# Pay Migration Tax to Homeland: Anchor-based Scalable Reference Counting for Multicores

**Seokyong Jung**, Jongbin Kim, Minsoo Ryu, Sooyong Kang, Hyungsoo Jung
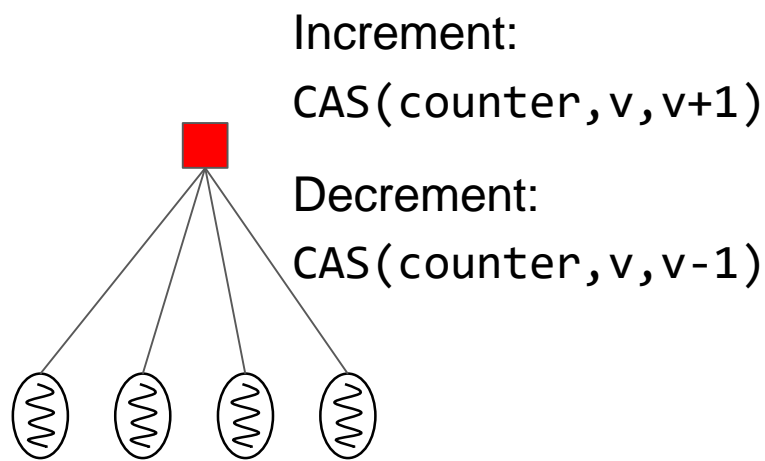
Hanyang University

# Reference counting

- It is
  - a general technique to manage the number of references for resources
  - mainly used to reclaim resources in timely manner
- Scalability is the most important challenge in multicore environment

1. REF (increase counter)
2. Use resource
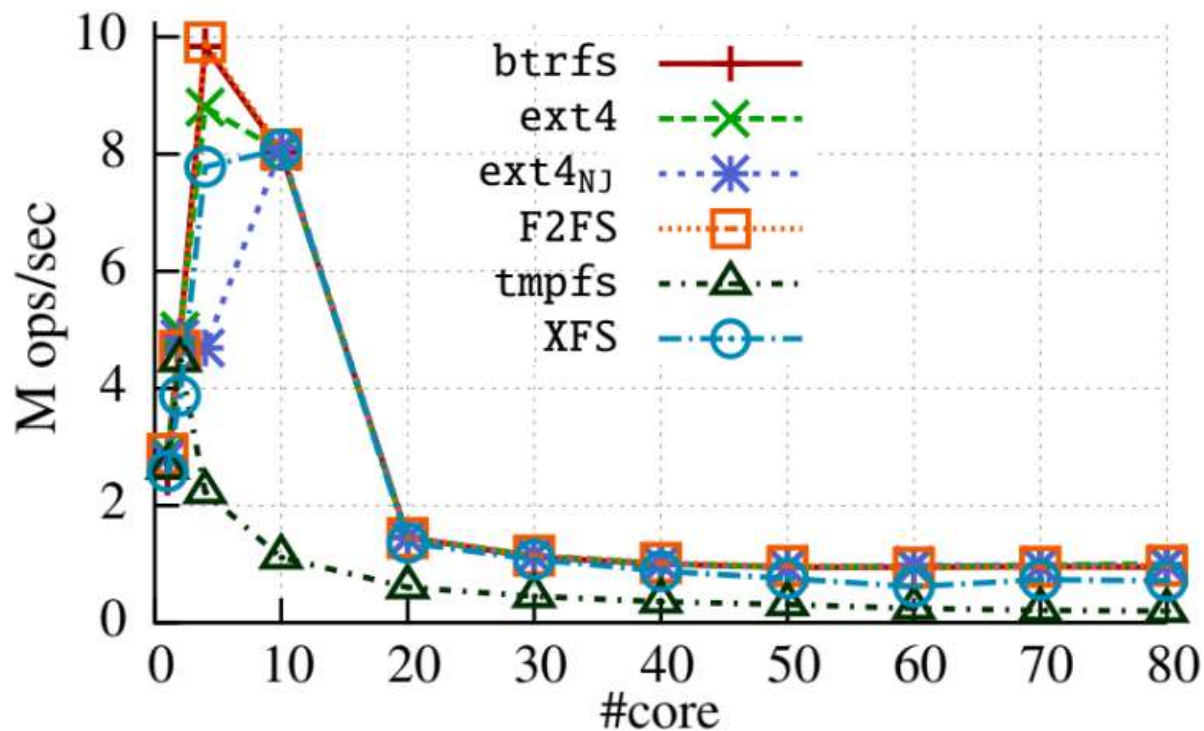3. UNREF (decrease counter)

# Known scalability issues of reference counting in Linux
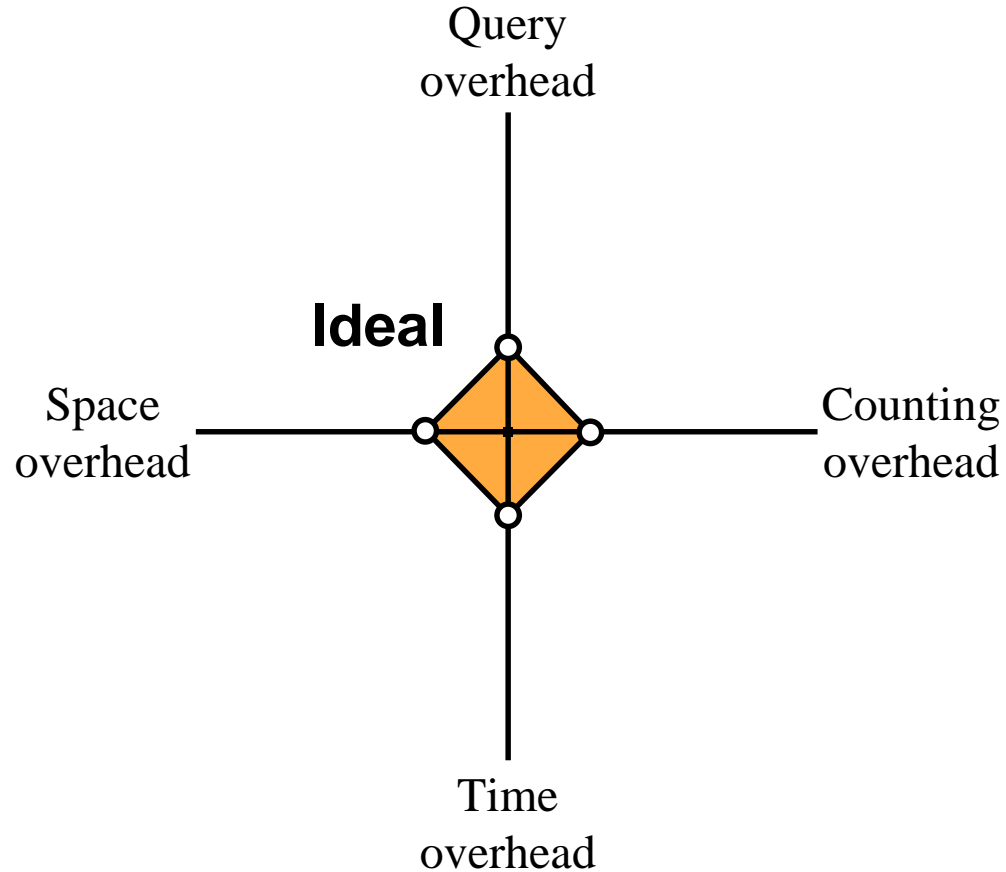
Thread  ■ Atomic counter  ☐ Counter

