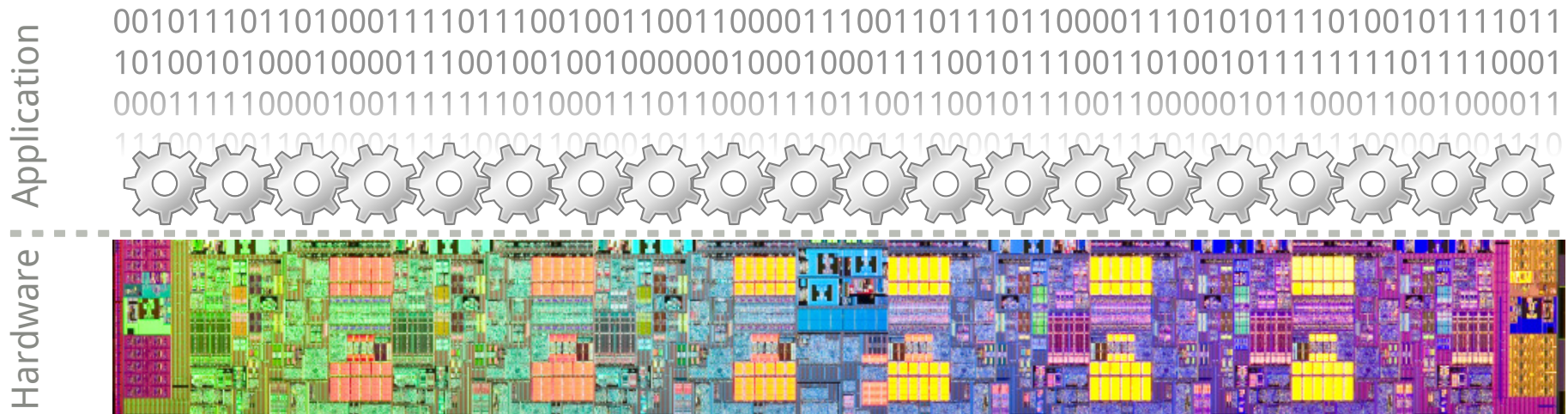# RadixVM: Scalable address spaces for multithreaded applications
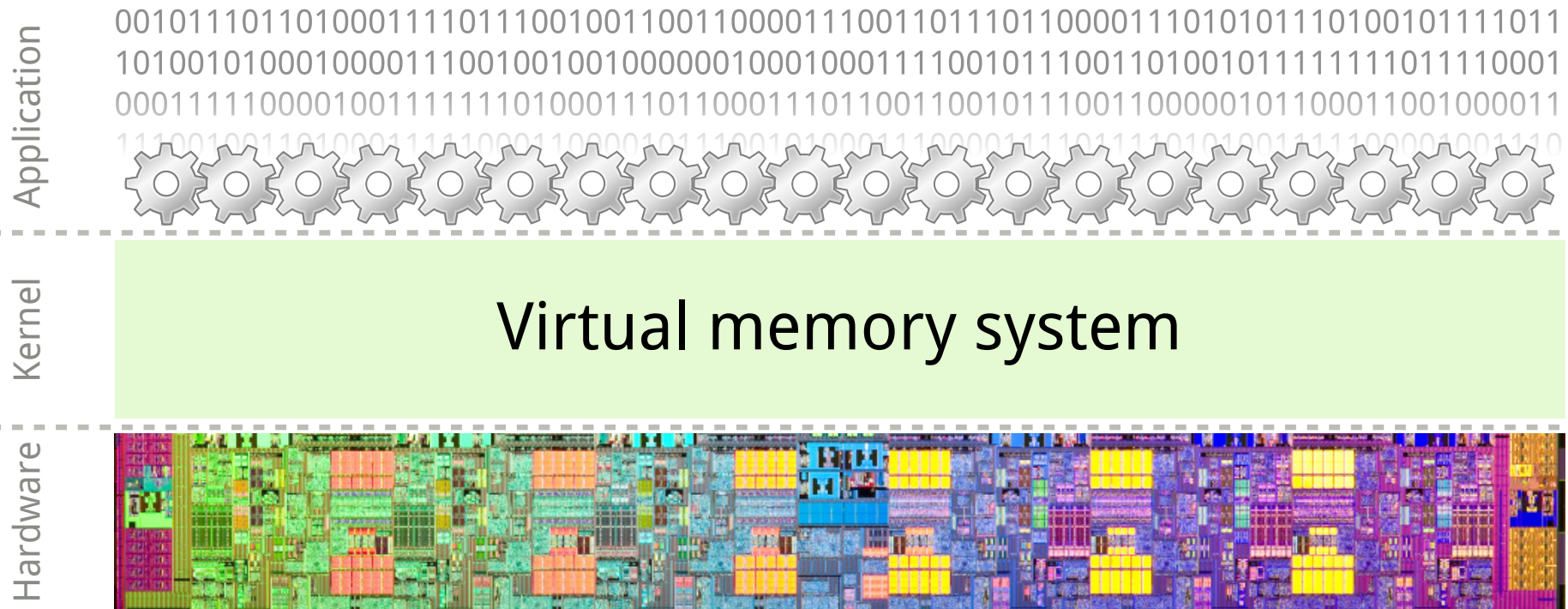
Austin T. Clements
M. Frans Kaashoek
Nickolai Zeldovich

MIT CSAIL

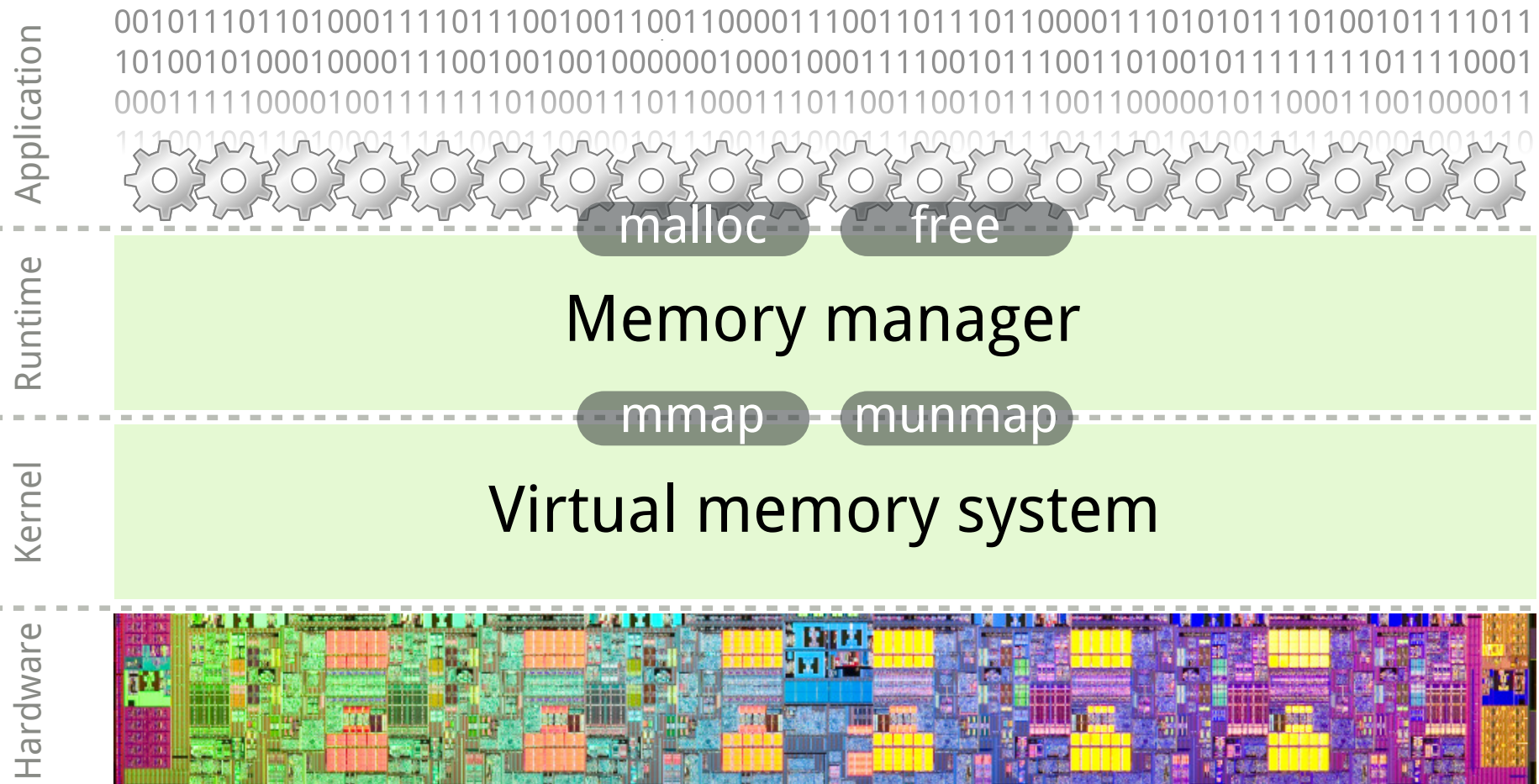# Parallel applications use VM intensively



