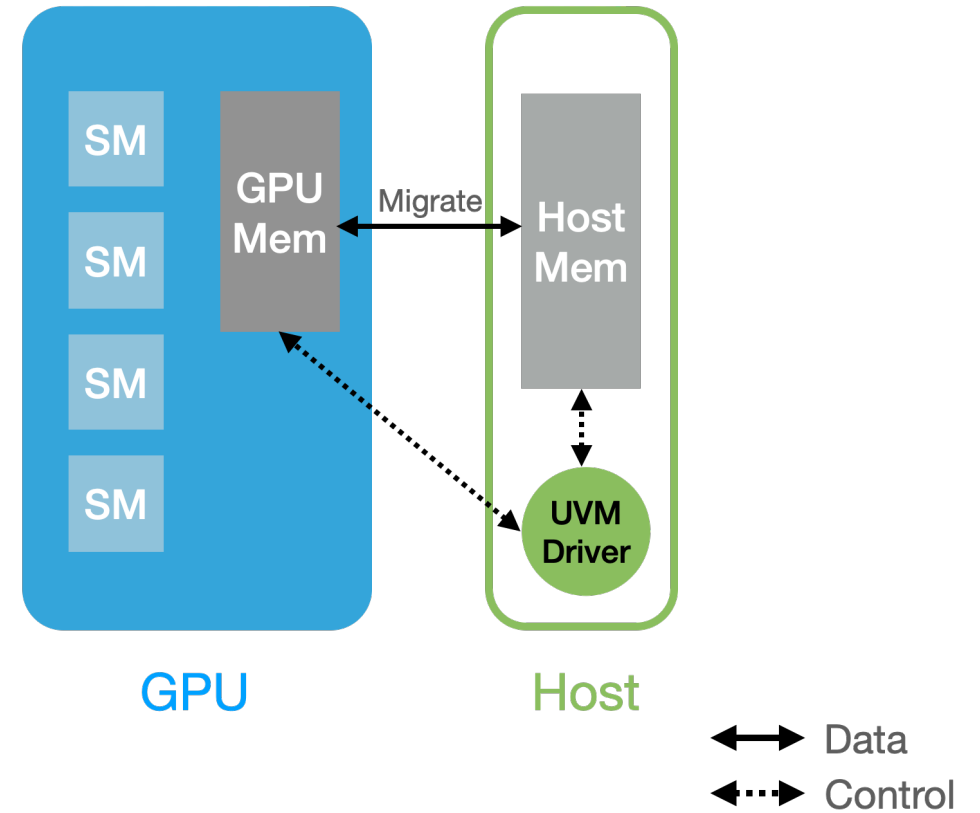
Memory space across GPUs is not fully utilized

- Each workload has highly varying memory demands
- GPU memory is not fully utilized in shared GPU clusters (e.g., Microsoft GPU cluster^[1])



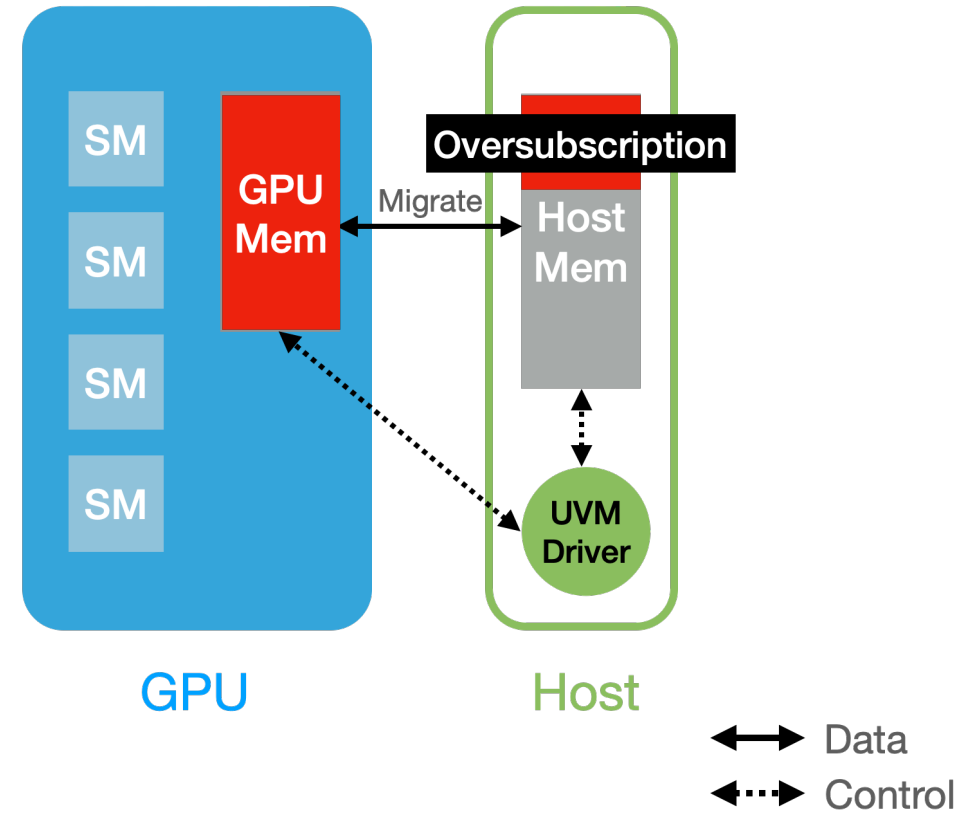
Performance overhead for memory oversubscription

- Unified Virtual Memory (UVM)
 - Swapping GPU memory to host memory

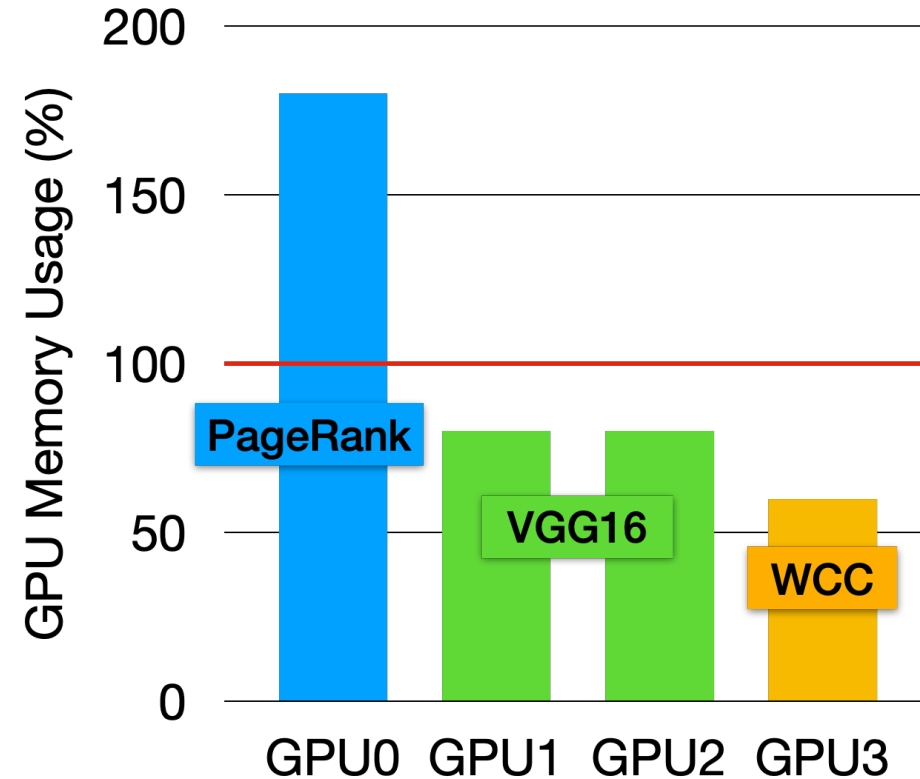


Performance overhead for memory oversubscription

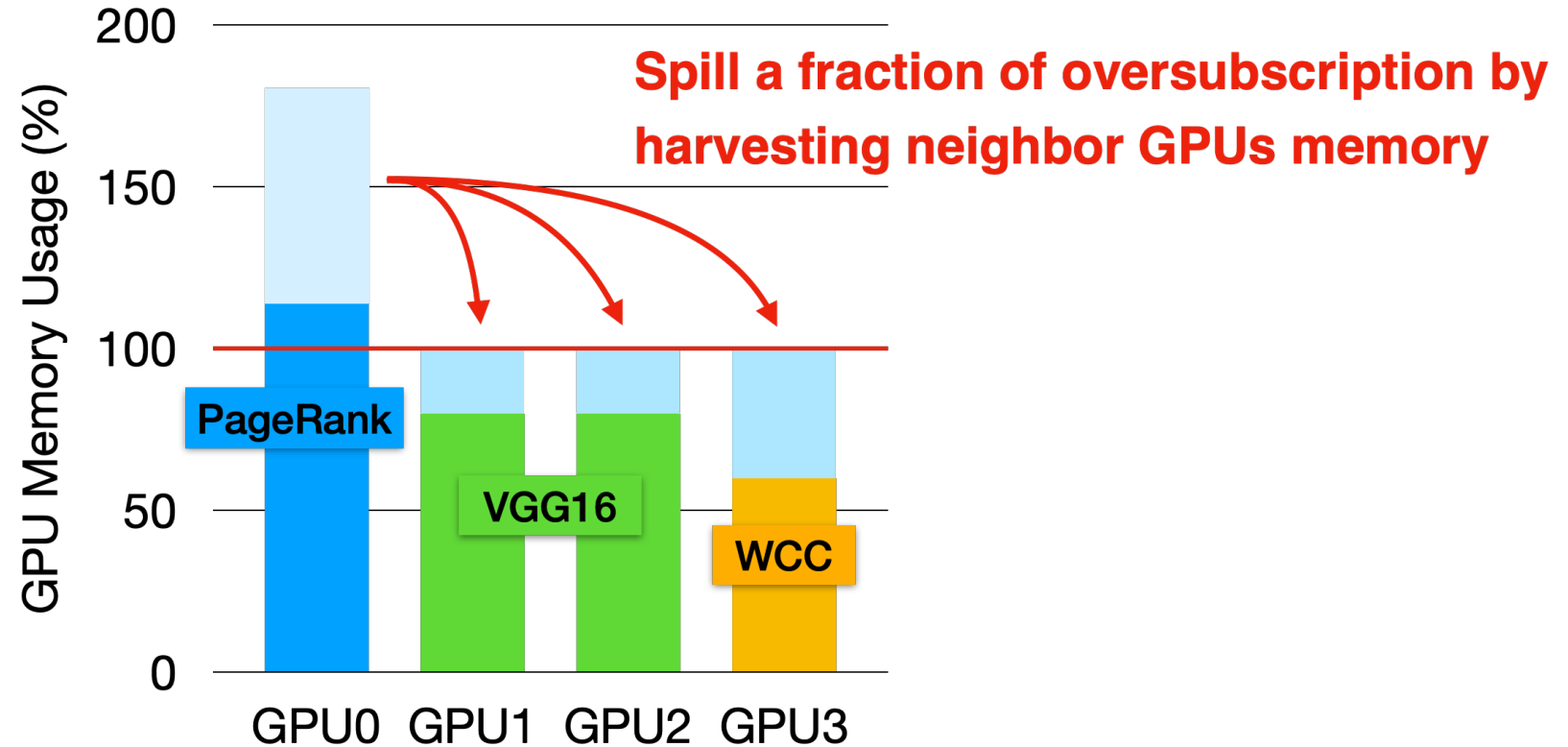
- Unified Virtual Memory (UVM)
 - Swapping GPU memory to host memory
- Cost for memory oversubscription is high
 - **2x ~ 64x** longer to complete with 40% oversubscription^[1]



Our approach: Harvesting neighbor GPU memory

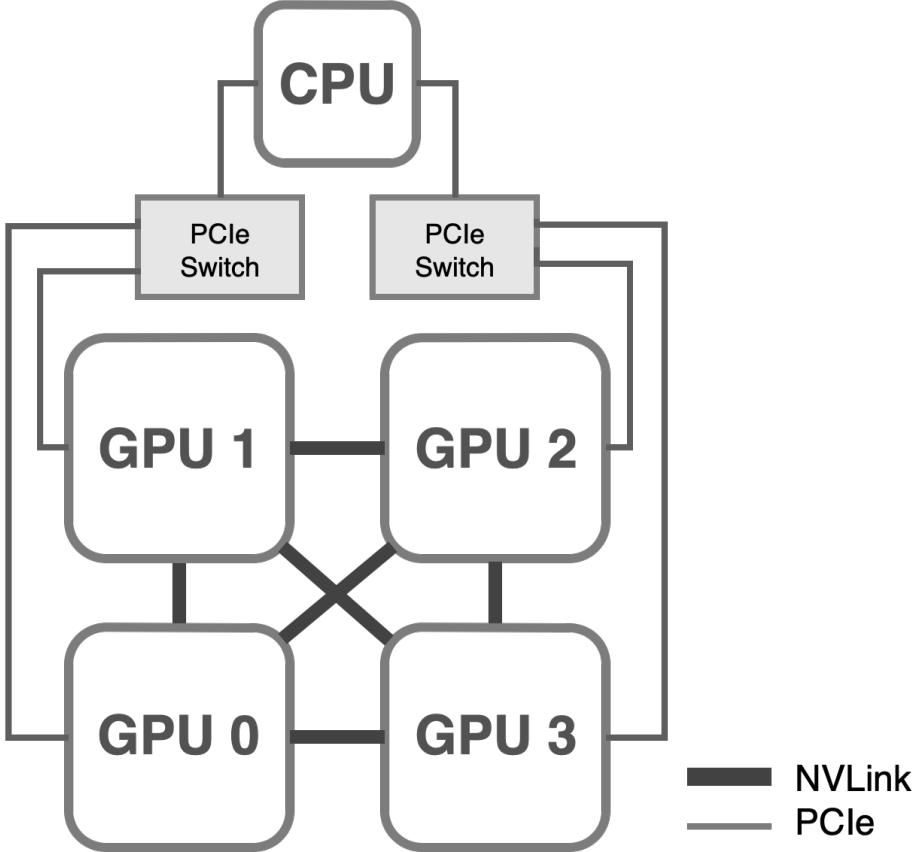


Our approach: Harvesting neighbor GPU memory



How to efficiently access neighbor GPU memory?

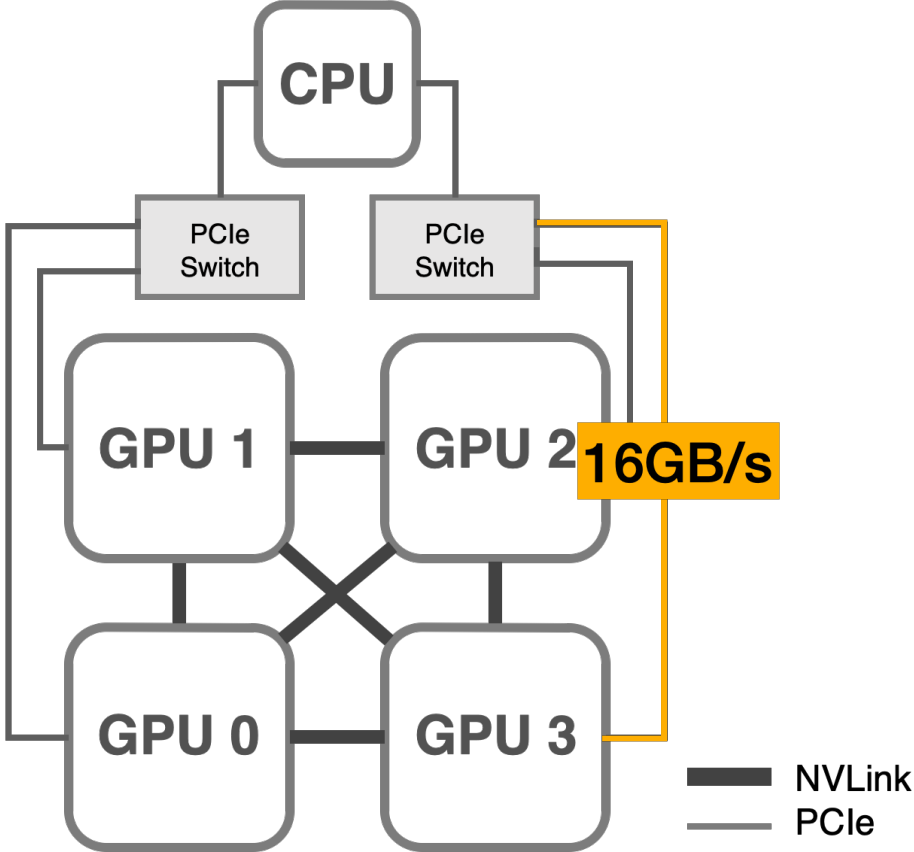
NVLink provides fast interconnect to neighbor GPU



	PCIe	NVLink
Throughput (GB/s)	12.3	40.1 (3.3x)
Latency (μ s)	16.7	5.1 (3.2x)

