# Contiguitas: The Pursuit of Physical Memory Contiguity in Datacenters

**Kaiyang Zhao**, Kaiwen Xue, Ziqi Wang, Dan Schatzberg, Leon Yang,
Antonis Manousis, Johannes Weiner, Rik van Riel, Bikash Sharma, Chunqiang Tang,
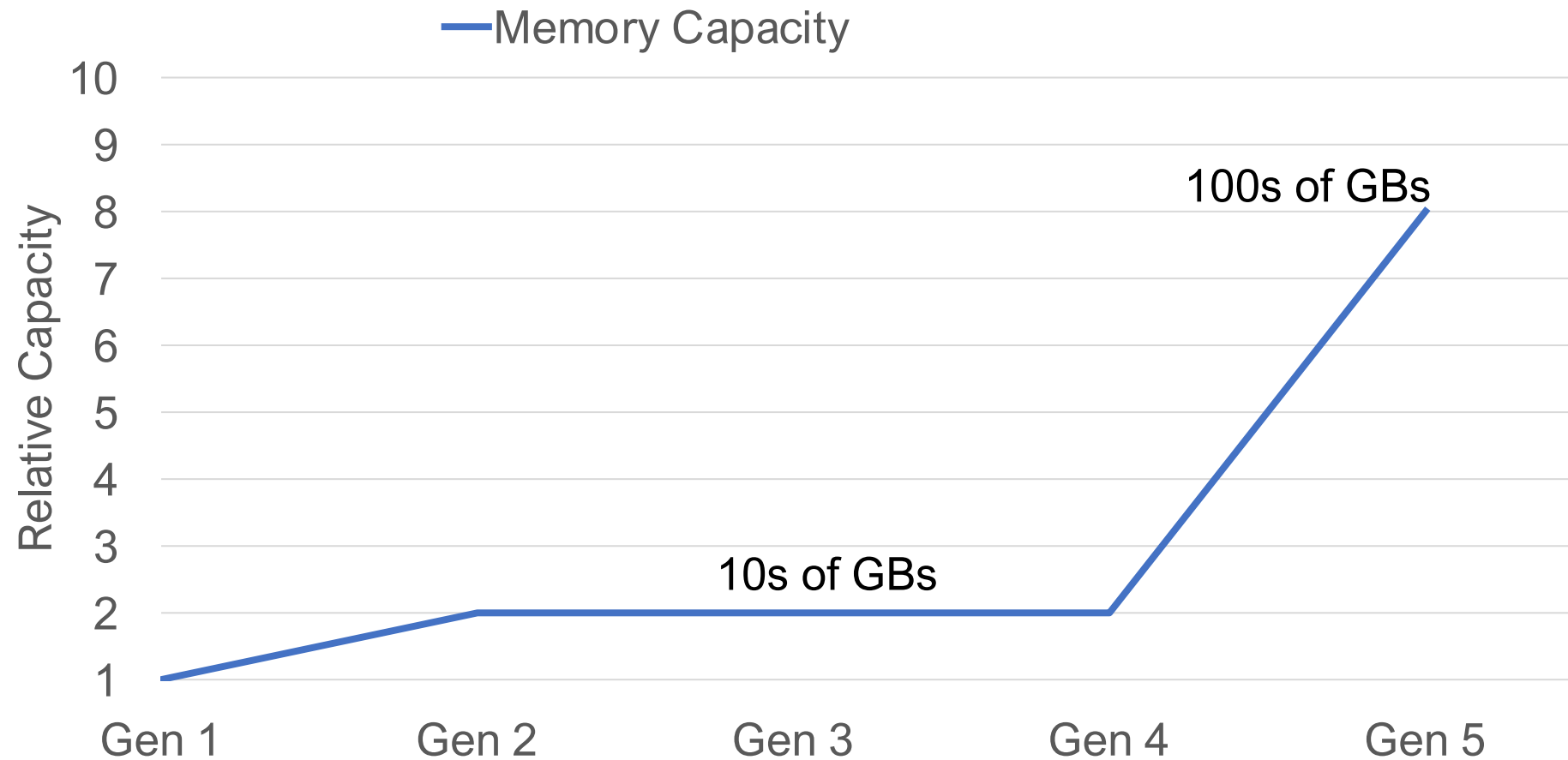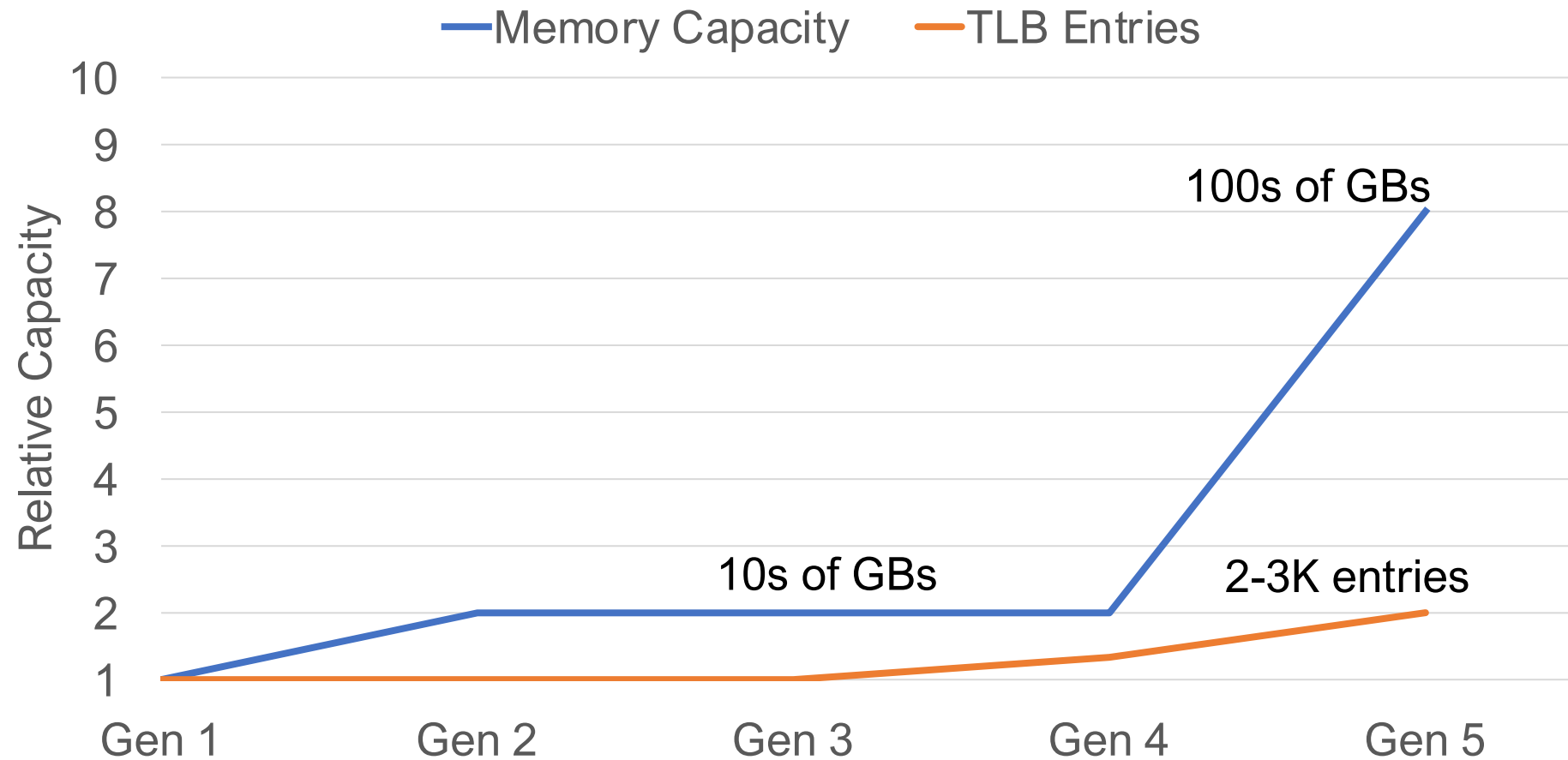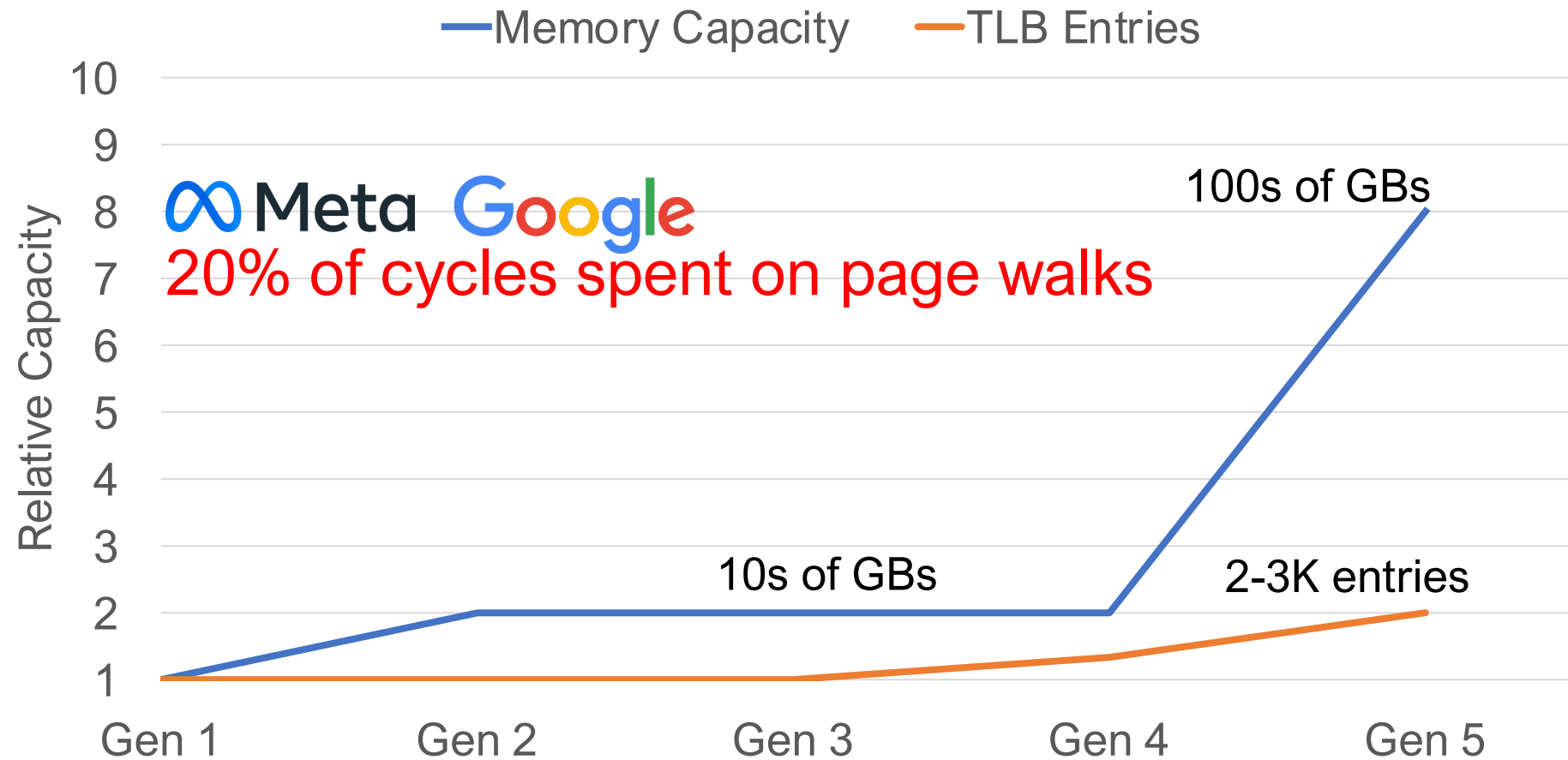Dimitrios Skarlatos

**Carnegie Mellon University**
Computer Science Department

Computer Architecture Operating Systems Group

# Memory Capacity Increases



Memory Capacity

Relative Capacity

100s of GBs

10s of GBs

Gen 1    Gen 2    Gen 3    Gen 4    Gen 5

Computer Architecture / Operating Systems Group

# TLB Does Not Scale

# The Virtual Memory Bottleneck



Meta Google
20% of cycles spent on page walks

100s of GBs

10s of GBs

2-3K entries

Relative Capacity

Memory Capacity — TLB Entries

Gen 1    Gen 2    Gen 3    Gen 4    Gen 5

Computer Architecture Operating Systems Group

# The Virtual Memory Bottleneck
# Will Only Get Worse



CXL Memory

Legend: Memory Capacity, TLB Entries

Y-axis: Relative Capacity (1–10)

X-axis: Gen 1, Gen 2, Gen 3, Gen 4, Gen 5, "Next Gen"

