# LATR: Lazy Translation Coherence

**Mohan Kumar**\*, Steffen Maass\*, Sanidhya Kashyap, Ján Veselý‡,
Zi Yan‡, Taesoo Kim, Abhishek Bhattacharjee‡, Tushar Krishna

Georgia Institute of Technology ‡Rutgers University

\* Co-First Authors

March 28, 2018

Supermicro Debuts 8-Socket Server for Intel Xeon Processors

By Sue Smith / NewsFactor Network

PUBLISHED:
OCTOBER

S upermicro just announced the latest addition to its line of
SuperServer systems, designed for data centers and
SuperServer 7089P-TR4T is an
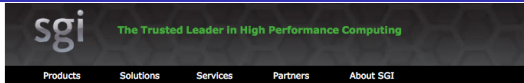ver for Intel Xeon scalable

Large NUMA machines

Large NUMA machines

Terabytes of memory

Large NUMA machines

Terabytes of memory

Microsecond latency

$\Rightarrow$ Problem of Microsecond Latency in System Services
$\Rightarrow$ TLB Coherence is Contributor in Important Subset

Large NUMA machines

Terabytes of memory

Microsecond latency

# Impact of TLB coherence on applications

- Multi-core MapReduce application
  - Prior research: **10x increase in shootdown time** with increasing core counts
- Web servers (e.g., Apache)
  - Prior research and our findings: $\approx$**35% of time spent in TLB shootdown**
- Die-stacked Memory
  - Swapping between on-chip and off-chip memory
- Disaggregated Memory
  - Swapping between local and remote memory

# Impact of TLB coherence on applications

- Multi-core MapReduce application
  - Prior research: **10x increase in shootdown time** with increasing core counts
- Web servers (e.g., Apache)
  - Prior research and our findings: ≈**35% of time spent in TLB shootdown**
- Die-stacked Memory
  - Swapping between on-chip and off-chip memory
- Disaggregated Memory
  - Swapping between local and remote memory

⇒ Can we mitigate this costly TLB shootdown?

# Table of contents

# Table of contents

# Translation lookaside buffer: Introduction

- Cache for virtual $\rightarrow$ physical mapping, per-core structures
- Accessed on every load/store
- Unlike data caches (L3, etc.), coherence managed by OS
- TLB coherence significantly impacts application performance

# TLB coherence: Background

- **Hardware-based Approaches**
  - Providing cache coherence to TLBs
  - ISA-level instruction support (ARM)
  - Microcode-based approaches

- **Software-based Approaches**
  - Current commodity OS design: Use **Inter-Processor Interrupts (IPI)**
  - Optimization: Reduce number of shootdowns, better tracking
  - Multikernel design: Use Message-Passing

# TLB coherence: Background

- **Hardware-based Approaches**
  - Providing cache coherence to TLBs

    $\Rightarrow$ More Hardware Complexity

- **Software-based Approaches**

    $\Rightarrow$ TLB shootdowns still significant

  - Optimization: Reduce number of shootdowns, better tracking
  - Multikernel design: Use Message-Passing

- munmap() on core 1, application running on cores 1, 2, and 5:



| | | | | | | | |
|---|---|---|---|---|---|---|---|
| | | | Application | | | | |
| App$_1$ | App$_2$ | Idle | Idle | App$_5$ | Idle | Idle | Idle |
| OS | OS | ⋯ | | OS | | | |
| | | | Operating System | | | | |

❶

| Core$_1$ | Core$_2$ | Core$_3$ | Core$_4$ |
|---|---|---|---|
| TLB | TLB | TLB | TLB |

| Core$_5$ | Core$_6$ | Core$_7$ | Core$_8$ |
|---|---|---|---|
| TLB | TLB | TLB | TLB |

Timeline:   ❶

# TLB shootdown internals in Linux

- `munmap()` on core 1, application running on cores 1, 2, and 5:

# TLB shootdown internals in Linux

- Context switch on core 1, local TLB shootdown:

# TLB shootdown internals in Linux

- Notify cores 2 and 5 via IPI, application blocked on core 1:

# TLB shootdown internals in Linux

- Execute context switch and TLB shootdown on cores 2 and 5:

# TLB shootdown internals in Linux

- Cores 2 and 5 respond ACK via shared memory:

# TLB shootdown internals in Linux

- Control is returned on all cores, TLB shootdown completed:



❶ munmap()
❷ Local Shootdown
❸ Send IPIs
❹ Remote Shootdown
❺ IPI ACK
❻ munmap() complete

Savings potential for asynchronous approach with LATR

# Observation

- **Synchronous TLB shootdown is expensive:**
  - Up to $6\,\mu$s delay with two sockets

- **Processing IPIs is expensive:**
  - Interrupt handler on remote core
  - Long wait time on initiating core

- **IPI send-and-wait delay:**
  - Unicast delivery of the IPIs (one at a time)

# TLB shootdown: A necessary evil

- Cost of a simple memory unmap operation (`munmap()`):
  - 1 page on 16 cores with 2 sockets: **up to 8 $\mu$s**
  - $\approx 70\%$ **from TLB shootdown alone**
- More expensive with more sockets:

# TLB shootdown: A necessary evil

- Cost of a simple memory unmap operation (`munmap()`):
    - 1 page on 16 cores with 2 sockets: **up to 8 $\mu$s**
    - $\approx 70\%$ **from TLB shootdown alone**
- More expensive with more sockets:

# TLB shootdown: A necessary evil

- Cost of a simple memory unmap operation (`munmap()`):
    - 1 page on 16 cores with 2 sockets: **up to 8 $\mu$s**
    - $\approx 70\%$ **from TLB shootdown alone**
- More expensive with more sockets:

# Table of contents

# In this talk: LATR

- LATR: **La**zy **Tr**anslation Coherence
- **Perform asynchronous TLB shootdown**
  - Remove remote shootdown from the critical path
  - Take advantage of change in ABI without affecting applications'
    correctness
- **Use shared memory instead of IPI**
  - Eliminate send-and-wait delay of IPIs
- **Scope:**
  - *free* operations (in this talk)
  - *migration* operations (see our paper)

# In this talk: LATR

- LATR: **La**zy **Tr**anslation Coherence
- **Perform asynchronous TLB shootdown**
    - Remove remote shootdown from the critical path
    - Take advantage of change in ABI without affecting applications' correctness
- **Use shared memory instead of IPI**
    - Eliminate send-and-wait delay of IPIs
- **Scope:**
    - *free* operations (in this talk)
    - *migration* operations (see our paper)

    ⇒ But: How to perform asynchronous shootdown?

# LATR States

- Store virtual addresses to be flushed
- Remote cores shootdown local TLB during
    - OS context switch
    - OS scheduler tick (**upper bound**: 1ms in Linux)

# LATR: Example

- munmap() initiated on core 1:



| | | Application | | | | | |
|---|---|---|---|---|---|---|---|
| App₁ | App₂ | Idle | Idle | App₅ | Idle | Idle | Idle |
| OS | OS | ⋯ | | OS | | | |
| | | Operating System | | | | | |

| Core₁ | Core₂ | Core₃ | Core₄ | | | Core₅ | Core₆ | Core₇ | Core₈ |
|---|---|---|---|---|---|---|---|---|---|
| LATR States | LATR States | LATR States | LATR States | | | LATR States | LATR States | LATR States | LATR States |

Timeline: ❶

# LATR: Example

- munmap() initiated on core 1:

- Set up LATR state (for cores 2 and 5), local shootdown:



| | App₁ | App₂ | Idle | Idle | App₅ | |
|---|---|---|---|---|---|---|
| | OS | OS | ··· | | OS | |

❶ munmap()
❷ Local Shootdown
❸ Create LATR State

| Core₁ | Core₂ | Core₃ | Core₄ |
|---|---|---|---|
| LATR States | LATR States | LATR States | LATR States |

Core₁, LATR State₁:

| start | end | mm | flags | Core list | active |
|---|---|---|---|---|---|
| 0x01 | 0x0F | 0x1234 | 0x1 | {2, 5} | True |