Increment:
CAS(counter,v,v+1)

Decrement:
CAS(counter,v,v-1)

Traditional Counting in Linux
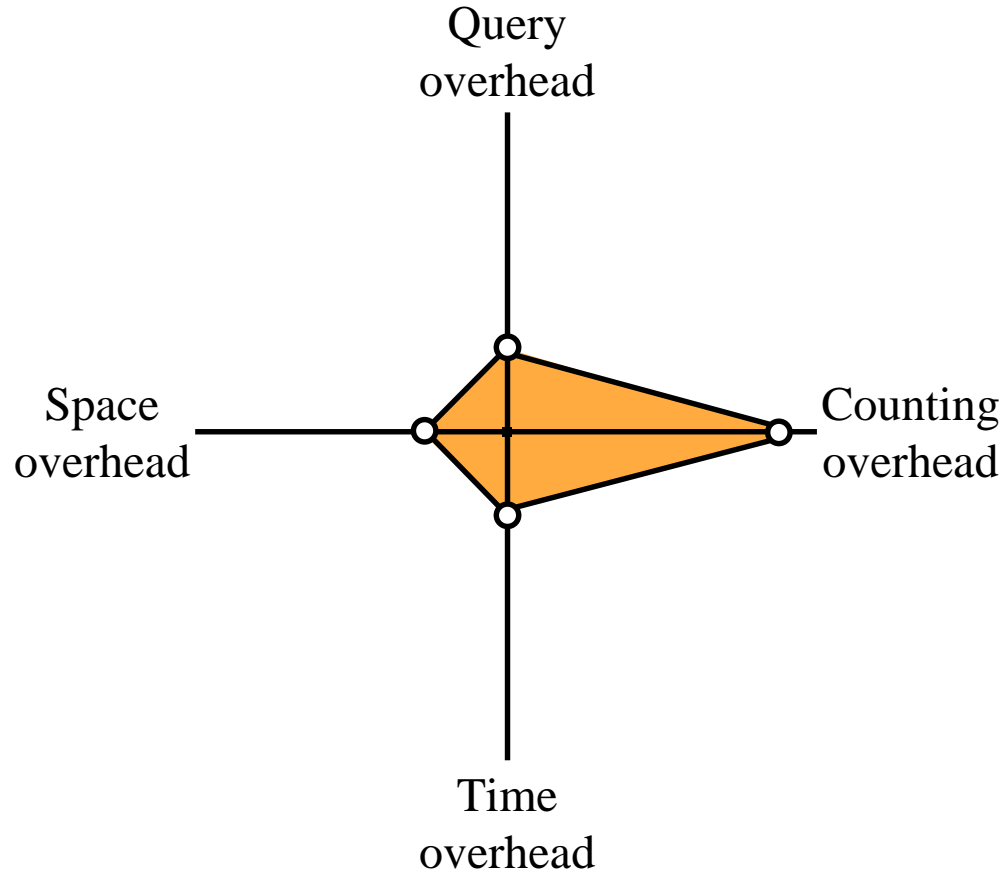


(c) DRBH

btrfs
ext4
$ext4_{NJ}$
F2FS
tmpfs
XFS

Read throughput for a same page in a single file
(Min et al. ATC'16)

3
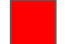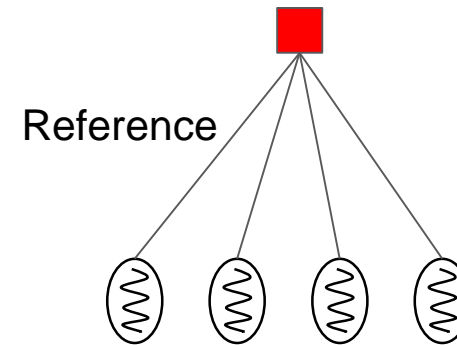
# Four performance metrics we established



- Counting Overhead
  - Cost for updating a reference counter
- Query Overhead
  - Cost for checking if a reference counter is zero
- Space Overhead
  - Space required for reference counter itself
- Time Overhead
  - Time for synchronizing between internal structures for reference counting
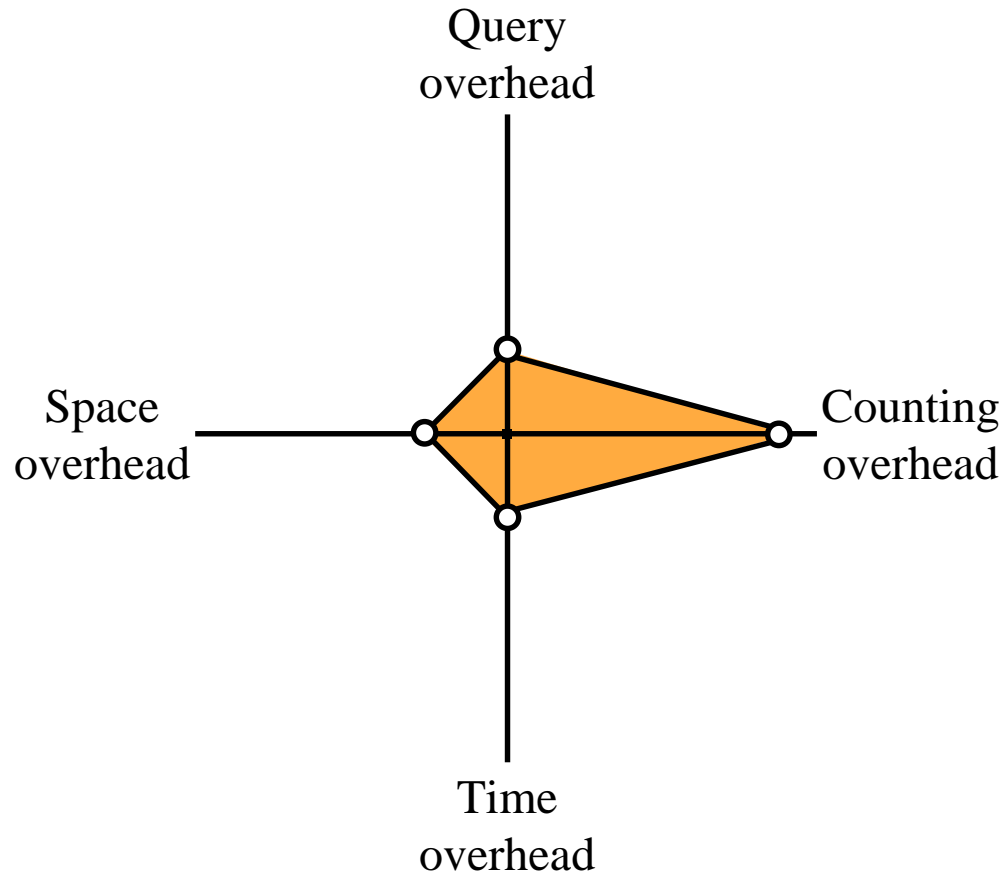
# Overhead analysis of prior proposals
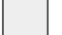