Application

0010111011010001110111001001100110000111001101110110000111010101110100101111011
1010010100010000111001001001000000100010001110010111001101001011111111011110001
0001111100001001111111010001110110001110110011001011100110000010110001100100011
1110010011010001111100011000010111001010001110000111101111010100111110000100110
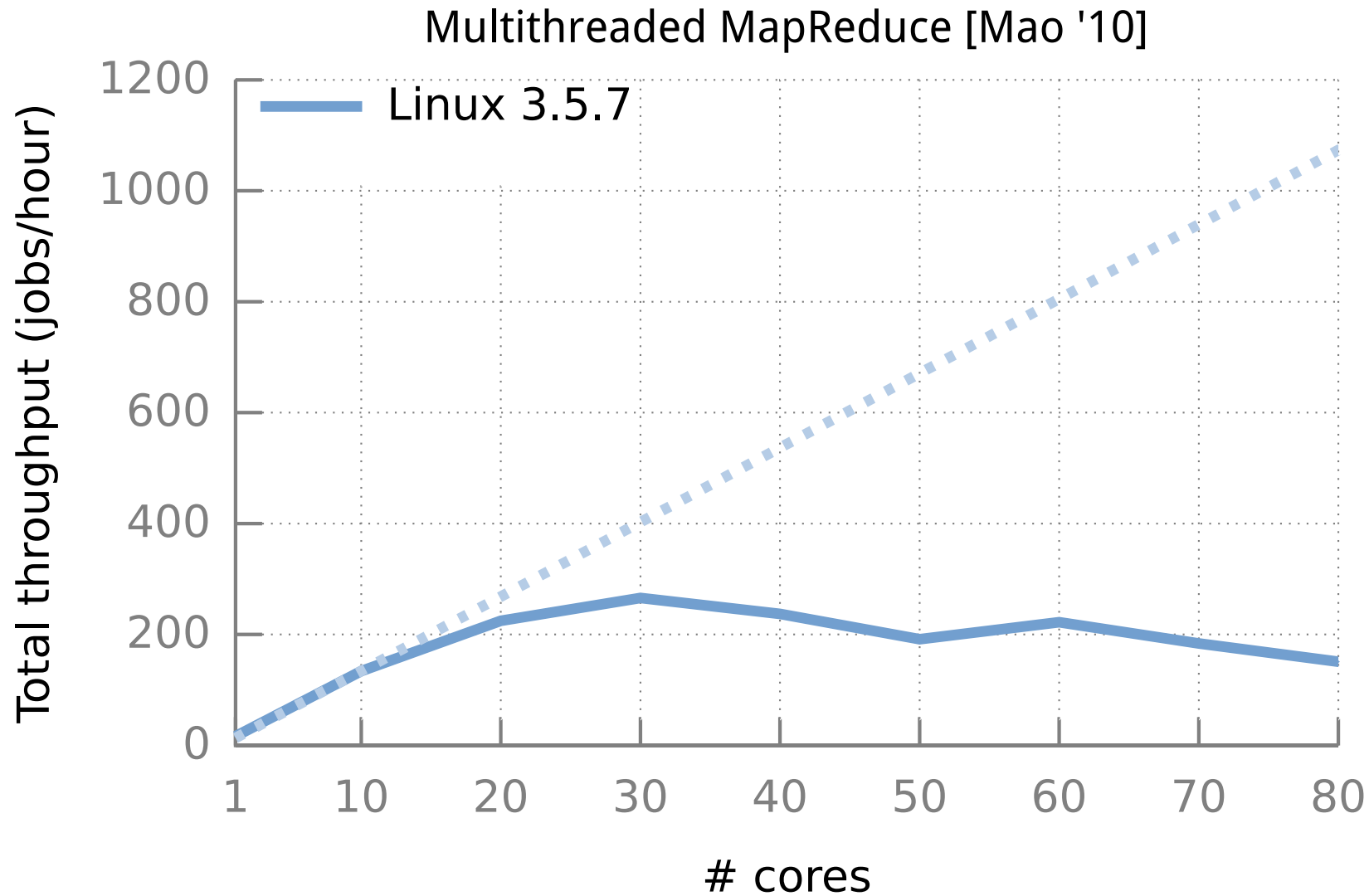
Hardware

# Parallel applications use VM intensively



Every popular operating system serializes
basic VM operations like mmap and munmap.
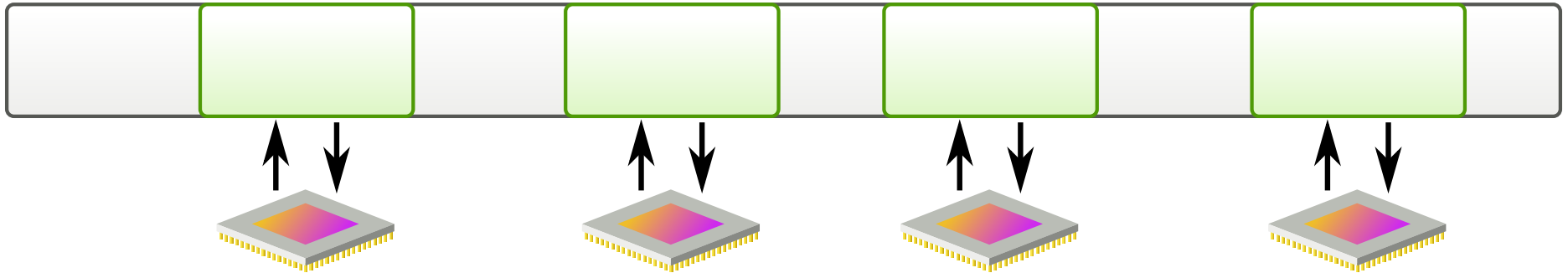
# Parallel applications use VM intensively



Every popular operating system serializes
basic VM operations like mmap and munmap.

# Application performance suffers



Multithreaded MapReduce [Mao '10]

Linux 3.5.7

Total throughput (jobs/hour) — vertical axis with values 0, 200, 400, 600, 800, 1000, 1200

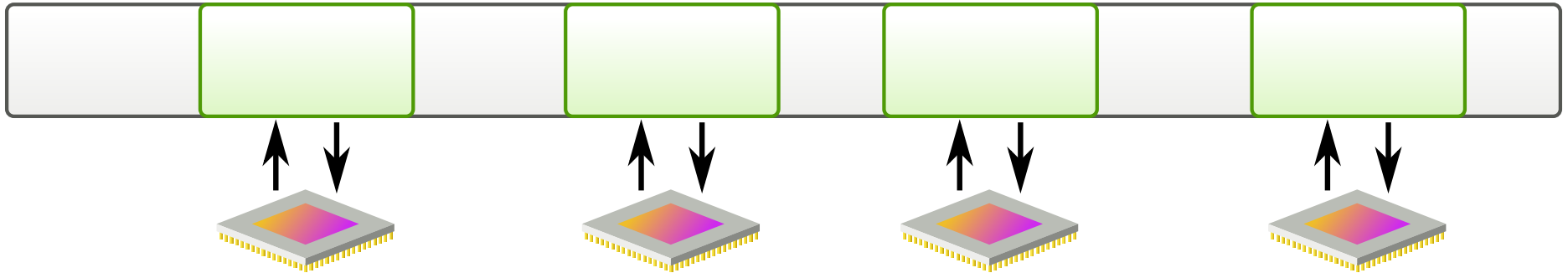# cores — horizontal axis with values 1, 10, 20, 30, 40, 50, 60, 70, 80

# Inside parallel applications



Independent VM operations on non-overlapping regions.

# Inside parallel applications



Independent VM operations on non-overlapping regions.
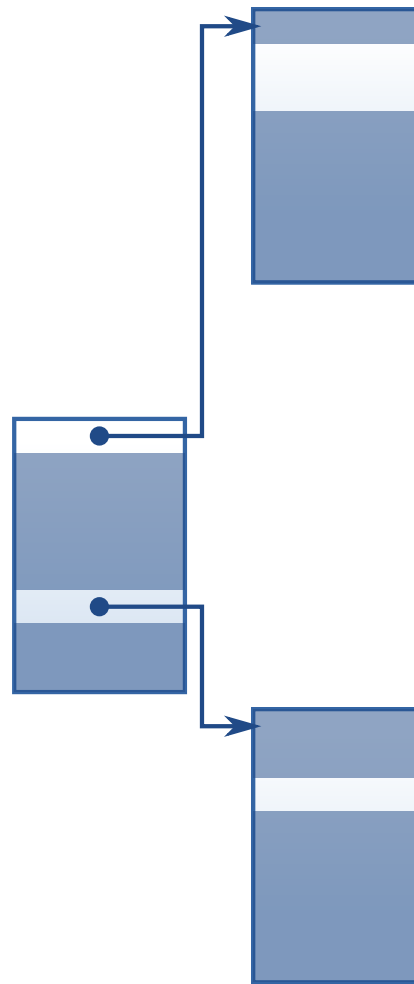
Common pattern for parallel applications.

Perfectly scalable mmap, munmap, and page fault operations on non-overlapping address space regions.
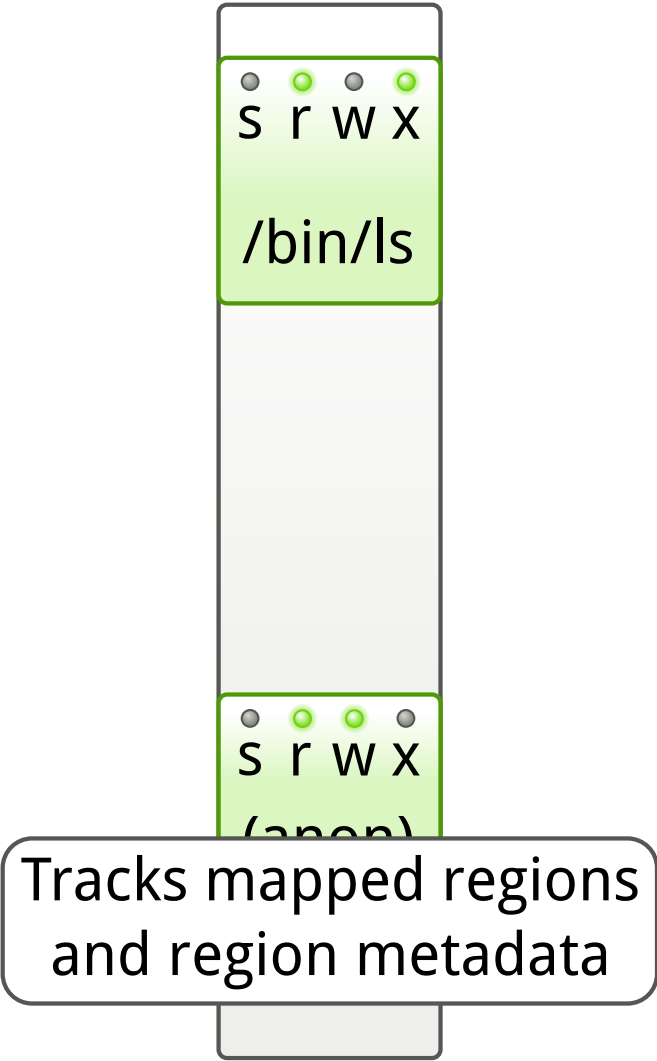
# Structure of a VM system