Computer Architecture / Operating Systems Group

# The Virtual Memory Bottleneck
## Will Only Get Worse



CXL Memory

Virtual Machines

Computer Architecture / Operating Systems Group

# The Virtual Memory Bottleneck
# Will Only Get Worse



CXL Memory · Virtual Machines · 5-level Page Tables

Computer Architecture Operating Systems Group

# Physical Memory Contiguity

Range of free space in memory that can back memory allocations

Available Contiguity

⬇

Larger Mappings

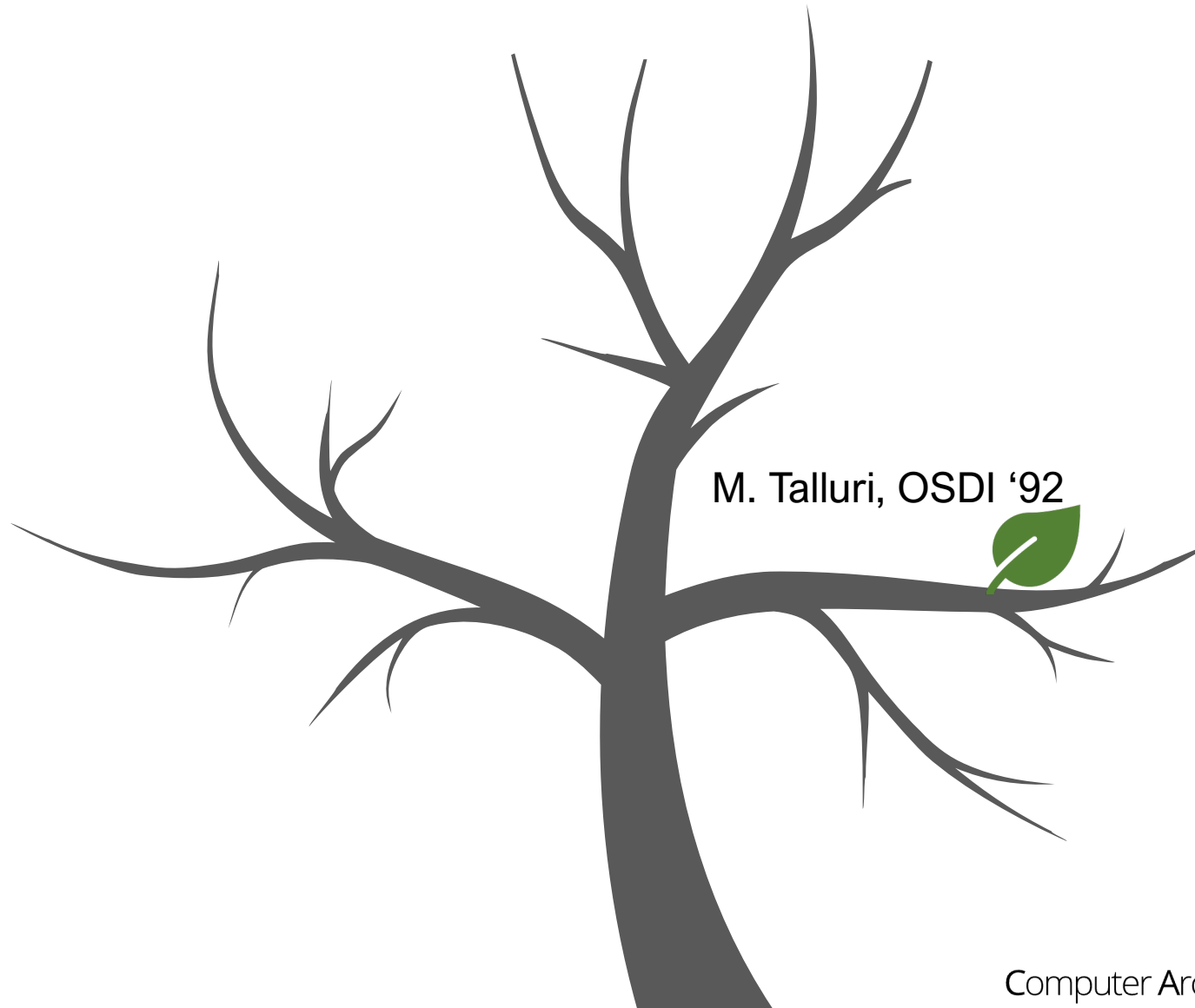Physical Address Space

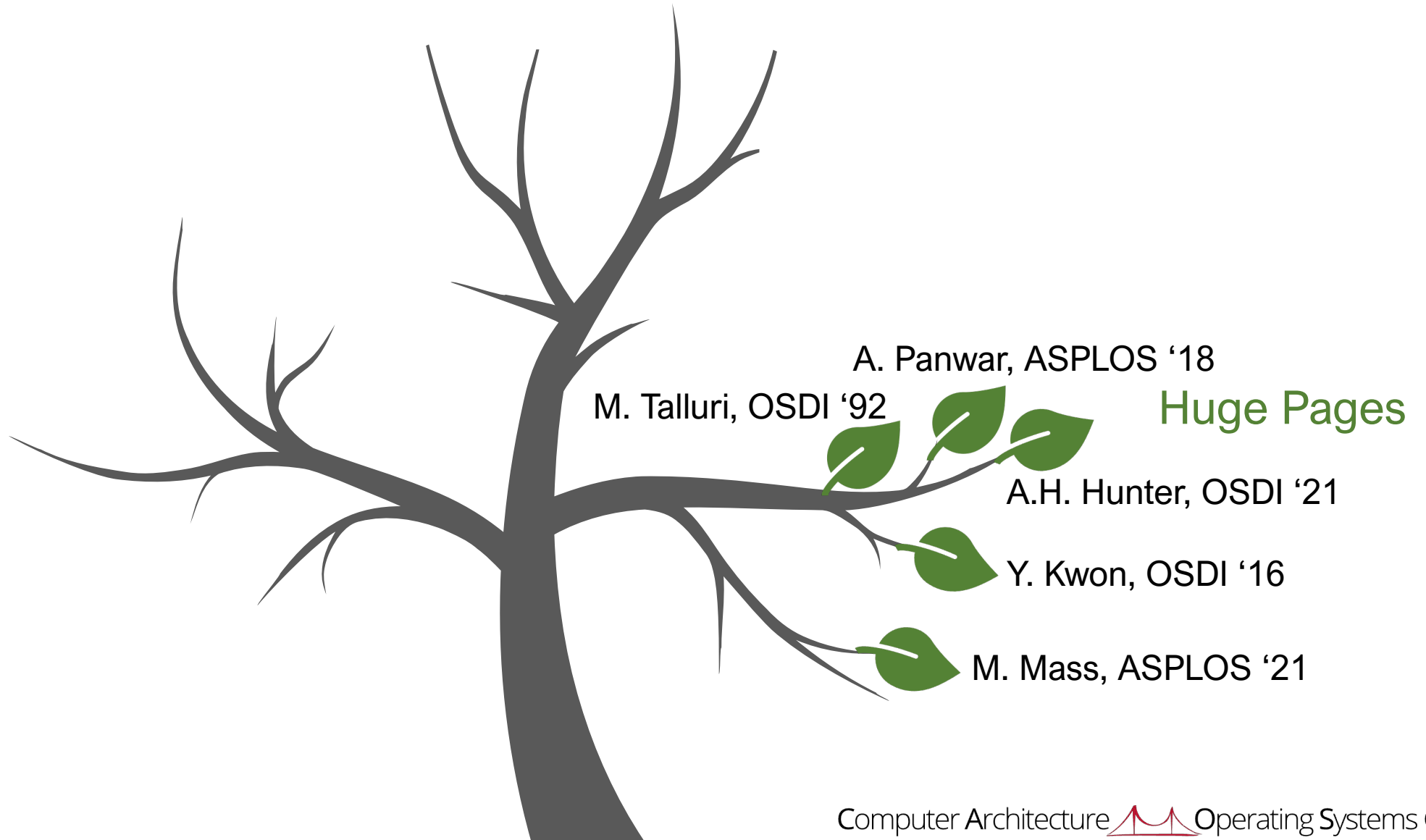Computer Architecture Operating Systems Group

# Physical Memory Contiguity

# Physical Contiguity for Huge Pages

M. Talluri, OSDI '92

Huge Pages

# Physical Contiguity for Huge Pages

A. Panwar, ASPLOS '18

M. Talluri, OSDI '92

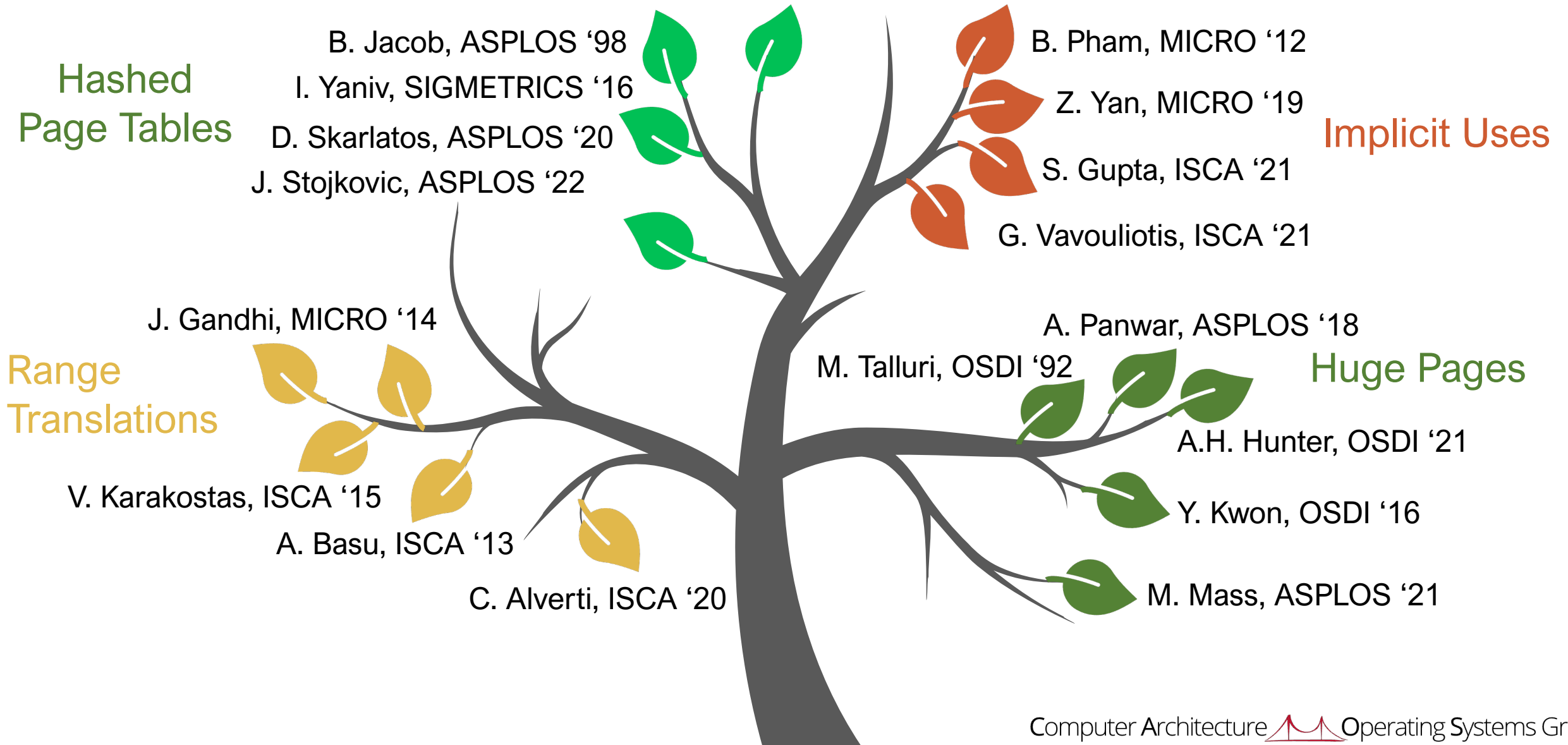Huge Pages

A.H. Hunter, OSDI '21

Y. Kwon, OSDI '16

M. Mass, ASPLOS '21

Computer Architecture Operating Systems Group

# Physical Contiguity for Range Translations

# Physical Contiguity for Hashed Page Tables



Hashed
Page Tables

B. Jacob, ASPLOS '98
I. Yaniv, SIGMETRICS '16
D. Skarlatos, ASPLOS '20
J. Stojkovic, ASPLOS '22

J. Gandhi, MICRO '14

Range
Translations

V. Karakostas, ISCA '15

A. Basu, ISCA '13

C. Alverti, ISCA '20

A. Panwar, ASPLOS '18

M. Talluri, OSDI '92

Huge Pages

A.H. Hunter, OSDI '21

Y. Kwon, OSDI '16

M. Mass, ASPLOS '21

Computer Architecture Operating Systems Group

# Physical Contiguity for Implicit Use Cases

# But is There Contiguity in Datacenters?

Computer Architecture Operating Systems Group

# But is There Contiguity in Datacenters?

Looking for contiguity across Meta's fleet



Datacenter
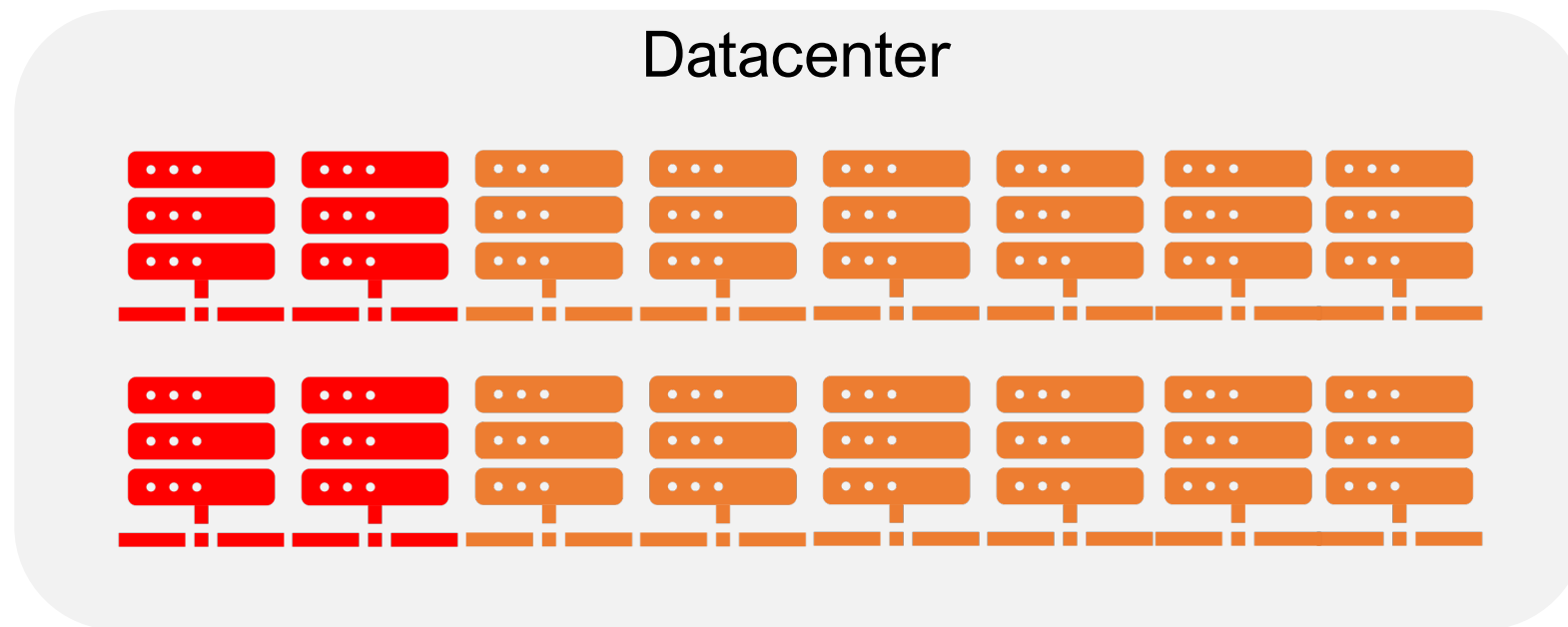
Computer Architecture Operating Systems Group

# Servers Are Highly Fragmented



Datacenter

# Servers Are Highly Fragmented

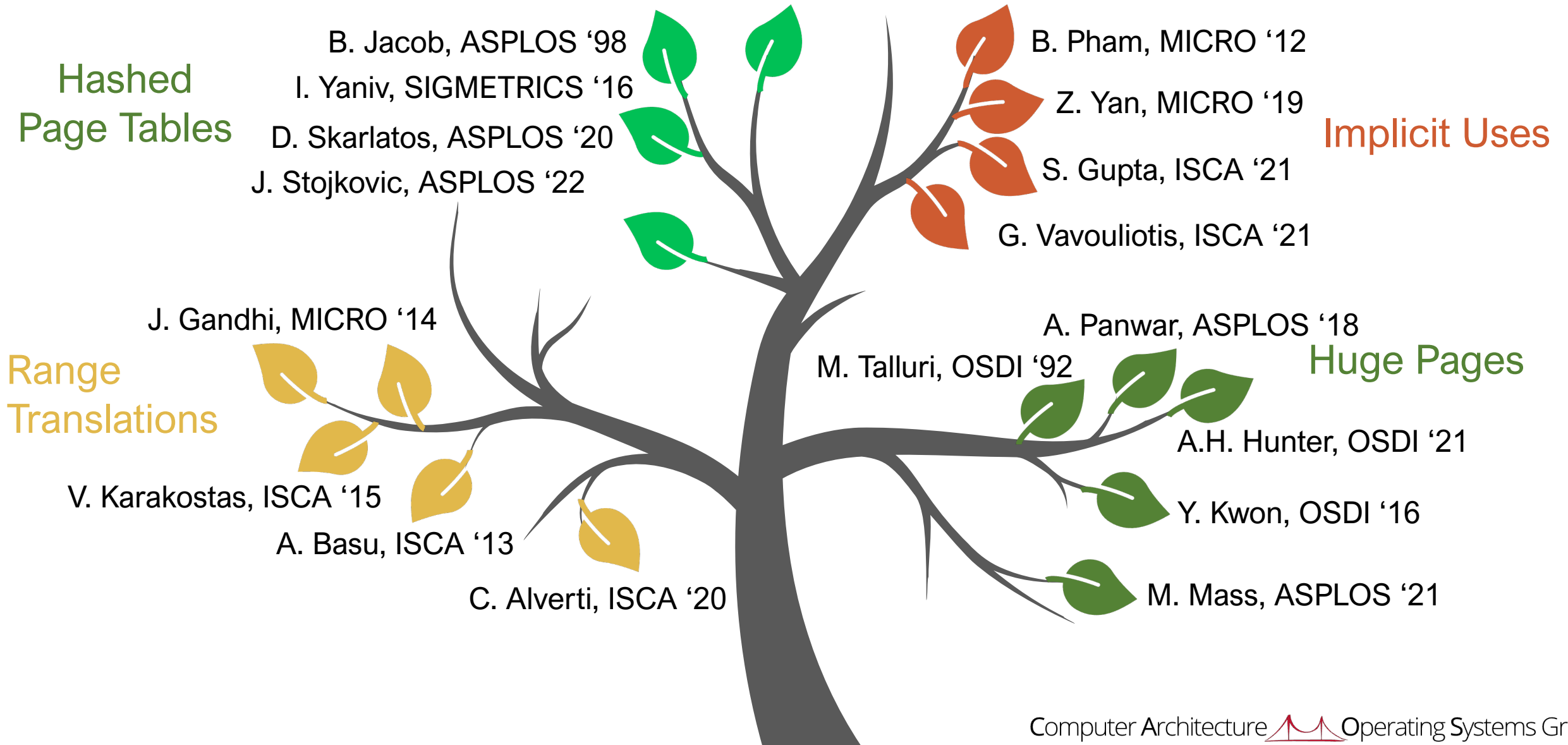A quarter of servers don't have contiguity for *even a single 2MB page*!



Datacenter

Computer Architecture Operating Systems Group

# Servers Are Highly Fragmented

Allocating 1GB pages is impossible in production!