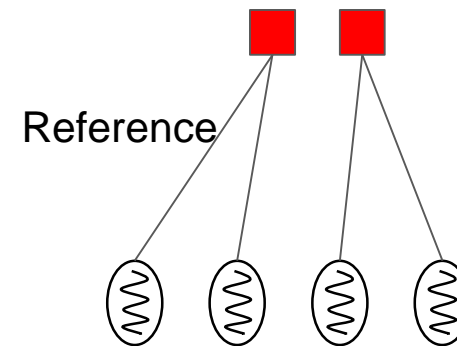Thread     Atomic counter     Counter

**Traditional Counting**

Reference

☑ Low query overhead

☑ Low space overhead

☑ Low time overhead

✗ High counting overhead
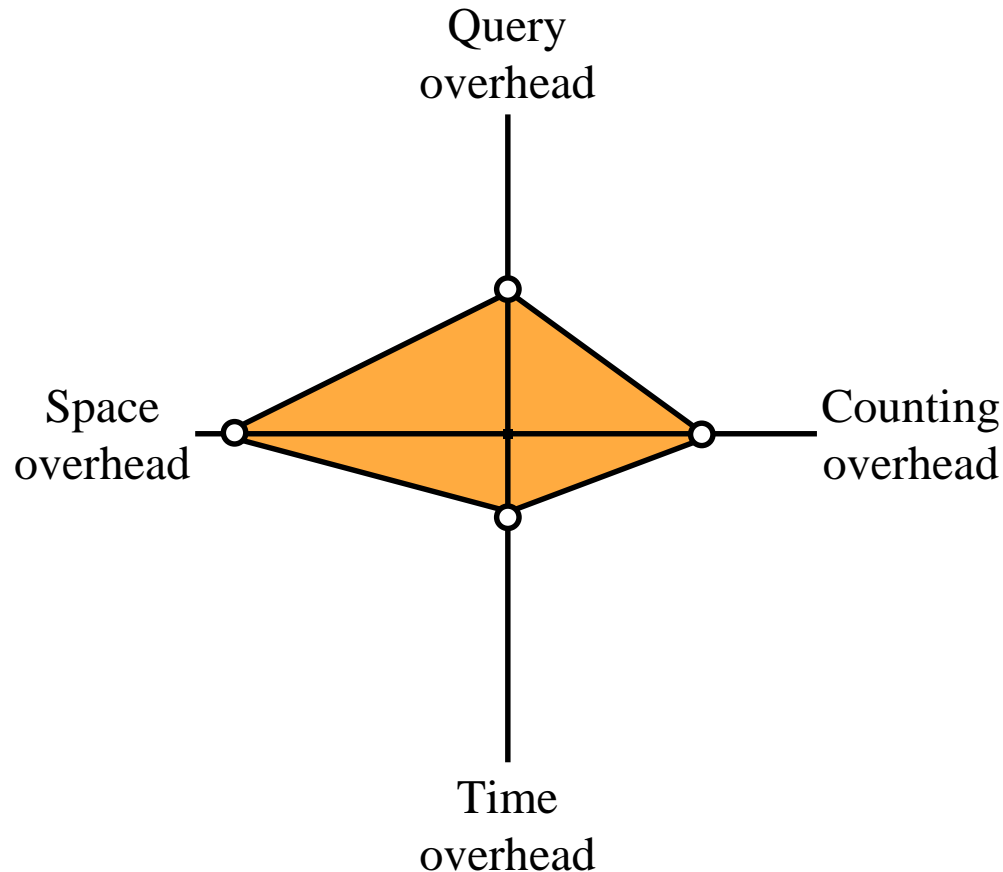
# Overhead analysis of prior proposals (cont.)