2MB Copy

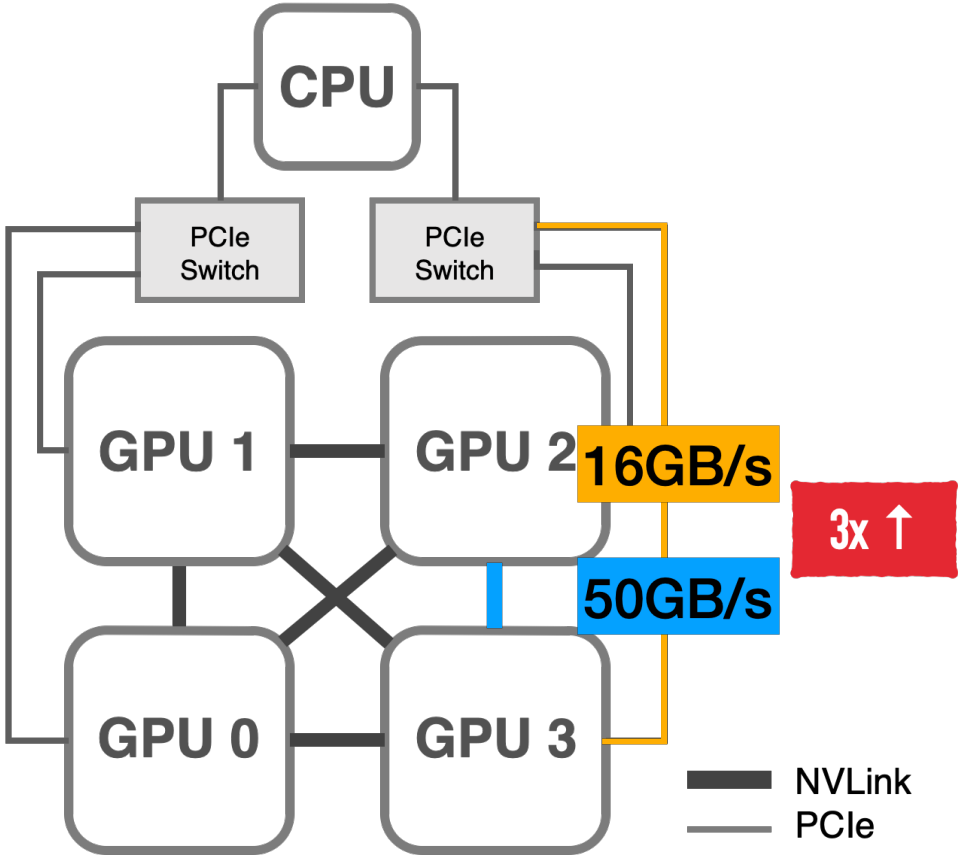
NVLink provides fast interconnect to neighbor GPU



	PCIe	NVLink
Throughput (GB/s)	12.3	40.1 (3.3x)
Latency (μs)	16.7	5.1 (3.2x)

2MB Copy

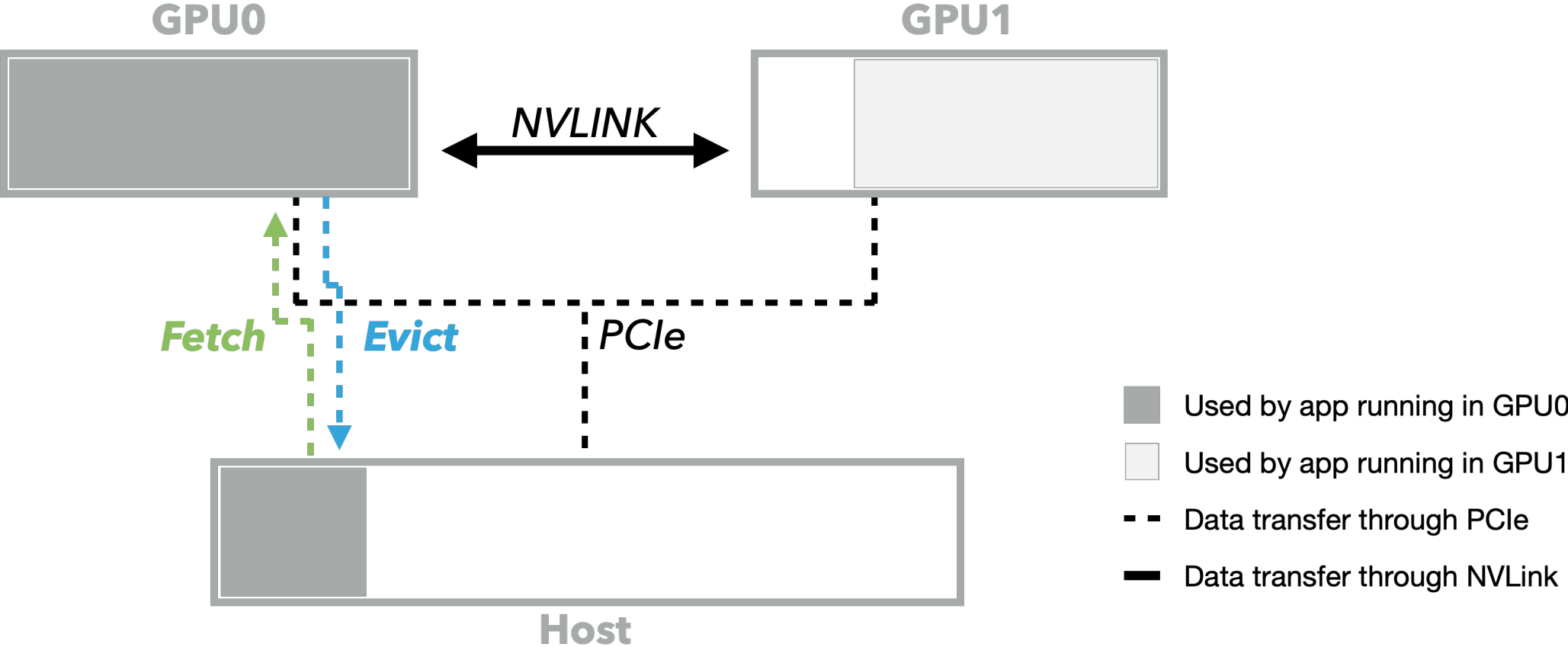
NVLink provides fast interconnect to neighbor GPU



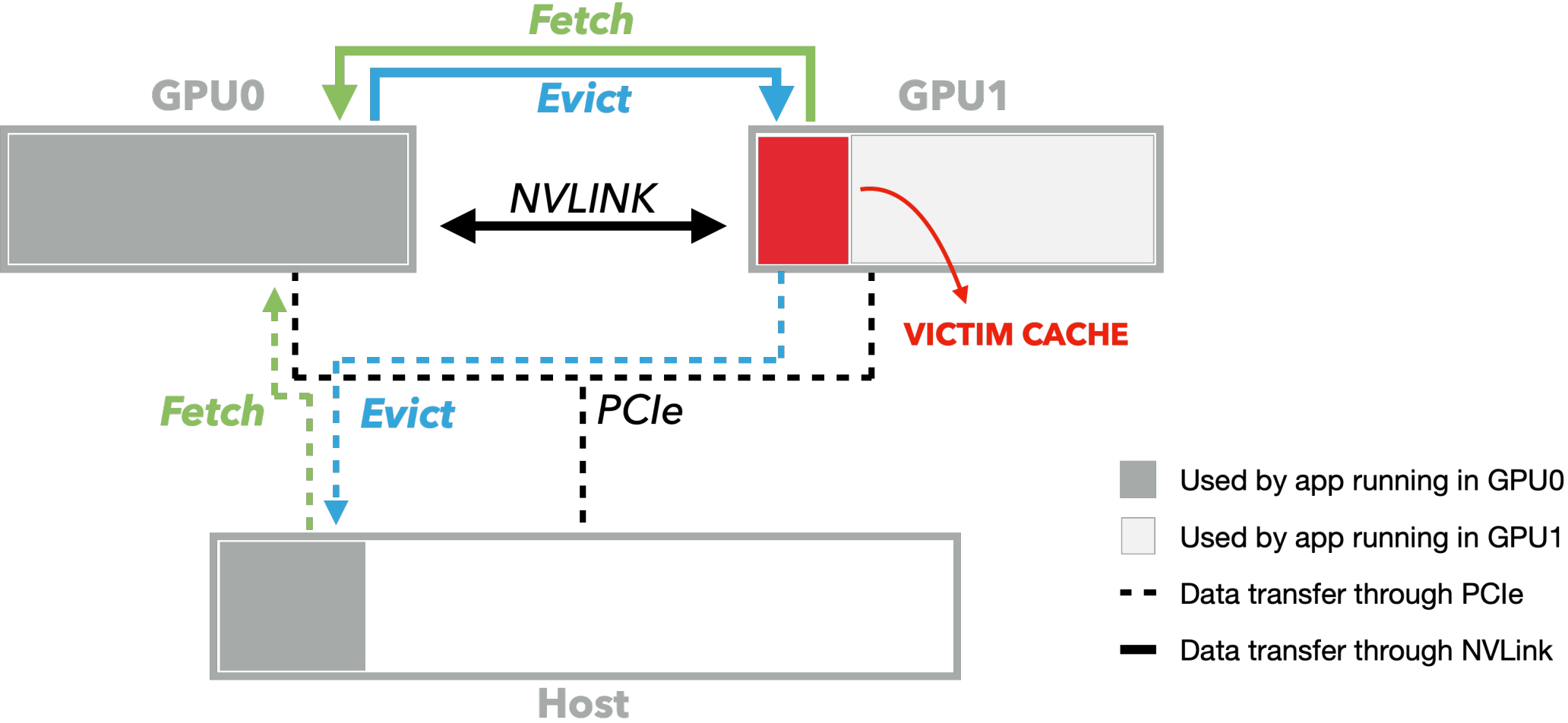
	PCIe	NVLink
Throughput (GB/s)	12.3	40.1 (3.3x)
Latency (μ s)	16.7	5.1 (3.2x)

2MB Copy

Unified Virtual Memory



Data-path: Hierarchical Unified Virtual Memory



Goals of Hierarchical Unified Virtual Memory

Reducing performance cost of memory oversubscription

Effective Harvesting

Harvest small and temporarily available spare memory of neighbor GPU

Reduce eviction/fetch latency with spare memory

Minimal Interference

Minimize performance impact of workloads running in neighbor GPU

Framework-agnostic

No modification of applications or frameworks

Memory manager for HUVIM: memHarvester

Centralized coordinator for data-path in HUVIM

1. Hiding eviction latency

2. Reducing fetch latency

3. Hiding fetch latency

Effective Harvesting

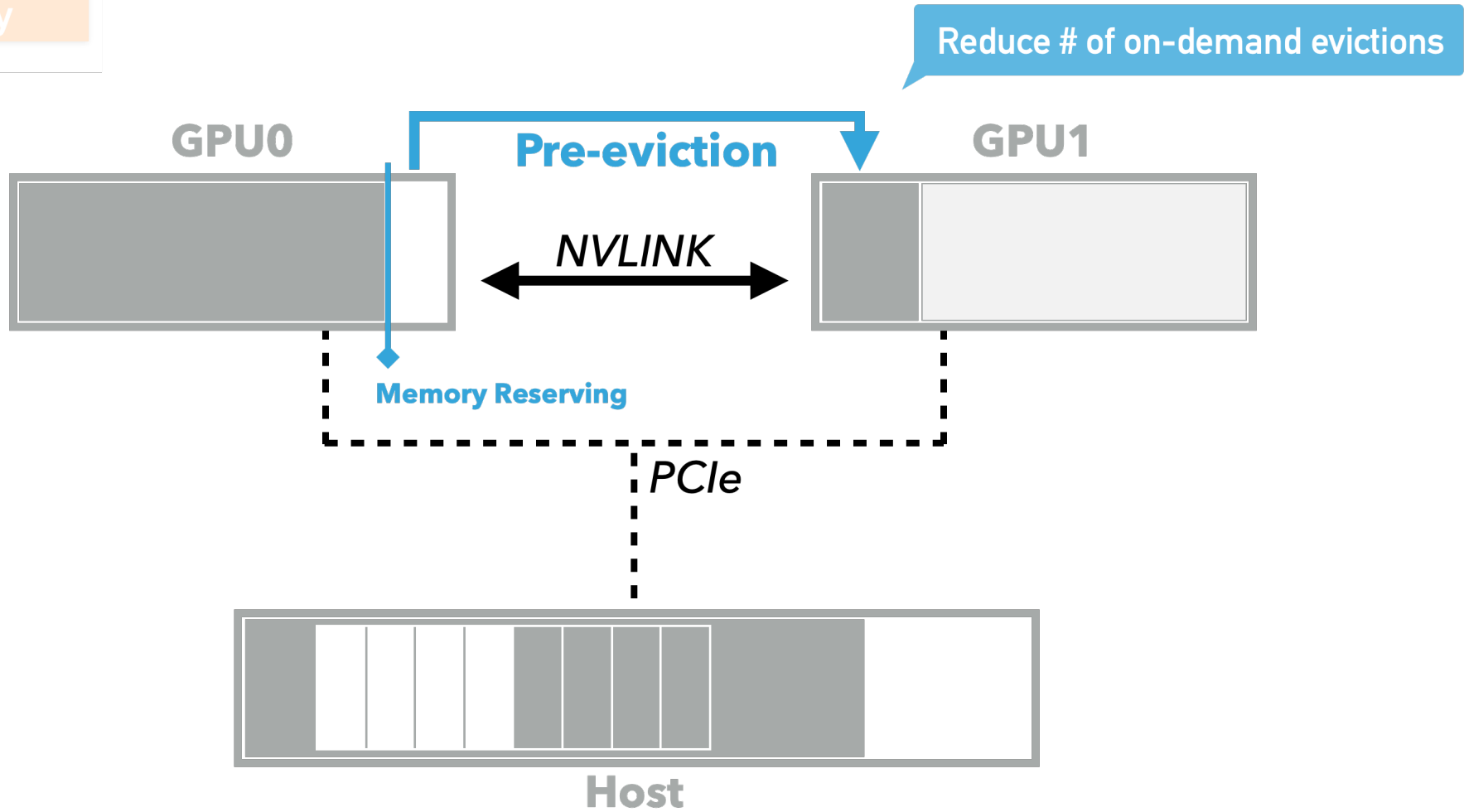
Minimal Interference

Framework-agnostic

1. Hiding eviction latency

- 2. Reducing fetch latency
- 3. Hiding fetch latency

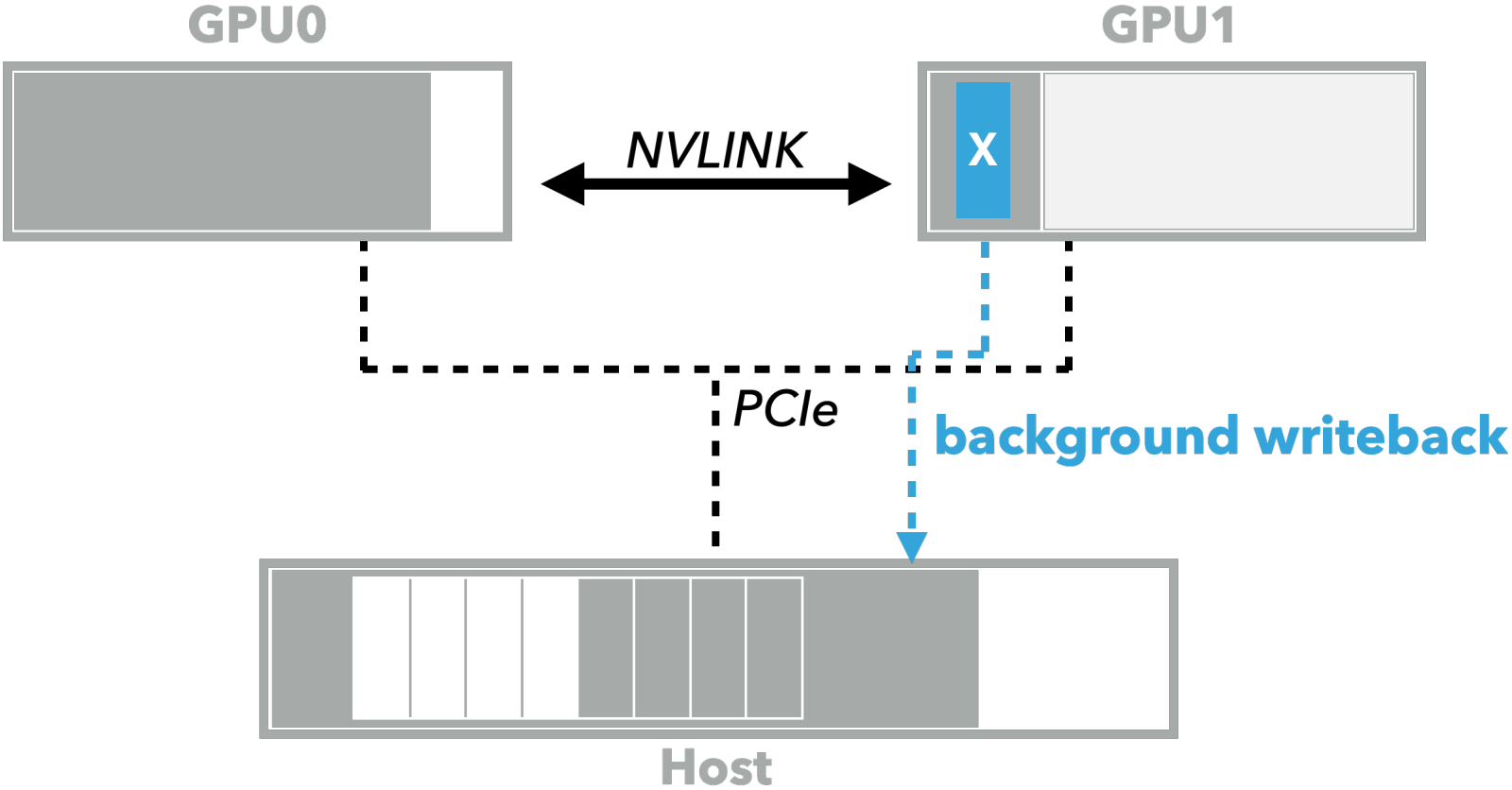
- ✓ Effective Harvesting
- Minimal Interference
- ✓ Framework-agnostic