Memory map

s r w x
/bin/ls

s r w x
(anon)

Hardware
page table

Per-CPU TLB

| Virt | Phys |
|------|------|
| 18bca | 00230 |
| 87c38 | 0049c |

| Virt | Phys |
|------|------|
| 8a4bd | 00382 |
| 87c38 | 0049c |

| Virt | Phys |
|------|------|
| b987a | 00520 |
| 8a4bd | 00382 |
| 87c38 | 0049c |

# Structure of a VM system

**s r w x**

/bin/ls

**s r w x**

(anon)

Tracks mapped regions
and region metadata

Memory map

Hardware
page table

| Virt | Phys |
|------|------|
| 18bca | 00230 |
| 87c38 | 0049c |

| Virt | Phys |
|------|------|
| 8a4bd | 00382 |
| 87c38 | 0049c |

| Virt | Phys |
|------|------|
| b987a | 00520 |
| 8a4bd | 00382 |
| 87c38 | 0049c |

Per-CPU TLB

# Structure of a VM system



Virt | Phys
--- | ---
18bca | 00230
87c38 | 0049c

Virt | Phys
--- | ---
8a4bd | 00382
87c38 | 0049c

Virt | Phys
--- | ---
b987a | 00520
8a4bd | 00382
87c38 | 0049c

s r w x
/bin/ls

s r w x
(anon)

Shared by OS and hardware.
Maps virtual to physical.

**Memory map**

**Hardware
page table**

**Per-CPU TLB**

# Structure of a VM system

| | Virt | Phys |
|---|---|---|
| | 18bca | 00230 |
| | 87c38 | 0049c |

s r w x
/bin/ls

s r w x
(anon)

| | Virt | Phys |
|---|---|---|
| | 8a4bd | 00382 |
| | 87c38 | 0049c |

| | Virt | Phys |
|---|---|---|
| | b987a | 00520 |
| | 8a4bd | 00382 |

Caches page tables.
Internal to CPU.

**Memory map**

**Hardware page table**

**Per-CPU TLB**

# Structure of a VM system

Application

mmap          munmap

load/store

s r w x

/bin/ls

s r w x

(anon)

| Virt | Phys |
|------|------|
| 18bca | 00230 |
| 87c38 | 0049c |

| Virt | Phys |
|------|------|
| 8a4bd | 00382 |
| 87c38 | 0049c |

| Virt | Phys |
|------|------|
| b987a | 00520 |
| 8a4bd | 00382 |
| 87c38 | 0049c |

# Structure of a VM system

Application

mmap    munmap                                        load/store

|         |         |                         | Virt  | Phys  |
| s r w x | /bin/ls |                         | 18bca | 00230 |
|         |         |                         | 87c38 | 0049c |
| s r w x | (anon)  |                         | Virt  | Phys  |
|         |         |                         | 8a4bd | 00382 |
|         |         |                         | 87c38 | 0049c |
| s r w x | (anon)  |                         | Virt  | Phys  |
|         |         |                         | b987a | 00520 |
|         |         |                         | 8a4bd | 00382 |
|         |         |                         | 87c38 | 0049c |

RadixVM: Scalable address spaces for multithreaded applications

# Structure of a VM system

Application

mmap          munmap                                    load/store

| Virt | Phys |
|------|-------|
| 18bca | 00230 |
| 87c38 | 0049c |

s r w x

/bin/ls

s r w x

(anon)

Page faults

| Virt | Phys |
|------|-------|
| 8a4bd | 00382 |
| 87c38 | 0049c |

TLB misses

s r w x

(anon)

| Virt | Phys |
|------|-------|
| b987a | 00520 |
| 8a4bd | 00382 |
| 87c38 | 0049c |

# Structure of a VM system

Application

mmap        munmap

load/store

s r w x

/bin/ls

s r w x

(anon)

Page faults

s r w x

(anon)

TLB misses

| Virt | Phys |
|------|------|
| 18bca | 00230 |
| 87c38 | 0049c |

| Virt | Phys |
|------|------|
| 8a4bd | 00382 |
| 87c38 | 0049c |

| Virt | Phys |
|------|------|
| b987a | 00520 |
| 8a4bd | 00382 |
| 87c38 | 0049c |

# Structure of a VM system

Application

mmap    munmap                                load/store

s r w x

/bin/ls

s r w x

(anon)

s r w x

Page faults

TLB misses

| Virt | Phys |
|------|------|
| 18bca | 00230 |
| 87c38 | 0049c |

| Virt | Phys |
|------|------|
| 8a4bd | 00382 |
| 87c38 | 0049c |

| Virt | Phys |
|------|------|
| b987a | 00520 |
| | 00382 |
| 87c38 | 0049c |

Seems reasonable.  Why doesn't it scale?

# Structure of a VM system

Application

mmap    munmap

load/store

| Virt | Phys |
|---|---|
| 18bca | 00230 |
| 87c38 | 0049c |

s r w x

/bin/ls

s r w x

(anon)

Page faults

TLB misses

| Virt | Phys |
|---|---|
| 8a4bd | 00382 |
| 87c38 | 0049c |

s r w x

Seems reasonable.  Why doesn't it scale?

Locking

| Virt | Phys |
|---|---|
| b987a | 00520 |
| | 00382 |
| 87c38 | 0049c |

# Structure of a VM system

Application

mmap    munmap                    load/store

| Virt | Phys |
|------|------|
| 18bca | 00230 |
| 87c38 | 0049c |

s r w x

/bin/ls

s r w x

(anon)

Page faults

| Virt | Phys |
|------|------|
| 8a4bd | 00382 |
| 87c38 | 0049c |

TLB misses

s r w x

| Virt | Phys |
|------|------|
| b987a | 00520 |
| | 00382 |
| 87c38 | 0049c |

Seems reasonable.  Why doesn't it scale?

Shootdown broadcast    Locking

# Structure of a VM system

Application

mmap          munmap                                    load/store

| Virt | Phys |
|------|------|
| 18bca | 00230 |
| 87c38 | 0049c |

s r w x

/bin/ls

| Virt | Phys |
|------|------|

s r w x

(anon)

| Virt | Phys |
|------|------|
| 8a4bd | 00382 |
| 87c38 | 0049c |

Page faults

TLB misses

s r w x

| Virt | Phys |
|------|------|
| b987a | 00520 |
| | 00382 |
| 87c38 | 0049c |

## Seems reasonable. Why doesn't it scale?

Shootdown broadcast          Locking          Cache contention

# Structure of a VM system

Application

mmap        munmap                                          load/store

| Virt | Phys |
|-------|-------|
| 18bca | 00230 |
| 87c38 | 0049c |

s r w x

/bin/ls

s r w x

(anon)

| Virt | Phys |
|-------|-------|
| 8a4bd | 00382 |
| 87c38 | 0049c |

Page faults

TLB misses

s r w x

| Virt | Phys |
|-------|-------|
| b987a | 00520 |
| | 00382 |
| 87c38 | 0049c |

(anon)

## Seems reasonable.  Why doesn't it scale?

Shootdown broadcast        Locking        Cache contention

Cross-core communication

# This talk: RadixVM

To achieve perfectly scalable non-overlapping operations, we eliminate communication between such operations.

Concurrent memory map representation

Method of targeting TLB shootdowns

Scalable, space-efficient reference counting

# Metadata management

Need to store OS-level memory mapping metadata

# Metadata management

Need to store OS-level memory mapping metadata

Popular operating systems use a balanced tree of region objects.

# Metadata management

Need to store OS-level memory mapping metadata



Popular operating systems use a balanced tree of region objects.

Memory-efficient

# Metadata management

Need to store OS-level memory mapping metadata



Popular operating systems use a balanced tree of region objects.

Unnecessary communication

Memory-efficient

# Metadata management

Need to store OS-level memory mapping metadata

Popular operating systems use a balanced tree of region objects.

Unnecessary communication

Memory-efficient

# Metadata management

Need to store OS-level memory mapping metadata



Popular operating systems use a balanced tree of region objects.

Unnecessary communication

Memory-efficient

~1000 cycles

# Metadata management

Need to store OS-level memory mapping metadata



Popular operating systems use a balanced tree of region objects.

Unnecessary communication

Memory-efficient

Most potential data structures (skip lists, B-trees, etc.) induce communication between disjoint operations.

~1... cycles

# Array-based memory map

# Array-based memory map

# Array-based memory map

0

$2^{35}$

# Array-based memory map

s r w x file

```
0
(anon)
(anon)
(anon)
(anon)
```

/lib/libc
/lib/libc
/lib/libc
/lib/libc
/lib/libc
/lib/libc

$2^{35}$

# Array-based memory map

s r w x file



0

Good: Operations on non-overlapping regions are concurrent and induce no communication.

(anon)
(anon)
(anon)
(anon)

/lib/libc
/lib/libc
/lib/libc
/lib/libc
/lib/libc
/lib/libc

$2^{35}$

# Array-based memory map

s  r  w  x  file

0

Good: Operations on non-overlapping regions are concurrent and induce no communication.