Query overhead

Space overhead

Counting overhead

Time overhead

Thread    Atomic counter    Counter

**Contention Distribution**

Reference

✓ Better counting overhead

✗ Worse space overhead

✗ Worse query overhead

# Overhead analysis of prior proposals (cont.)

Query
overhead

Space
overhead

Counting
overhead

Time
overhead
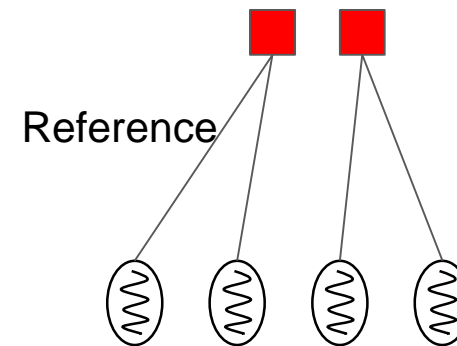
**Prior proposals:**
SNZI (Ellen et al., SOSP'07)
Carrefour (Dashti et al., SIGPLAN Notices (2013))
Doppel (Nurula et al., OSDI'14)
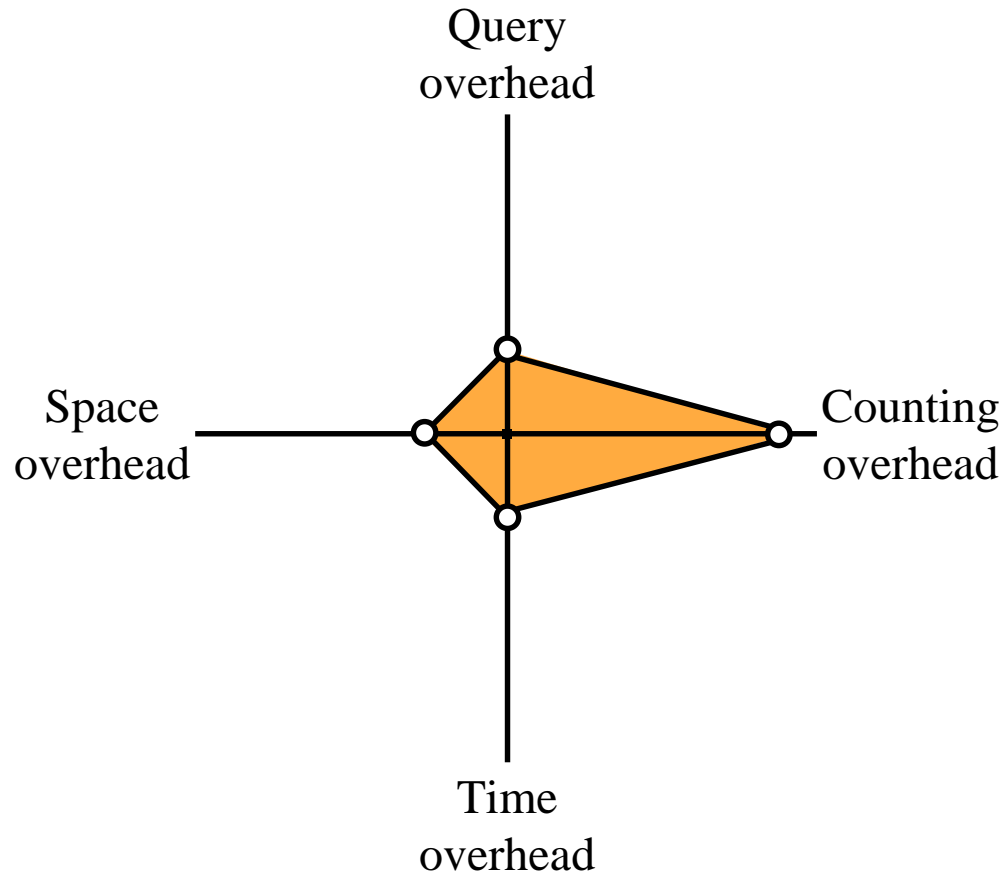Dynamic SNZI (Acar et al., PPoPP'17)

Thread    Atomic counter    Counter

**Contention Distribution**

Reference

☑  Better counting overhead

✗  Worse space overhead

✗  Worse query overhead

# Overhead analysis of prior proposals (cont.)



Query overhead

Space overhead

Counting overhead

Time overhead

Thread    Atomic counter    Counter

**Cache Affinity**

Reference

core    core    core    core
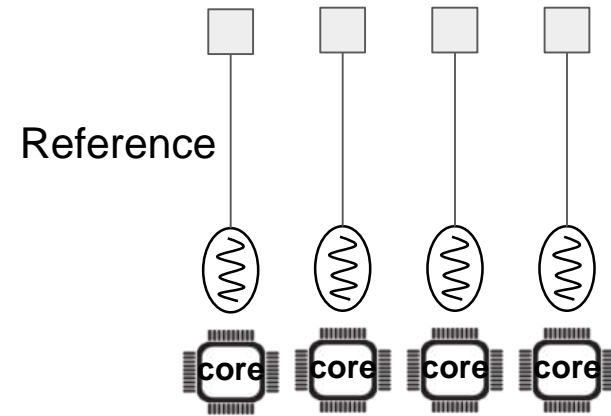
✓ Better counting overhead

✗ Worse space overhead

✗ Worse query overhead

# Overhead analysis of prior proposals (cont.)



Query overhead

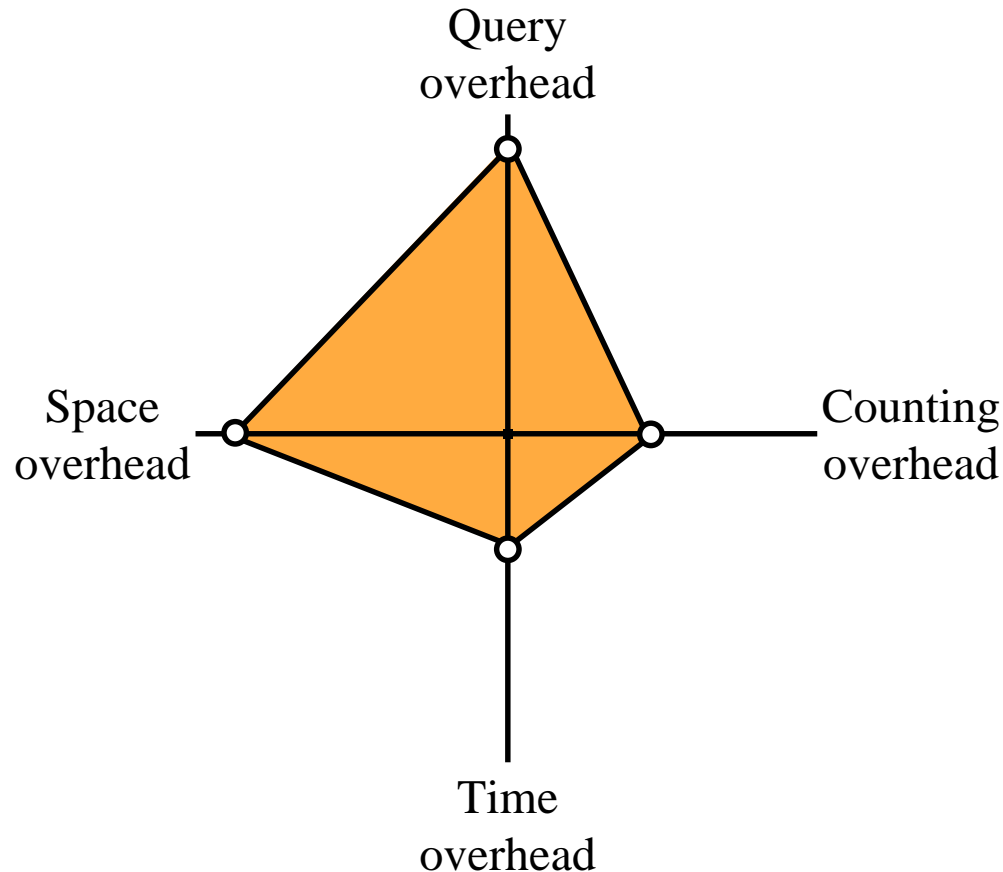Space overhead

Counting overhead

Time overhead

**Prior proposals:**
`percpu_counter` structure in Linux (2006)
Sloppy counter (Boyd-Wickizer et al., OSDI'10)
`percpu_ref` structure in Linux (2013)