Datacenter

Computer Architecture Operating Systems Group

# Without Physical Contiguity

Hashed
Page Tables

B. Jacob, ASPLOS '98

I. Yaniv, SIGMETRICS '16

D. Skarlatos, ASPLOS '20

J. Stojkovic, ASPLOS '22

B. Pham, MICRO '12

Z. Yan, MICRO '19

S. Gupta, ISCA '21

G. Vavouliotis, ISCA '21

Implicit Uses

J. Gandhi, MICRO '14

A. Panwar, ASPLOS '18

M. Talluri, OSDI '92

Huge Pages

Range
Translations

A.H. Hunter, OSDI '21

V. Karakostas, ISCA '15

Y. Kwon, OSDI '16

A. Basu, ISCA '13

C. Alverti, ISCA '20

M. Mass, ASPLOS '21

Computer Architecture Operating Systems Group

# Without Physical Contiguity

Address translation overheads will remain high!

# Contribution: Contiguitas

Computer Architecture Operating Systems Group

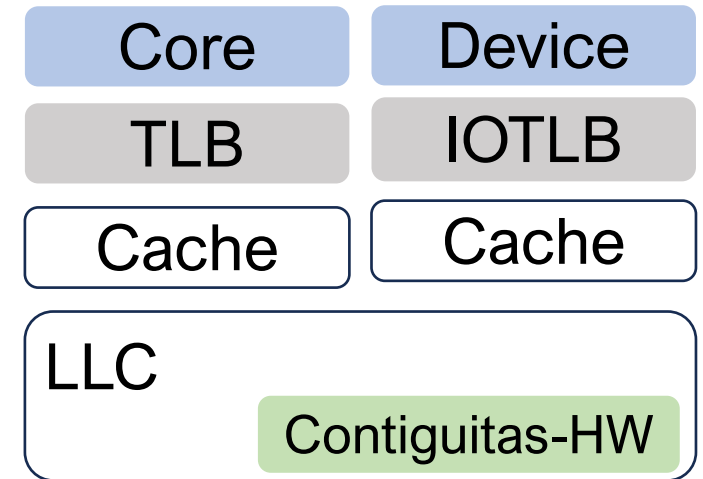# Contribution: Contiguitas

Unmovable pages are detrimental to contiguity

- I/O pages → Networking, RDMA, GPUs, Accelerators

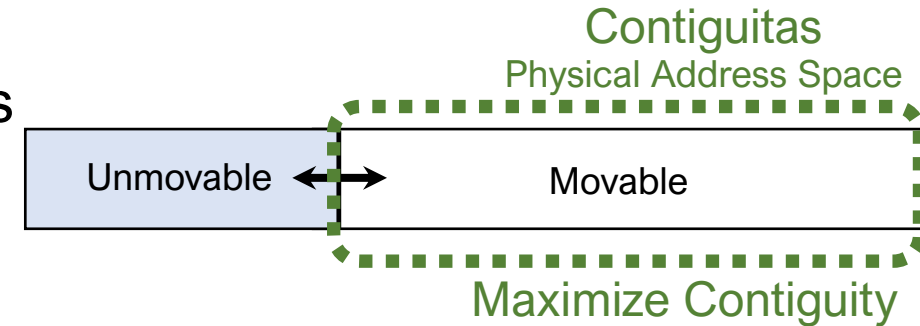Computer Architecture Operating Systems Group

# Contribution: Contiguitas

Unmovable pages are detrimental to contiguity

- I/O pages → Networking, RDMA, GPUs, Accelerators

Redesign OS memory management

Contiguitas
Physical Address Space

| Unmovable | ⟷ | Movable |

Maximize Contiguity

Computer Architecture Operating Systems Group

# Contribution: Contiguitas

Unmovable pages are detrimental to contiguity

- I/O pages → Networking, RDMA, GPUs, Accelerators

Redesign OS memory management
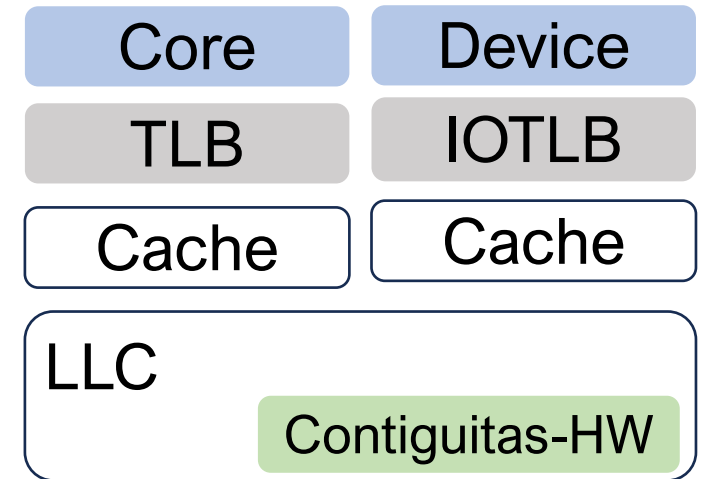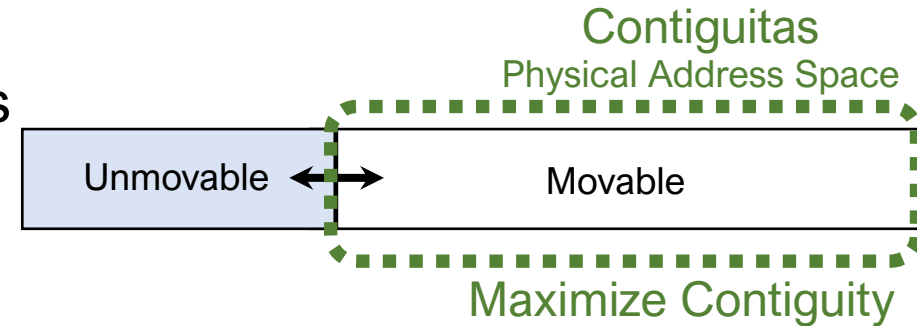
HW support for unmovable page migration

**Contiguitas**
Physical Address Space

| Unmovable | ← → | Movable |

Maximize Contiguity

| Core | Device |
| TLB | IOTLB |
| Cache | Cache |

LLC

Contiguitas-HW

Computer Architecture / Operating Systems Group

# Contribution: Contiguitas

Unmovable pages are detrimental to contiguity

- I/O pages → Networking, RDMA, GPUs, Accelerators

Redesign OS memory management

HW support for unmovable page migration

Ample physical memory contiguity

- ~90% of total memory

Performance gains 2-18% with production workloads

Efficiently reduce unmovable pages with HW

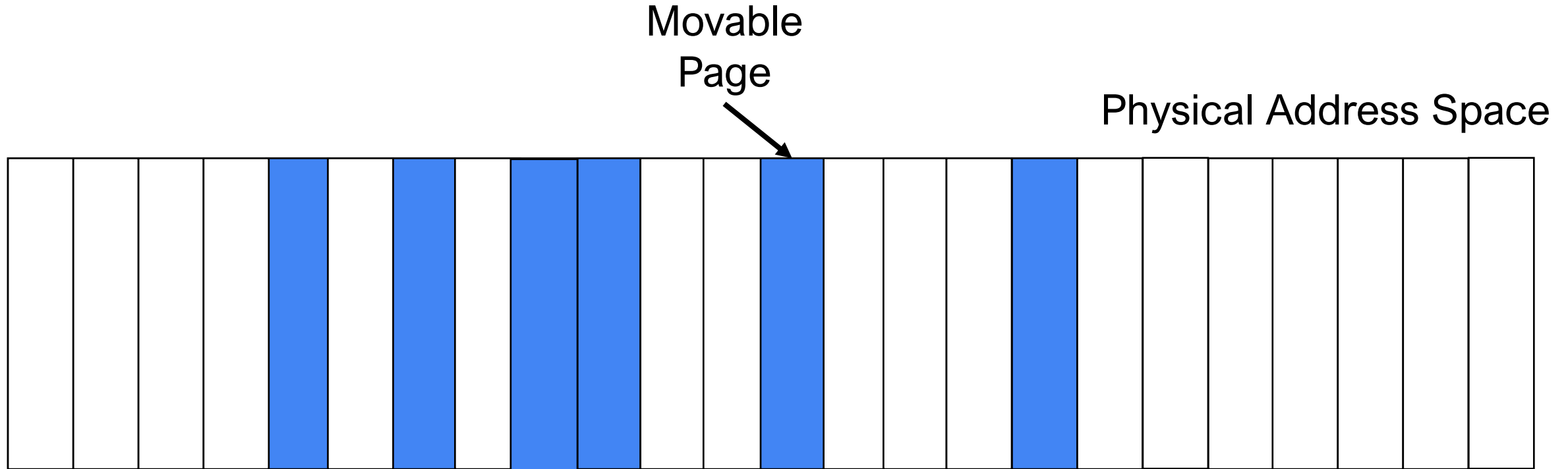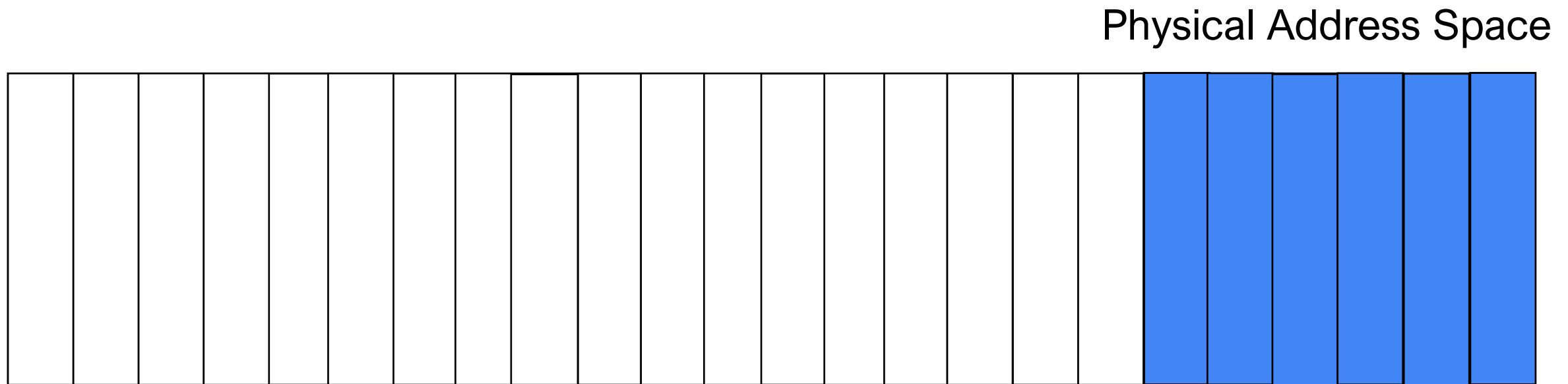In the process of upstreaming to Linux



Contiguitas
Physical Address Space

| Unmovable | Movable |

Maximize Contiguity

| Core | Device |
| TLB | IOTLB |
| Cache | Cache |

LLC

Contiguitas-HW

Computer Architecture ∧∧∧ Operating Systems Group

# Current Physical Address Space

Physical Address Space

Computer Architecture Operating Systems Group

# Current Physical Address Space

Movable
Page

Physical Address Space



Movable: Can be defragmented

Computer Architecture Operating Systems Group
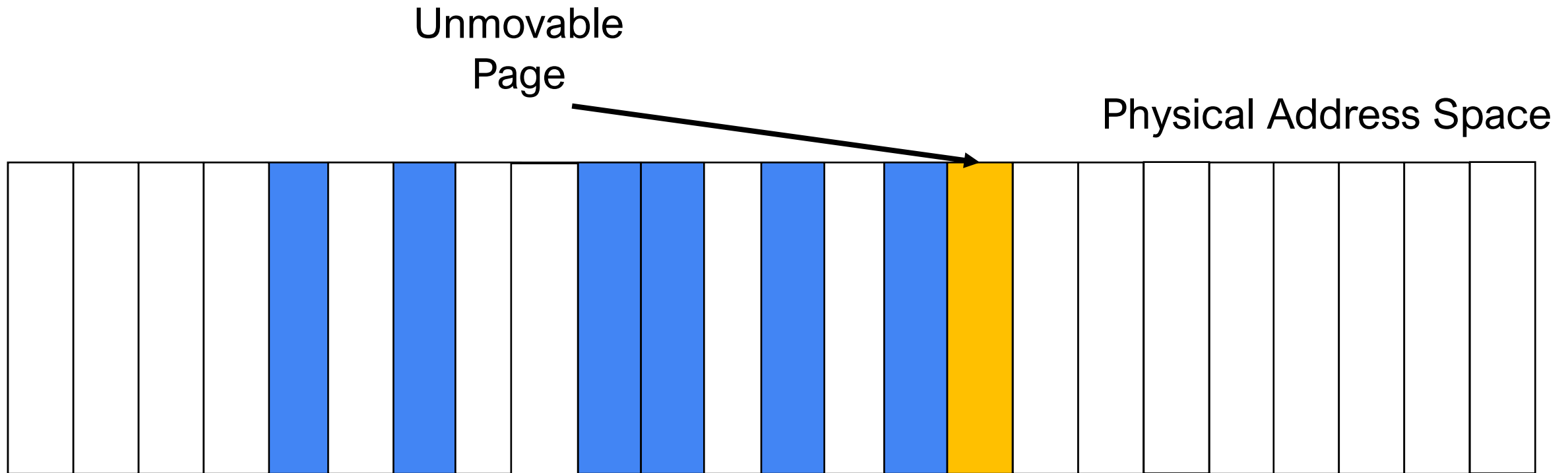
# Memory Defragmentation

Physical Address Space

Physical memory contiguity can be formed by defragmentation

Movable: Can be defragmented

Computer Architecture Operating Systems Group

# Unmovable Pages Block Contiguity

Unmovable
Page

Physical Address Space



Unmovable: Cannot be moved by the OS

Movable: Can be defragmented

Computer Architecture Operating Systems Group
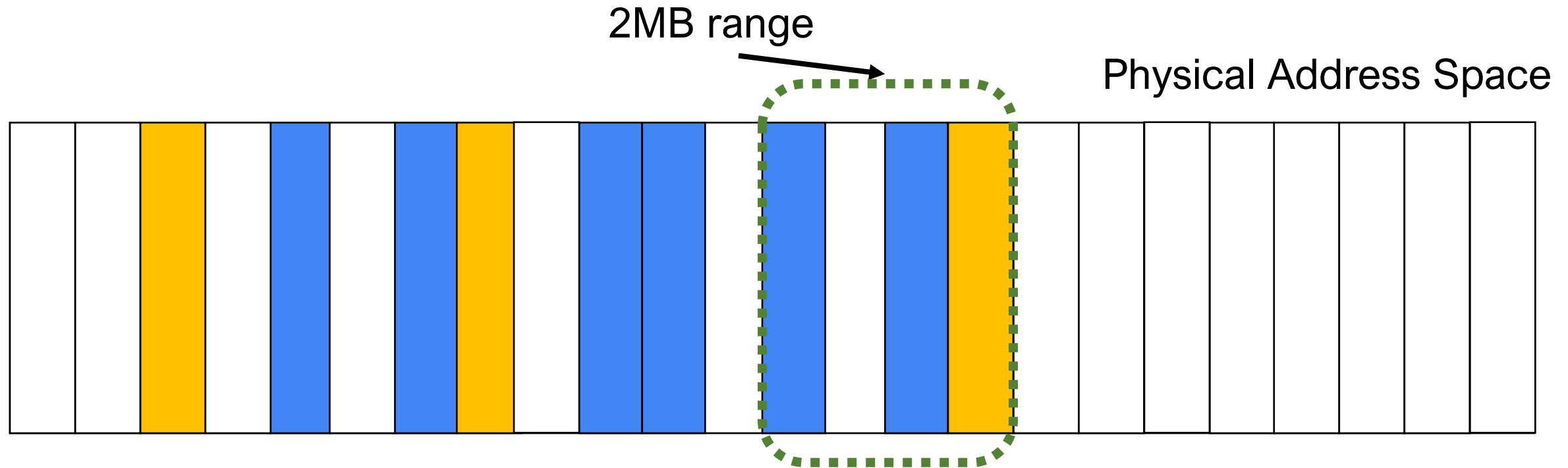
# Unmovable Pages Block Contiguity

Unmovable Page

Physical Address Space

Unmovable pages cannot be moved by the OS → Block contiguity!