1. Hiding eviction latency

- 2. Reducing fetch latency
- 3. Hiding fetch latency

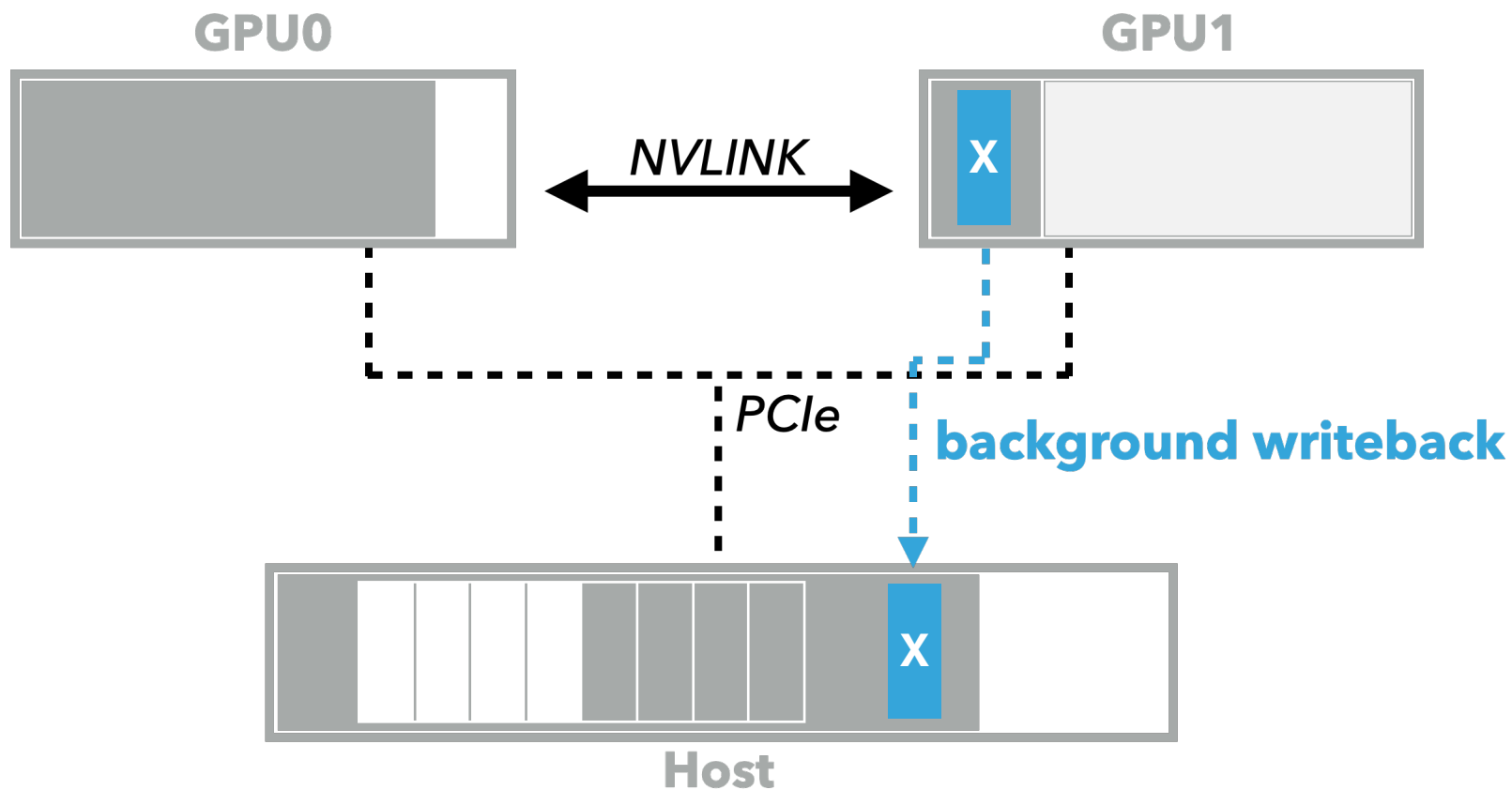
- Effective Harvesting
- ✓ Minimal Interference
- ✓ Framework-agnostic



1. Hiding eviction latency

- 2. Reducing fetch latency
- 3. Hiding fetch latency

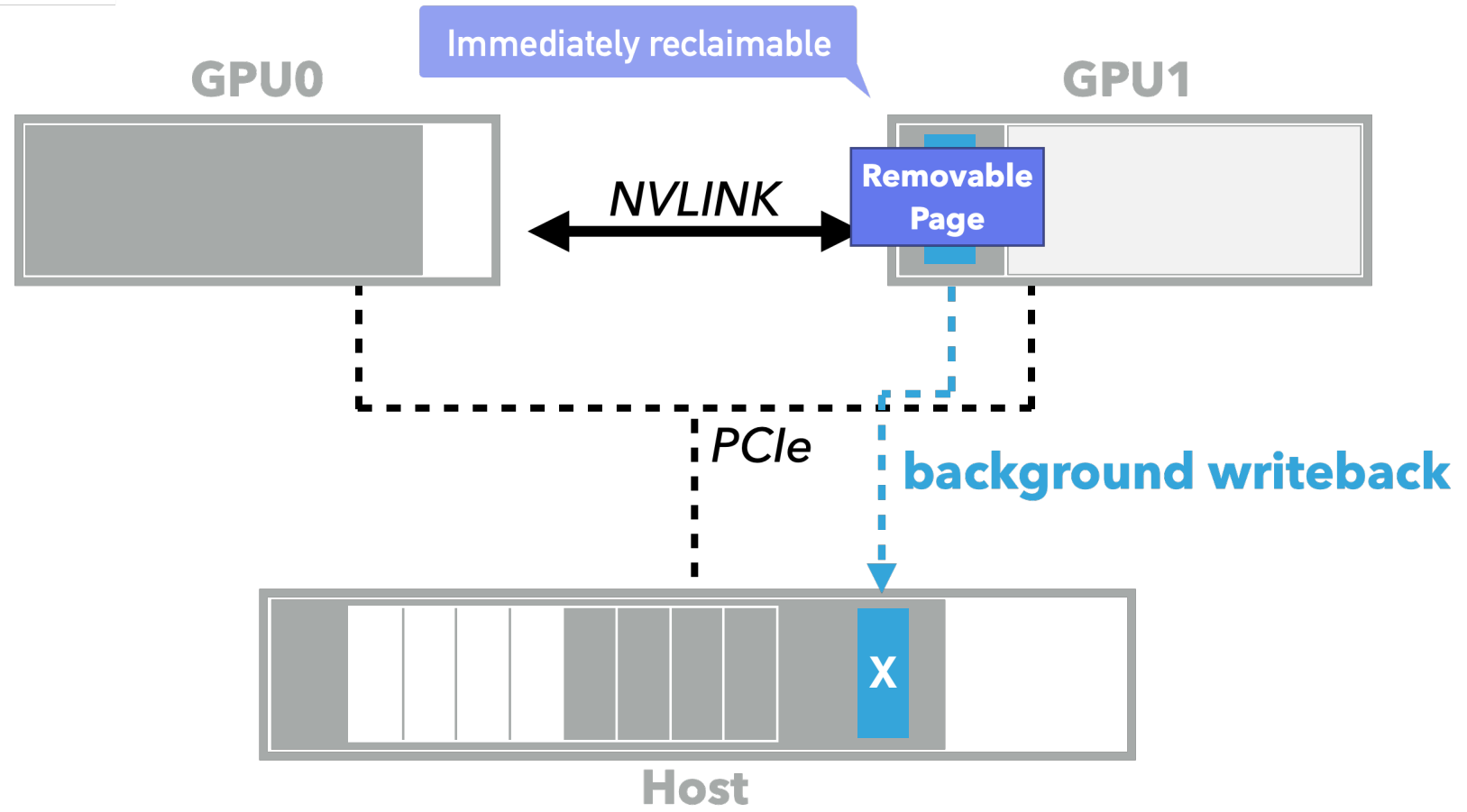
- Effective Harvesting
- ✓ Minimal Interference
- ✓ Framework-agnostic



1. Hiding eviction latency

- 2. Reducing fetch latency
- 3. Hiding fetch latency

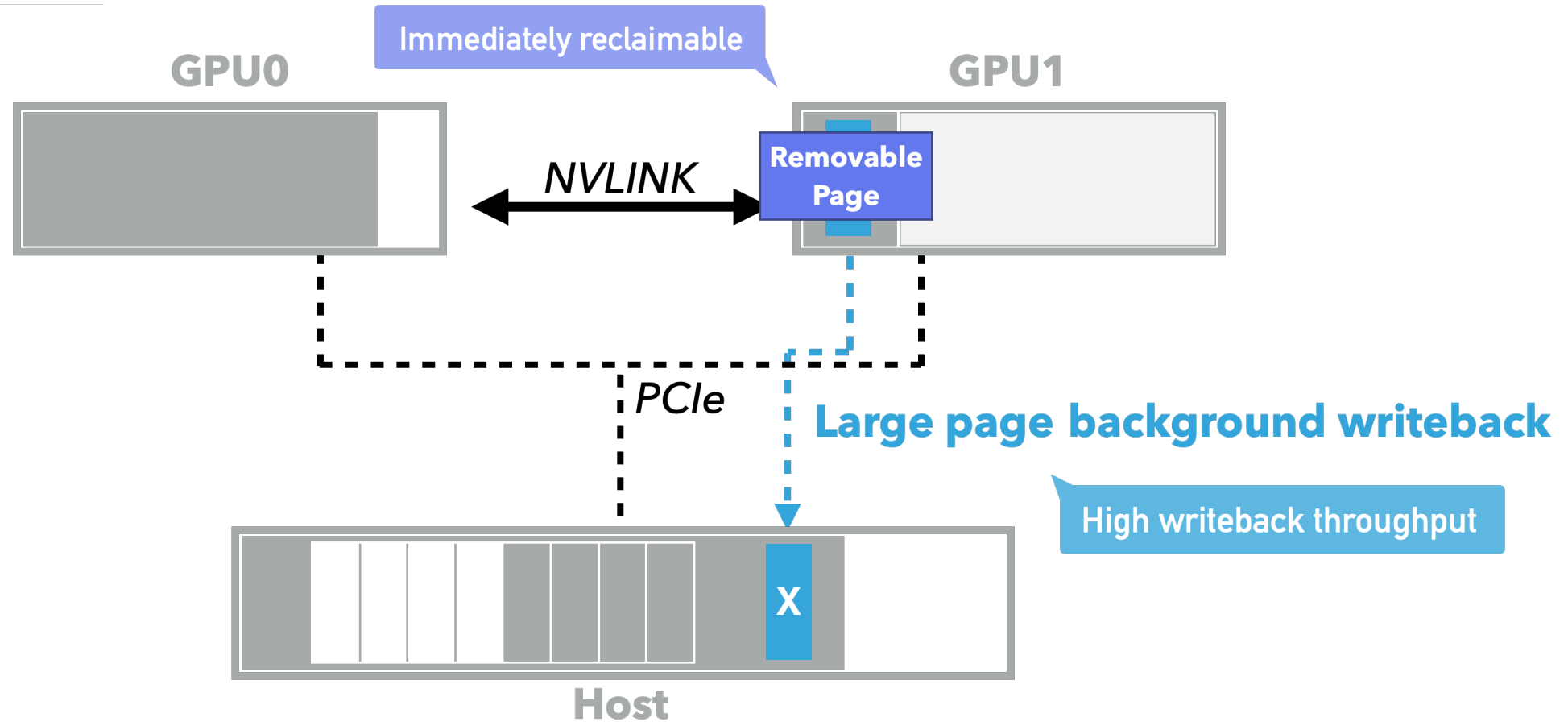
- Effective Harvesting
- Minimal Interference
- Framework-agnostic