Thread   Atomic counter   Counter

**Cache Affinity**

Reference

core   core   core   core

☑ Better counting overhead

✗ Worse space overhead

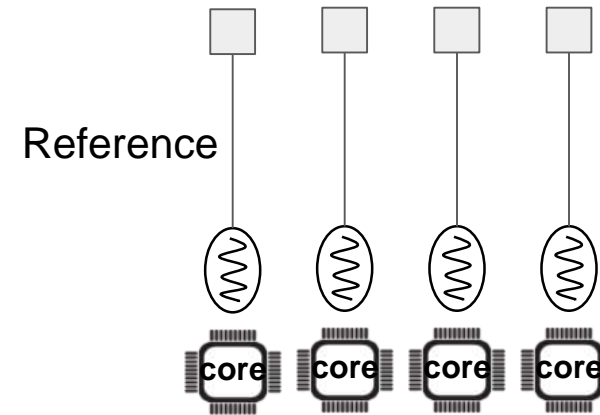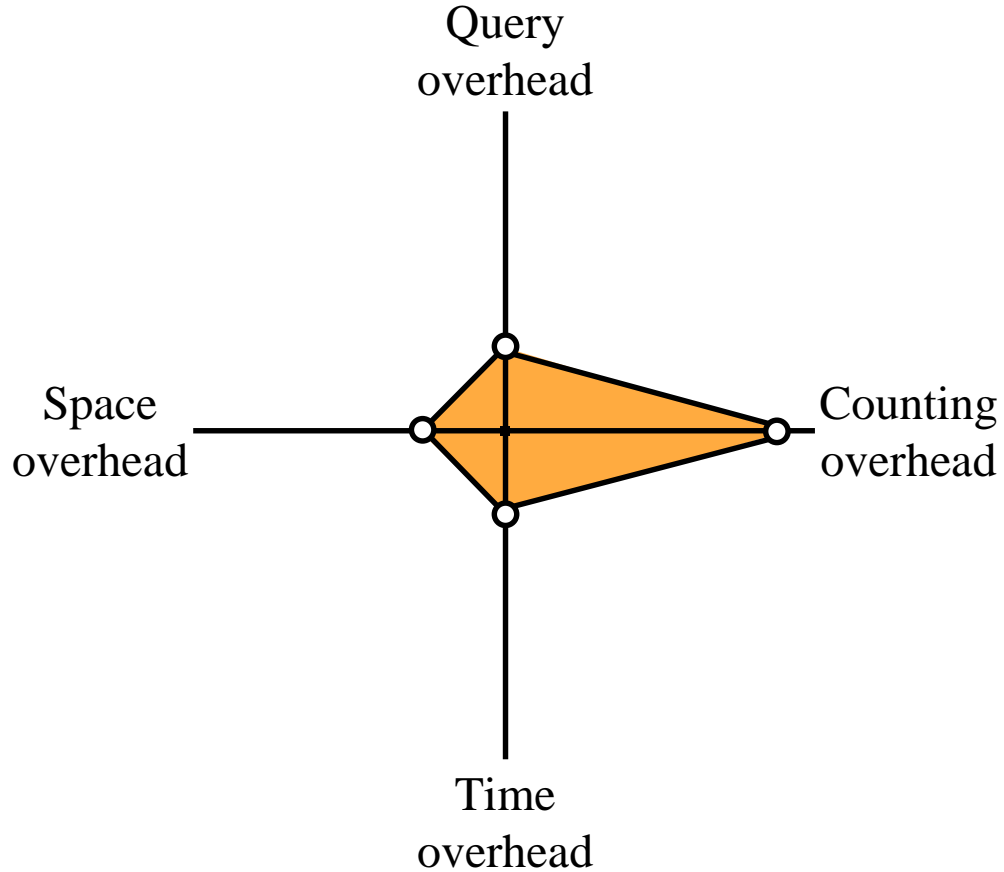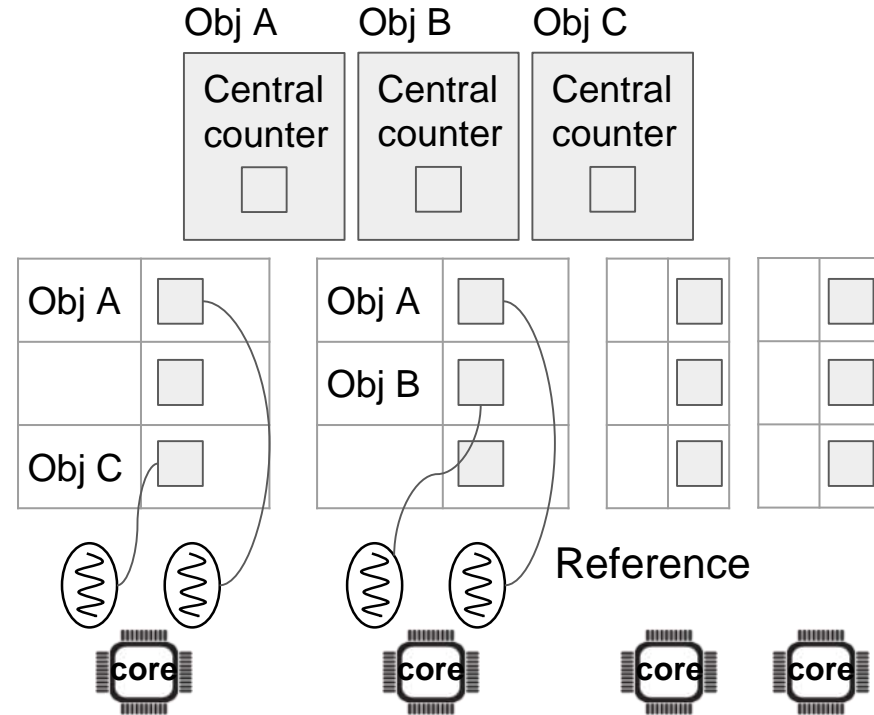✗ Worse query overhead

