

Exploiting superpages in a nonvolatile memory file system

Sheng Qiu, Narasimha Reddy
Texas A&M University

What is SCM?

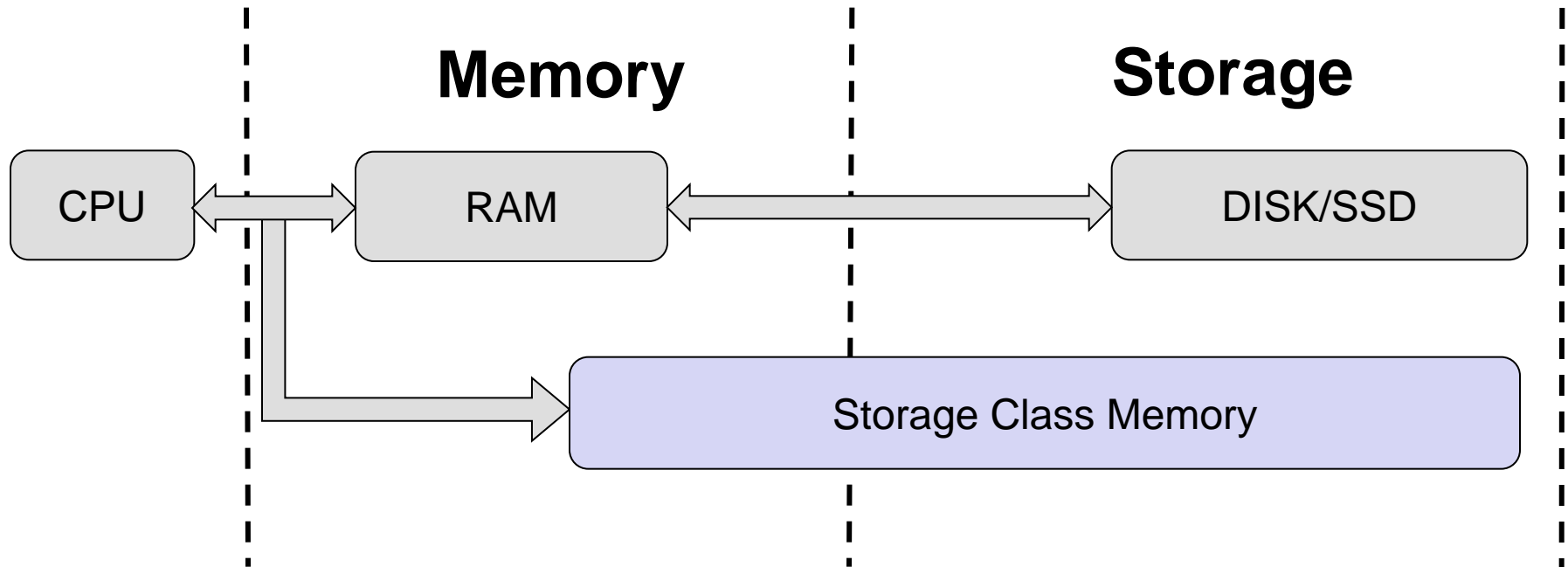
- **Storage Class Memory**
 - Byte-addressable, like DRAM
 - Non-volatile, persistent storage
- **Example: Phase Change Memory**

PCM Attributes

Attributes	PCM	DRAM	NAND	NOR	EEPROM
Bit Alterable	Green	Green	Red	Red	Green
Non-volatile	Green	Red	Green	Green	Green
Read Speed	Yellow	Green	Red	Yellow	Yellow
Write Speed	Yellow	Green	Yellow	Yellow	Red
Scaling	Green	Yellow	Yellow	Yellow	Red