1. Hiding eviction latency

- 2. Reducing fetch latency
- 3. Hiding fetch latency

- Effective Harvesting
- ✓ Minimal Interference
- ✓ Framework-agnostic

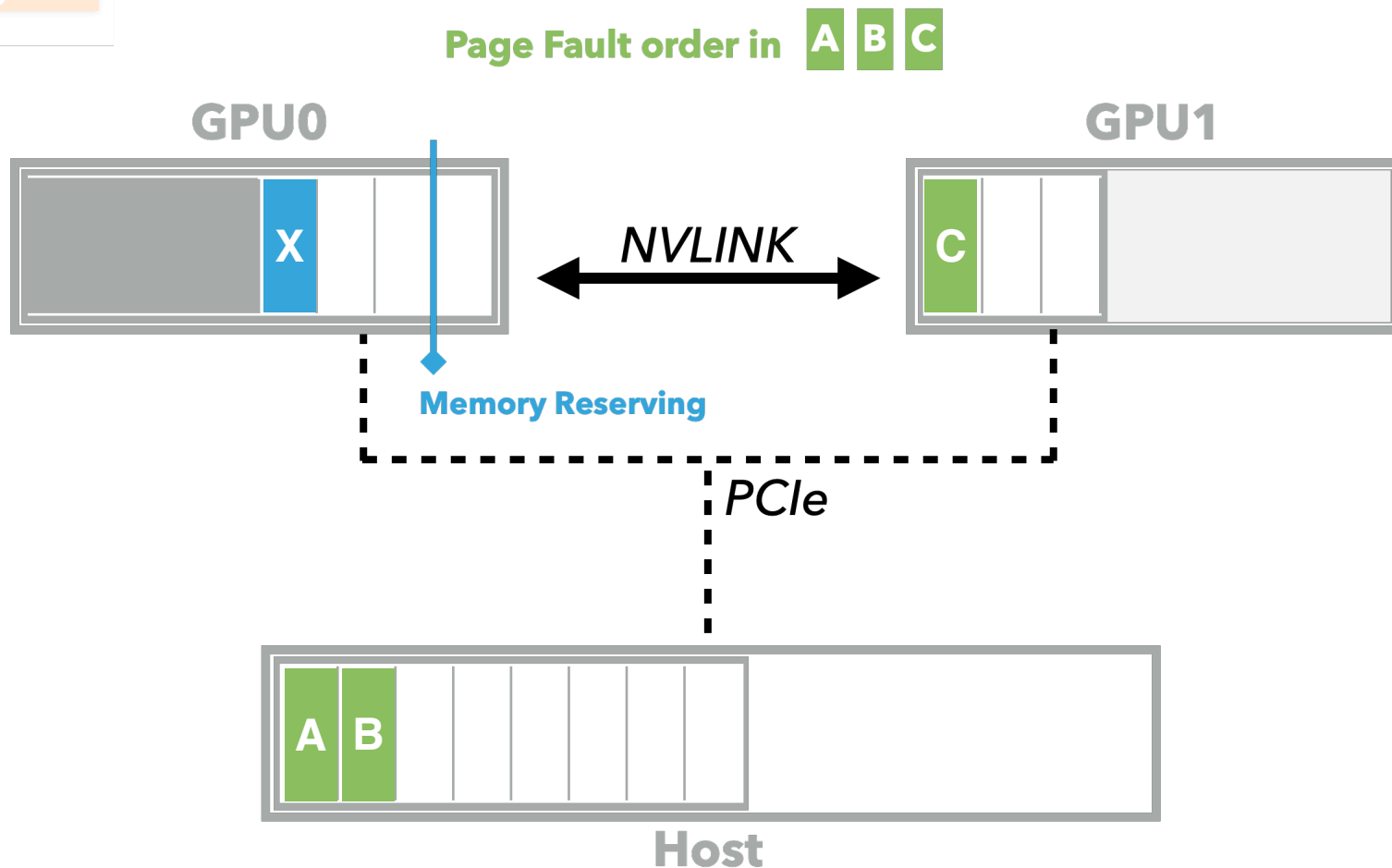


1. Hiding eviction latency

2. Reducing fetch latency

3. Hiding fetch latency

- ✓ Effective Harvesting
- ✓ Minimal Interference
- ✓ Framework-agnostic

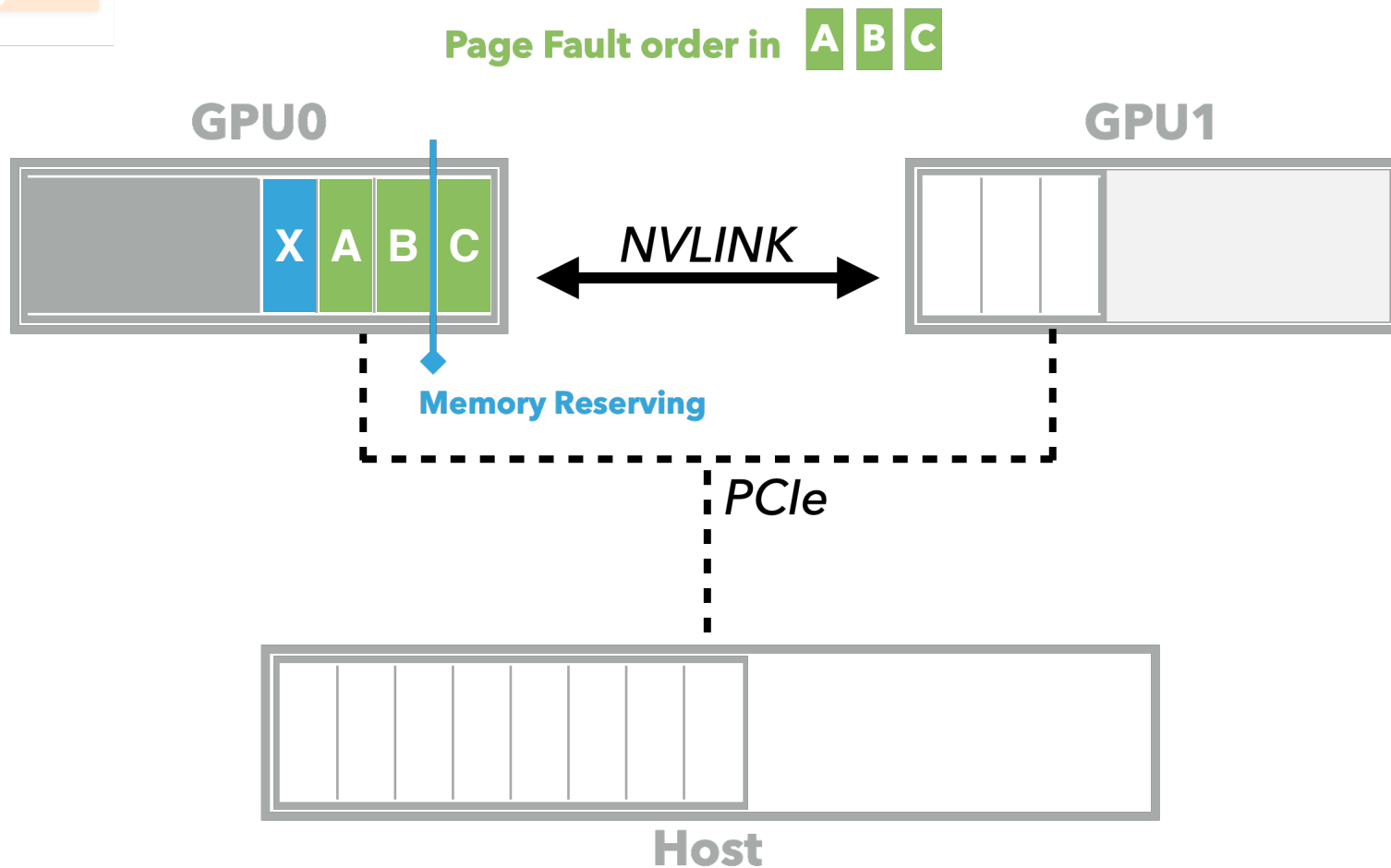


1. Hiding eviction latency

2. Reducing fetch latency

3. Hiding fetch latency

- ✓ Effective Harvesting
- ✓ Minimal Interference
- ✓ Framework-agnostic

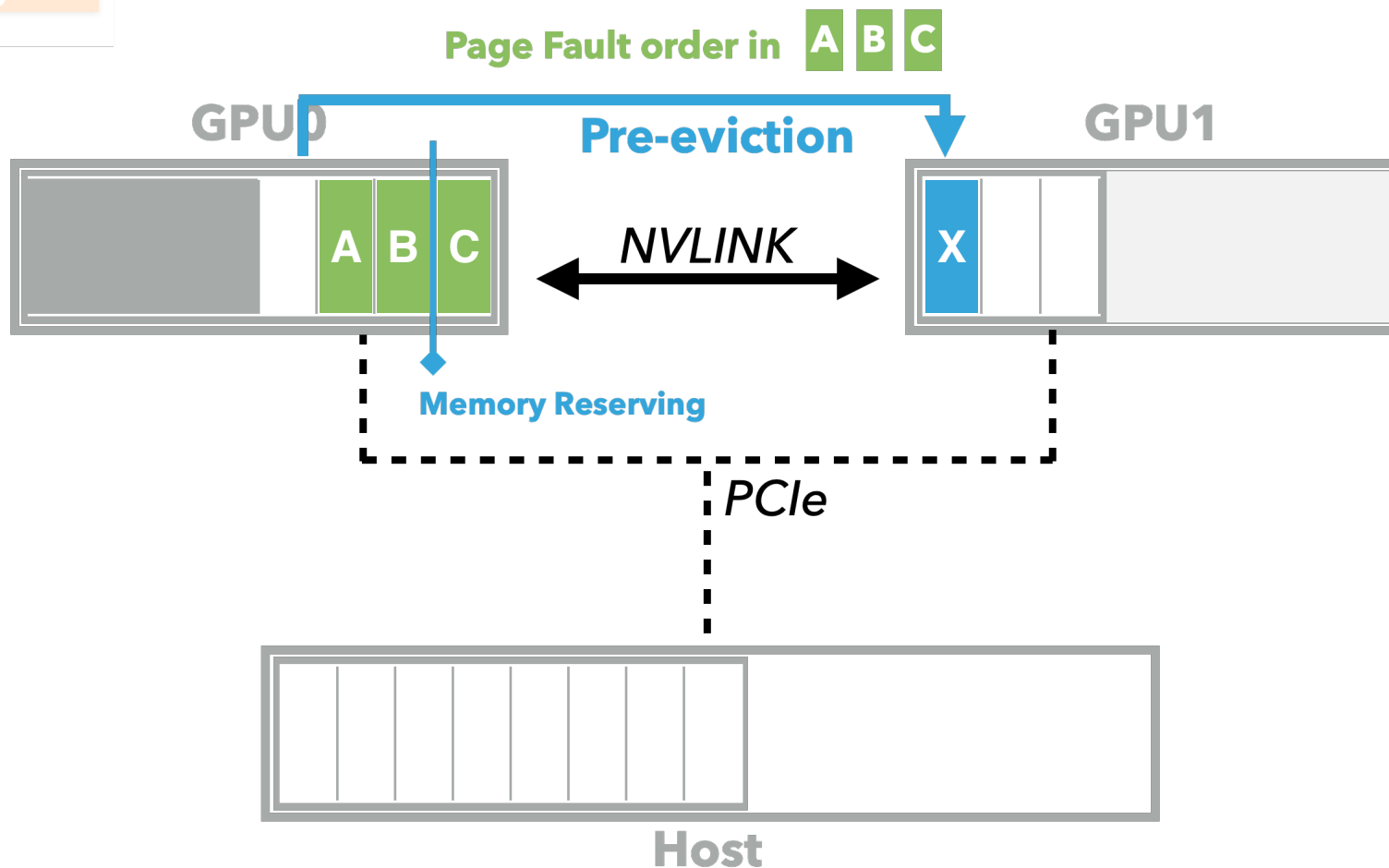


1. Hiding eviction latency

2. Reducing fetch latency

3. Hiding fetch latency

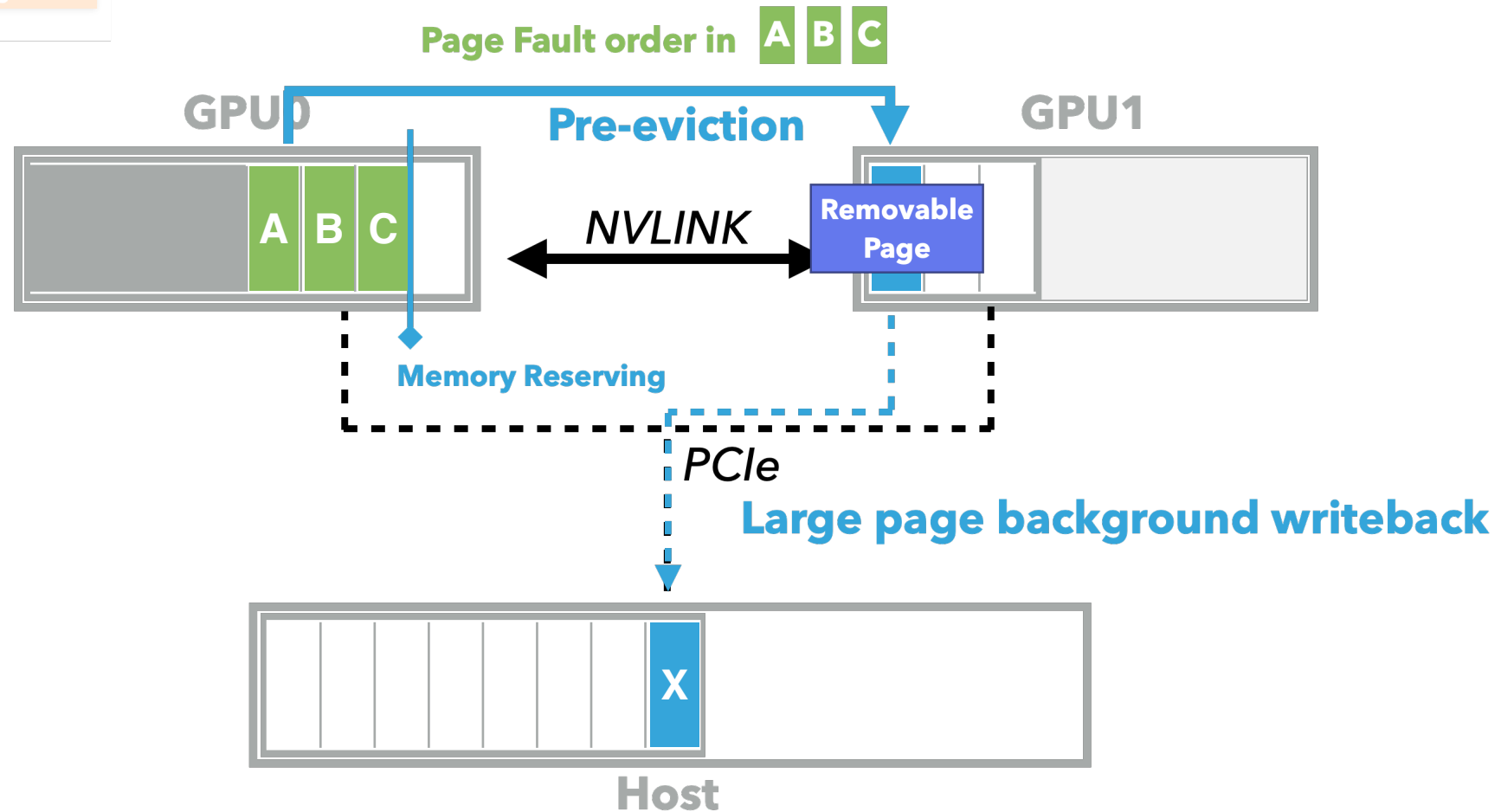
- ✓ Effective Harvesting
- ✓ Minimal Interference
- ✓ Framework-agnostic