# Overhead analysis of prior proposals (cont.)

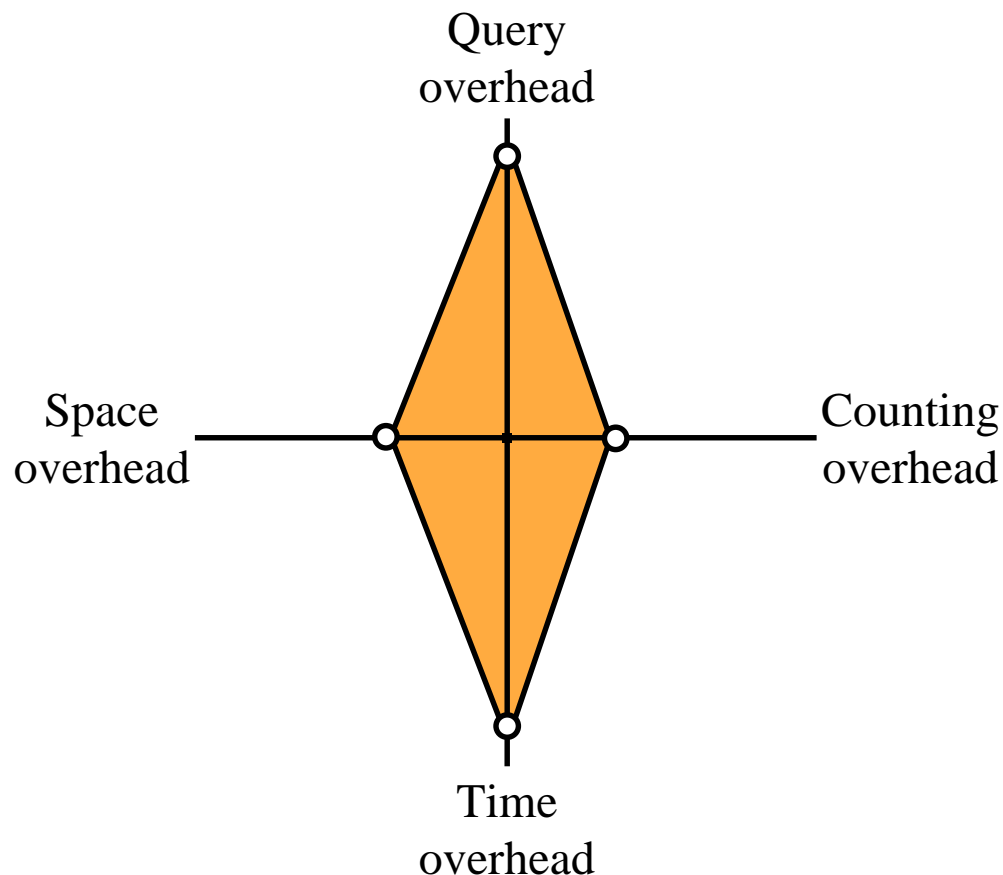# Overhead analysis of prior proposals (cont.)

Query
overhead

Space
overhead

Counting
overhead

Time
overhead

**Prior proposals:**
Reference counting using Linux RCU (Mckenny et al., TR (2013))
RadixVM (Clements et al., EuroSys'13)
OpLog (Boyd-Wickizer, PhD thesis (2014))
ScaleFS (Bhat et al., SOSP'17)

Thread    Atomic counter    Counter

**Per-core Hash**

Obj A    Obj B    Obj C

Central
counter

Central
counter

Central
counter

Obj A

Obj C

Obj A

Obj B

Reference

core    core    core    core
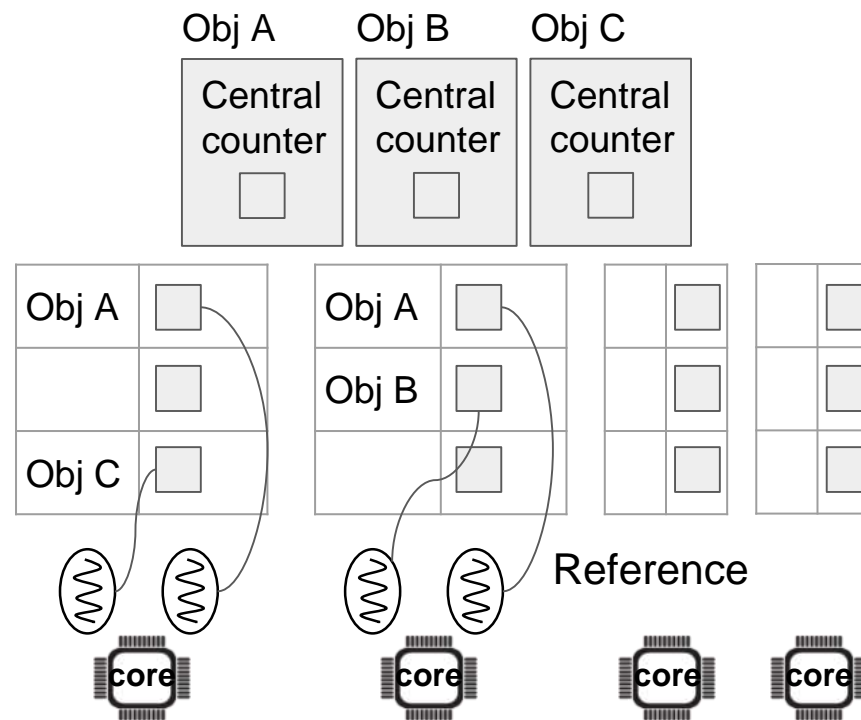
☑ Better counting overhead

✕ Worse query overhead

✕ Worse time overhead

# Overhead analysis of prior proposals (cont.)



**Per-core Hash**

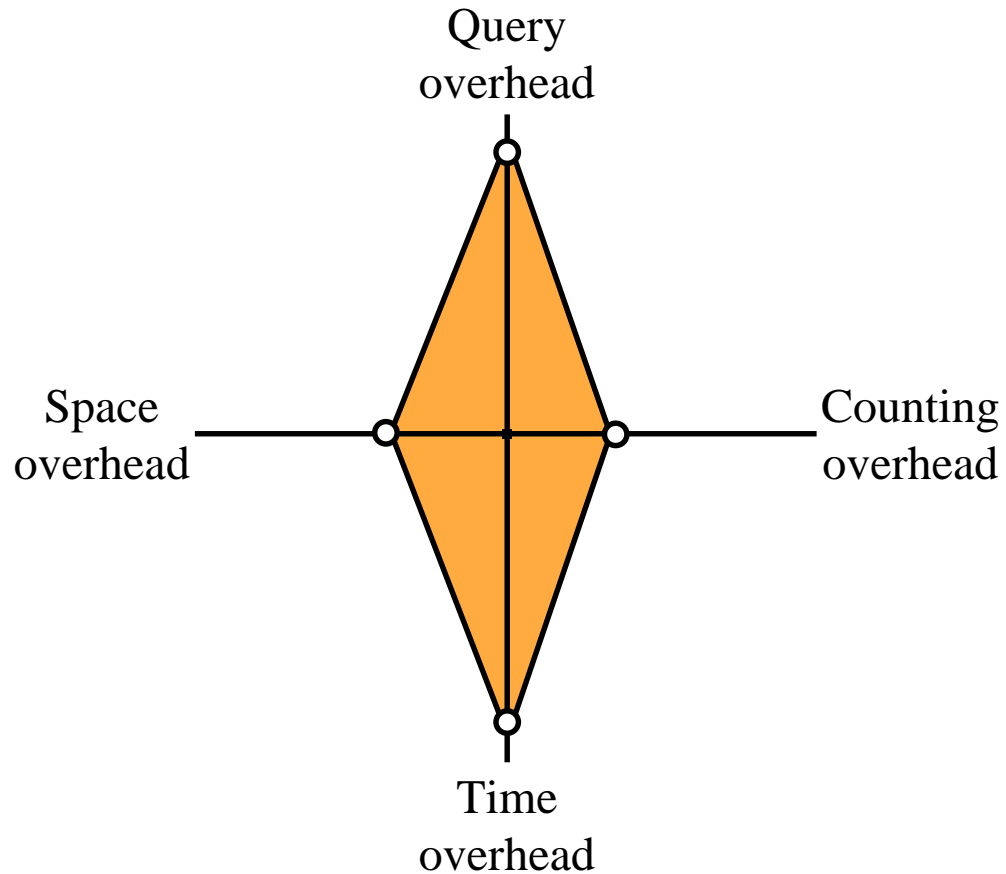# Overhead analysis of prior proposals (cont.)