Bad: Space use is obscene,
time is proportional to region size

(anon)
(anon)
(anon)
(anon)

/lib/libc
/lib/libc
/lib/libc
/lib/libc
/lib/libc
/lib/libc

$2^{35}$

# Array-based memory map

s r w x file

0



$2^{35}$

Good: Operations on non-overlapping regions are concurrent and induce no communication.

Bad: Space use is obscene,
time is proportional to region size

How can we achieve good concurrency while keeping space and time under control?

# Radix tree

s r w x file

Solution: Range-oriented radix tree



0

2^{35}

# Radix tree

s r w x file



**Solution: Range-oriented radix tree**

# Radix tree

s r w x file

Solution: Range-oriented radix tree

# Radix tree

s r w x file



Solution: Range-oriented radix tree

(anon)
(anon)
(anon)
(anon)

...

/lib/libc
/lib/libc
/lib/libc
/lib/libc
/lib/libc
/lib/libc

# Radix tree

s r w x file



Solution: Range-oriented radix tree

Fold constant-valued chunks into parent, recursively.

/lib/libc
/lib/libc
/lib/libc
/lib/libc
/lib/libc
/lib/libc

(anon)
(anon)
(anon)
(anon)

# Radix tree

s r w x file



Solution: Range-oriented radix tree

Fold constant-valued chunks into parent, recursively.

# Radix tree

s r w x file

Solution: Range-oriented radix tree

Fold constant-valued chunks into parent, recursively.

(anon)
(anon)
(anon)
(anon)

...

...

/lib/libc

# Radix tree

s r w x file

Solution: Range-oriented radix tree

Fold constant-valued chunks into parent, recursively.

(anon)
(anon)
(anon)
(anon)

/lib/libc

...

...

2-3x the size of the balanced region tree

# Radix tree

s r w x file

Solution: Range-oriented radix tree

Fold constant-valued chunks into parent, recursively.

(anon)
(anon)
(anon)
(anon)

...

/lib/libc

...

2-3x the size of the balanced region tree

We can achieve array-like concurrency
with time and space similar to the balanced tree.

# TLB shootdown

munmap must notify cores of changes to cached mappings

# TLB shootdown

munmap must notify cores of changes to cached mappings

Which cores have a mapping cached?  Who knows?!

munmap must notify cores of changes to cached mappings

Which cores have a mapping cached?  Who knows?!

munmap must notify cores of changes to cached mappings

Which cores have a mapping cached?  Who knows?!

munmap must notify cores of changes to cached mappings

Which cores have a mapping cached?  Who knows?!

# TLB shootdown

munmap must notify cores of changes to cached mappings

Which cores have a mapping cached?  Who knows?!

In the common case, there is little or no sharing.

# TLB tracking

A software-managed TLB would make this easy.

# TLB tracking

A software-managed TLB would make this easy.



| Virt | Phys |
|------|------|
| 18bca | 00230 |
| 87c38 | 0049c |

?

| Virt | Phys |
|------|------|
| 8a4bd | 00382 |
| 87c38 | 0049c |

?

| Virt | Phys |
|------|------|
| b987a | 00520 |
| 8a4bd | 00382 |
| 87c38 | 0049c |

?

Page faults

TLB misses

Trap and track

# TLB tracking

A software-managed TLB would make this easy.



Page faults

TLB misses

| Virt | Phys |
|------|------|
| 18bca | 00230 |
| 87c38 | 0049c |

?

| Virt | Phys |
|------|------|
| 8a4bd | 00382 |
| 87c38 | 0049c |

?

| Virt | Phys |
|------|------|
| b987a | 00520 |
| 8a4bd | 00382 |
| 87c38 | 0049c |

?

~~Trap and track~~

# Soft TLBs, the hard way

## Solution: Per-core page tables for precise TLB tracking



0

s r w x

/bin/ls

Page faults

s r w x

(anon)

$2^{35}$

| Virt | Phys |
|---|---|
| 18bca | 00230 |
| 87c38 | 0049c |

TLB misses

| Virt | Phys |
|---|---|
| 8a4bd | 00382 |
| 87c38 | 0049c |

| Virt | Phys |
|---|---|
| b987a | 00520 |
| 8a4bd | 00382 |
| 87c38 | 0049c |

# Soft TLBs, the hard way

## Solution: Per-core page tables for precise TLB tracking



| Virt | Phys |
|------|------|
| 18bca | 00230 |
| 87c38 | 0049c |

| Virt | Phys |
|------|------|
| 8a4bd | 00382 |
| 87c38 | 0049c |

| Virt | Phys |
|------|------|
| b987a | 00520 |
| 8a4bd | 00382 |
| 87c38 | 0049c |

s r w x

/bin/ls

s r w x

(anon)

0

$2^{35}$

Page faults

TLB misses

# Soft TLBs, the hard way

Solution: Per-core page tables for precise TLB tracking



0

s r w x

/bin/ls

Page faults

s r w x

(anon)

$2^{35}$

Trap and track

Virt | Phys
--- | ---
18bca | 00230
87c38 | 0049c

Virt | Phys
--- | ---
8a4bd | 00382
87c38 | 0049c

Virt | Phys
--- | ---
b987a | 00520
8a4bd | 00382
87c38 | 0049c

TLB misses

Solution: Per-core page tables for precise TLB tracking



| Virt | Phys |
|------|------|
| 18bca | 00230 |
| 87c38 | 0049c |

| Virt | Phys |
|------|------|
| 8a4bd | 00382 |
| 87c38 | 0049c |

| Virt | Phys |
|------|------|
| b987a | 00520 |
| 8a4bd | 00382 |

s r w x

/bin/ls

s r w x

(anon)

0

Page faults

TLB misses

TLB tracking allows us to target TLB shootdowns, eliminating unnecessary shootdown communication.

Trap and track

# Reference counting

Reference counting for physical pages and radix nodes

# Reference counting

Reference counting for physical pages and radix nodes

Shared
counters



Scalable inc/dec | N

# Reference counting

## Reference counting for physical pages and radix nodes

|  | Shared counters | Distributed counters |
|---|---|---|
| Scalable inc/dec | N | Y |
| Zero-detection cost | O(1) | O(objs*cpus) |
| Space | O(1) | O(cpus) |

# Reference counting

## Reference counting for physical pages and radix nodes

|  | Shared counters | Distributed counters | SNZIs [Ellen '07] |
|---|---|---|---|
| |  |  |  |
| Scalable inc/dec | N | Y | Mostly |
| Zero-detection cost | O(1) | O(objs*cpus) | O(1) |
| Space | O(1) | O(cpus) | O(cpus) |

# Reference counting

## Reference counting for physical pages and radix nodes

|  | Shared counters | Distributed counters | SNZIs [Ellen '07] | Refcache |
|---|---|---|---|---|
| Scalable inc/dec | N | Y | Mostly | Y |
| Zero-detection cost | O(1) | O(objs*cpus) | O(1) | O(1) |
| Space | O(1) | O(cpus) | O(cpus) | O(1) |

# Reference counting

## Reference counting for physical pages and radix nodes

|  | Shared counters | Distributed counters | SNZIs [Ellen '07] | Refcache |
|---|---|---|---|---|
| Scalable inc/dec | N | Y | Mostly | Y |
| Zero-detection cost | O(1) | O(objs*cpus) | O(1) | O(1) |
| Space | O(1) | O(cpus) | O(cpus) | O(1) |
| Immediate zero detection | Y | N | Y | N |