Unmovable: Cannot be moved by the OS

Movable: Can be defragmented

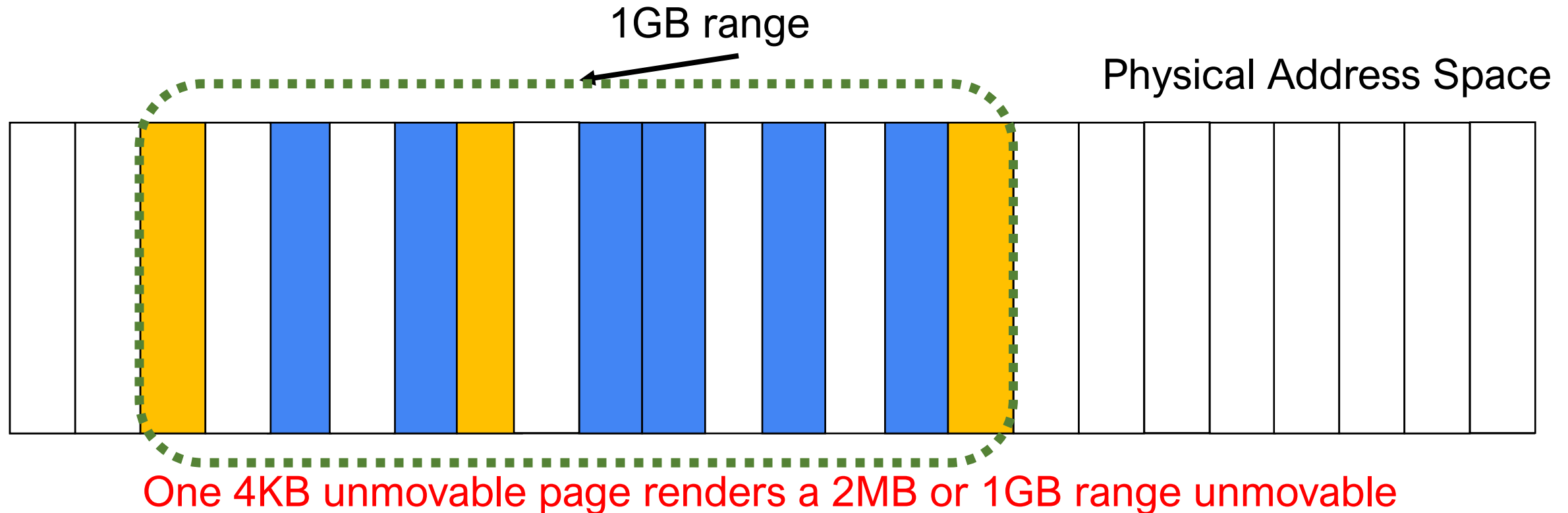Computer Architecture Operating Systems Group

# Unmovable Pages Block Contiguity

2MB range

Physical Address Space

One 4KB unmovable page renders a 2MB or 1GB range unmovable

Unmovable: Cannot be moved by the OS

Movable: Can be defragmented

Computer Architecture Operating Systems Group

# Unmovable Pages Block Contiguity

1GB range

Physical Address Space

One 4KB unmovable page renders a 2MB or 1GB range unmovable

☐ Unmovable: Cannot be moved by the OS

☐ Movable: Can be defragmented

Computer Architecture Operating Systems Group

# Unmovable Pages Block Contiguity

1GB range

Physical Address Space
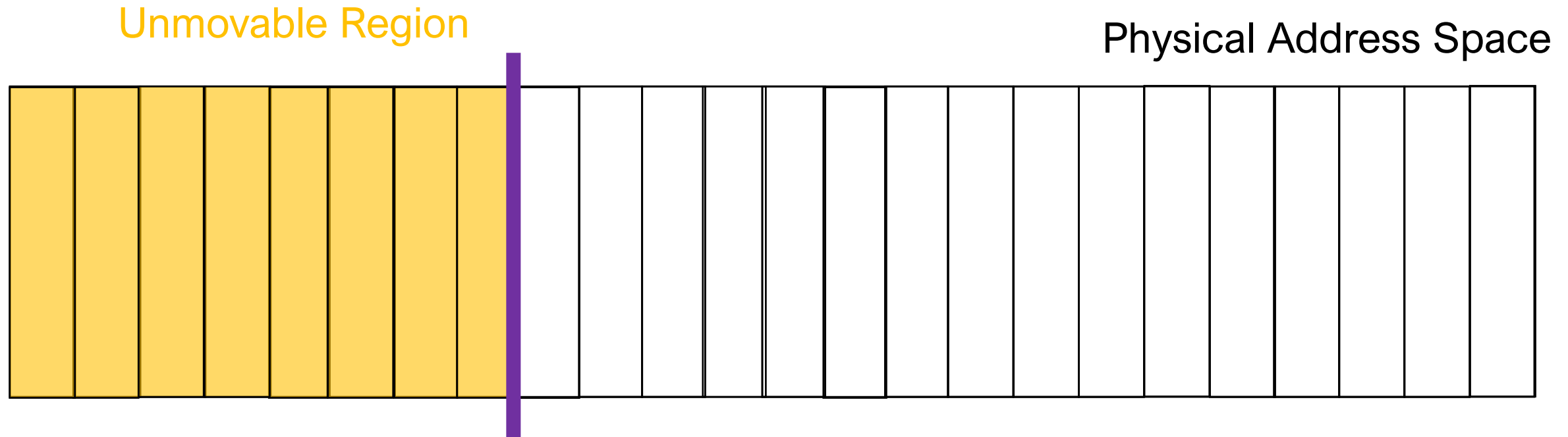


One 4KB unmovable page renders a 2MB or 1GB range unmovable

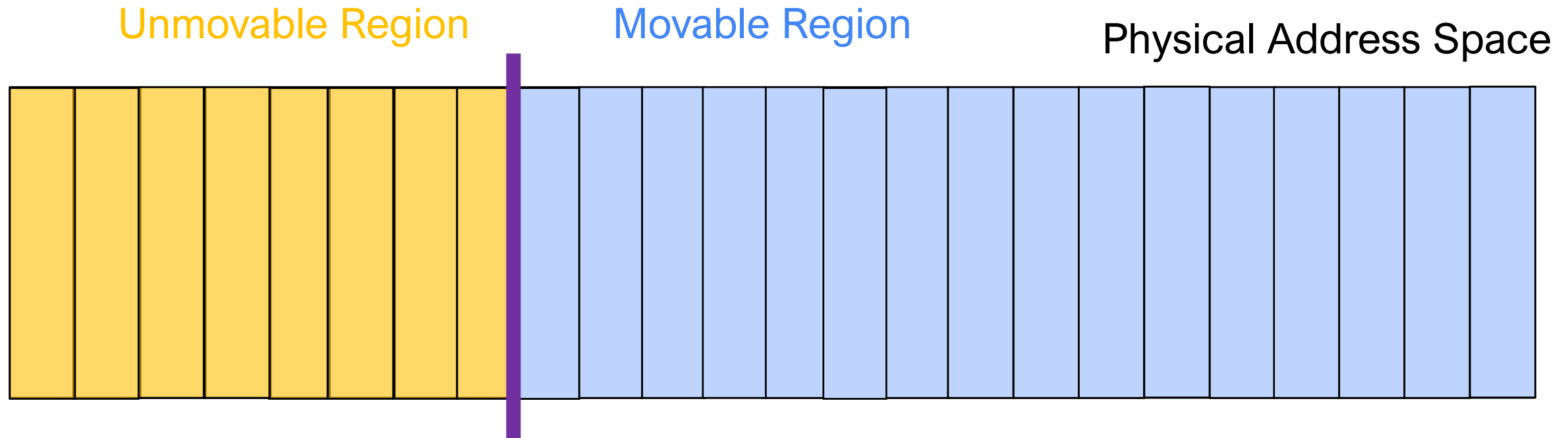Only 0.19% of bad allocations can fragment the whole memory irrecoverably

Computer Architecture Operating Systems Group

# Contiguitas Software:
# Ample Physical Contiguity by Design

Computer Architecture Operating Systems Group

# Separate the Address Space Into Regions

Unmovable Region

Physical Address Space

Confines unmovable allocations to the unmovable region

Computer Architecture Operating Systems Group

# Separate the Address Space Into Regions

# Challenge 1: Resizing Two Regions



Unmovable Region

Physical Address Space

Computer Architecture Operating Systems Group

# Challenge 1: Resizing Two Regions

Unmovable Region

Physical Address Space

Computer Architecture Operating Systems Group

# Challenge 1: Resizing Two Regions

Allocate

Unmovable Region

Physical Address Space

If either region runs out of memory → allocations block

Computer Architecture Operating Systems Group

# Solution: New Pressure Metric

Pressure: Time spent freeing up memory

Physical Address Space

Computer Architecture Operating Systems Group

# Solution: New Pressure Metric

Physical Address Space

# Solution: New Pressure Metric

Physical Address Space

Computer Architecture Operating Systems Group

# Solution: New Pressure Metric

Physical Address Space

# Solution: Resize in the Background



Pressure high!

Physical Address Space

A background thread resizes proactively off the critical path

Computer Architecture Operating Systems Group

# Solution: Resize in the Background

Physical Address Space

A background thread resizes proactively off the critical path

Computer Architecture Operating Systems Group

# Challenge 2: Fragmentation in Unmovable

# Solution: Allocation Policy

Kernel code,
long-lived kernel data
structures

Physical Address Space

More challenges are discussed in the paper:

Resizing Algorithm

Pressure Tracking

Prefers allocations far from the boundary

Computer Architecture ⋀⋀ Operating Systems Group

# The Need to Reduce Unmovable Pages

Most unmovable pages come from I/O

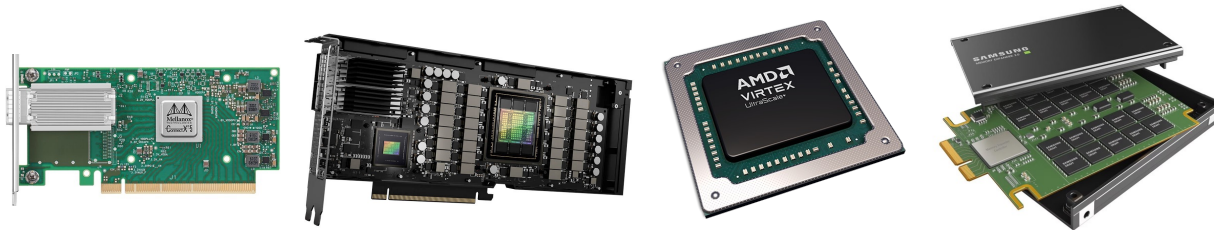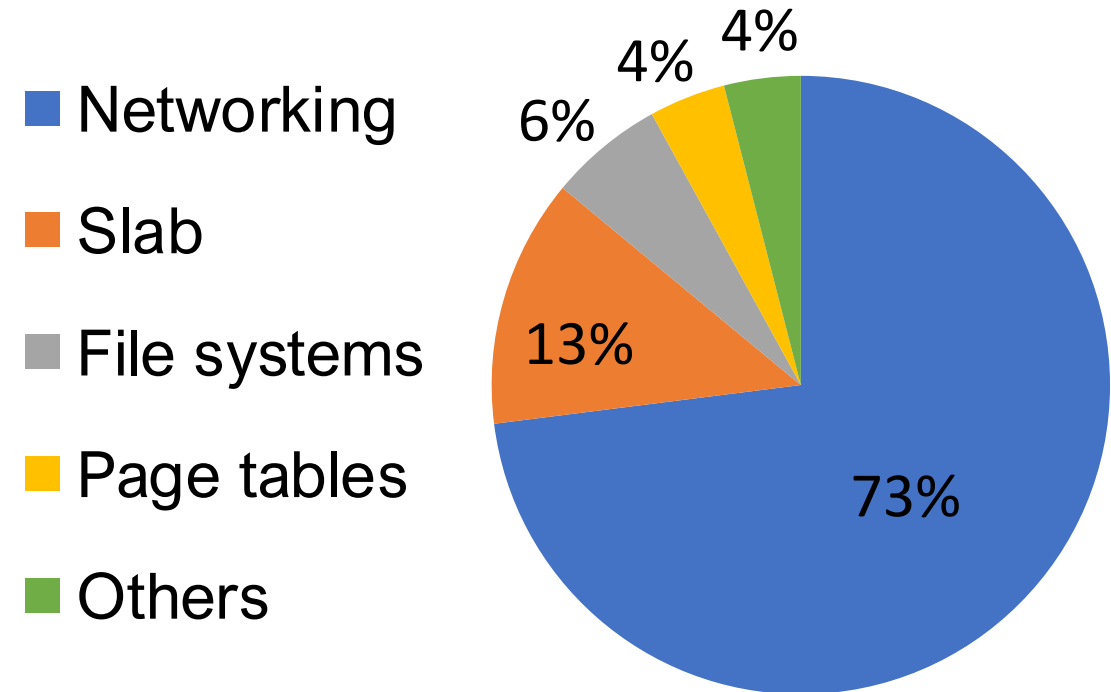More in the future from:

RDMA, GPU, accelerators, CXL



## Sources of Unmovable Pages at Meta



- Networking
- Slab
- File systems
- Page tables
- Others

4%
4%
6%
13%
73%

Unmovable allocations will get worse!

Computer Architecture Operating Systems Group

# Why I/O Pages are Unmovable

Computer Architecture Operating Systems Group

# Why I/O Pages are Unmovable

Core    Core    Device

Computer Architecture    Operating Systems Group

# Why I/O Pages are Unmovable

| Core | Core | Device |
|------|------|--------|
| **TLB** | **TLB** | **IOTLB** |
| VA PA | VA PA | VA PA |

# Why I/O Pages are Unmovable

| Core | Core | Device |
|------|------|--------|
| **TLB** VA PA | **TLB** VA PA | **IOTLB** VA PA |
| Cache | Cache | Cache |

LLC

Computer Architecture · Operating Systems Group

# Why I/O Pages are Unmovable

Software Page Migration
* Access must be blocked during migration

| Core | Core | Device |
|------|------|--------|
| TLB | TLB | IOTLB |
| VA PA | VA PA | VA PA |

| Cache | Cache | Cache |
|-------|-------|-------|

| LLC |
|-----|

Computer Architecture ⋀⋀ Operating Systems Group
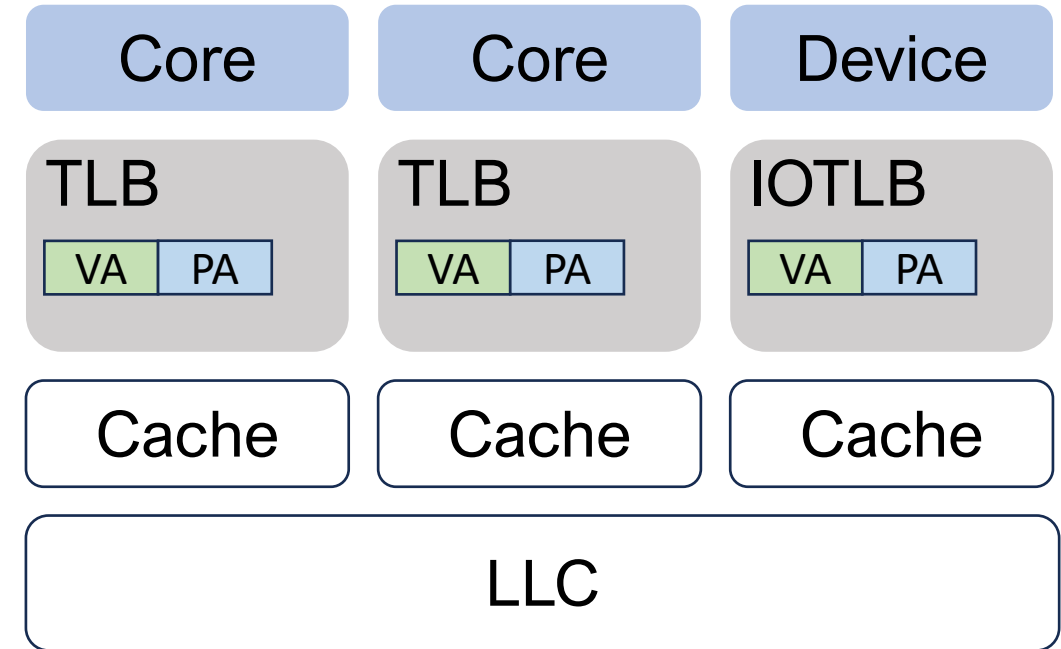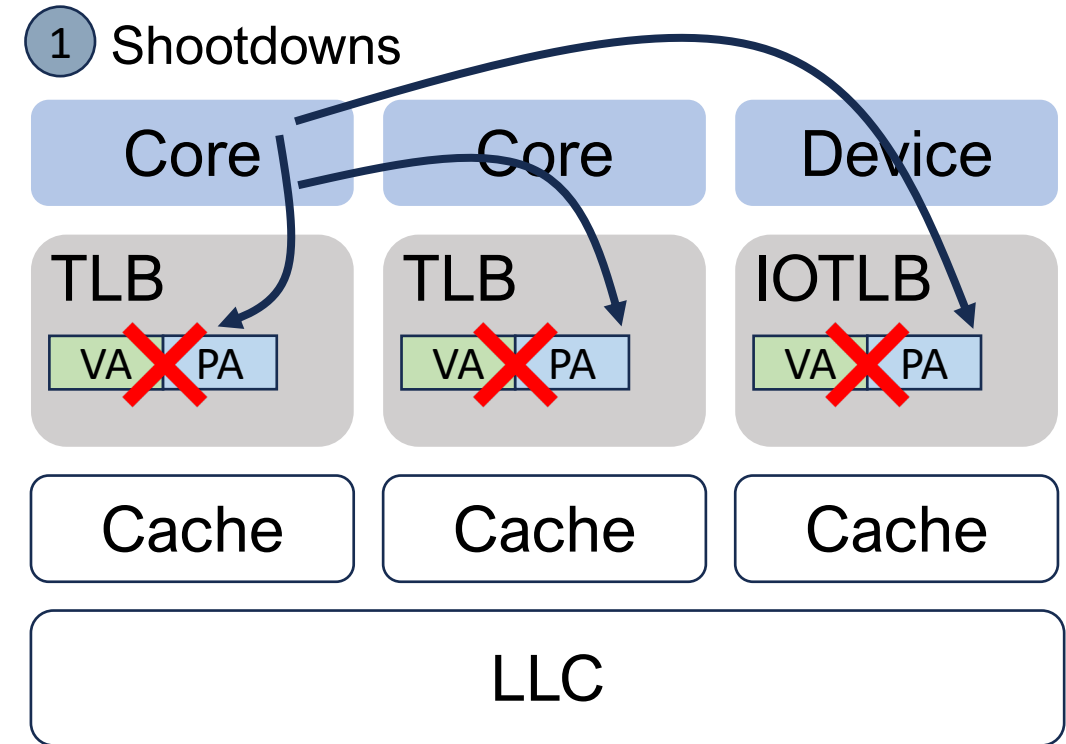
# Why I/O Pages are Unmovable