**Per-core Hash**

# Overhead analysis of prior proposals (cont.)



Query overhead

Space overhead

Counting overhead

Time overhead

**Per-core Hash (RefCache)**



Throughput (Mops/sec)

400

300

200

100

0

256    1024    4096    16384

Number of Shared Objects (Pages)

# Summarizing all these ...

# Summarizing all these ...

# Challenges for the space-time tradeoff

Obj A

Obj B

Obj C

Obj D

Obj E

Obj F

Obj G

Obj H

**core**   **core**

**Cache Affinity**

Obj A  Obj B  Obj C  Obj D  Obj E  Obj F  Obj G  Obj H

Central counter

0

**core**                    **core**

**Per-core Hash**

# Challenges for the space-time tradeoff (cont.)

# Our solution to this issue ...

**Anchoring**



Obj A

Central counter

0

Obj A    0

⚓ Core 0

**core**    **core**

REF A    UNREF A

# Our solution to this issue ...

## **PayGo**
## Pay Migration Tax to Homeland

# Issue for a single local counter

# Anchoring in action

# Paygo
## (**Pay** migration tax as you **go** to other core)

- Low counting overhead

  - Scalable for local counters

- Low space overhead (per core hash)

  - Proportional to the number of CPU cores

- Query overhead is still high

  - Escaping the counting-query tradeoff is beyond the scope of this work.

# Overhead Analysis for a Reference Counter

# Data structures in PayGo

(per-core)

PAYGO ← core core

(per-process)

task struct ←

anchor info

anchor core IDs

PAYGO entry

| object pointer | local counter | anchor counter | |
|---|---|---|---|
| 63      56 | 55      52 | 51      48 | 47      0 |

cache line size (byte)

- Local counter increases the local count by REF operation
- A process records core IDs along with object pointer when REF operation is performed

# Extending PayGo design to support user-level objects

user threads ·········· object

user mode

**sys_ref**(void* *obj*) or **sys_unref**(void* *obj*)

REF(*obj*, *pid*)
UNREF(*obj*, *pid*)

**preempt_disable()**          **preempt_enable()**

**referencing** or **unreferencing**

kernel mode

# Experimental Setup

Kernel: Linux 4.12.5

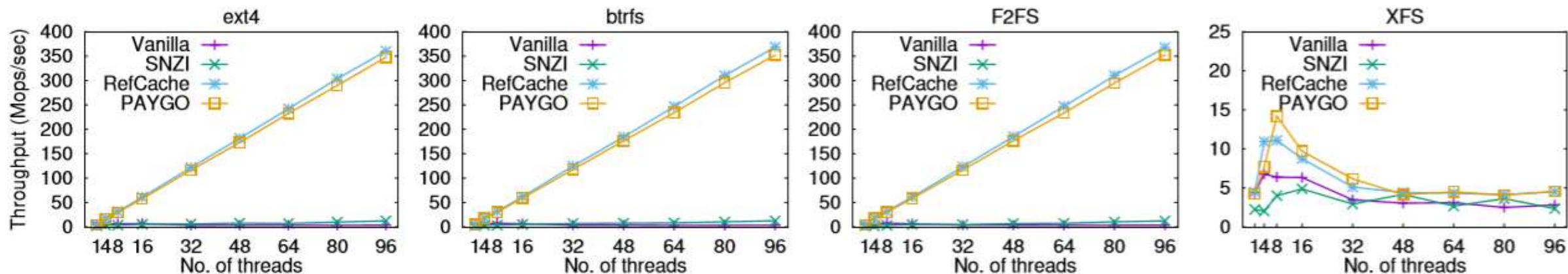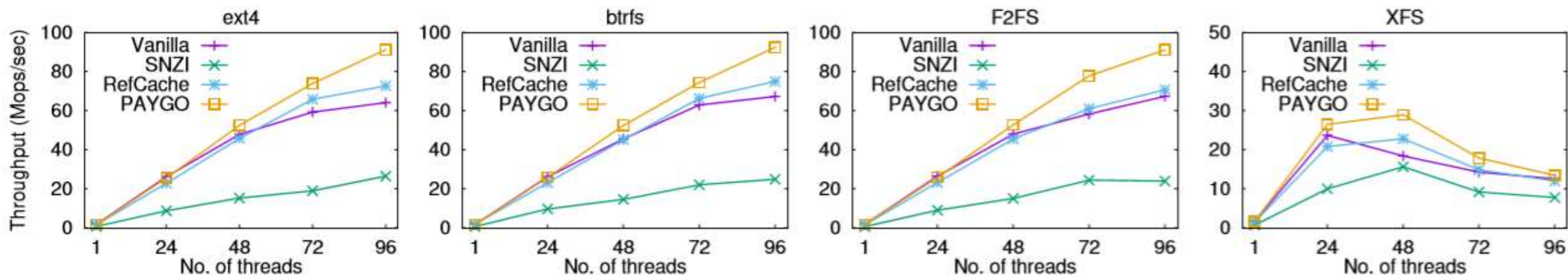CPU: four 24-core Intel Xeon E7-8890 v4 CPUs