Approach: Shared counters with per-core delta caches

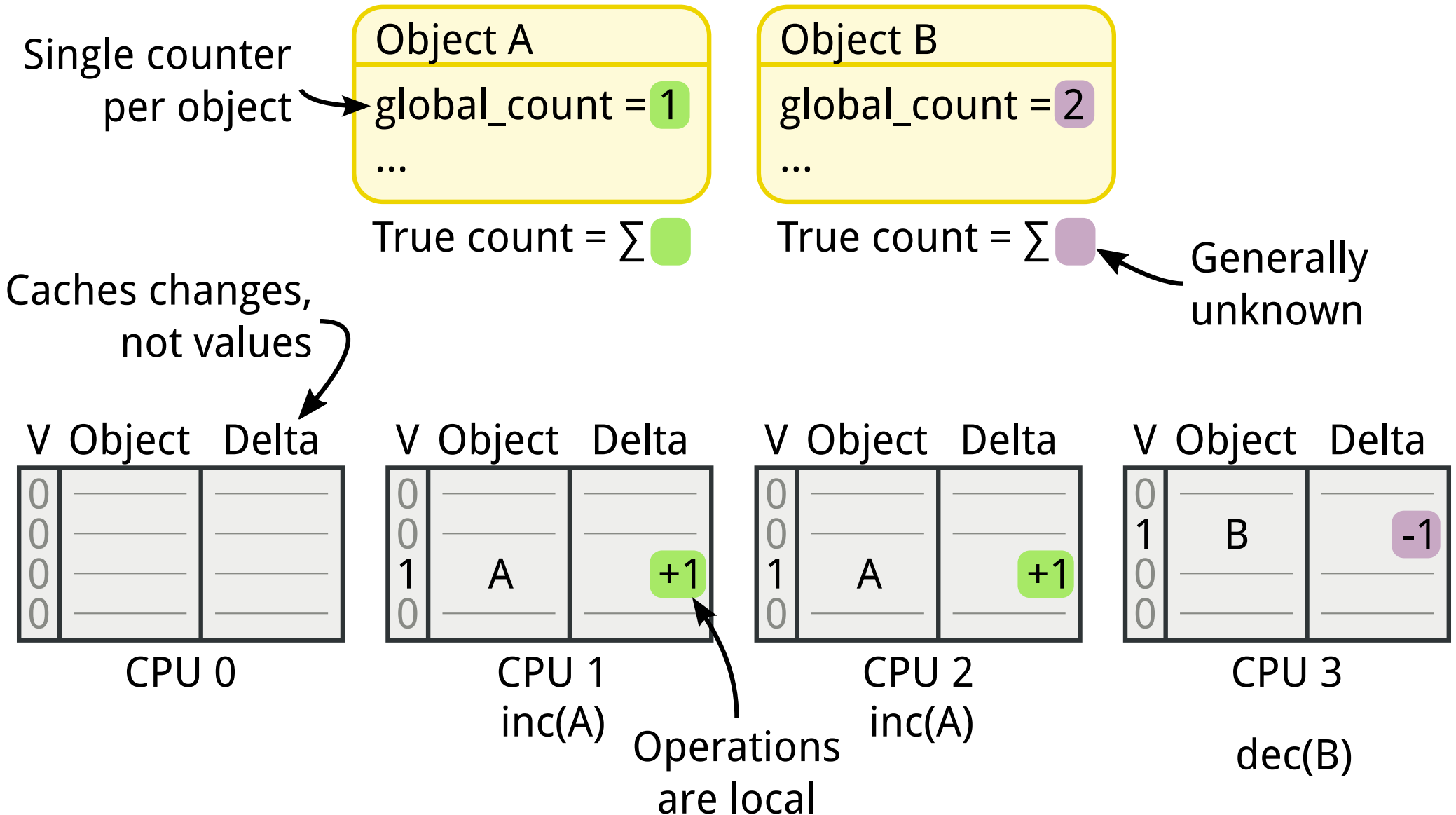# Refcache

Approach: Shared counters with per-core delta caches

Single counter
per object

| Object A | Object B |
|---|---|
| global_count = 1 ... | global_count = 2 ... |

# Refcache

Approach: Shared counters with per-core delta caches

Single counter
per object

| Object A | | Object B |
|---|---|---|
| global_count = 1 | | global_count = 2 |
| ... | | ... |

Caches changes,
not values

| V | Object | Delta | | V | Object | Delta | | V | Object | Delta | | V | Object | Delta |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | | | | 0 | | | | 0 | | | | 0 | | |
| 0 | | | | 0 | | | | 0 | | | | 0 | | |
| 0 | | | | 0 | | | | 0 | | | | 0 | | |
| 0 | | | | 0 | | | | 0 | | | | 0 | | |

CPU 0     CPU 1     CPU 2     CPU 3

# Refcache

Approach: Shared counters with per-core delta caches

Single counter
per object →

| Object A | Object B |
|---|---|
| global_count = 1 | global_count = 2 |
| ... | ... |

Caches changes,
not values →

| V | Object | Delta |
|---|---|---|
| 0 | | |
| 0 | | |
| 0 | | |
| 0 | | |

CPU 0

| V | Object | Delta |
|---|---|---|
| 0 | | |
| 0 | | |
| 1 | A | +1 |
| 0 | | |

CPU 1
inc(A)

| V | Object | Delta |
|---|---|---|
| 0 | | |
| 0 | | |
| 1 | A | +1 |
| 0 | | |

CPU 2
inc(A)

| V | Object | Delta |
|---|---|---|
| 0 | | |
| 0 | | |
| 0 | | |
| 0 | | |

CPU 3

Operations
are local

# Refcache

## Approach: Shared counters with per-core delta caches

Single counter
per object →

| Object A | Object B |
|---|---|
| global_count = 1 | global_count = 2 |
| ... | ... |

Caches changes,
not values →

| V | Object | Delta |
|---|---|---|
| 0 | | |
| 0 | | |
| 0 | | |
| 0 | | |

CPU 0

| V | Object | Delta |
|---|---|---|
| 0 | | |
| 0 | | |
| 1 | A | +1 |
| 0 | | |

CPU 1
inc(A)

| V | Object | Delta |
|---|---|---|
| 0 | | |
| 0 | | |
| 1 | A | +1 |
| 0 | | |

CPU 2
inc(A)

| V | Object | Delta |
|---|---|---|
| 0 | | |
| 1 | B | -1 |
| 0 | | |
| 0 | | |

CPU 3

dec(B)

Operations
are local

# Refcache

Approach: Shared counters with per-core delta caches

Single counter
per object →

| Object A | Object B |
|---|---|
| global_count = 1 | global_count = 2 |
| ... | ... |

True count = ∑ 🟩        True count = ∑ 🟪 ← Generally unknown

Caches changes,
not values →

| V | Object | Delta |
|---|---|---|
| 0 | | |
| 0 | | |
| 0 | | |
| 0 | | |

CPU 0

| V | Object | Delta |
|---|---|---|
| 0 | | |
| 0 | | |
| 1 | A | +1 |
| 0 | | |

CPU 1
inc(A)

Operations
are local

| V | Object | Delta |
|---|---|---|
| 0 | | |
| 0 | | |
| 1 | A | +1 |
| 0 | | |

CPU 2
inc(A)

| V | Object | Delta |
|---|---|---|
| 0 | | |
| 1 | B | -1 |
| 0 | | |
| 0 | | |

CPU 3

dec(B)

# Refcache

When is the true count zero?

# Refcache

When is the true count zero?

Assumption: When the true count is zero, it will stay zero.

When is the true count zero?

Assumption: When the true count is zero, it will stay zero.

Divide time in to epochs.  Each epoch, all CPUs flush their delta caches.  If an object's global count stays zero for a whole epoch, then its true count is zero.

t

# Refcache

When is the true count zero?

Assumption: When the true count is zero, it will stay zero.

Divide time in to epochs.  Each epoch, all CPUs flush their delta caches.  If an object's global count stays zero for a whole epoch, then its true count is zero.

# Refcache

When is the true count zero?

Assumption: When the true count is zero, it will stay zero.

Divide time in to epochs.  Each epoch, all CPUs flush their delta caches.  If an object's global count stays zero for a whole epoch, then its true count is zero.

# Refcache example

Initially: Global count is 1, no cached deltas (so true count is 1)

global

CPU 0
1
2
3
t

true

# Refcache example

Initially: Global count is 1, no cached deltas (so true count is 1)

global

CPU 0
1
2
3
t

true

# Refcache example

Initially: Global count is 1, no cached deltas (so true count is 1)

CPU 0 decrements and flushes; global count is now 0.  What about true count?
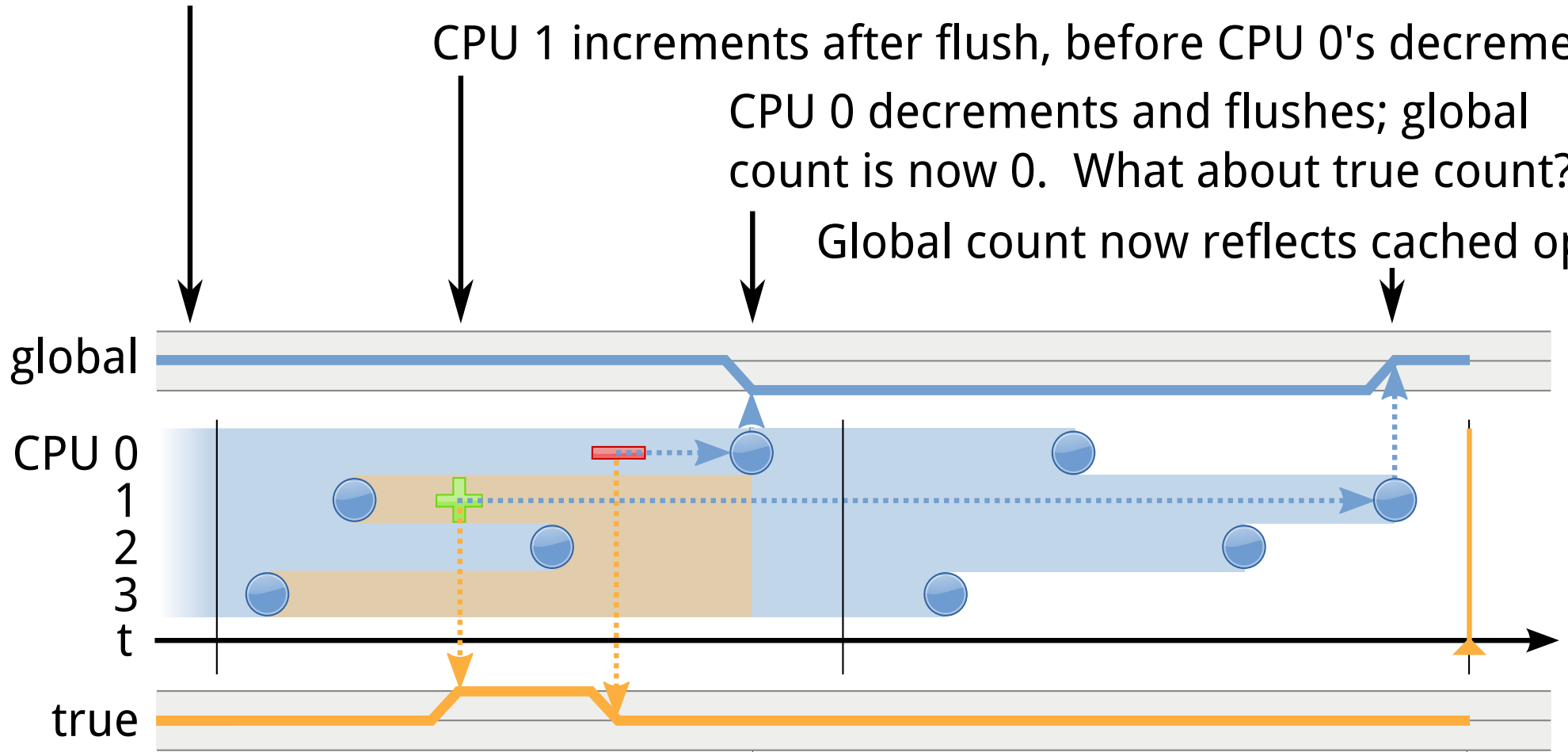
global