Software Page Migration
- Access must be blocked during migration

# Why I/O Pages are Unmovable

Software Page Migration
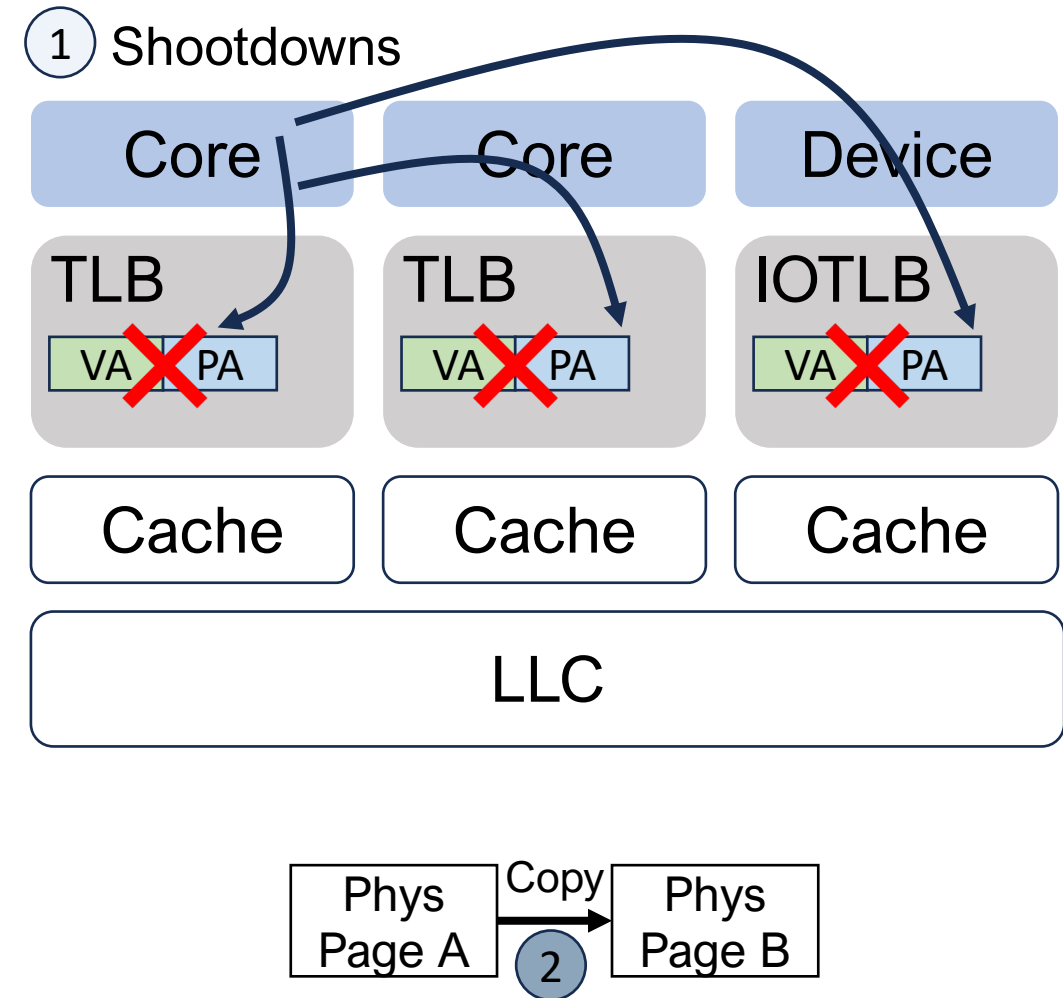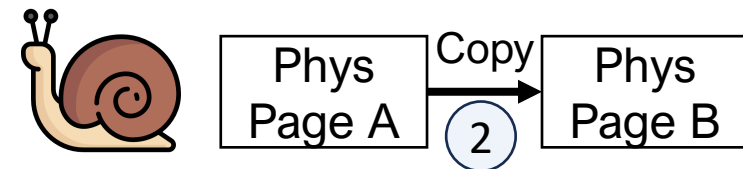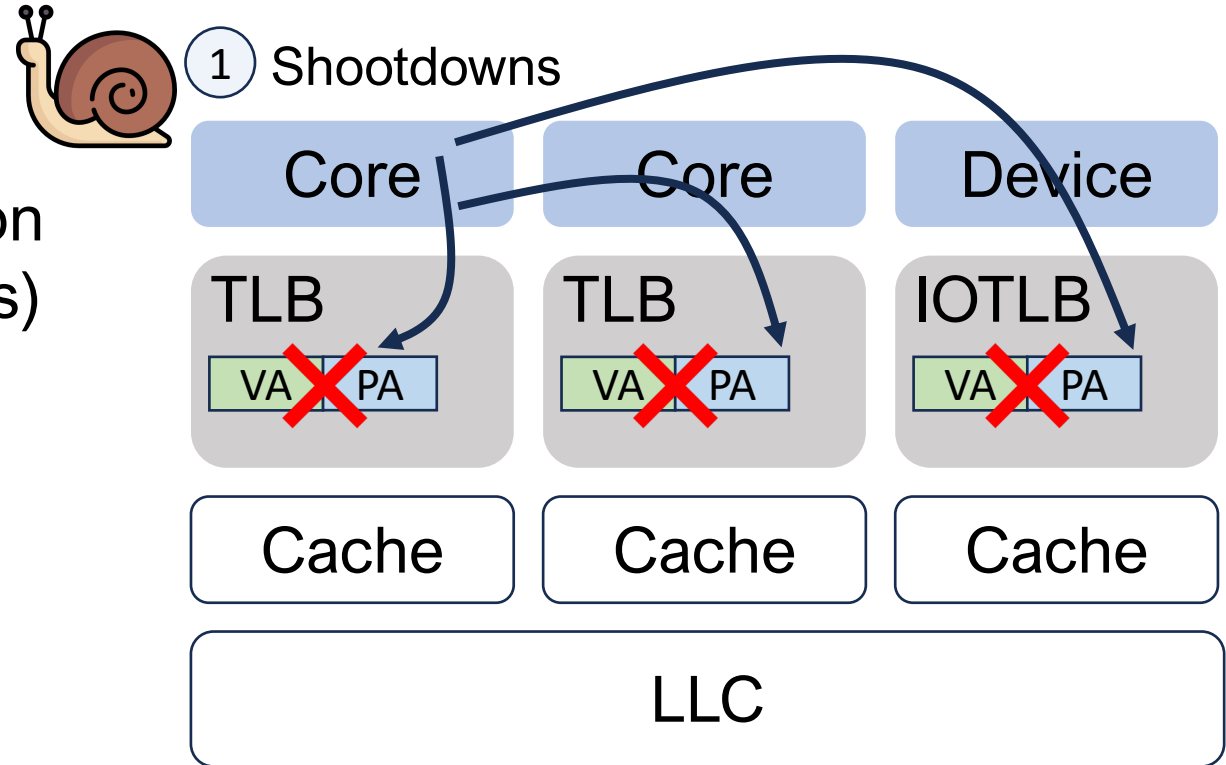- Access must be blocked during migration

① Shootdowns

| Core | Core | Device |
|------|------|--------|
| TLB | TLB | IOTLB |
| VA ✗ PA | VA ✗ PA | VA ✗ PA |
| Cache | Cache | Cache |

LLC

Phys Page A → Copy → Phys Page B ②

Computer Architecture ⩙ Operating Systems Group

# Why I/O Pages are Unmovable

Software Page Migration
- Access must be blocked during migration
- **Impossible** for many devices (page faults)

Computer Architecture / Operating Systems Group

# Why I/O Pages are Unmovable

Software Page Migration
- Access must be blocked during migration
- Impossible for many devices (page faults)
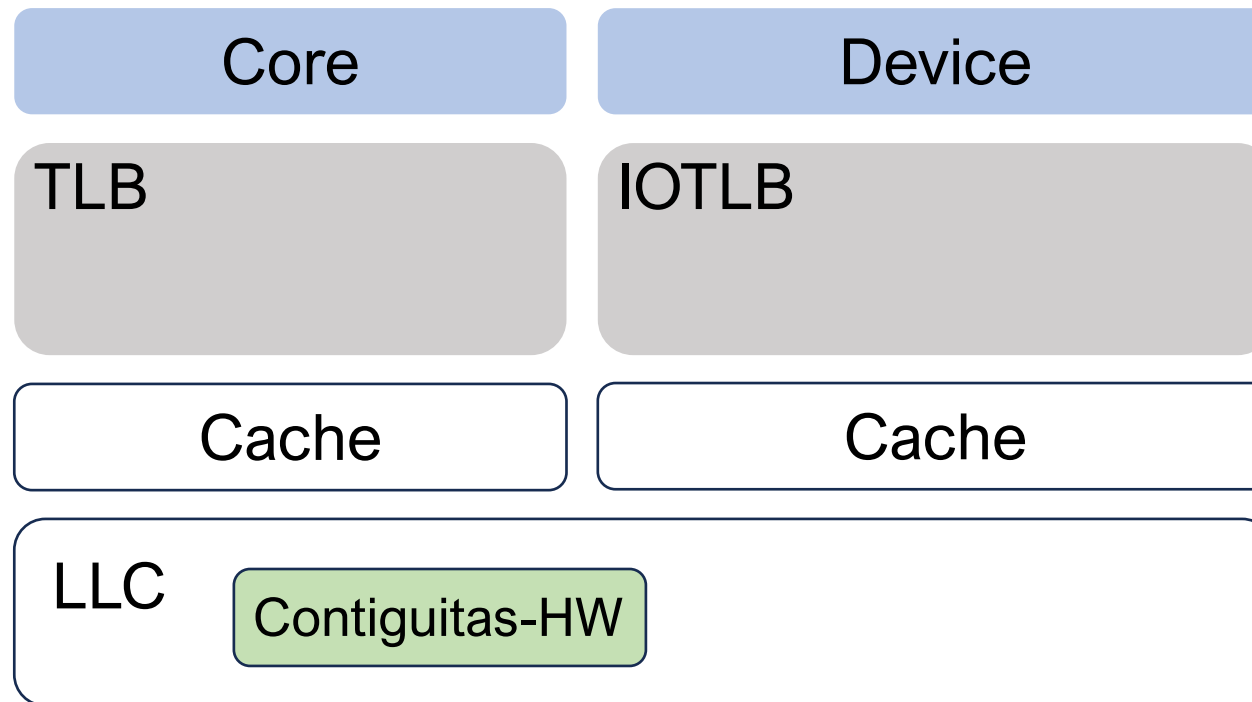- Long page unavailable time

Computer Architecture Operating Systems Group

# Contiguitas Hardware:
# Enabling Migration of I/O Pages

Computer Architecture Operating Systems Group

# Contiguitas Hardware

Contiguitas-HW in LLC provides transparent page migration

| Core | Device |
|------|--------|
| TLB | IOTLB |
| Cache | Cache |

LLC   Contiguitas-HW

Computer Architecture / Operating Systems Group

# Transparent Page Migration

Computer Architecture ⌂⌂ Operating Systems Group

# Transparent Page Migration

Core

Device

TLB

VA | PA

IOTLB

VA | PA

Cache

Cache

LLC

Contiguitas-HW

Line copied?