1. Hiding eviction latency

- 2. Reducing fetch latency
- 3. Hiding fetch latency

- ✓ Effective Harvesting
- ✓ Minimal Interference
- ✓ Framework-agnostic

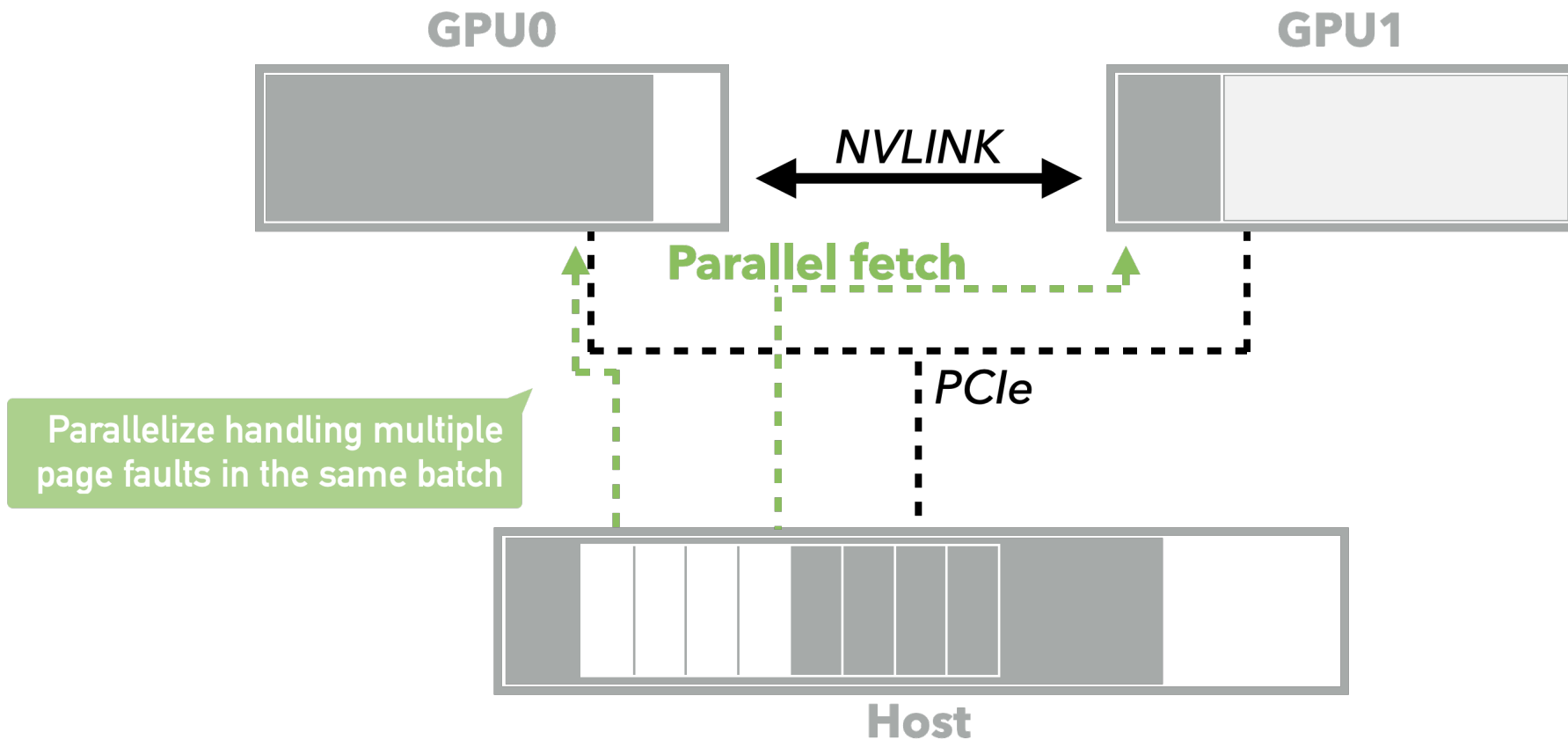


1. Hiding eviction latency

2. Reducing fetch latency

3. Hiding fetch latency

- ✓ Effective Harvesting
- Minimal Interference
- ✓ Framework-agnostic

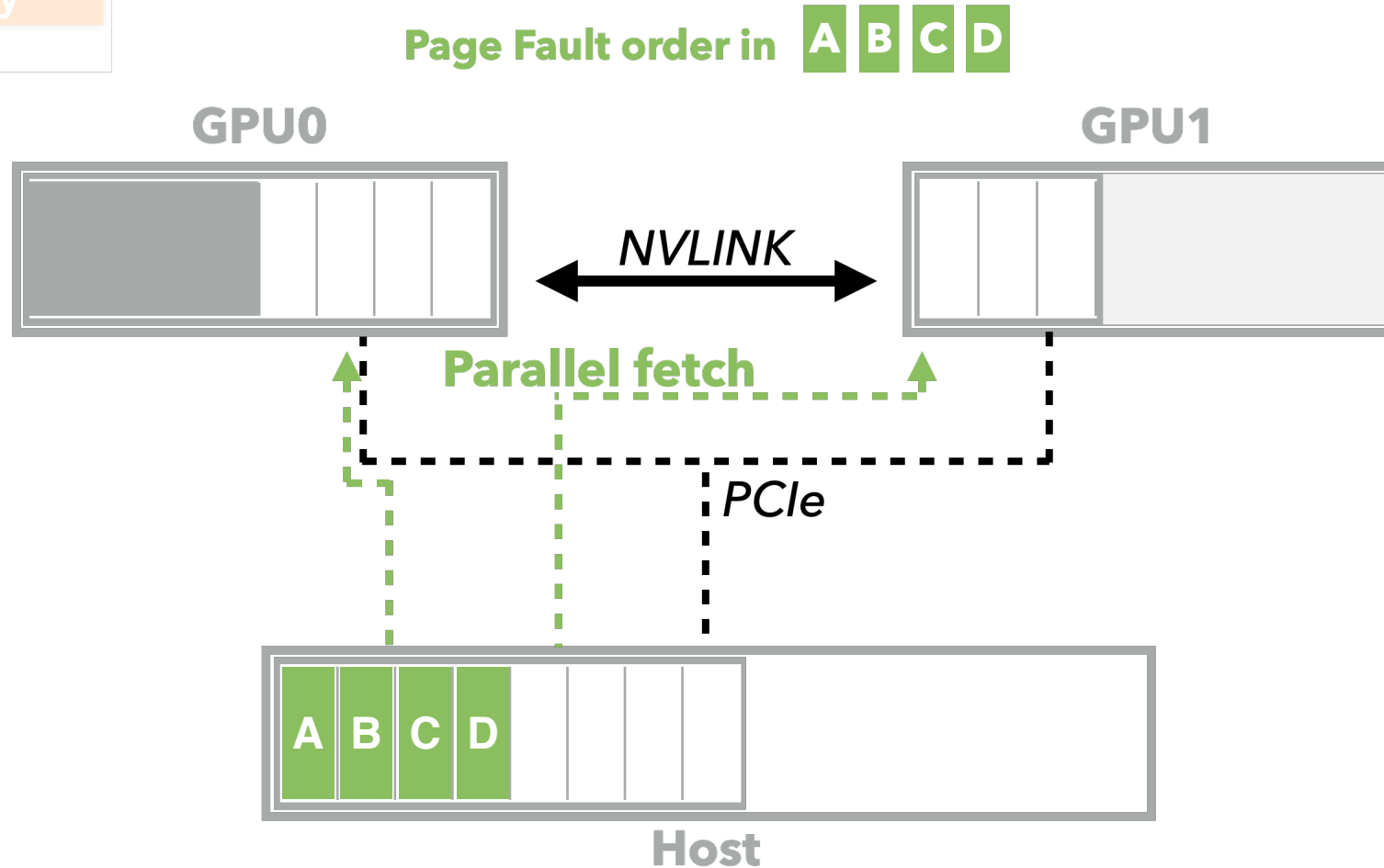


1. Hiding eviction latency

2. Reducing fetch latency

3. Hiding fetch latency

- ✓ Effective Harvesting
- Minimal Interference
- ✓ Framework-agnostic

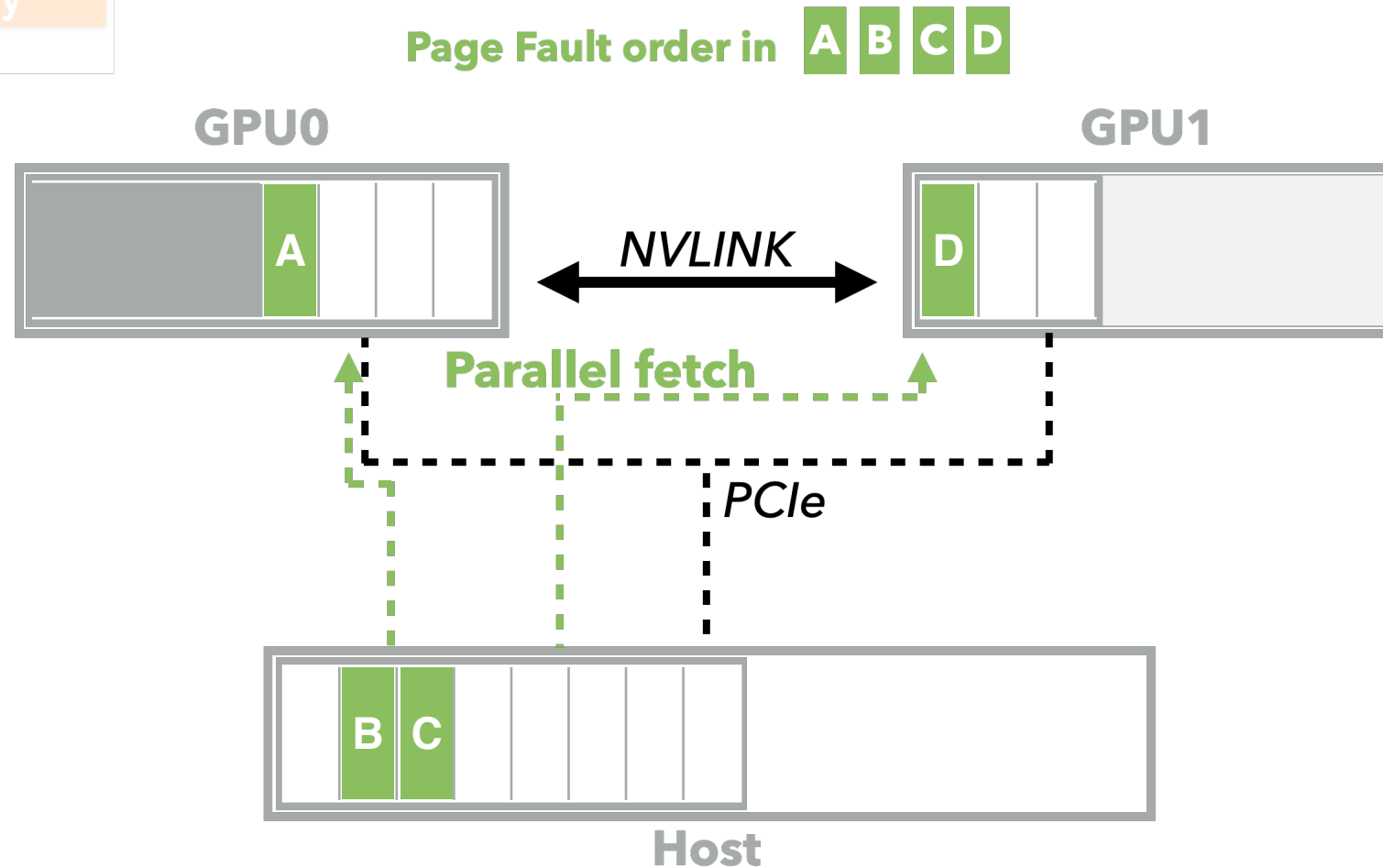


1. Hiding eviction latency

2. Reducing fetch latency

3. Hiding fetch latency

- ✓ Effective Harvesting
- Minimal Interference
- ✓ Framework-agnostic



1. Hiding eviction latency

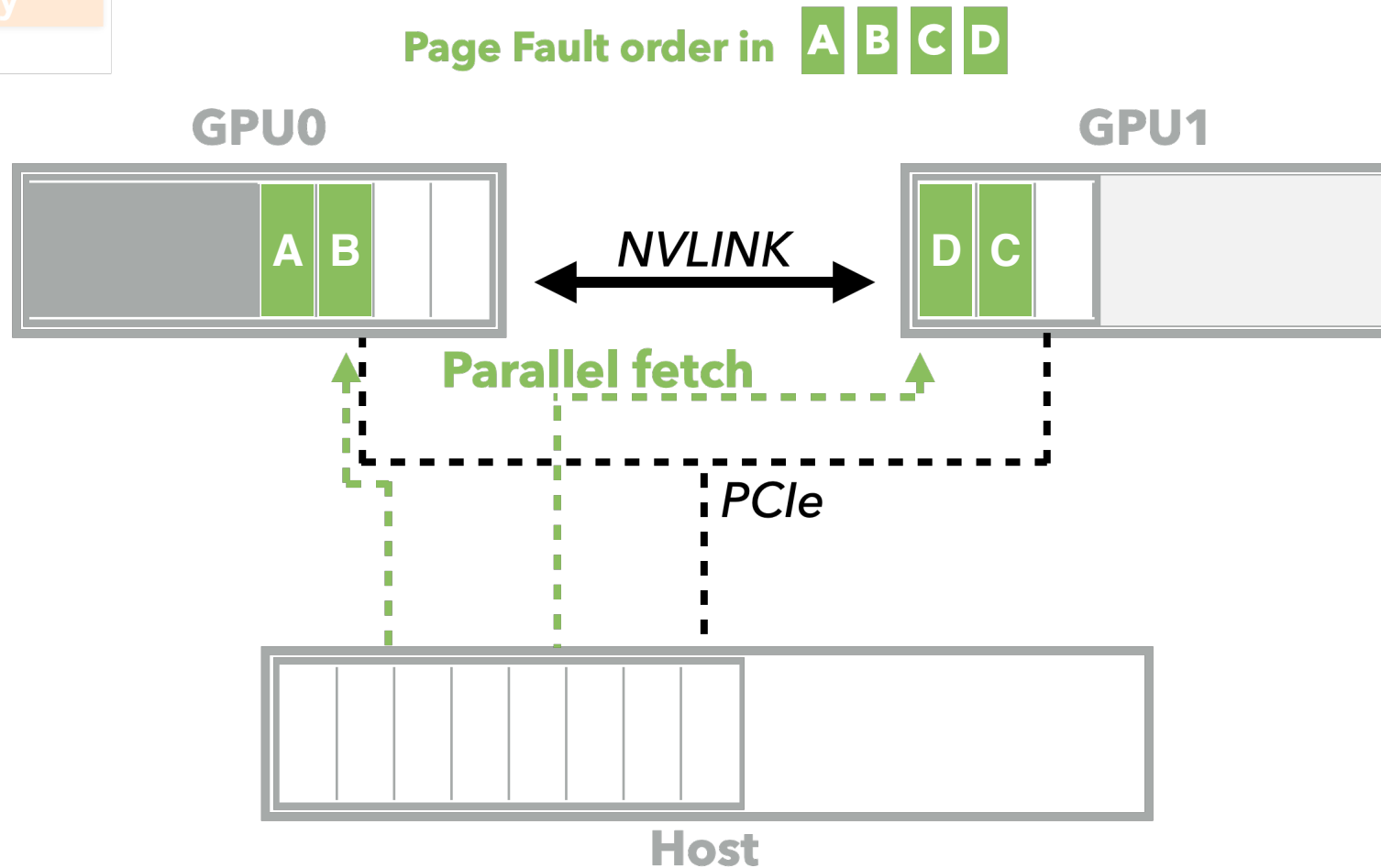
2. Reducing fetch latency

3. Hiding fetch latency

✓ Effective Harvesting

Minimal Interference

✓ Framework-agnostic

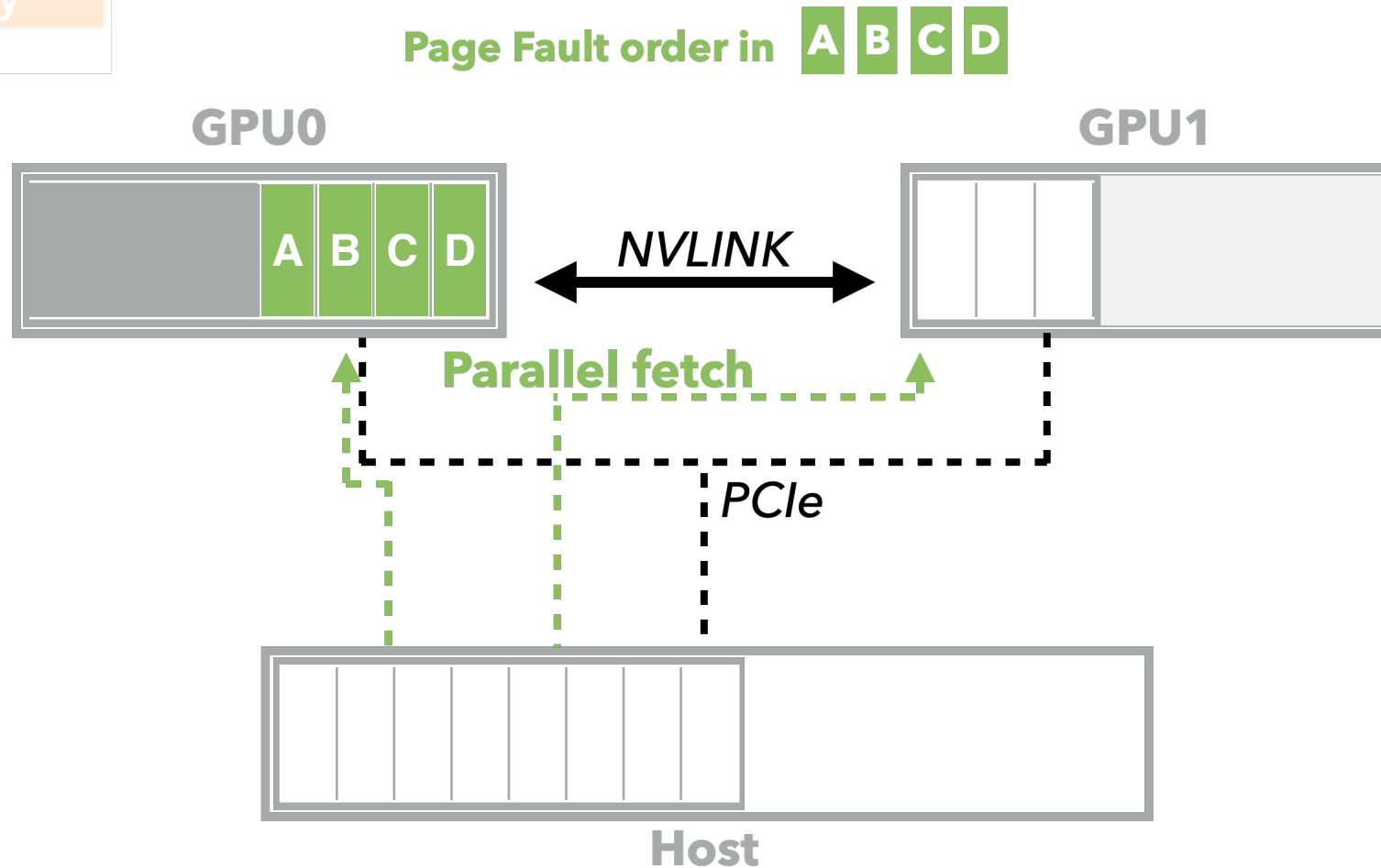


1. Hiding eviction latency

2. Reducing fetch latency

3. Hiding fetch latency

- ✓ Effective Harvesting
- Minimal Interference
- ✓ Framework-agnostic



1. Hiding eviction latency

2. Reducing fetch latency

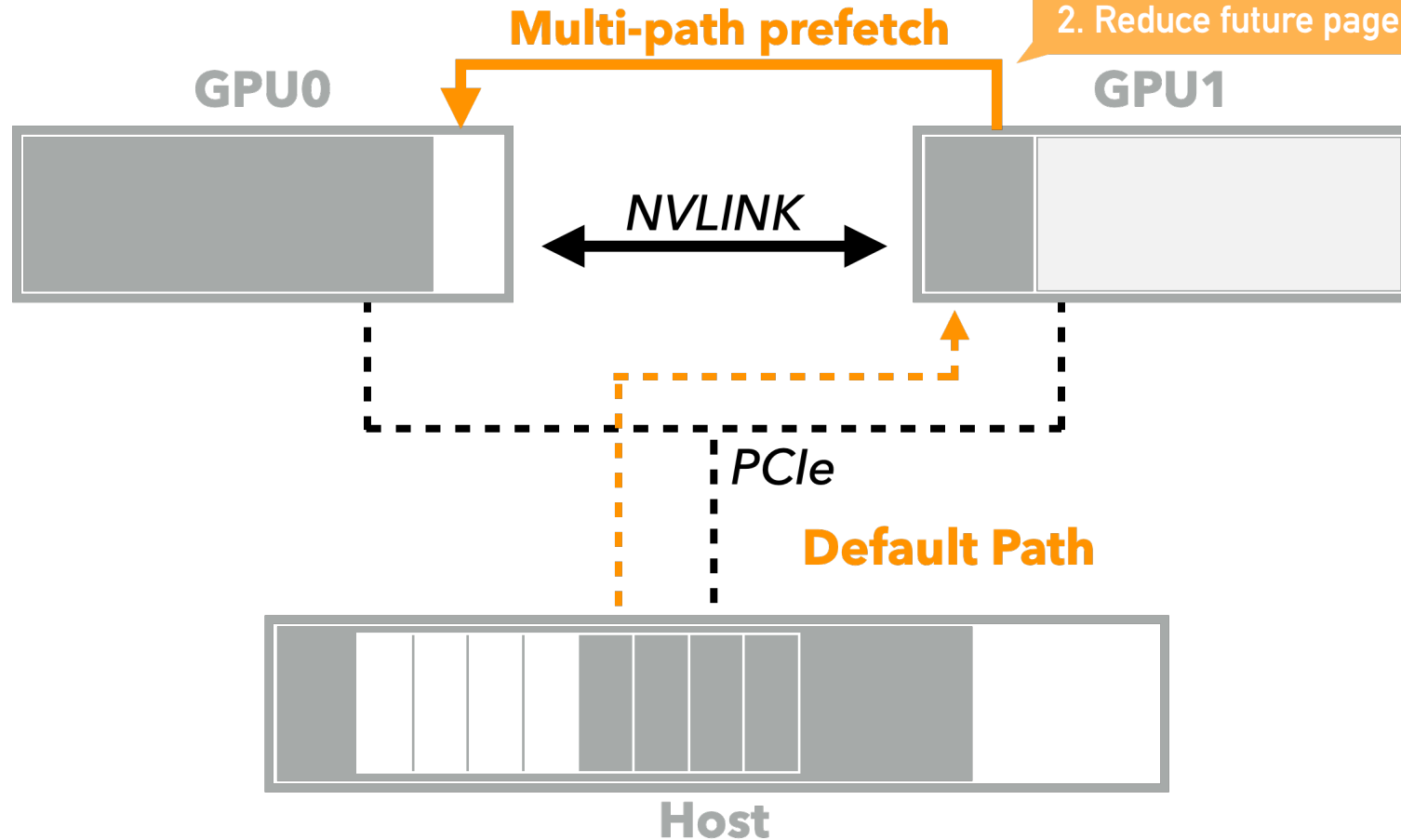
3. Hiding fetch latency

✓ Effective Harvesting

Minimal Interference

✓ Framework-agnostic

- 1. Reduce page faults by proactive prefetching
- 2. Reduce future page fault latency