RAM: 1 TiB DDR4 DRAM

Storage: Samsung SM1725 NVMe SSD

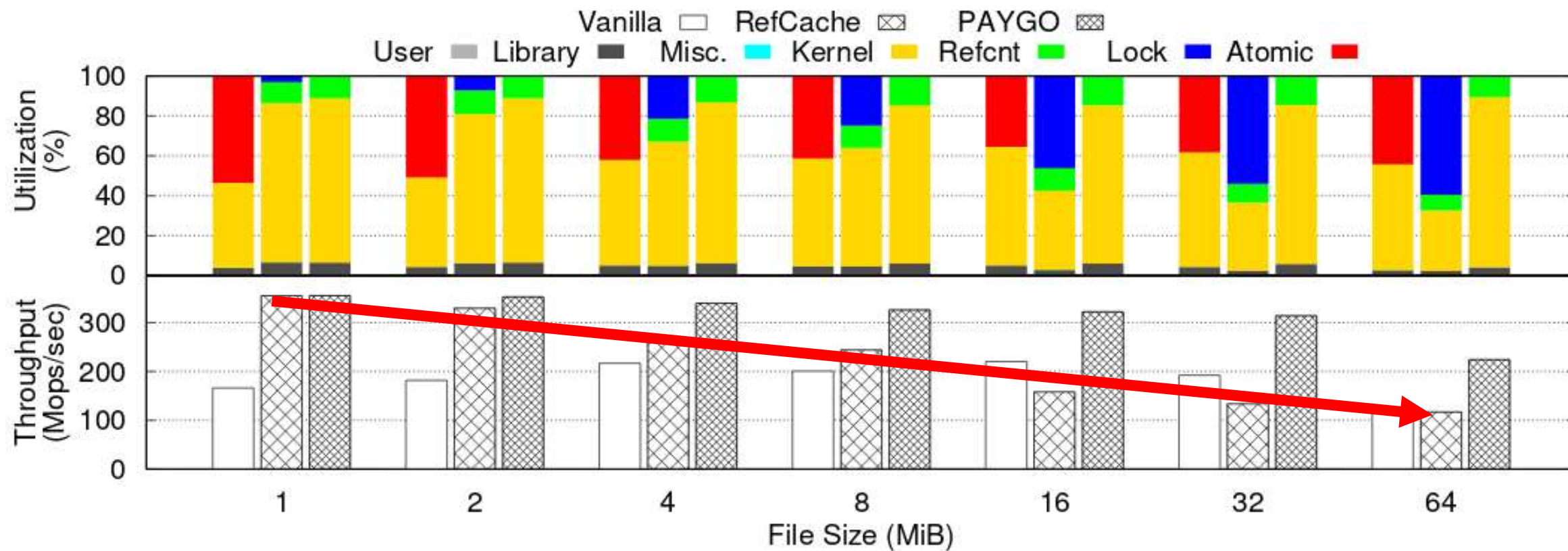# Scalability Comparison of the Linux Page Cache



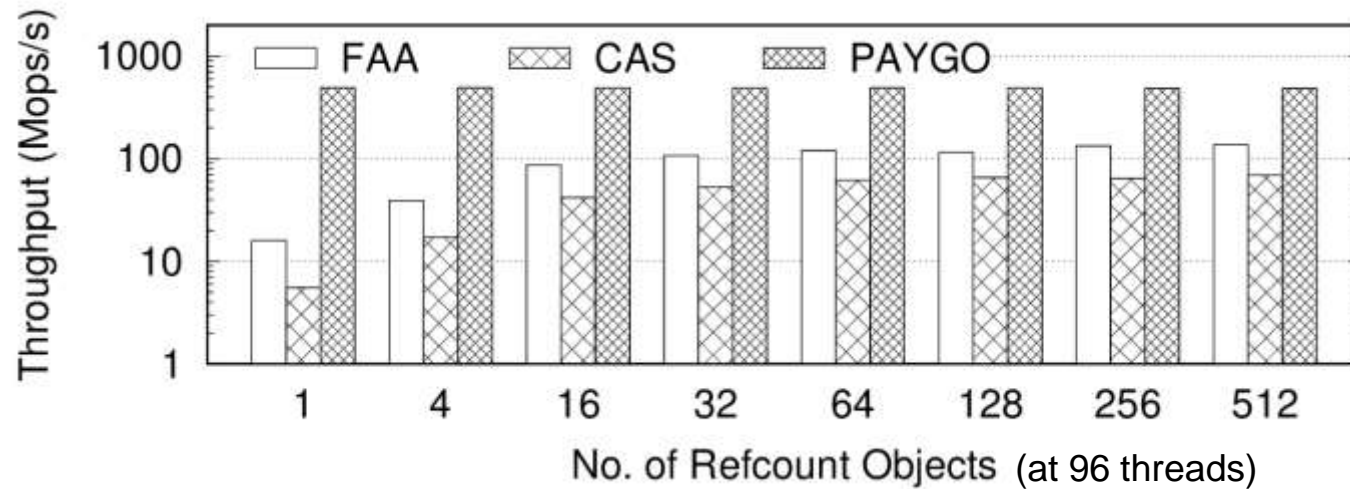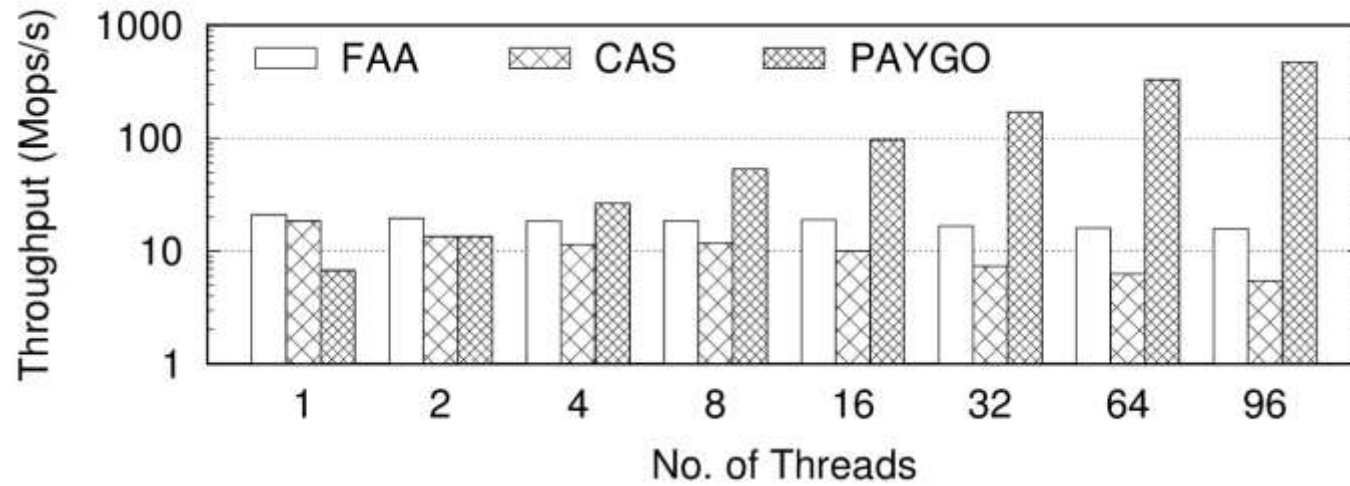Strongly contending workloads: FxMark DRBH workload



Weakly contending workloads: filebench modified fileserver workload

# Performance Spectrum on Varying Contention Levels

# Scalability of User-level Paygo



No. of Threads

No. of Refcount Objects (at 96 threads)

# Conclusion

- Designing scalable reference counting techniques should consider space-time tradeoff as well as counting-query tradeoff.

- PayGo escapes the space-time tradeoff by using anchoring technique.

- PayGo provides scalable counting and space efficiency with negligible time delay for reclaiming obsolete hash entries.