No

Yes

Copy

Src Page

Dst Page

Migrated Line

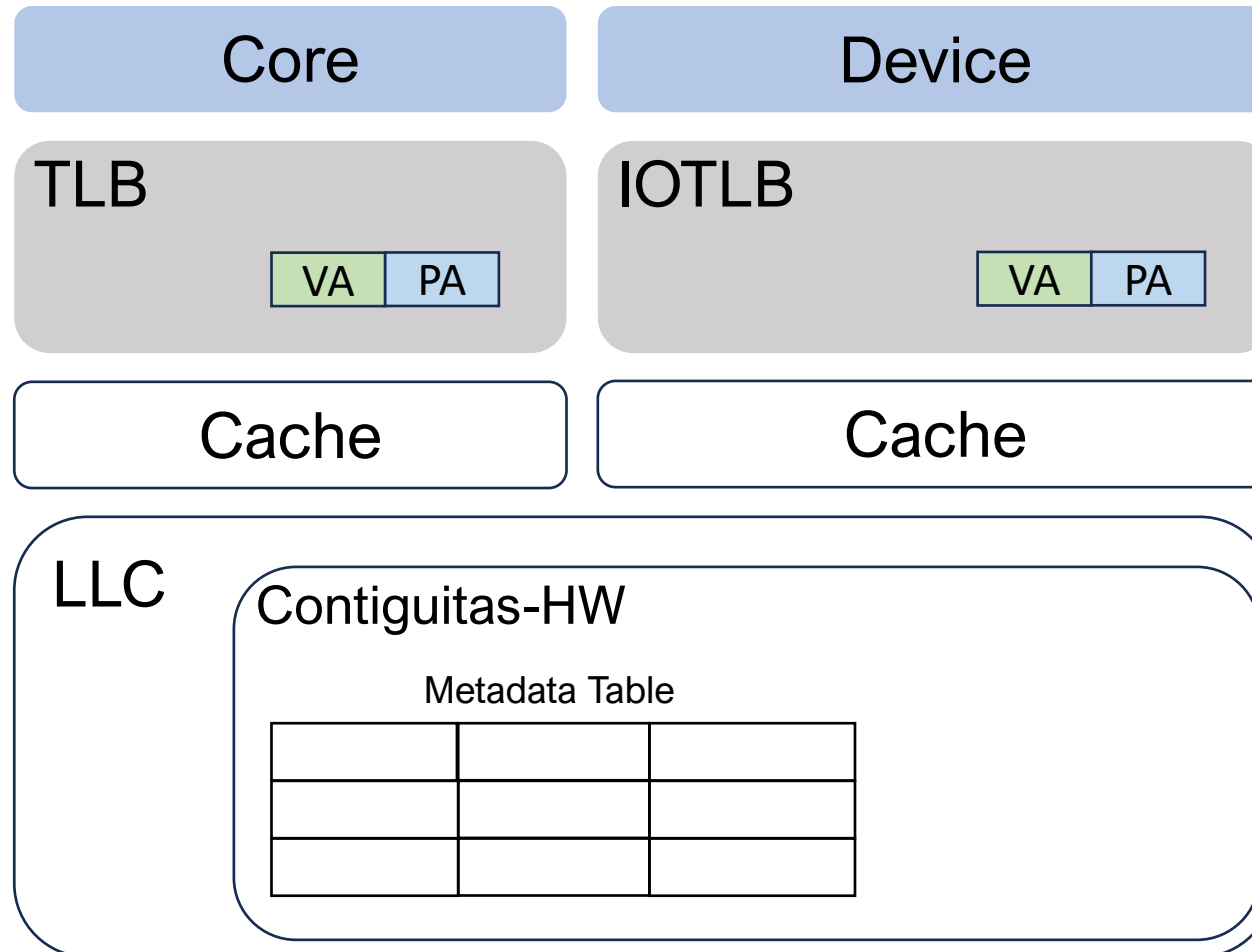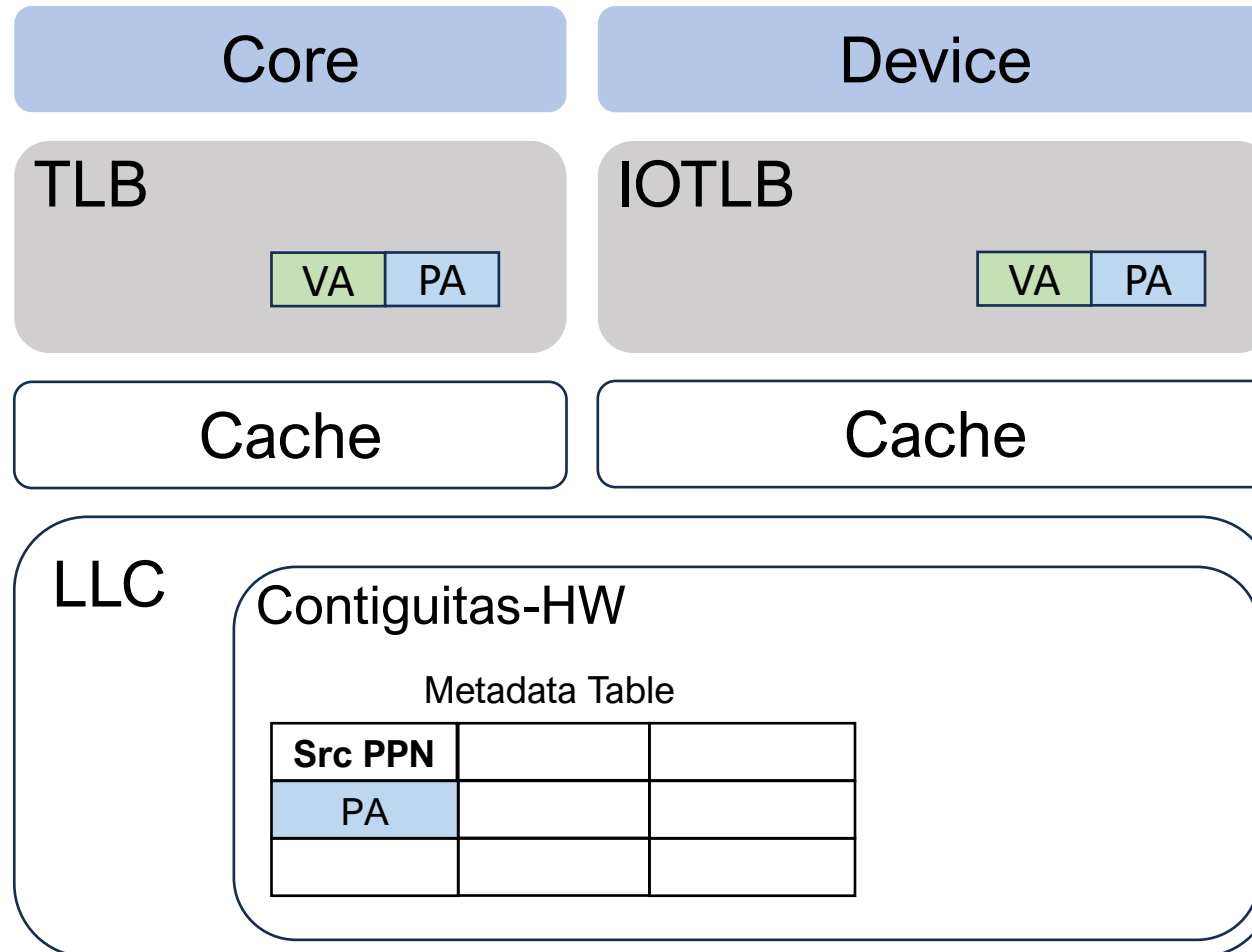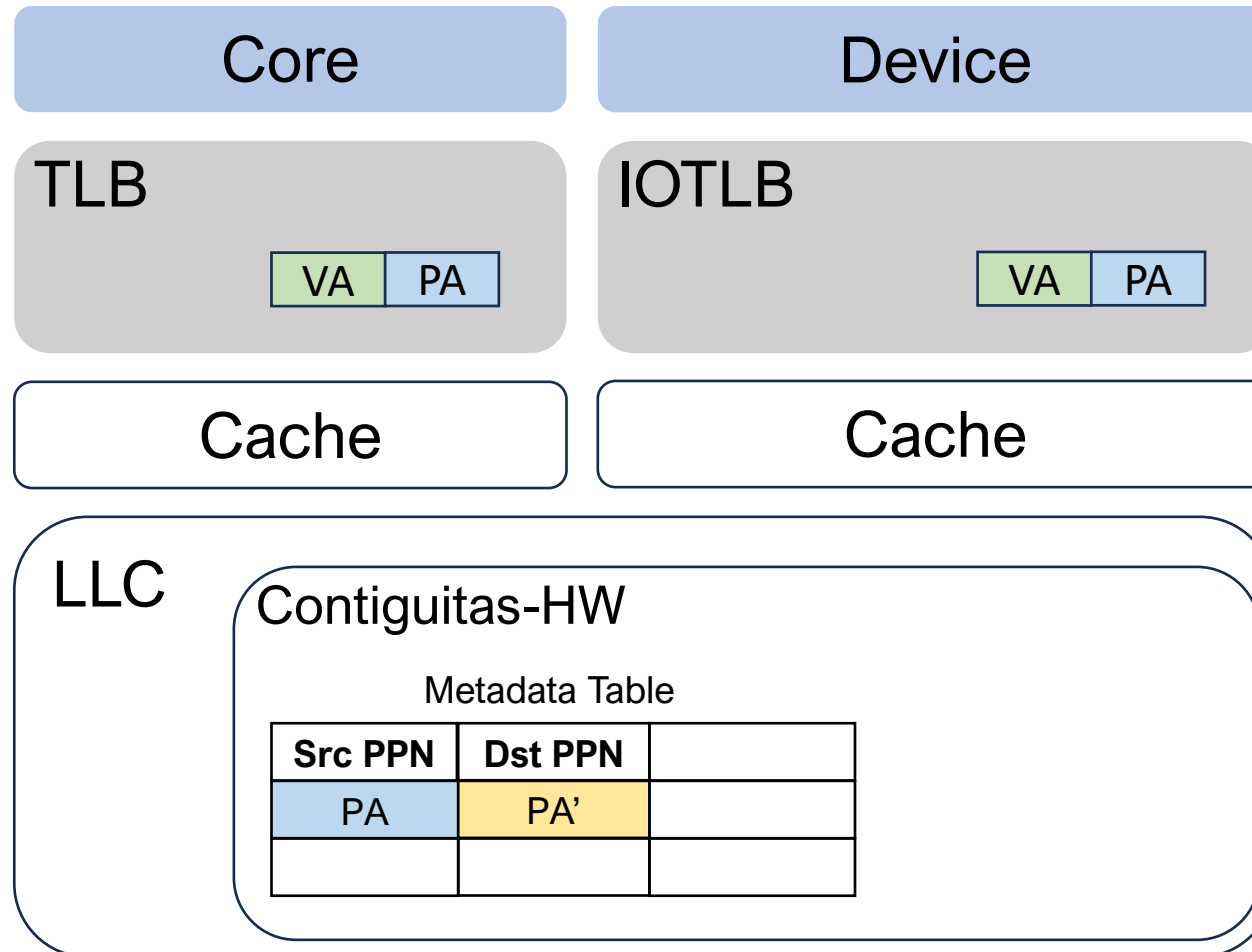Computer Architecture ⌒⌒ Operating Systems Group

# Transparent Page Migration

# Transparent Page Migration

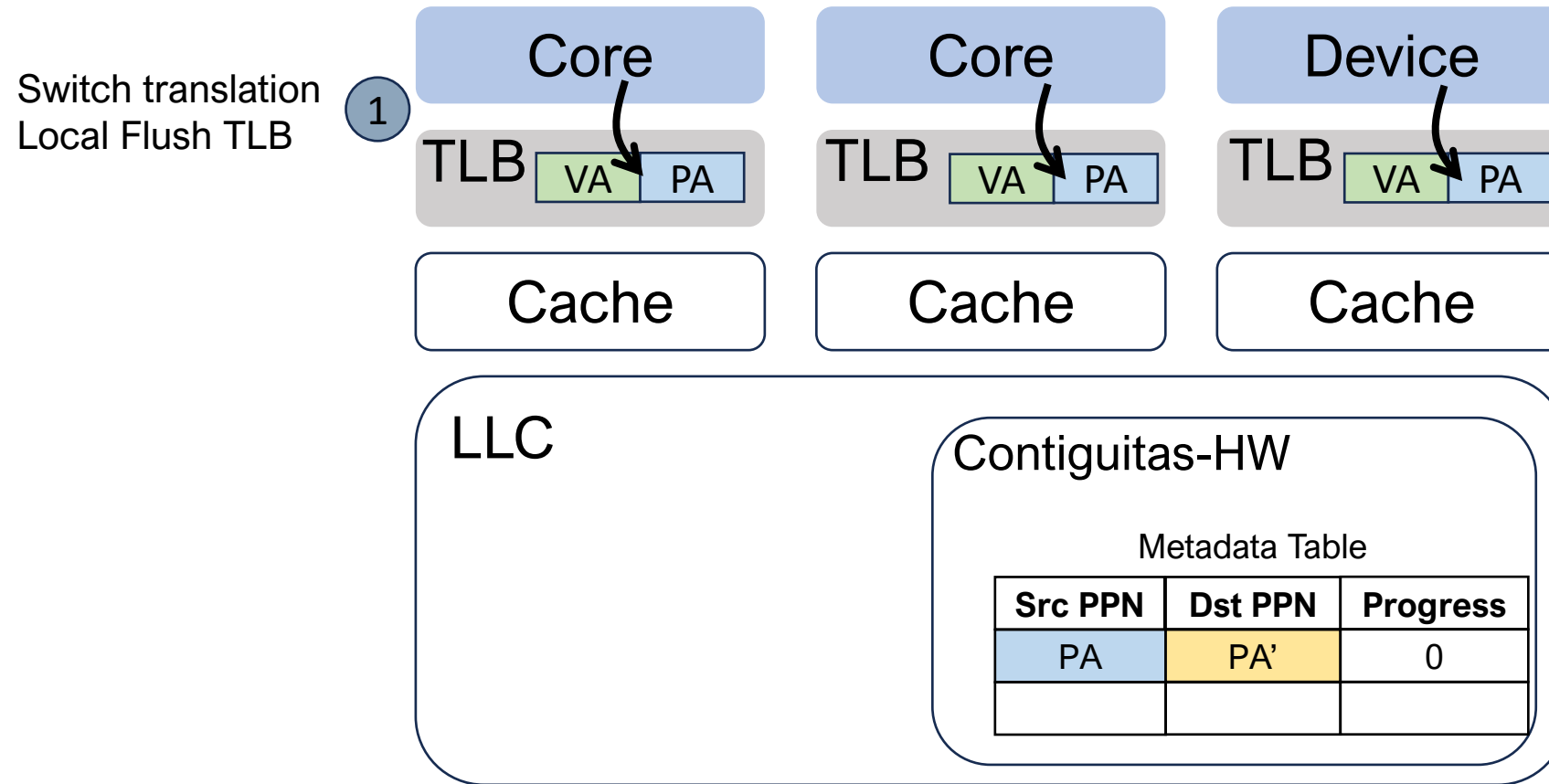# Transparent Page Migration

# Transparent Page Migration

# Scalable Local TLB Invalidations

Switch translation
Local Flush TLB ①

| Core | Core | Device |
|------|------|--------|
| TLB [VA → PA] | TLB [VA → PA] | TLB [VA → PA] |
| Cache | Cache | Cache |

**LLC**

**Contiguitas-HW**

Metadata Table

| Src PPN | Dst PPN | Progress |
|---------|---------|----------|
| PA | PA' | 0 |
| | | |

# Moving One Cacheline

| Core | Core | Device |
|---|---|---|

| TLB | VA | PA' | | TLB | VA | PA' | | TLB | VA | PA' |

| Cache | Cache | Cache |
|---|---|---|

**LLC**

**Contiguitas-HW**

Line $PA_0$

Line $PA'_0$

Metadata Table

| Src PPN | Dst PPN | Progress |
|---|---|---|
| PA | PA' | 0 |
| | | |

# Moving One Cacheline

# Moving One Cacheline

# Redirecting Access to Moved Lines

**(1)** Load $VA_5$

| Core | Core | Device |
|---|---|---|

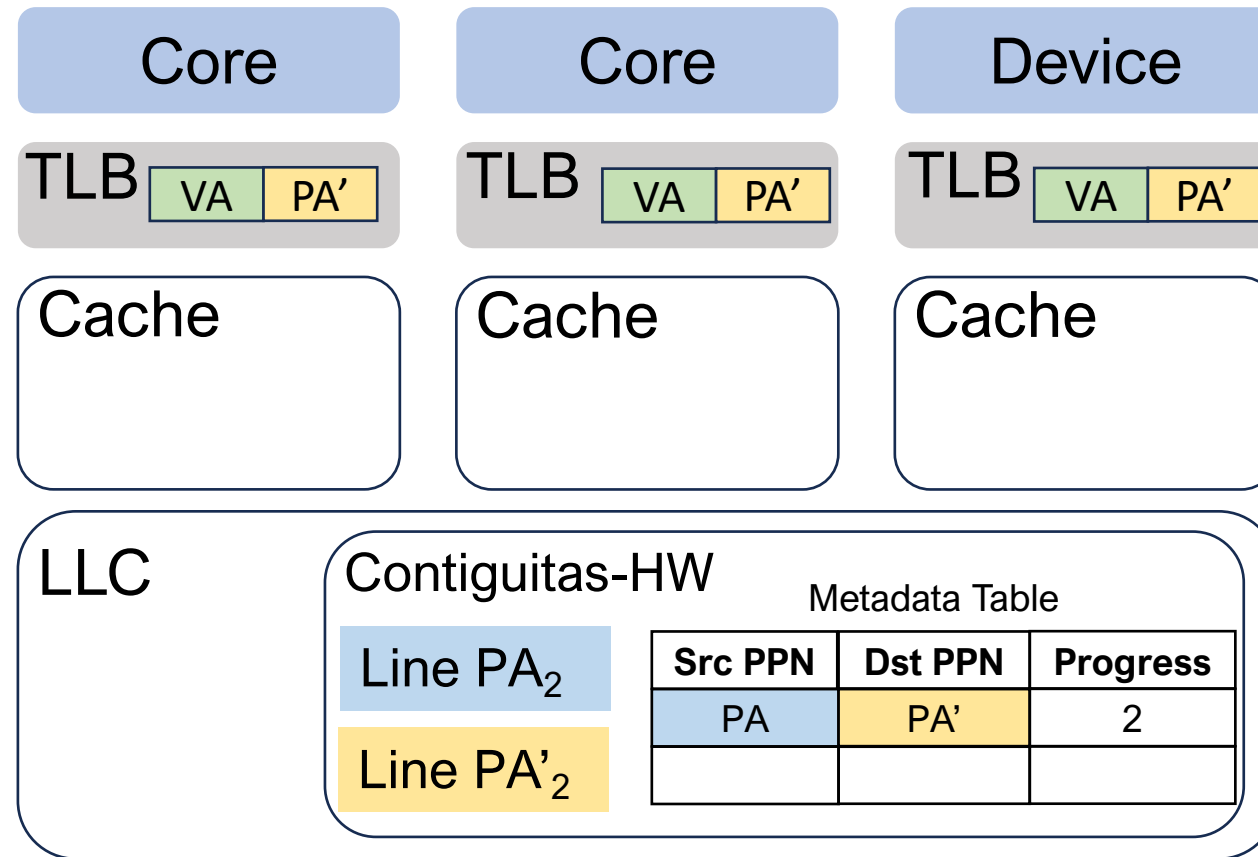| TLB | VA | PA' | | TLB | VA | PA' | | TLB | VA | PA' |

| Cache | Cache | Cache |

**LLC**

**Contiguitas-HW**

Metadata Table

| Src PPN | Dst PPN | Progress |
|---|---|---|
| PA | PA' | 3 |
| | | |

# Redirecting Access to Moved Lines

# Challenge: Concurrent Active Physical Mappings

(1) Load $VA_5$

| Core | Core | Device |
|------|------|--------|

**TLB** | VA | PA' |

**TLB**

**TLB** | VA | PA |

**Cache**

**Cache**

**Cache**

Line $PA_5$

(3)