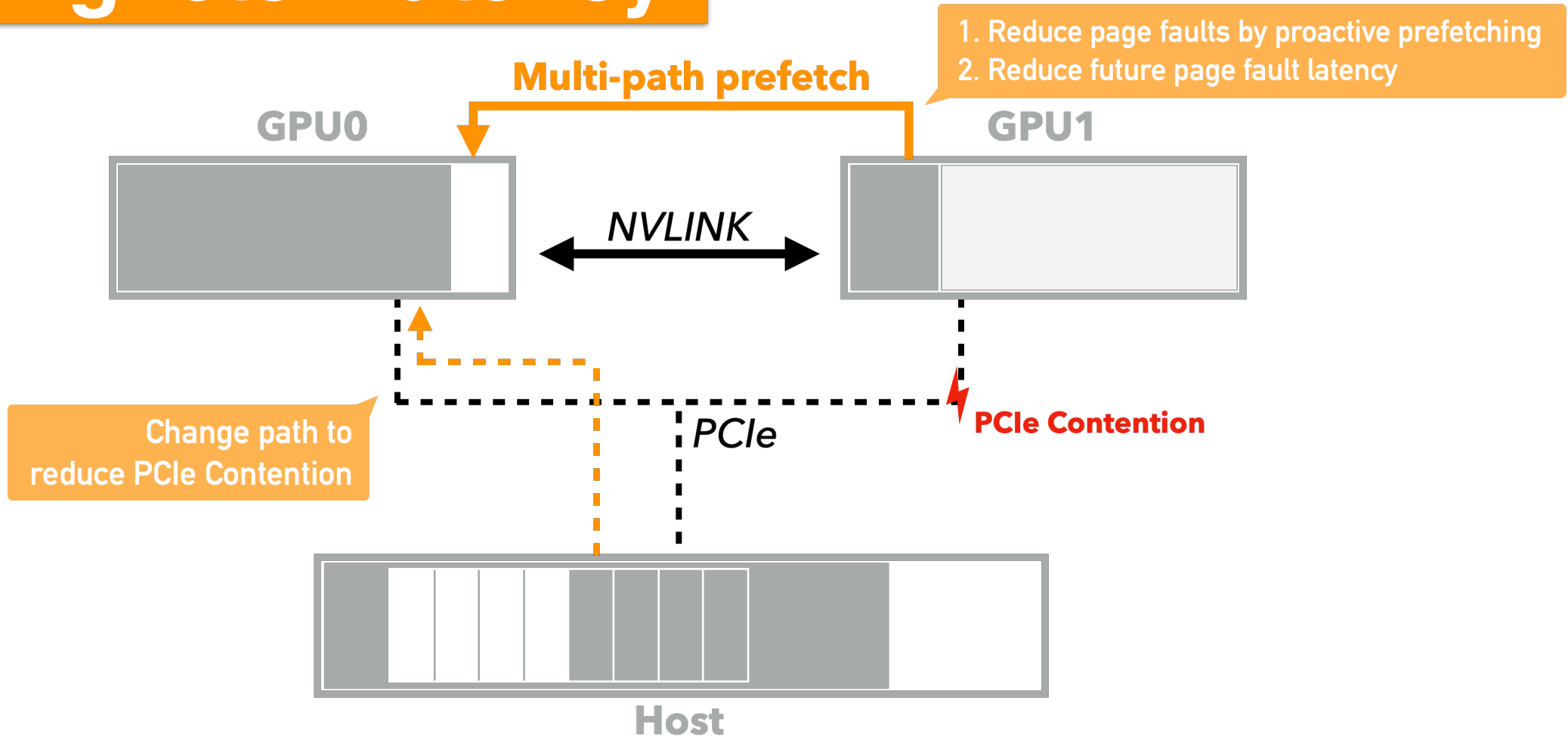


1. Hiding eviction latency

2. Reducing fetch latency

3. Hiding fetch latency

- ✓ Effective Harvesting
- ✓ Minimal Interference
- ✓ Framework-agnostic



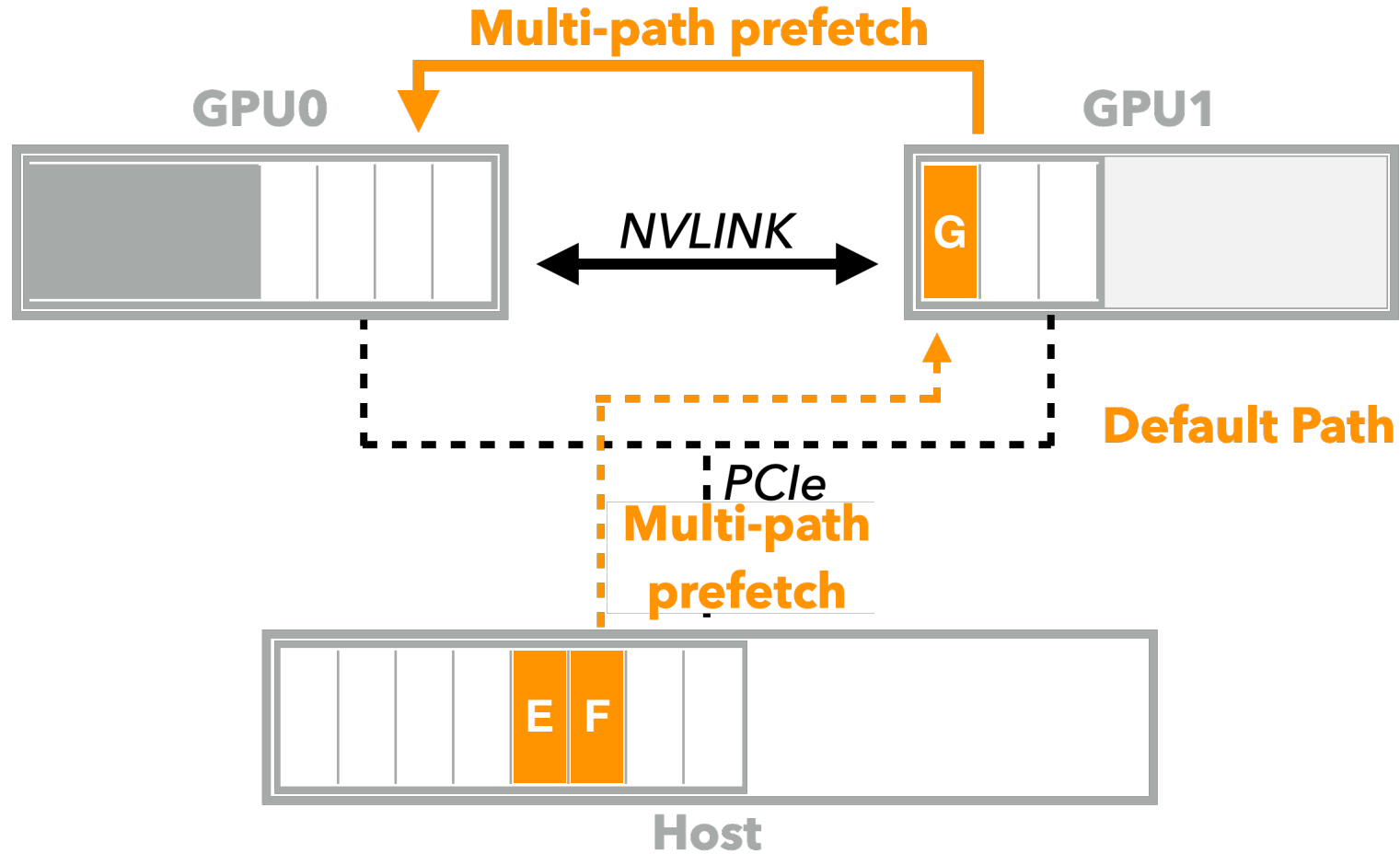
1. Hiding eviction latency

2. Reducing fetch latency

3. Hiding fetch latency

- ✓ Effective Harvesting
- Minimal Interference
- ✓ Framework-agnostic

Prefetch **E** **F** **G**



1. Hiding eviction latency

2. Reducing fetch latency

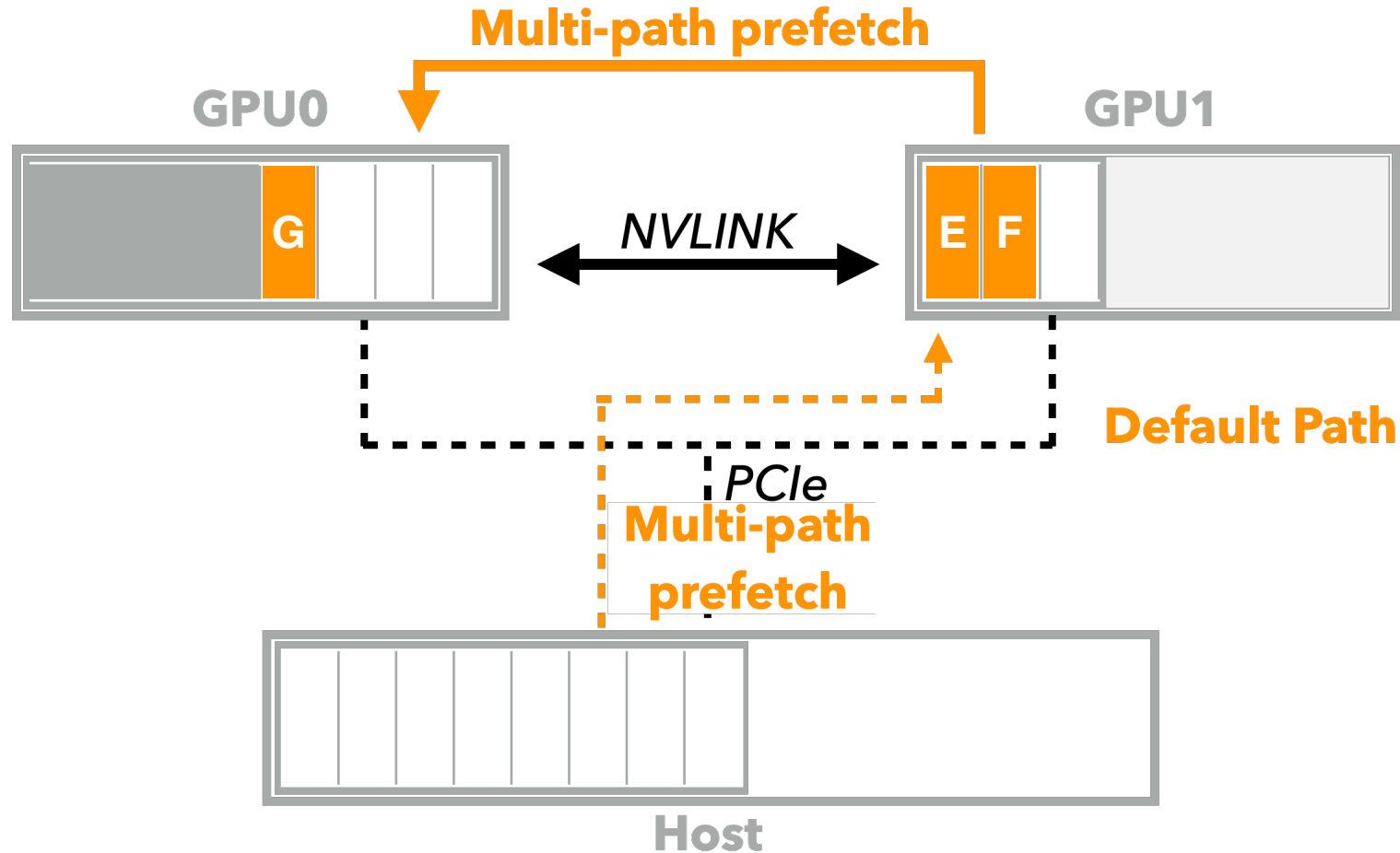
3. Hiding fetch latency

✓ Effective Harvesting

Minimal Interference

✓ Framework-agnostic

Prefetch **E** **F** **G**



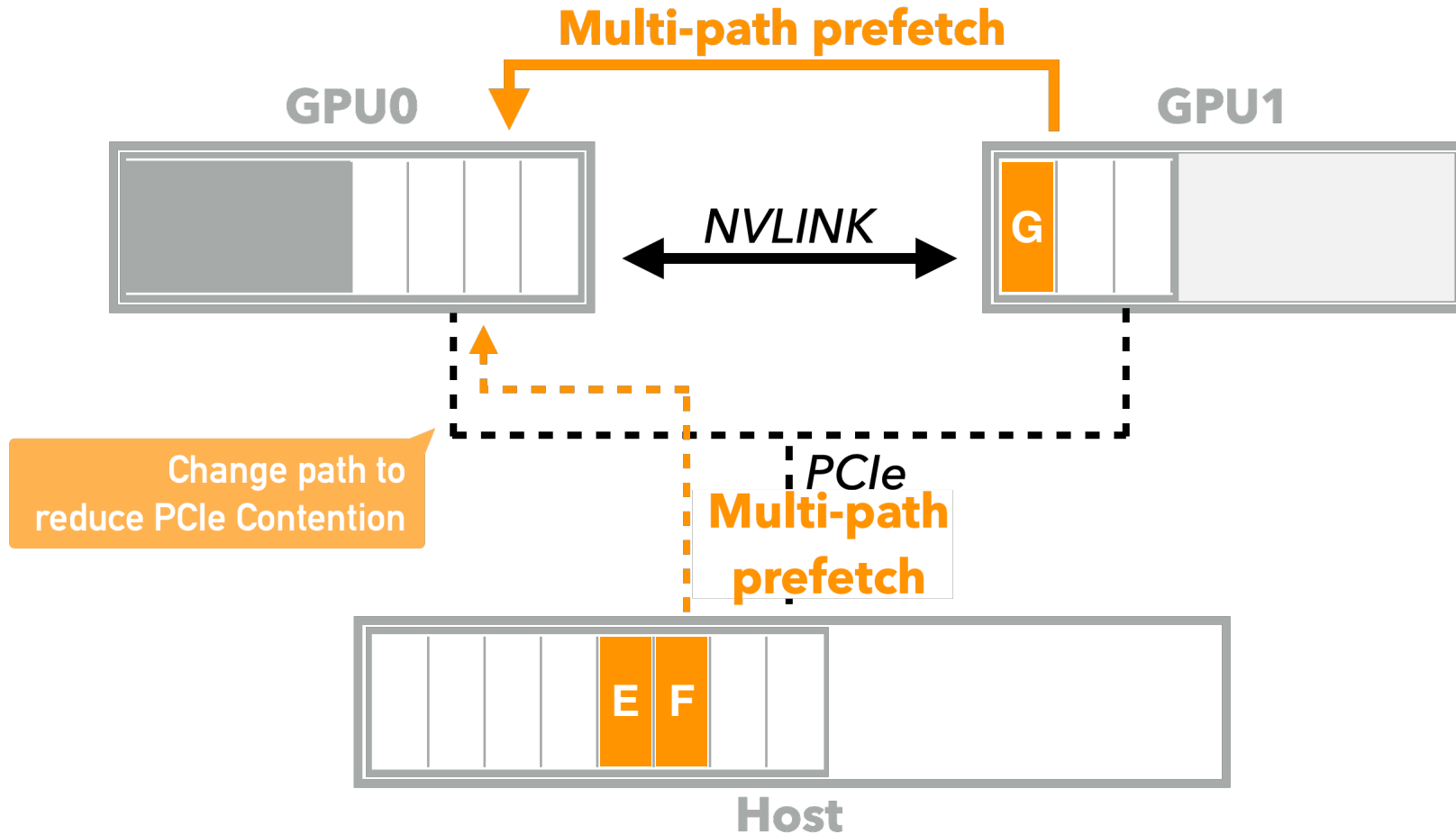
1. Hiding eviction latency

2. Reducing fetch latency

3. Hiding fetch latency

- ✓ Effective Harvesting
- ✓ Minimal Interference
- ✓ Framework-agnostic

Prefetch **E** **F** **G**



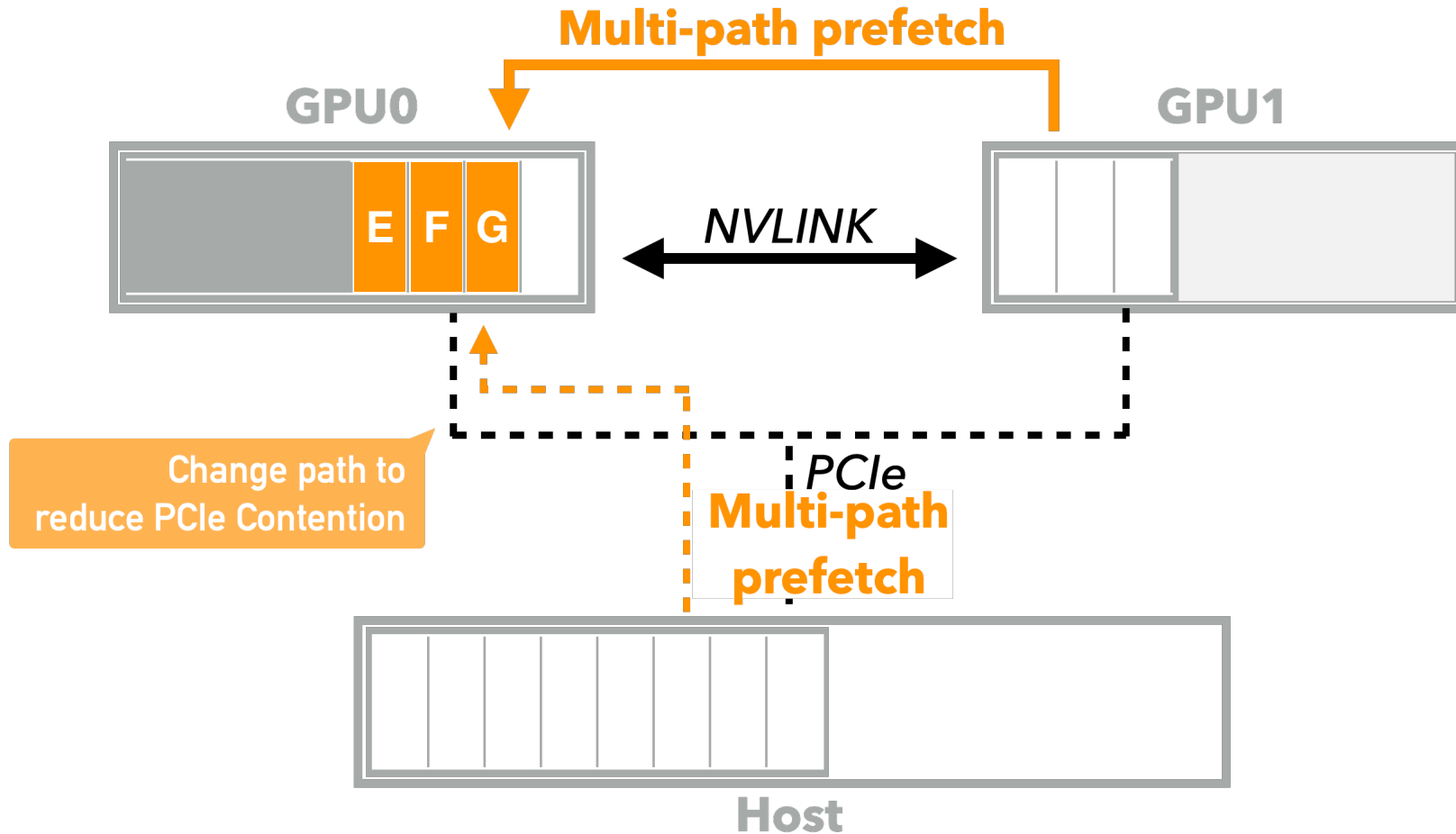
1. Hiding eviction latency

2. Reducing fetch latency

3. Hiding fetch latency

- ✓ Effective Harvesting
- ✓ Minimal Interference
- ✓ Framework-agnostic

Prefetch **E** **F** **G**



Goals of Hierarchical Unified Virtual Memory

Reducing performance cost of memory oversubscription

Effective Harvesting

Harvest small and temporarily available spare memory of neighbor GPU

Reduce eviction/fetch latency with spare memory

Minimal Interference

Minimize performance impact of workloads running in neighbor GPU

Framework-agnostic

No modification of applications or frameworks

Memory manager for HUVIM: memHarvester

Centralized coordinator for data-path in HUVIM

Effective Harvesting

Harvest small and temporarily available spare memory or neighbor
Pre-emption
Parallel Fetch
Multi-path Prefetch
Reduce eviction/fetch latency with spare memory