Timeline: ❶······❷··❸

- Return control on core 1. Time taken: $2.3\,\mu s$, 70% reduction:



❶ munmap()
❷ Local Shootdown
❸ Create LATR State
❹ munmap() complete

| App$_1$ | App$_2$ | Idle | Idle | App$_5$ |
| OS | OS | $\cdots$ | | OS |

| Core$_1$ | Core$_2$ | Core$_3$ | Core$_4$ |
| LATR States | LATR States | LATR States | LATR States |

| Cor | | | |
| LATR States | States | States | States |

Core$_1$, LATR State$_1$:

| start | end | mm | flags | Core list | active |
|-------|------|--------|-------|-----------|--------|
| 0x01  | 0x0F | 0x1234 | 0x1   | {2, 5}    | True   |

Timeline: ❶ ❷ ❸ ❹
$\leftarrow$ 2.3µs $\rightarrow$

# LATR: Example

- Scheduler tick on core 2, local shootdown, reset state:



**Core₁, LATR State₁:**

| start | end | mm | flags | Core list | active |
|-------|-------|--------|-------|-----------|--------|
| 0x01 | 0x0F | 0x1234 | 0x1 | {5} | True |

- Scheduler tick on core 5, local shootdown, reset state:



❶ munmap()
❷ Local Shootdown
❸ Create LATR State
❹ munmap() complete
❺ Shootdown Core₂
❻ Shootdown Core₅

| App₁ | App₂ | Idle | Idle | App₅ |
|------|------|------|------|------|
| OS | OS | ⋯ | | OS |

| Core₁ | Core₂ | Core₃ | Core₄ |
|-------|-------|-------|-------|
| LATR States | LATR States | LATR States | LATR States |

Core₁, LATR State₁:

| start | end | mm | flags | Core list | active |
|-------|-----|------|-------|-----------|--------|
| 0x01 | 0x0F | 0x1234 | 0x1 | {} | False |

# LATR: Example

- Shootdown complete, LATR entry can be reused:



| App₁ | App₂ | Idle | Idle | App₅ |
|------|------|------|------|------|
| OS | OS | ⋯ | | OS |

❶ munmap()
❷ Local Shootdown
❸ Create LATR State
❹ munmap() complete
❺ Shootdown Core₂
❻ Shootdown Core₅
❼ Shootdown complete

Core₁, LATR State₁:

| start | end | mm | flags | Core list | active |
|-------|-----|-----|-------|-----------|--------|
| 0x01 | 0x0F | 0x1234 | 0x1 | {} | False |

# Lazy TLB shootdown: Correctness

- Same physical memory or virtual memory is reused
  - Leads to memory corruption
- $\Rightarrow$ Avoid same physical/virtual page reuse
  - Upper bound for TLB shootdown with LATR is 1ms
  - OS physical/virtual memory reclamation delayed by two scheduler ticks (2ms)
  - Memory overhead is bounded by 21 MB

# Lazy TLB shootdown: Incorrect accesses

- Memory accesses before LATR shootdown:
    - Consequence of incorrect application: Use After Free
    - Before LATR shootdown, access (reads and writes) allowed
    - Exists in the current OS implementation
    - After LATR shootdown, access results in segmentation fault

# Scope of LATR

- ABI change for *free* operations
- Support for operations limited to few, frequently used operations:

| Classification | Operations | Lazy operation possible |
|---|---|---|
| Free | munmap(): unmap address range | ✓ |
| | madvise(): free memory range | ✓ |
| Migration | AutoNUMA page migration ($\Rightarrow$ See paper) | ✓ |
| | Page swap: swap page to disk | ✓ |
| Permission | mprotect(): change page permission | - |
| Ownership | CoW: Copy on Write | - |
| Remap | mremap(): change physical address | - |