**LLC**

**Contiguitas-HW**

Metadata Table

(2)

Check progress

| Src PPN | Dst PPN | Progress |
|---------|---------|----------|
| PA | PA' | 0 |
|  |  |  |
|  |  |  |

Invalidate:
Only one of src or dst mapping
can exist in the private caches

Computer Architecture / Operating Systems Group

# More in the Paper

① Load VA$_5$

Core   Core   Device

Sliced LLC

Non-cacheable access

Software interface

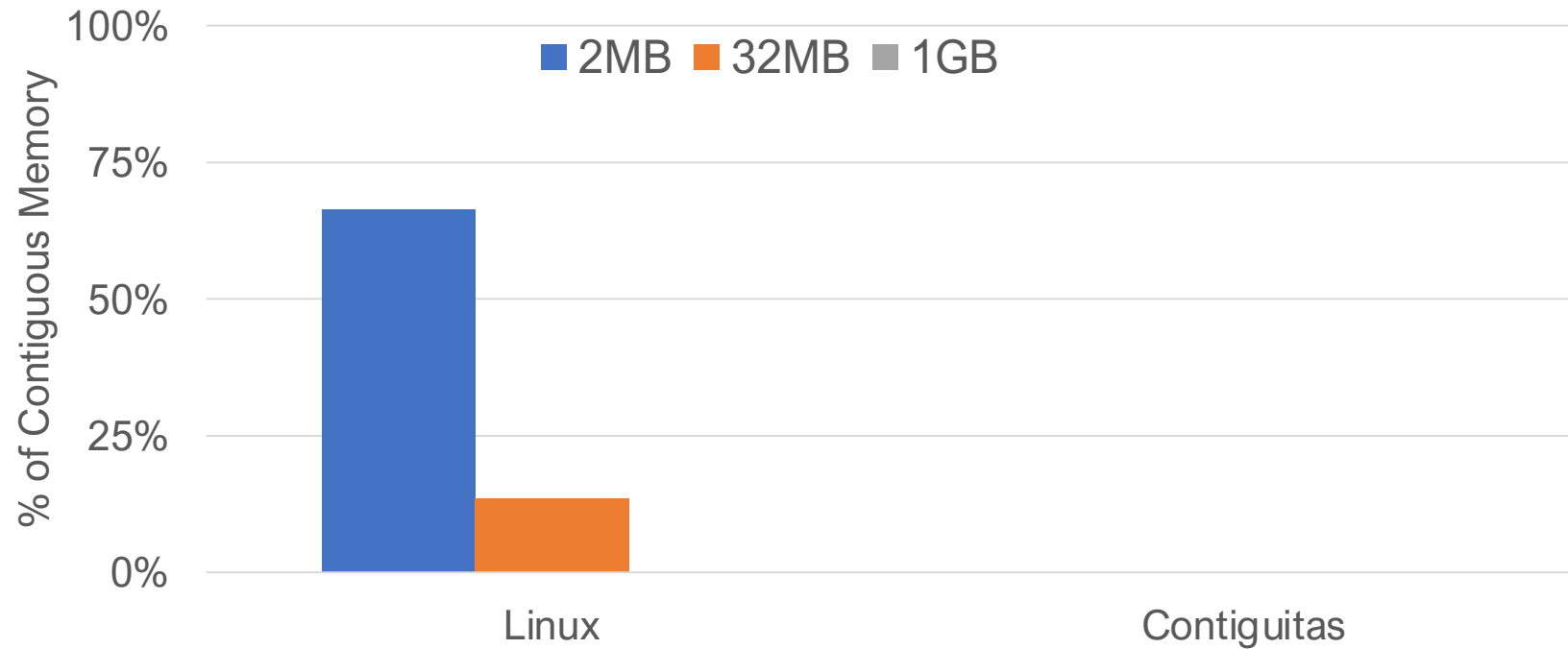Area and power

ping
hes

② Check progress

# Evaluation Overview

Live production traffic at Meta for Contiguitas OS

Major workloads at Meta → Web and two caching services

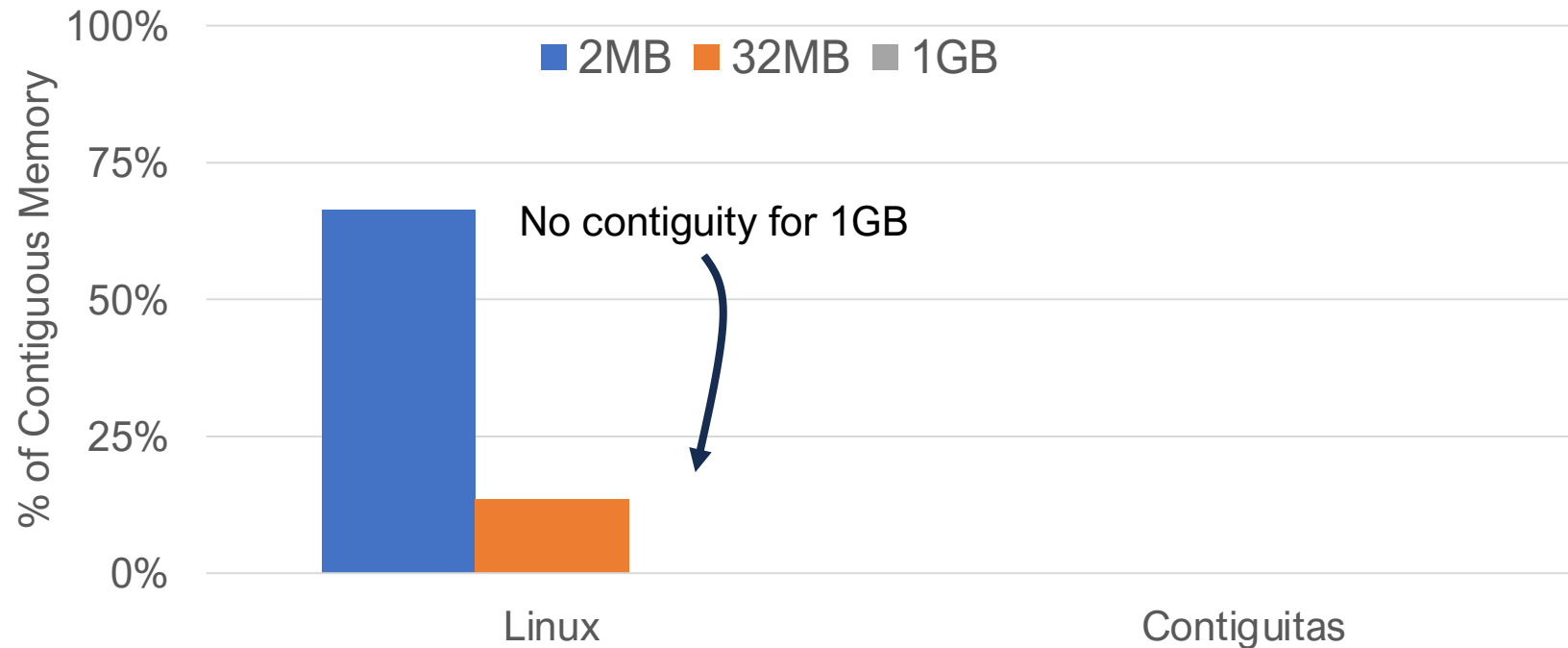End-to-end requests per second (RPS)

Computer Architecture Operating Systems Group

# Potential Memory Contiguity

Computer Architecture ⚊⚊ Operating Systems Group

# Potential Memory Contiguity

Computer Architecture ΛＪＪＬ Operating Systems Group

# Potential Memory Contiguity



100%

75%

50%

25%

0%

% of Contiguous Memory

■ 2MB  ■ 32MB  ■ 1GB

No contiguity for 1GB

Linux

Contiguitas

Linux struggles to produce contiguous memory regions

Computer Architecture  Operating Systems Group

# Potential Memory Contiguity



Linux struggles to produce contiguous memory regions
Contiguitas can use the entire movable region for contiguity
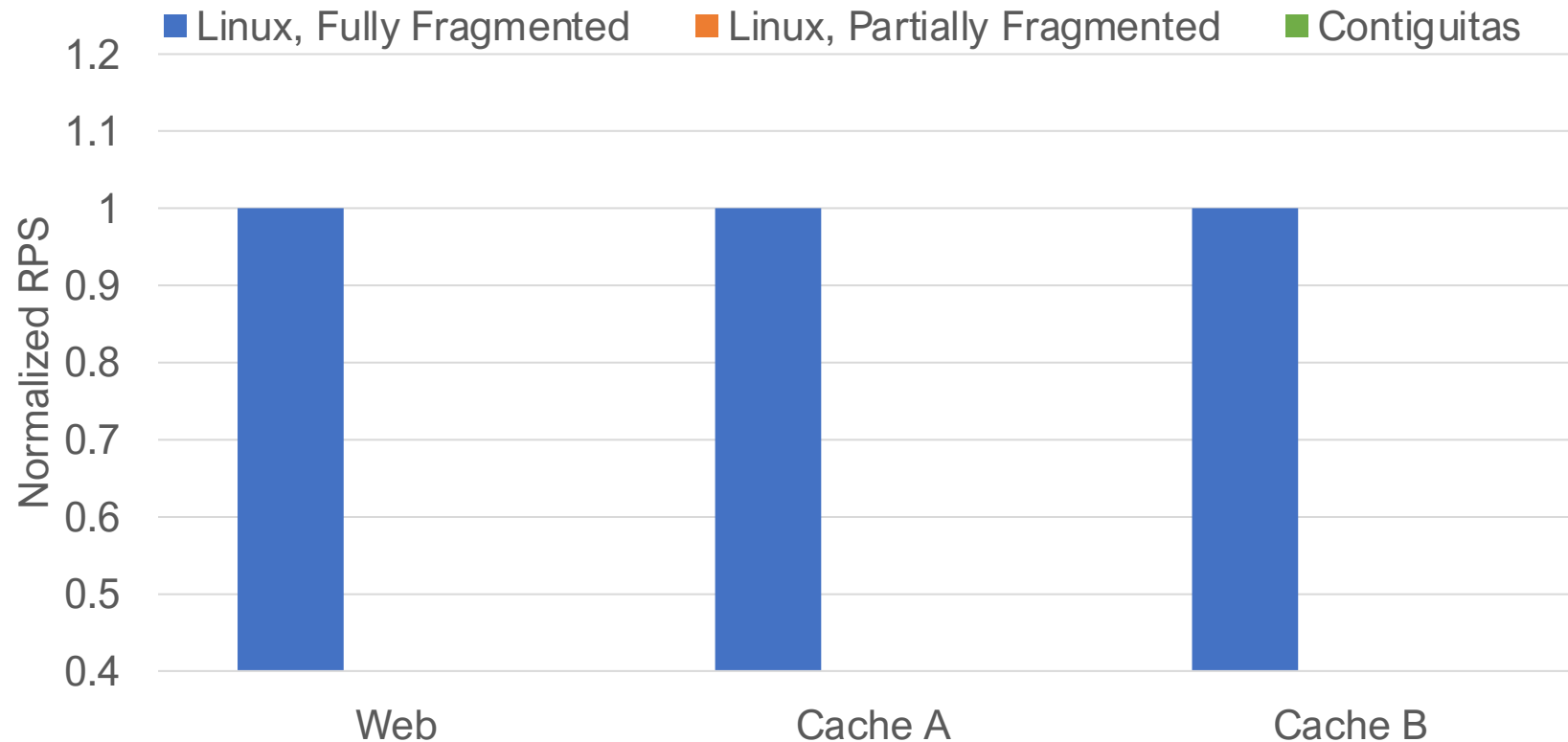
Computer Architecture Operating Systems Group
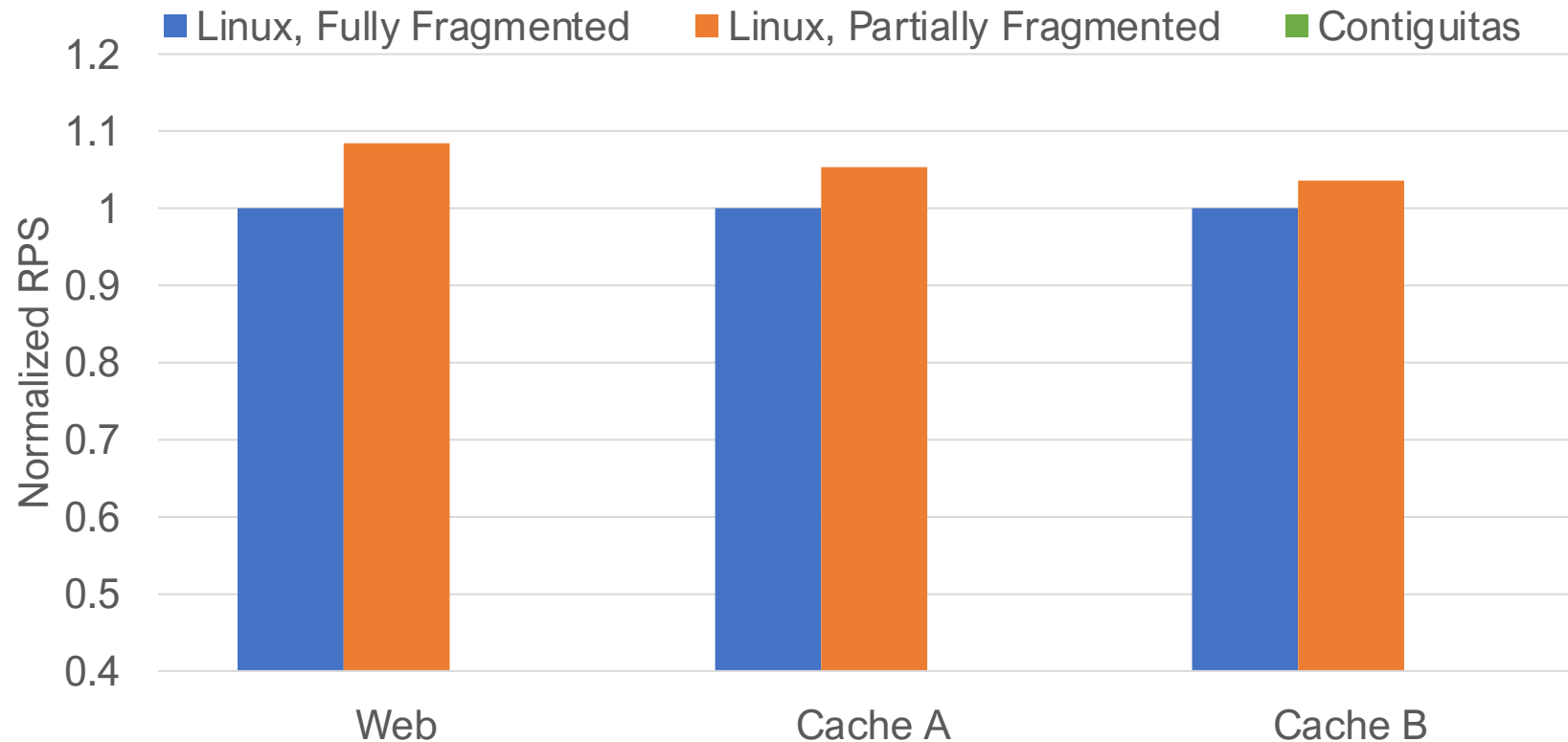
# Potential Memory Contiguity



Linux struggles to produce contiguous memory regions
Contiguitas can use the entire movable region for contiguity
Contiguitas can dynamically allocate 1GB pages