Minimal Interference

Minimize performance impact of workloads running on GPU
Removable Page

Framework-agnostic

No modification of applications or frameworks
Leverage UVM

Evaluation

- System configuration (AWS p3.8xlarge)

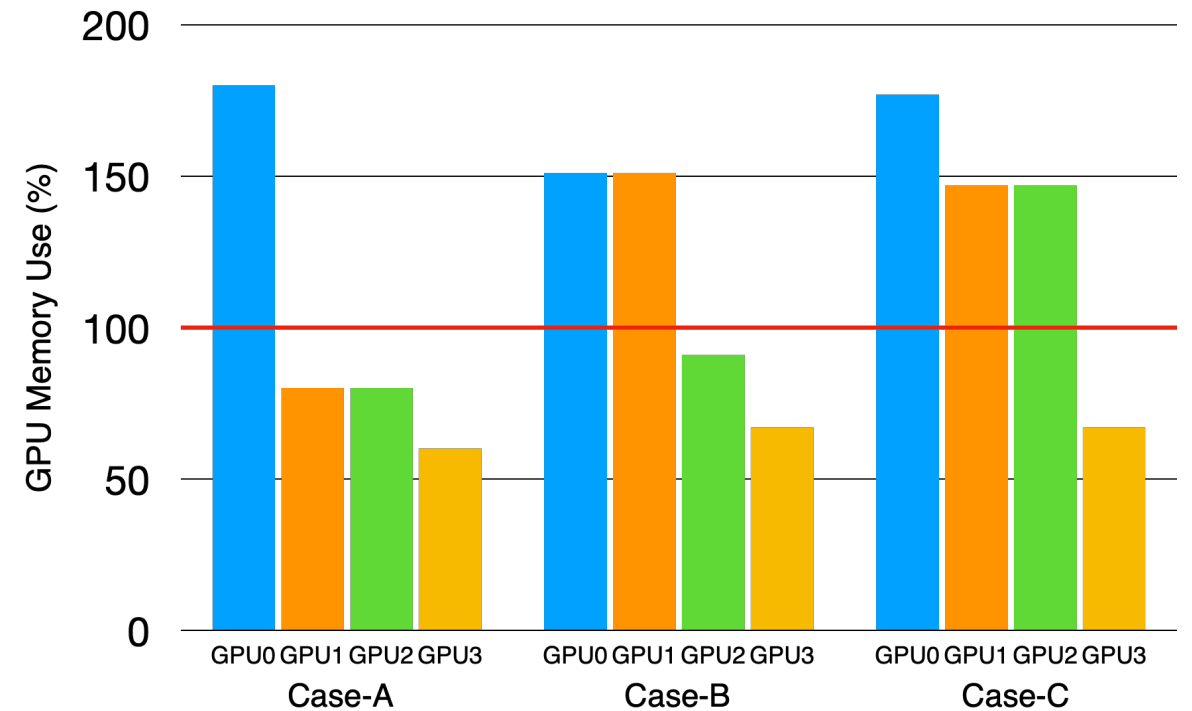
GPU	NVIDIA V100 (x4) (16 GB each) (Connected through NVLink bridge)
Processors	Intel Xeon Skylake vCPU (x32)
Memory	224GB DDR4
Driver Version	NVIDIA Driver 460.67

- Benchmarks
 - ▶ cuGraph (version 21.12)
 - ▶ PyTorch (version 1.10.1)

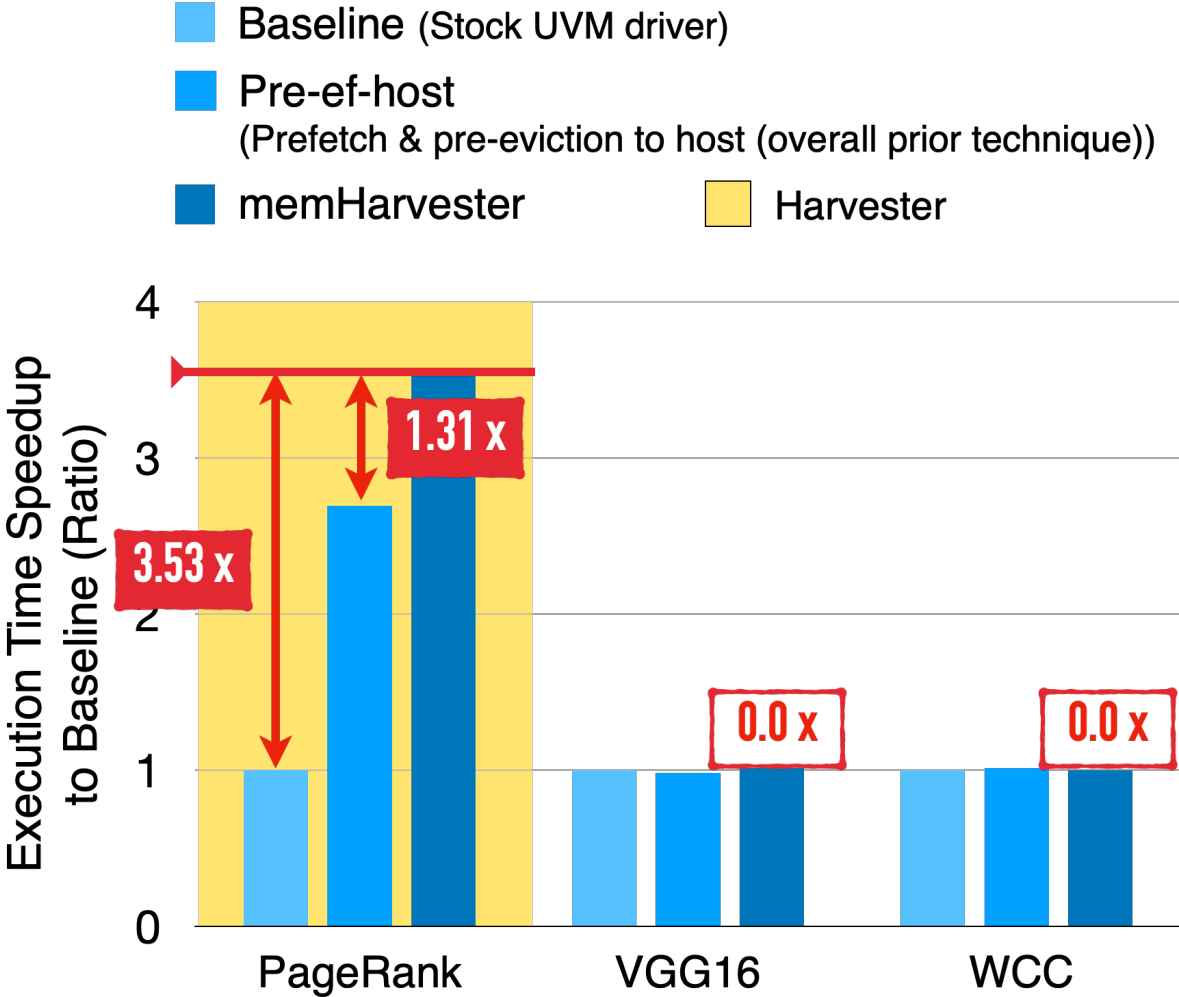
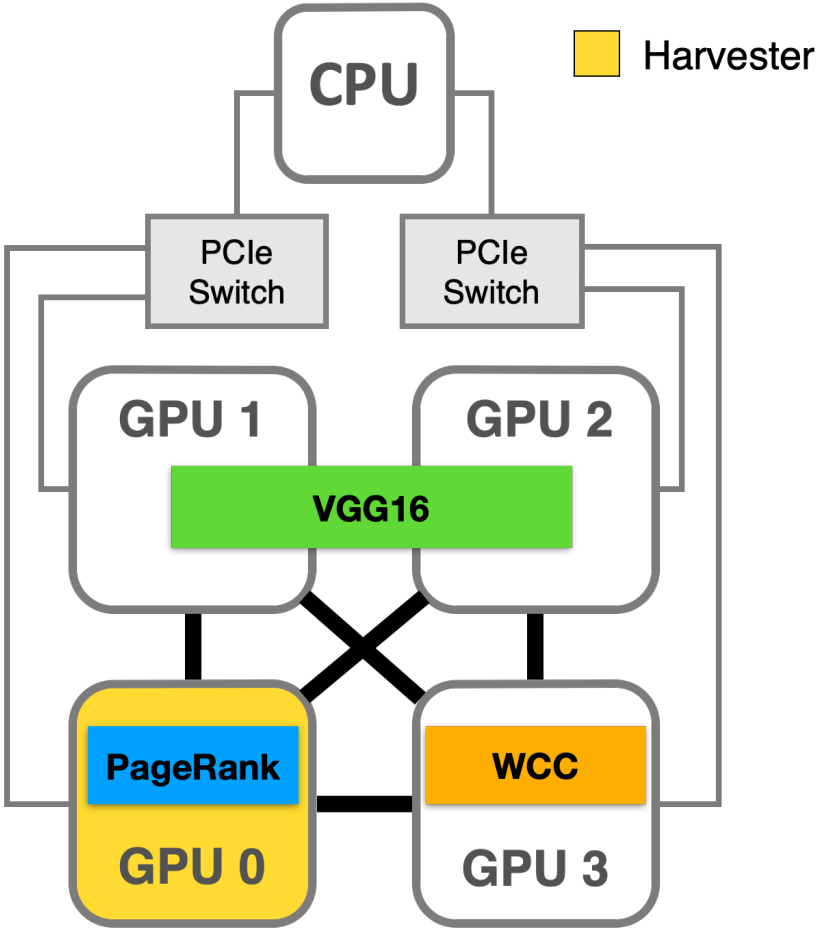
Workload running scenarios

■ Harvester

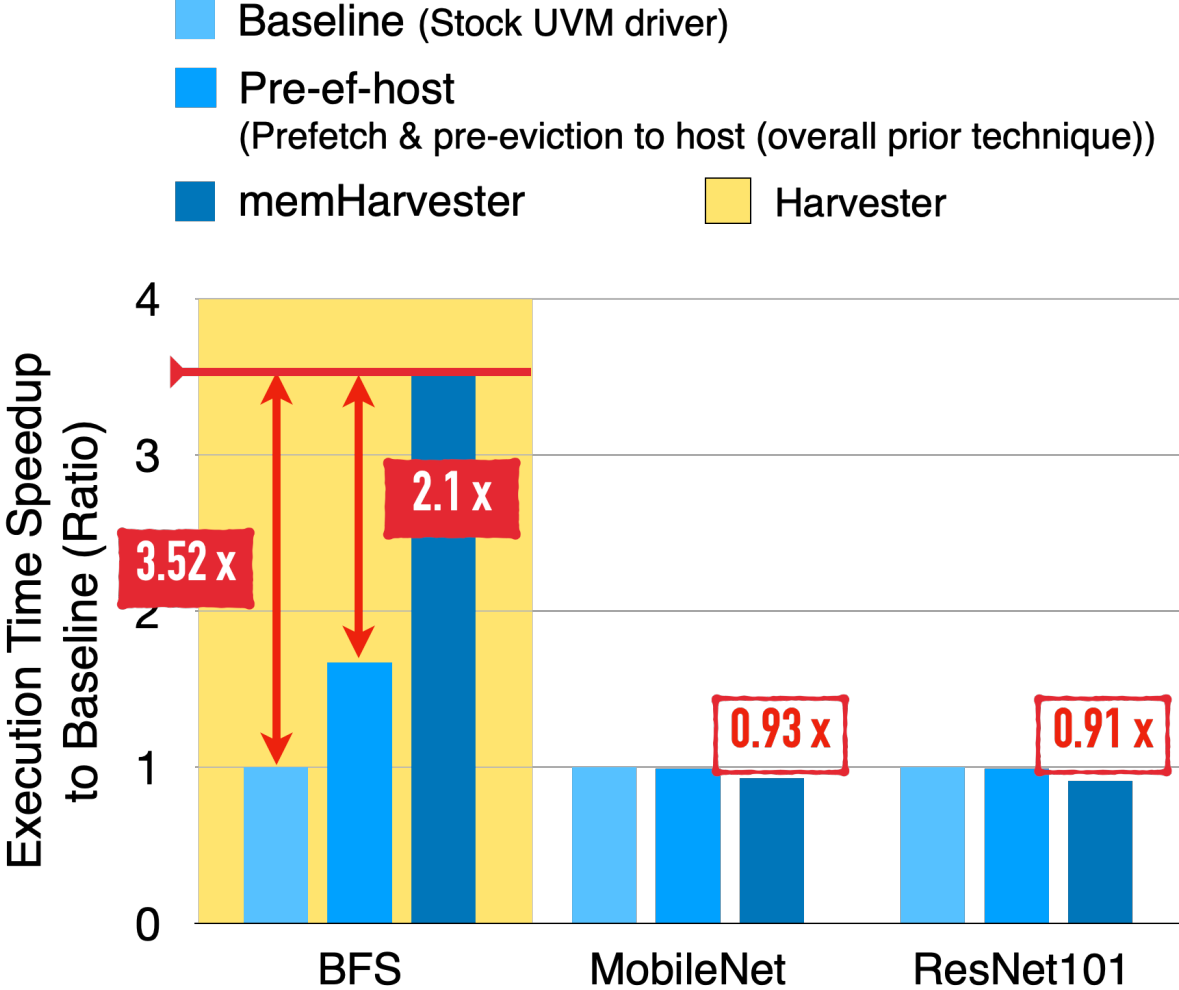
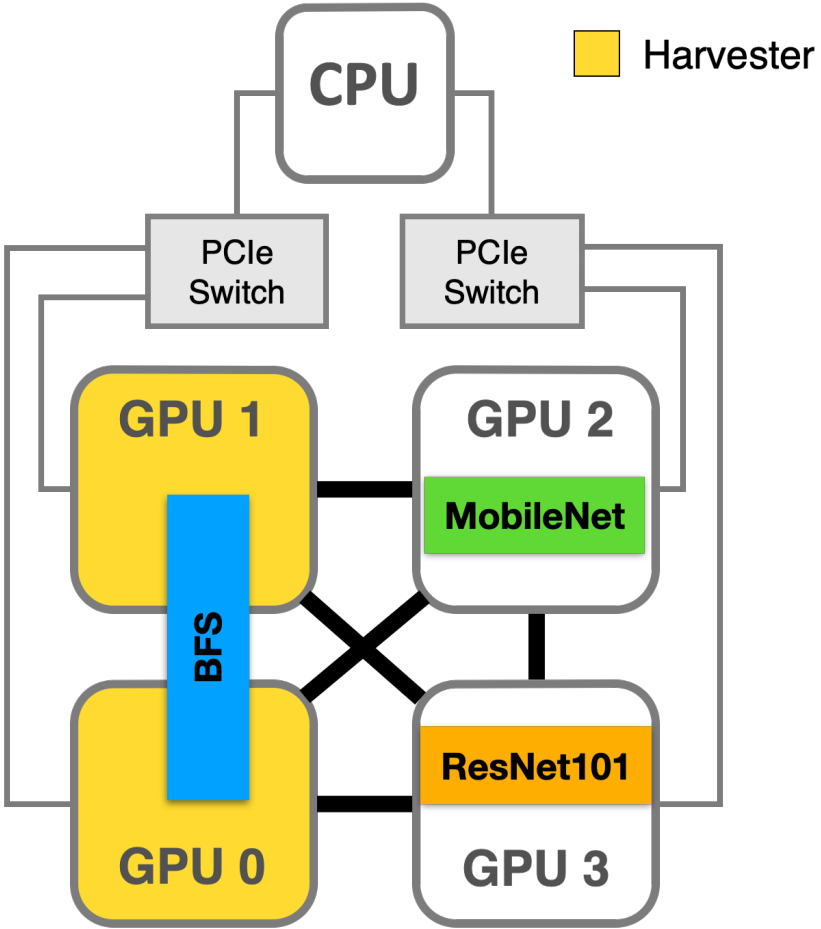
	GPU-0	GPU-1	GPU-2	GPU-3
Case-A	PageRank soc-twitter-2010	VGG16 Batch 256		WCC soc-sinaweibo
Case-B	BFS web-uk-2005		MobileNet Batch 256	ResNet101 Batch 64
Case-C	WCC soc-twitter-2010	Louvain web-uk-2005		ResNet101 Batch 64