# Table of contents

- LATR prototype developed for Linux 4.10
- Evaluation questions
    - What are LATR's benefits with microbenchmarks?
    - What are LATR's benefits with real-world applications exhibiting many TLB shootdowns?
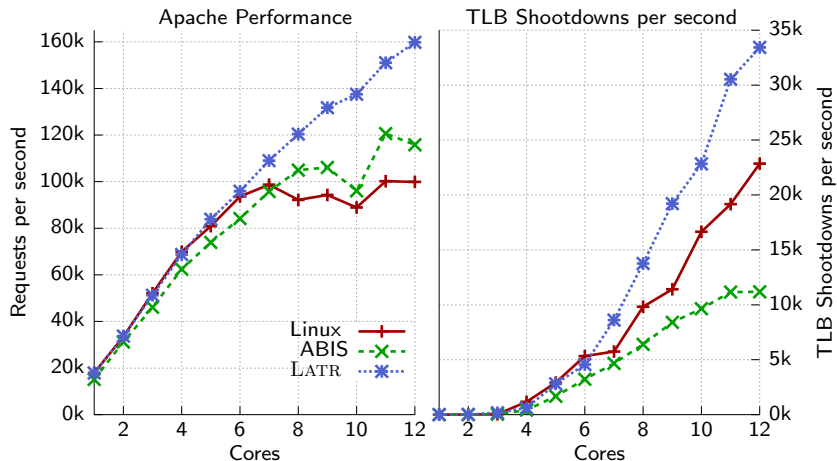    - What is the cost for LATR?

- Linux and LATR calling `munmap()` with one page on 120 cores:



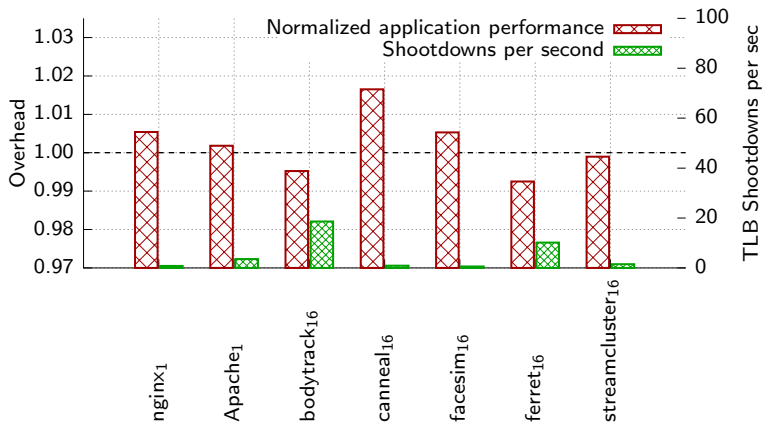$\Rightarrow$ **Up to 66.7% reduction for `munmap()`**

# Serving files with Apache

- Linux, ABIS [ATC17], and LATR on 2 sockets:



$\Rightarrow$ **Up to 59.9% more $\frac{requests}{second}$ than Linux, 37.9% higher than ABIS.**

# Cost of LATR

- Memory overhead is bounded by 21 MB
- Performance overheads for applications with few TLB shootdowns:



$\Rightarrow$ LATR shows small performance overheads of up to 1.7% due to added operations during scheduling.

# Future work

Further applications of LATR in:

- Disaggregated data centers
- Heterogeneous memory
- Applicability to PCID/ASID-based approaches
- Impact on new features such as KPTI, . . . ?

# LATR: Takeaways

- The **synchronous** TLB shootdown is expensive
- We propose a software-based **asynchronous** shootdown mechanism
- Significant improvement in application performance with LATR
  - **70%** reduction for `munmap()`, for 16-core and 120-core machines
  - Improves Apache's throughput by **60%**
- Asynchronous mechanism applicable to other services:
  - AutoNUMA (see our paper)

- The **synchronous** TLB shootdown is expensive
- We propose a software-based **asynchronous** shootdown mechanism
- Significant improvement in application performance with LATR
  - **70%** reduction for `munmap()`, for 16-core and 120-core machines
  - Improves Apache's throughput by **60%**
- Asynchronous mechanism applicable to other services:
  - AutoNUMA (see our paper)

**Thanks!**