CPU 0
1
2
3
t

true

# Refcache example

Initially: Global count is 1, no cached deltas (so true count is 1)

CPU 1 increments after flush, before CPU 0's decrement

CPU 0 decrements and flushes; global count is now 0. What about true count?
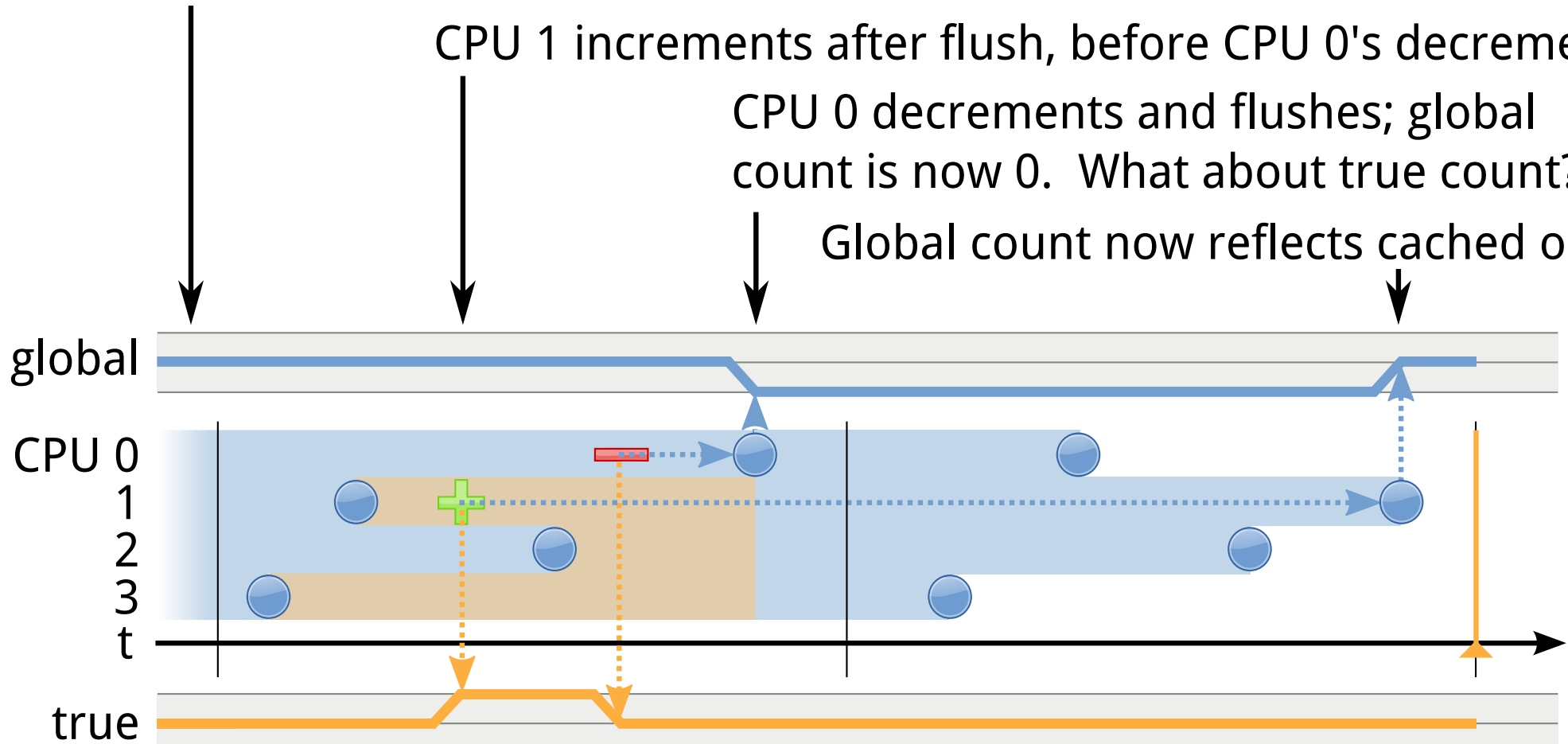
global

CPU 0

1

2

3

t

true

# Refcache example

Initially: Global count is 1, no cached deltas (so true count is 1)

CPU 1 increments after flush, before CPU 0's decrement

CPU 0 decrements and flushes; global count is now 0.  What about true count?

global

CPU 0
1
2
3
t

true

The true count is the sum of everything up to right now.

# Refcache example

Initially: Global count is 1, no cached deltas (so true count is 1)

CPU 1 increments after flush, before CPU 0's decrement

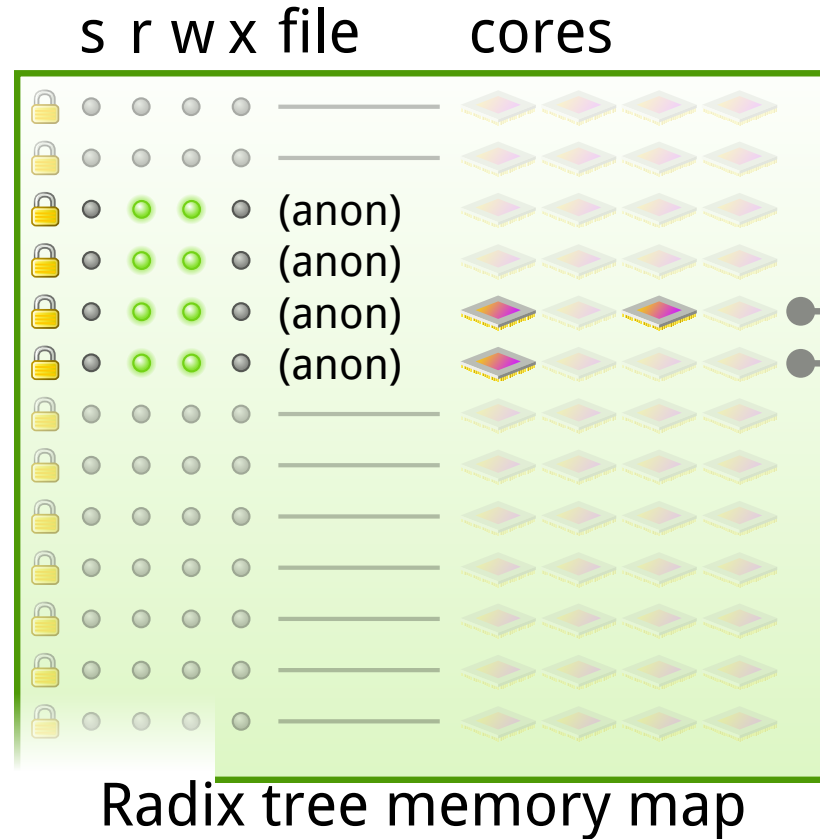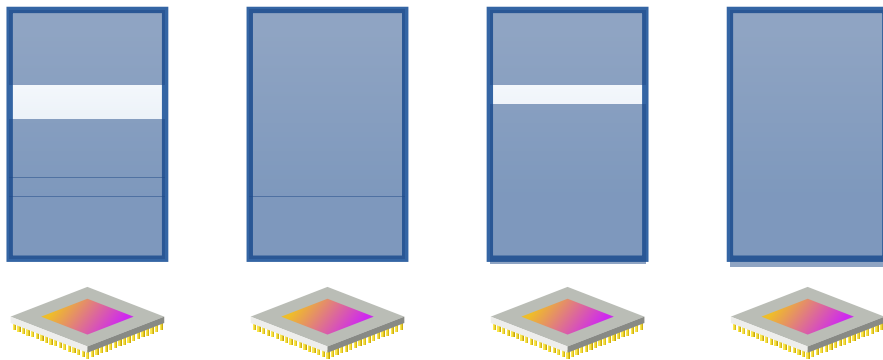CPU 0 decrements and flushes; global count is now 0. What about true count?

global

CPU 0
1
2
3
t

true

The true count is the sum of everything up to right now.
But the global count only reflects the blue region.
Operations in the orange region are still cached.

# Refcache example

Initially: Global count is 1, no cached deltas (so true count is 1)

CPU 1 increments after flush, before CPU 0's decrement

CPU 0 decrements and flushes; global count is now 0.  What about true count?

global

CPU 0
1
2
3
t

true

The true count is the sum of everything up to right now.
But the global count only reflects the blue region.
Operations in the orange region are still cached.

# Refcache example

Initially: Global count is 1, no cached deltas (so true count is 1)

CPU 1 increments after flush, before CPU 0's decrement

CPU 0 decrements and flushes; global count is now 0.  What about true count?

Global count now reflects cached ops

global

CPU 0
1
2
3
t

true

The true count is the sum of everything up to right now.
But the global count only reflects the blue region.
Operations in the orange region are still cached.

# Refcache example



Initially: Global count is 1, no cached deltas (so true count is 1)

CPU 1 increments after flush, before CPU 0's decrement

CPU 0 decrements and flushes; global count is now 0.  What about true count?

Global count now reflects cached ops

global

CPU 0
1
2
3
t

true

The true count is the sum of everything up to right now.
But the global count only reflects the blue region.
Operations in the orange region are still cached.

Abort delete

# Refcache example

Initially: Global count is 1, no cached deltas (so true count is 1)

CPU 1 increments after flush, before CPU 0's decrement

CPU 0 decrements and flushes; global count is now 0.  What about true count?

Global count now reflects cached ops

global

CPU 0
1
2
3
t

true

Refcache enables time- and space-efficient scalable reference counting with minimal latency.

Operations in the orange region are still cached.