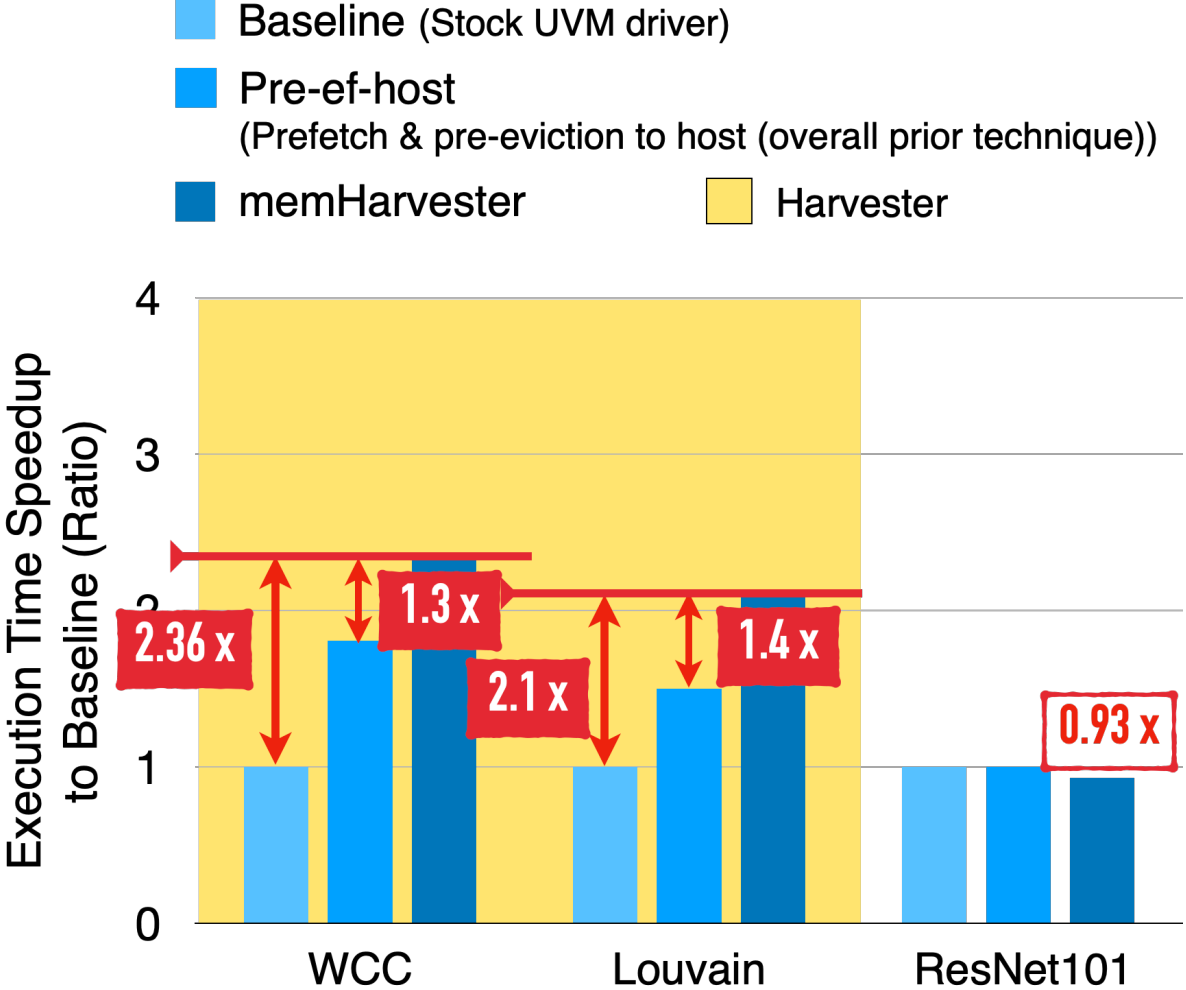
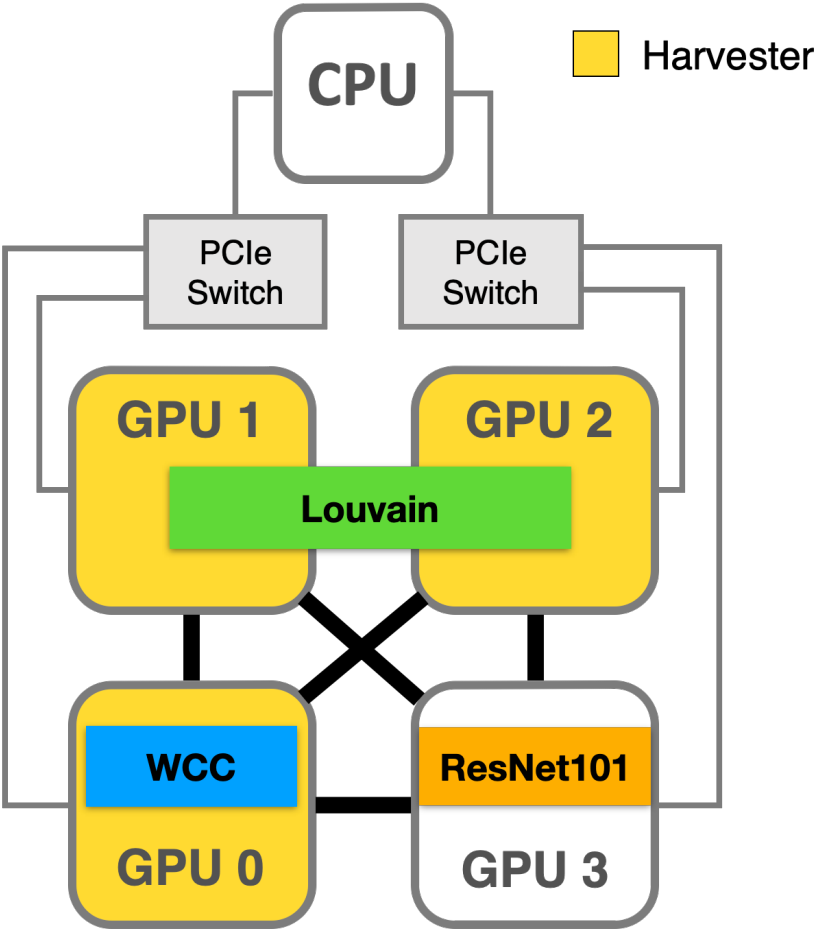
Case-A (1 harvests 3)



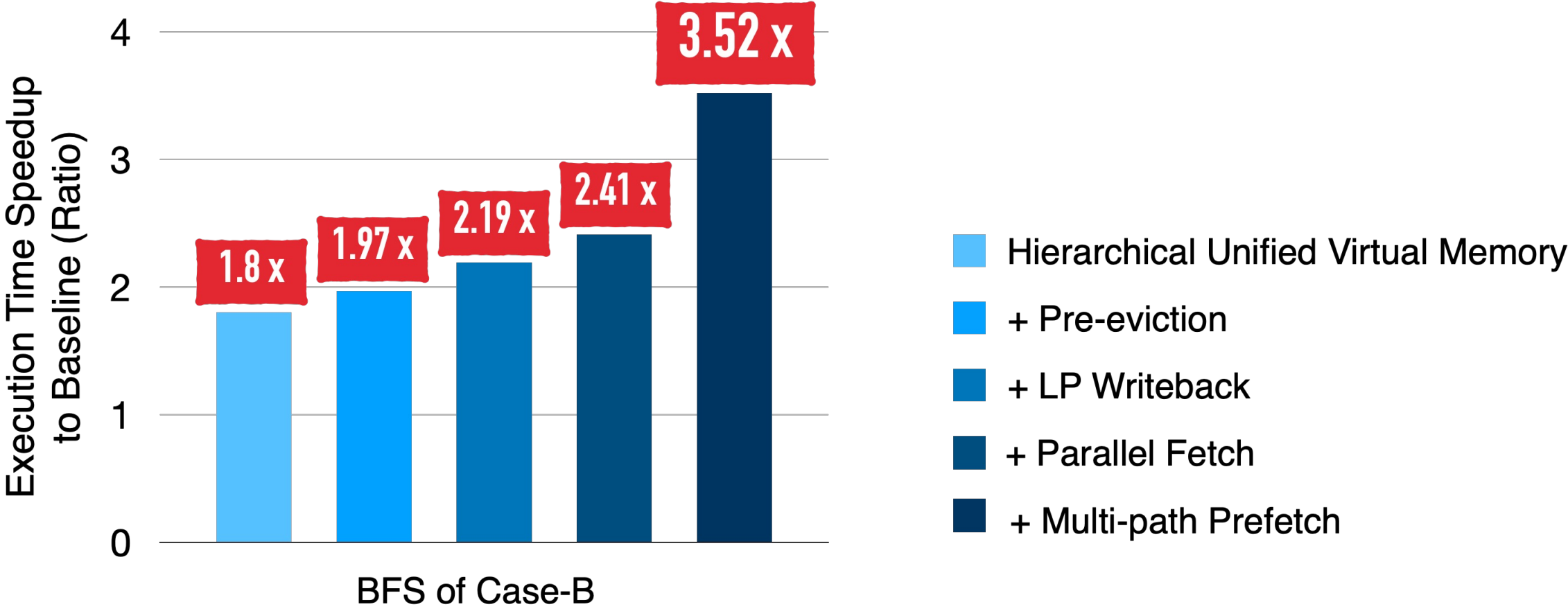
Case-B (2 harvests 2)



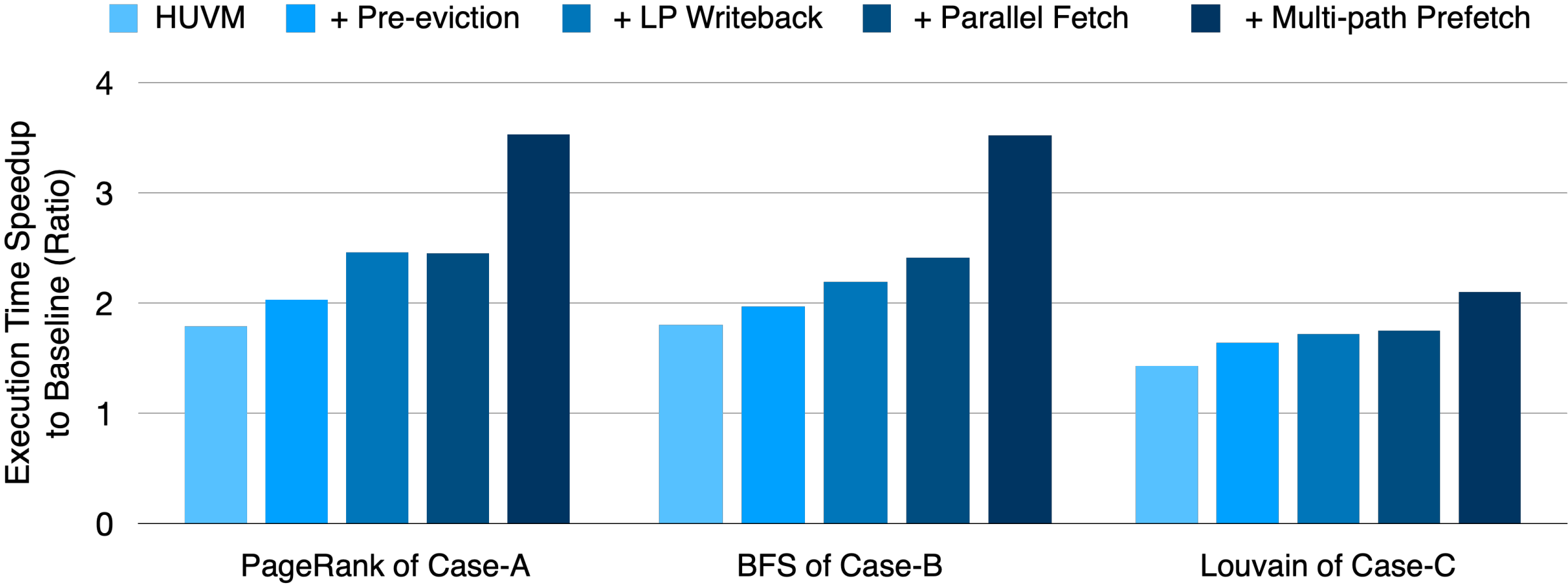
Case-C (3 harvests 1)



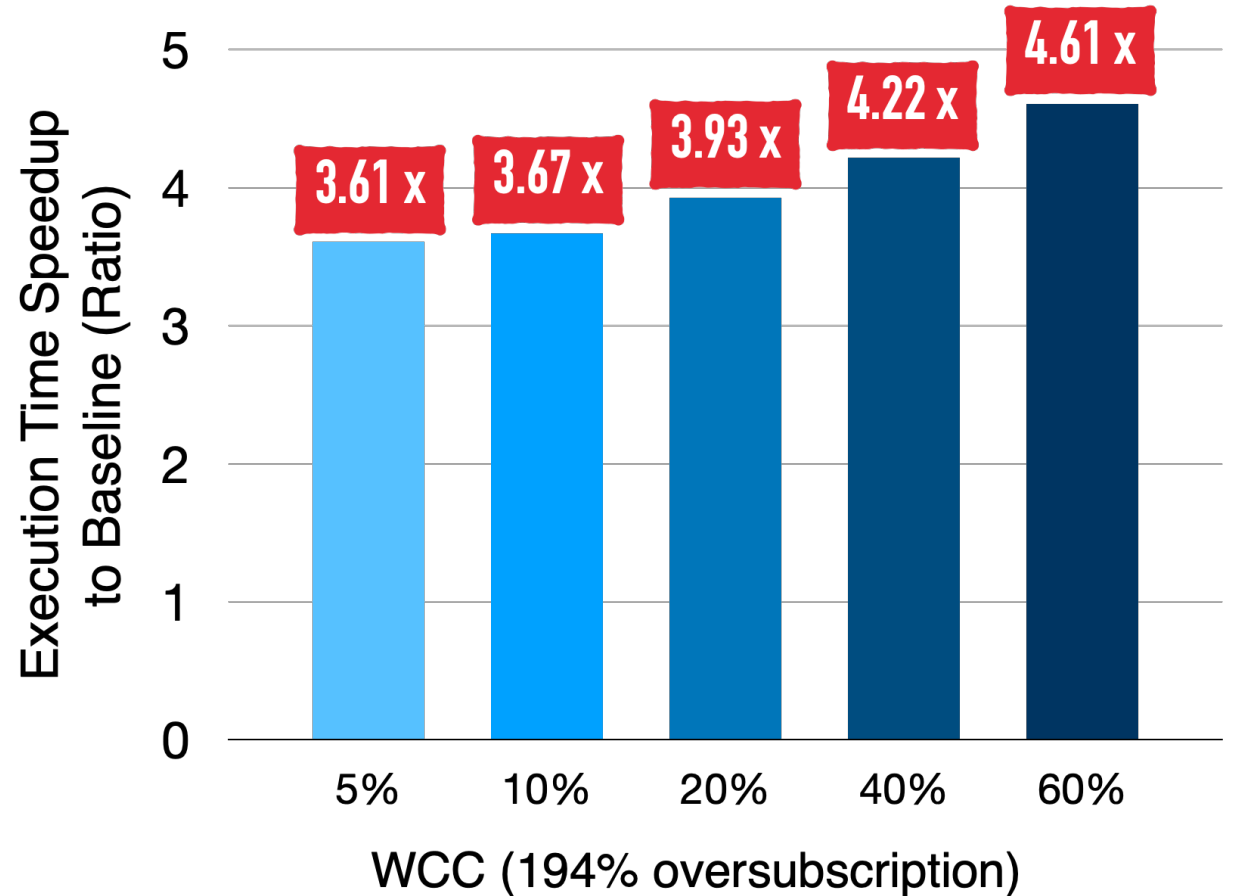
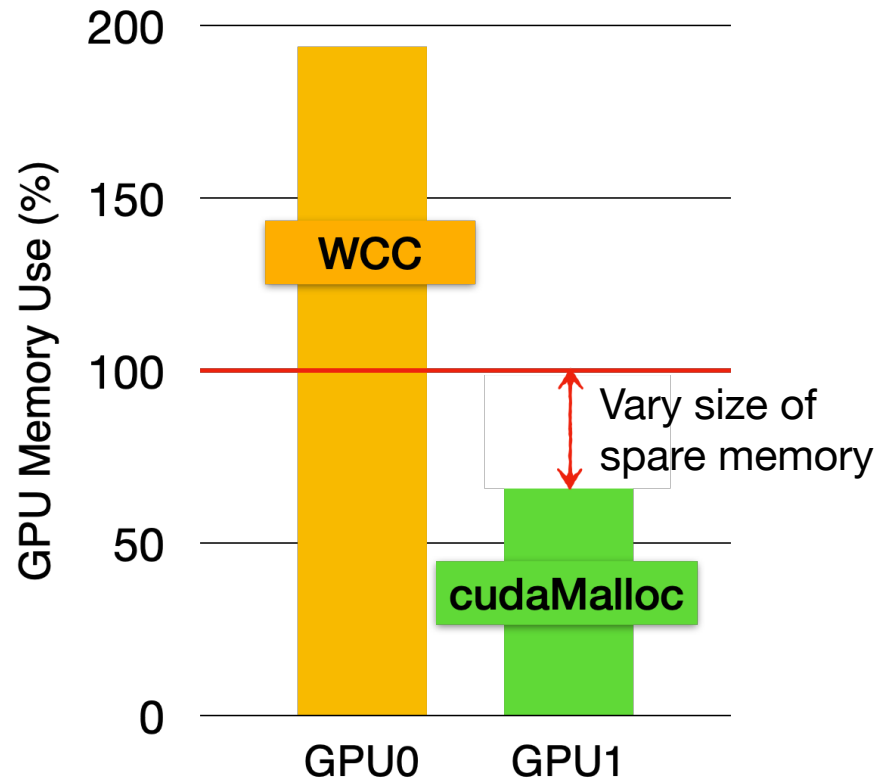
Individual technique performance breakdown



Individual technique performance breakdown



Sensitivity to size of spare memory



Conclusion

Problem

Resource imbalance in shared multi-GPU server
Significant performance overhead for memory oversubscription

Approach

1. Spill fraction of oversubscription to neighbor GPUs
2. Build new data-path with fast interconnect (NVLink)
3. Propose memory manager for HUVM

Result

Hierarchically unified virtual memory (w/ memory manager)

Performance improvement compared to naive UVM

3x ↑

W/ small fraction of spare memory

W/ Less performance interference

< 10%