Computer Architecture ⋀⋀ Operating Systems Group

# End-to-End Application Performance

Computer Architecture ⩕⩕ Operating Systems Group

# End-to-End Application Performance

Computer Architecture ⋀⋀⋀ Operating Systems Group

# End-to-End Application Performance

Computer Architecture ⏞ Operating Systems Group

# End-to-End Application Performance



Additional +7.5% from 1GB pages

Legend: Linux, Fully Fragmented | Linux, Partially Fragmented | Contiguitas

Y-axis: Normalized RPS (0.4 to 1.2)

X-axis categories: Web, Cache A, Cache B

**Contiguitas achieves 2% - 18% speedup**

Computer Architecture ∧∧∧ Operating Systems Group

# Evaluation Overview

Live production traffic at Meta for Contiguitas OS
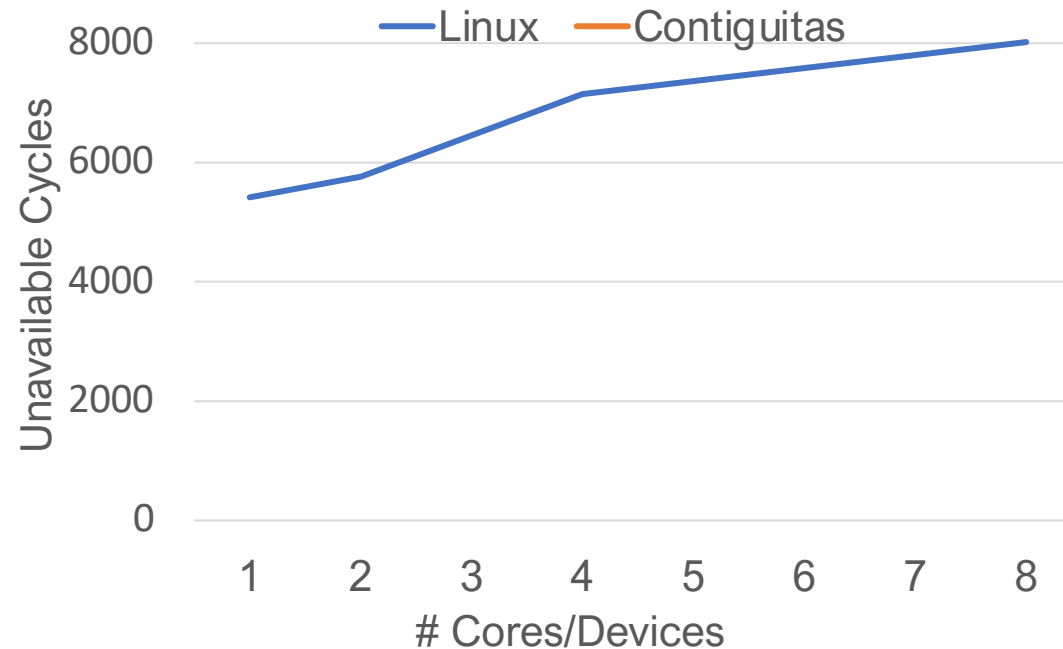  Major workloads at Meta: web and two caching services
  End-to-end requests per second (RPS)

Full-system cycle-accurate simulation for hardware
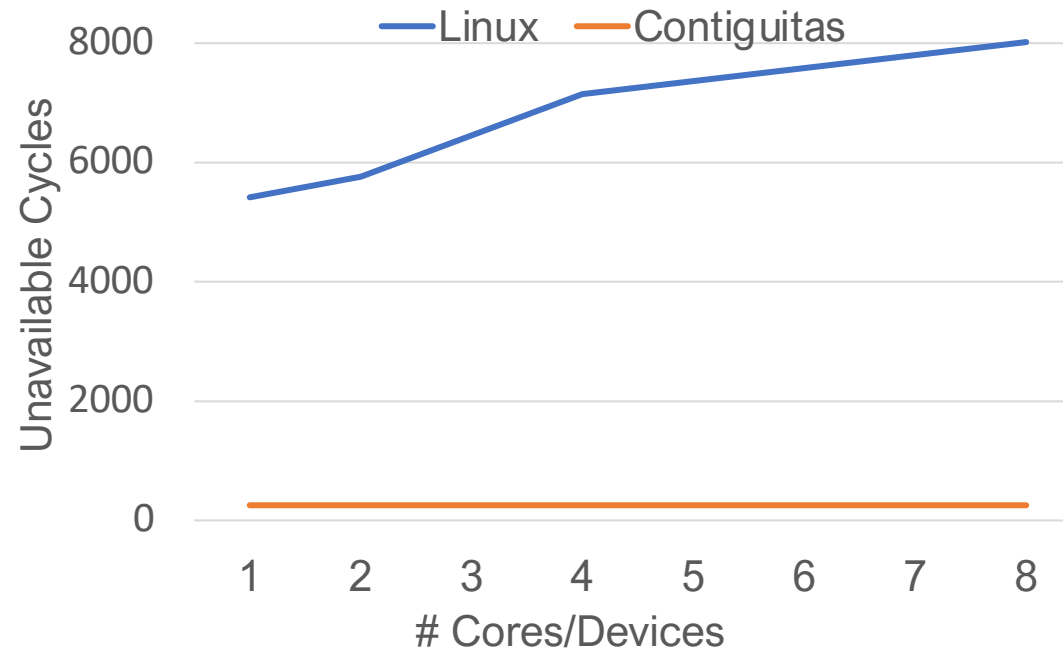  Microbenchmark for page unavailable time
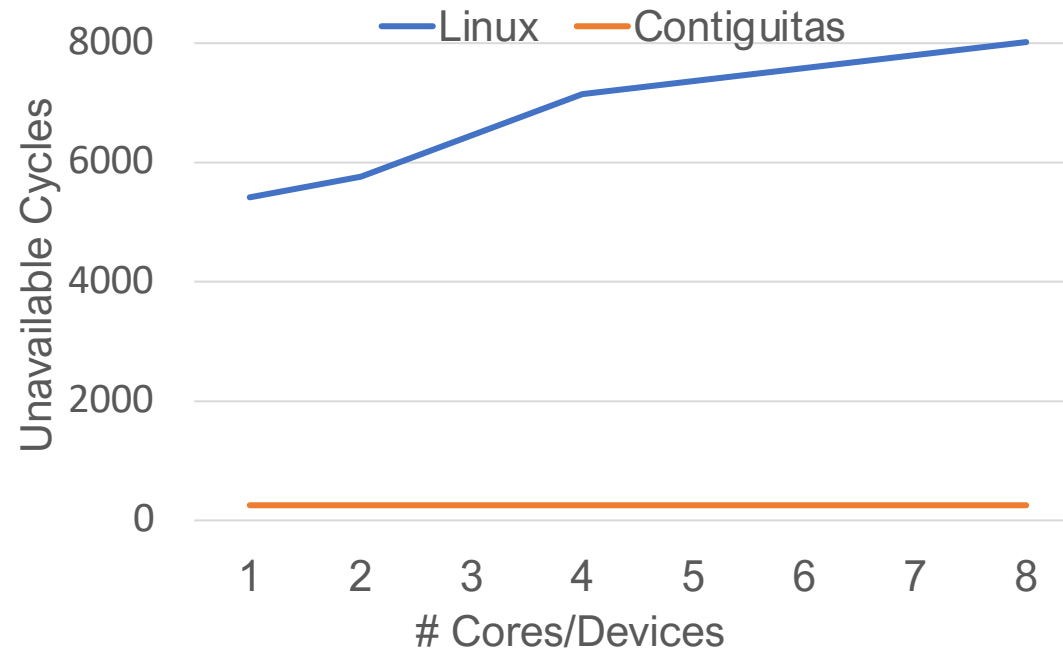  Application performance impacts of page migration

Computer Architecture Operating Systems Group

# Scalability of Contiguitas-HW



Migrating a page in software blocks for linear time w.r.t. TLBs

Computer Architecture Operating Systems Group

# Scalability of Contiguitas-HW



Migrating a page in software blocks for linear time w.r.t. TLBs
Contiguitas-HW blocks for only local TLB invalidation

Computer Architecture ⩘⩘ Operating Systems Group

# Scalability of Contiguitas-HW



Migrating a page in software blocks for linear time w.r.t. TLBs
Contiguitas-HW blocks for only local TLB invalidation
Negligible performance impacts to applications

Computer Architecture ⋀⋀ Operating Systems Group

# Takeaways: <span style="color:red">Contiguitas</span>

Unmovable pages are detrimental to contiguity

- They will only get worse

Contiguitas is a holistic solution across OS and hardware

- Unmovable page confinement + transparent page migration
- Reduces memory fragmentation due to unmovable allocations

Ample physical memory contiguity
Performance gains of 2-18% with production traffic
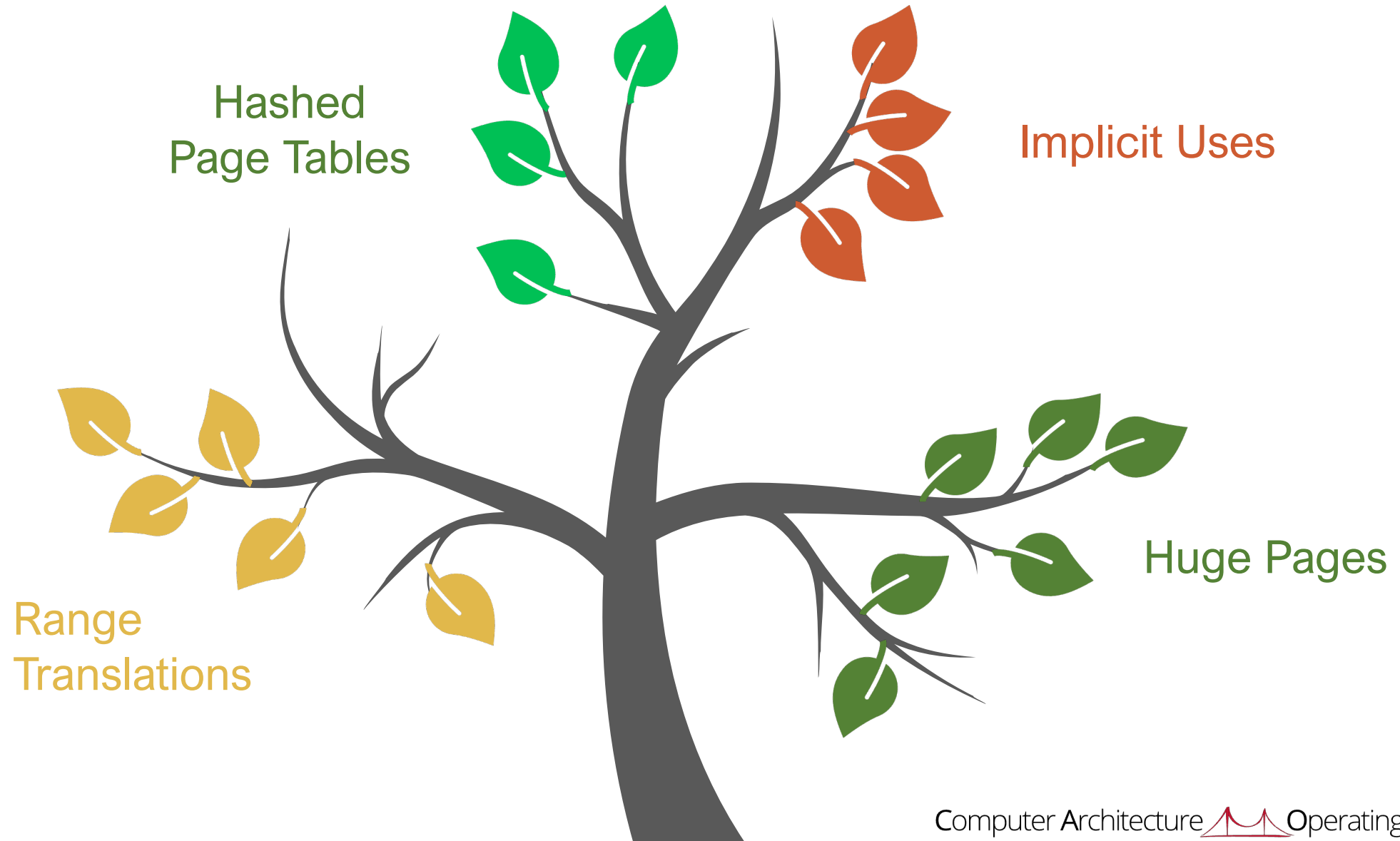Reduce unmovable pages with HW support

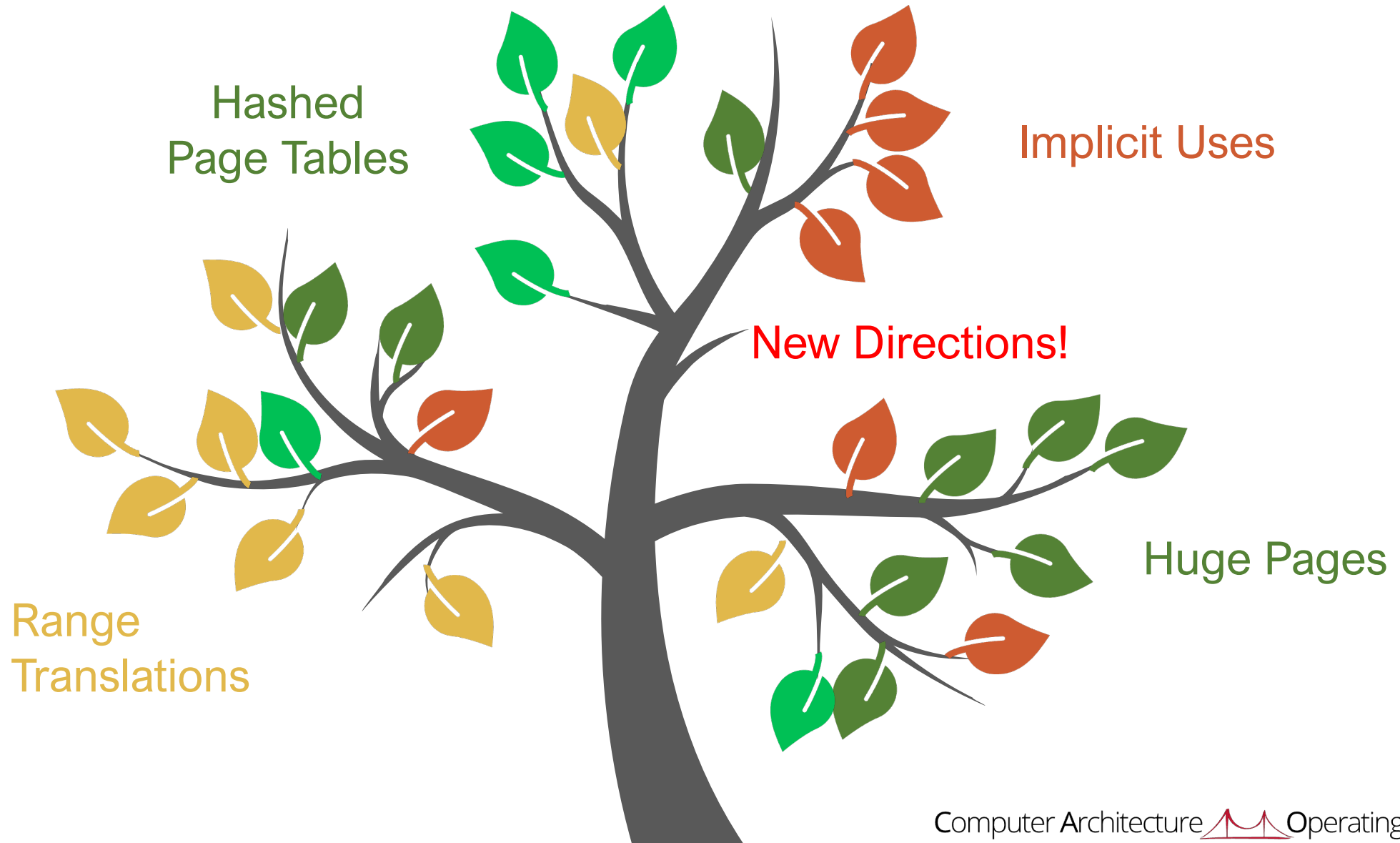In the process of upstreaming to Linux

# With Ample Physical Contiguity..

# Contiguitas: The Pursuit of Physical Memory Contiguity in Datacenters

**Kaiyang Zhao**, Kaiwen Xue, Ziqi Wang, Dan Schatzberg, Leon Yang, Antonis Manousis, Johannes Weiner, Rik van Riel, Bikash Sharma, Chunqiang Tang, Dimitrios Skarlatos

Computer Architecture Operating Systems Group