Abort delete

# Bringing it all together

s r w x file          cores
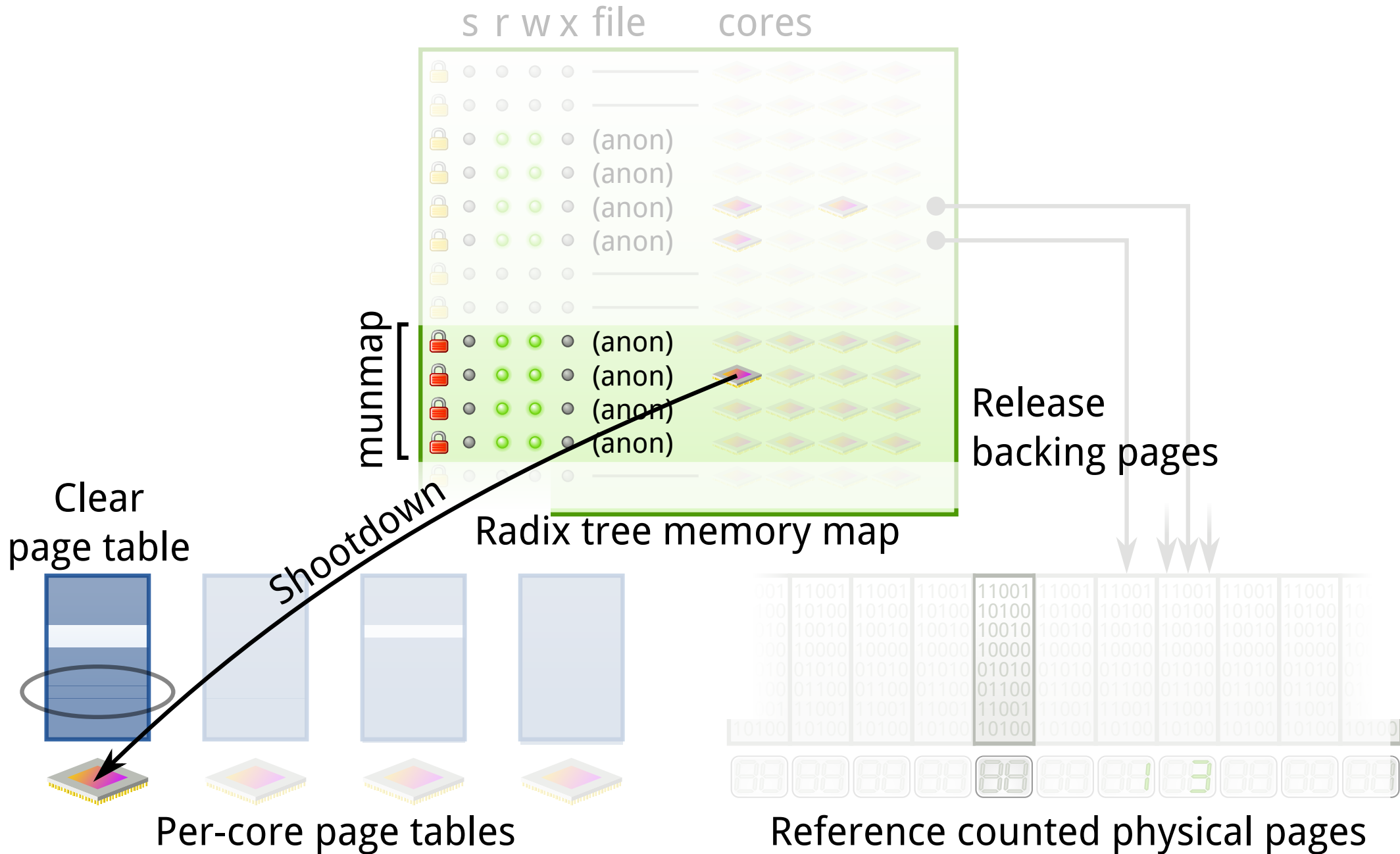
(anon)
(anon)
(anon)
(anon)

Radix tree memory map

Per-core page tables

Reference counted physical pages

# Bringing it all together

s r w x file        cores

(anon)
(anon)
(anon)
(anon)

mmap
(anon)
(anon)
(anon)
(anon)

**Radix tree memory map**

Per-core page tables

Reference counted physical pages

# Bringing it all together

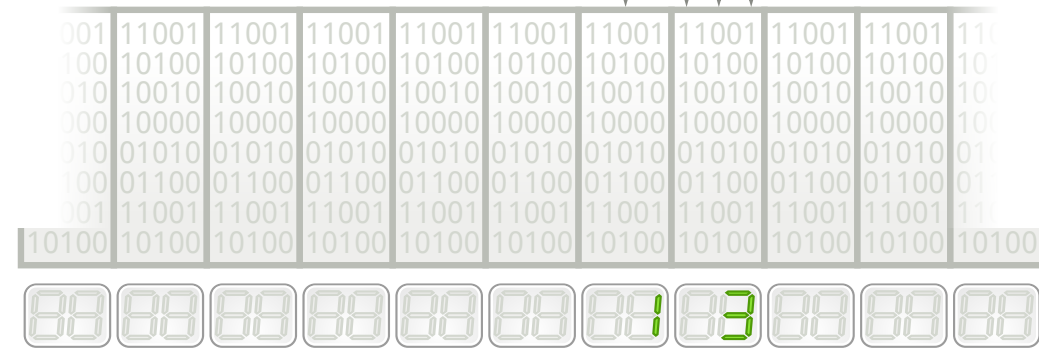s r w x file    cores

Radix tree memory map

Per-core page tables

Reference counted physical pages

# Bringing it all together

s r w x file          cores



Radix tree memory map

Page fault

Per-core page tables

Reference counted physical pages

# Bringing it all together

s r w x file          cores

Record faulting CPU

Page fault

(anon)

Allocate
backing page

Radix tree memory map

Install in local page table

Per-core page tables

Reference counted physical pages

# Bringing it all together

s r w x file          cores



Radix tree memory map

Per-core page tables

Reference counted physical pages

RadixVM: Scalable address spaces for multithreaded applications

# Bringing it all together



s r w x file        cores

Release
backing pages

munmap

Shootdown

Radix tree memory map

Clear
page table

Per-core page tables

Reference counted physical pages

# Bringing it all together

s r w x file        cores



Radix tree memory map

Per-core page tables

Reference counted physical pages

# Implementation

We built RadixVM in a custom research kernel.

We believe RadixVM could be built in a mainstream kernel.

All benchmarks are source-compatible with Linux.

# The other 99% is perspiration

Booting 80 cores (ACPI, x2APIC, IOMMU, oh my!)

NUMA-aware everything (memory allocation, per-CPU data, etc)

Performance analysis tools (NMI profiling, PEBS, load latency profiling, statistics counters)

Hardware curve balls (false sharing, bad prefetch behavior, etc)

All necessary for good results; all standard engineering.

Does parallel mmap/munmap matter to applications?

Are all of RadixVM's components necessary for scalability?

# RadixVM improves application scalability

Metis multicore MapReduce [Mao '10], inverse indexing application

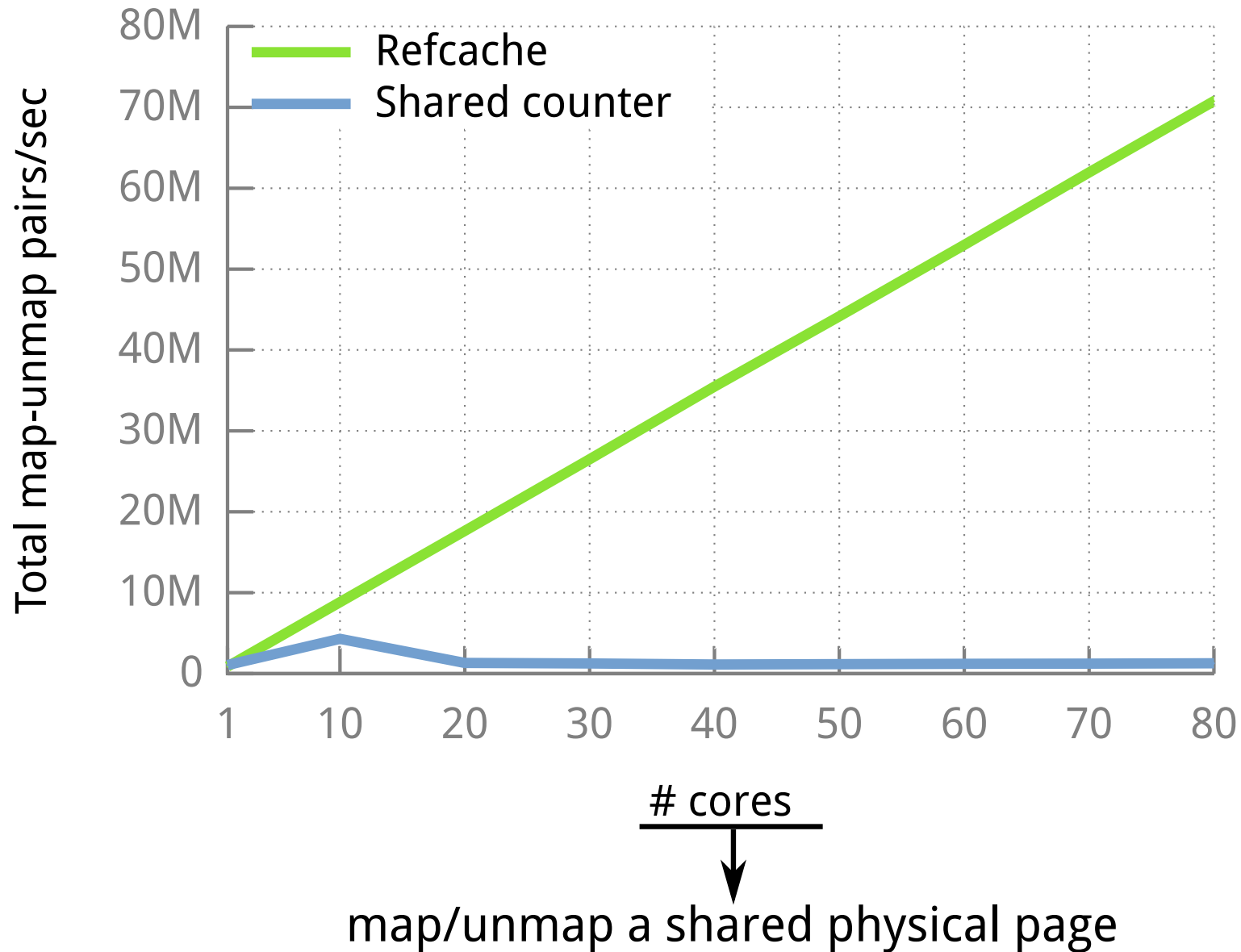# RadixVM improves application scalability

Metis multicore MapReduce [Mao '10], inverse indexing application
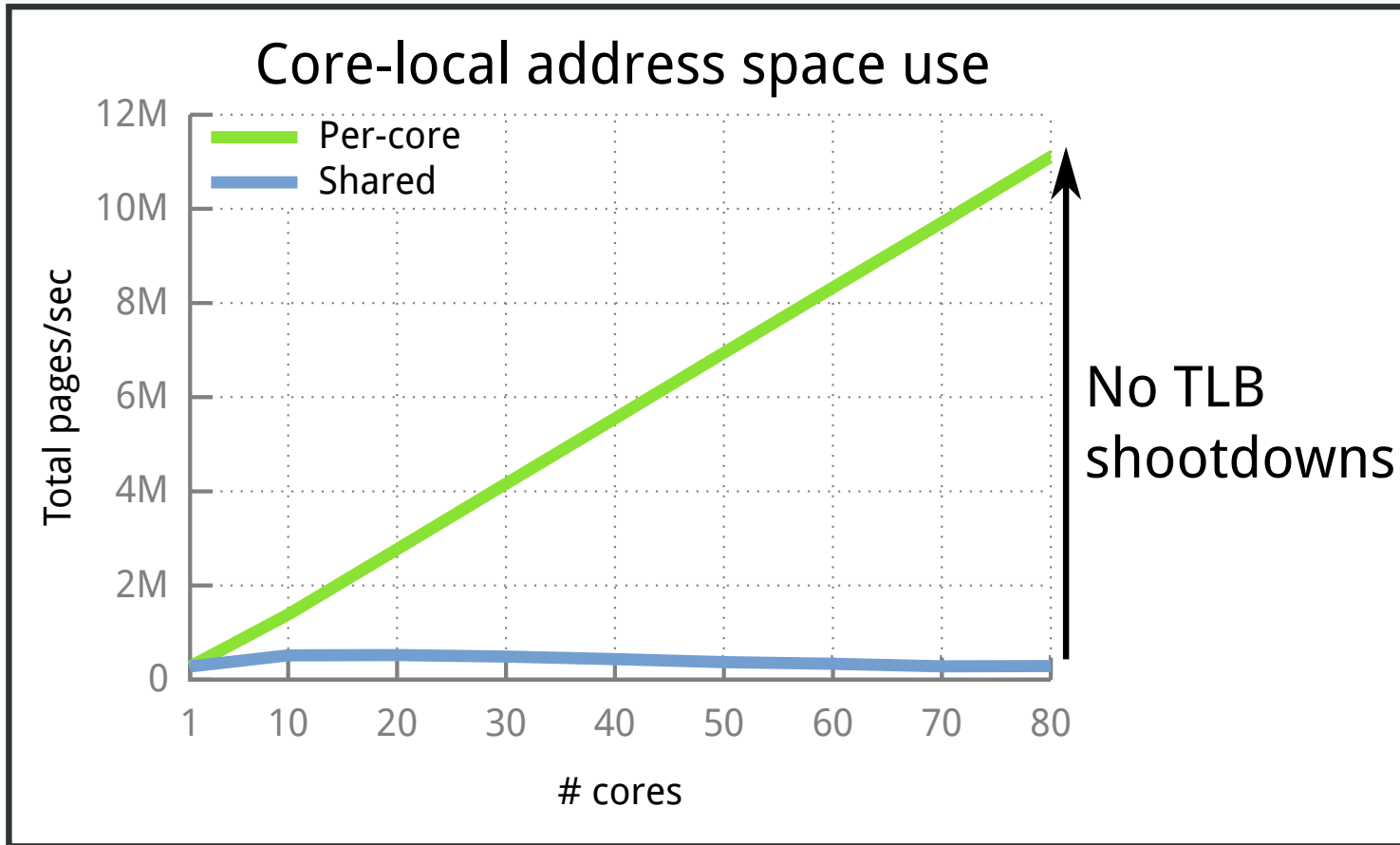
# RadixVM improves application scalability

Metis multicore MapReduce [Mao '10], inverse indexing application



RadixVM: Scalable address spaces for multithreaded applications

# Radix trees avoid communication



~No communication
Linear scalability

lookup existing keys    insert/delete random keys

# Refcache avoids cache line sharing



map/unmap a shared physical page

# Targeted TLB shootdown improves scalability

# Targeted TLB shootdown improves scalability

## Core-local address space use



Total pages/sec (y-axis): 0, 2M, 4M, 6M, 8M, 10M, 12M

# cores (x-axis): 1, 10, 20, 30, 40, 50, 60, 70, 80

Legend:
- Per-core
- Shared

No TLB shootdowns

# Targeted TLB shootdown improves scalability

## Core-local address space use

Total pages/sec

12M
10M
8M
6M
4M
2M
0

— Per-core
— Shared

No TLB ↑

1   10   20

## Global address space use

Total pages/sec

600k
500k
400k
300k
200k
100k
0

— Per-core
— Shared

Per-core overhead ↑

Page table contention ↓

1   10   20   30   40   50   60   70   80

# cores

RadixVM: Scalable address spaces for multithreaded applications

# Related work

Scalable VM systems
- K42 [Krieger '06]
- Corey [Boyd-Wickizer '08]
- Bonsai [Clements '12]

Scalable reference counters
- Modula-2+ local refs [DeTreville '90]
- Distributed counters [Appavoo '07]
- Scalable non-zero indicators [Ellen '07]
- Sloppy counters [Boyd-Wickizer '10]

# Conclusion



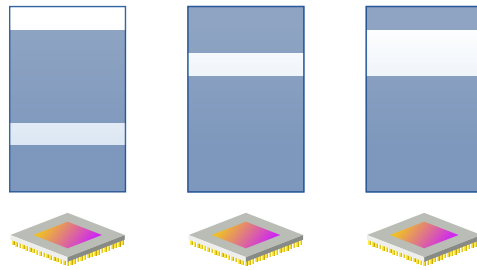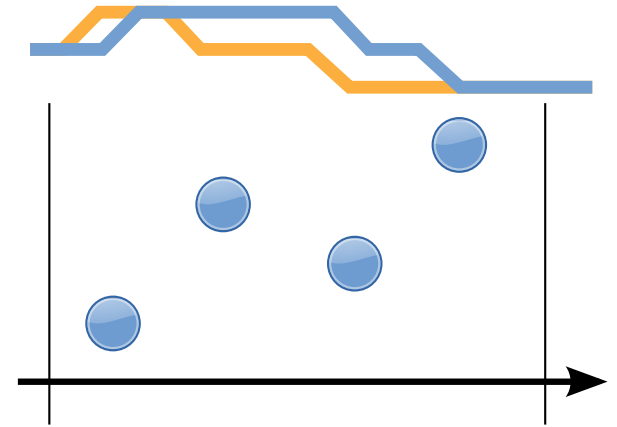Radix trees



Per-core page tables



Refcache

# Conclusion
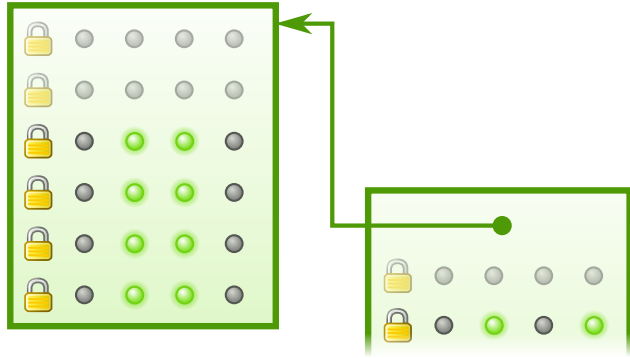


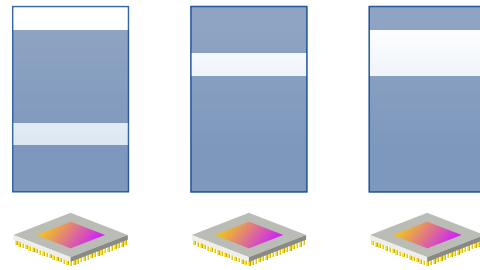Radix trees · Per-core page tables · Refcache

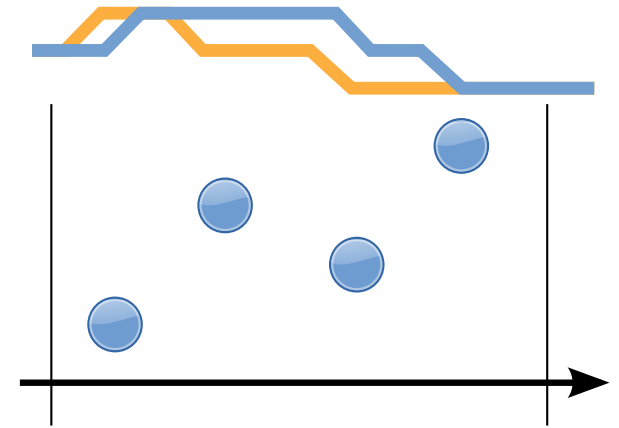## Perfect scalability for non-overlapping VM operations

# Conclusion



Radix trees

Per-core page tables

Refcache

## Perfect scalability for non-overlapping VM operations

## Check it out: http://pdos.csail.mit.edu/multicore