From Numonyx

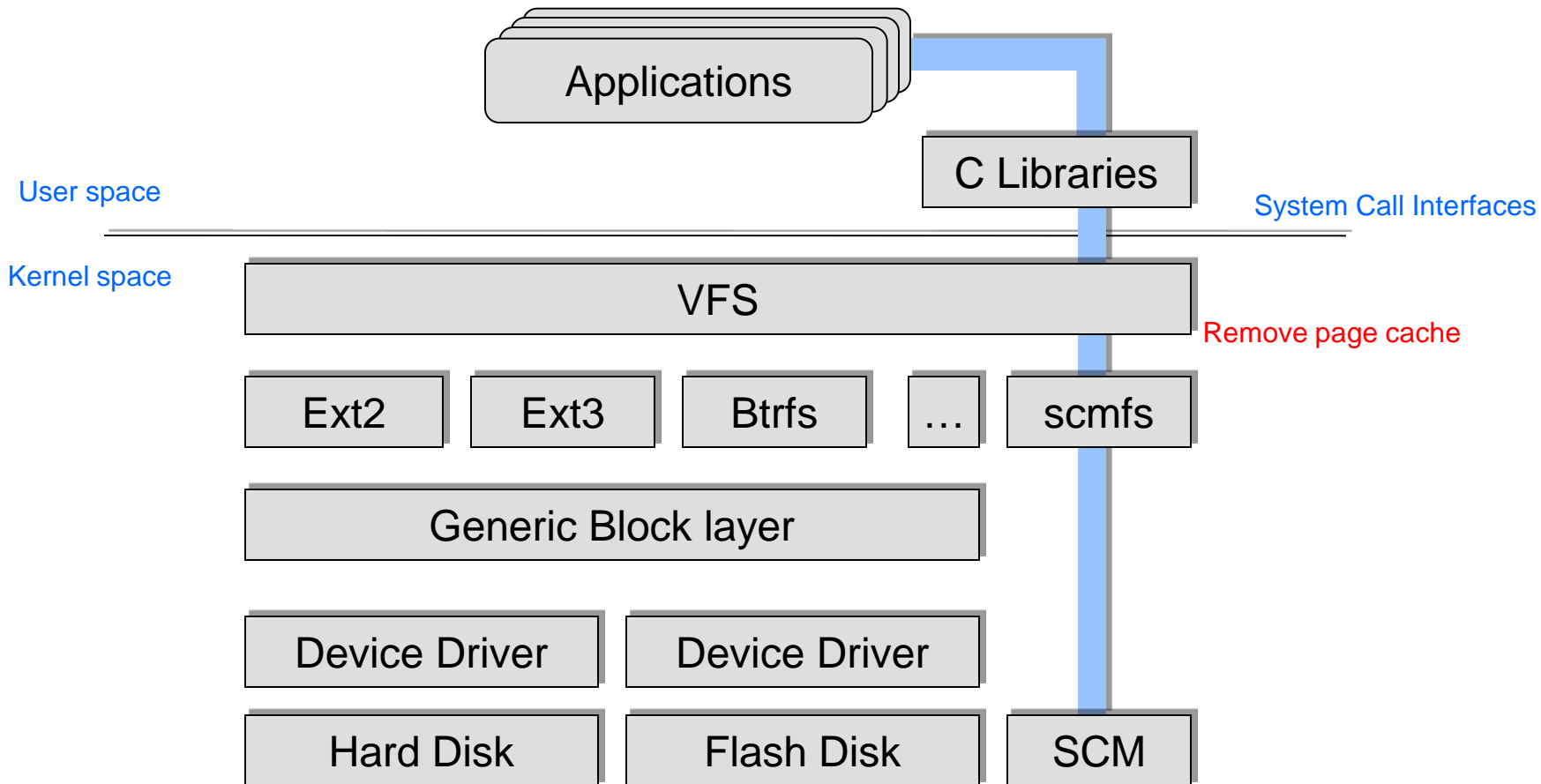
Hardware hierarchy



How to use SCM as storage?

- **Device level**
 - Use existing file system on RamDisk
- **File system level**
 - Design a new file system

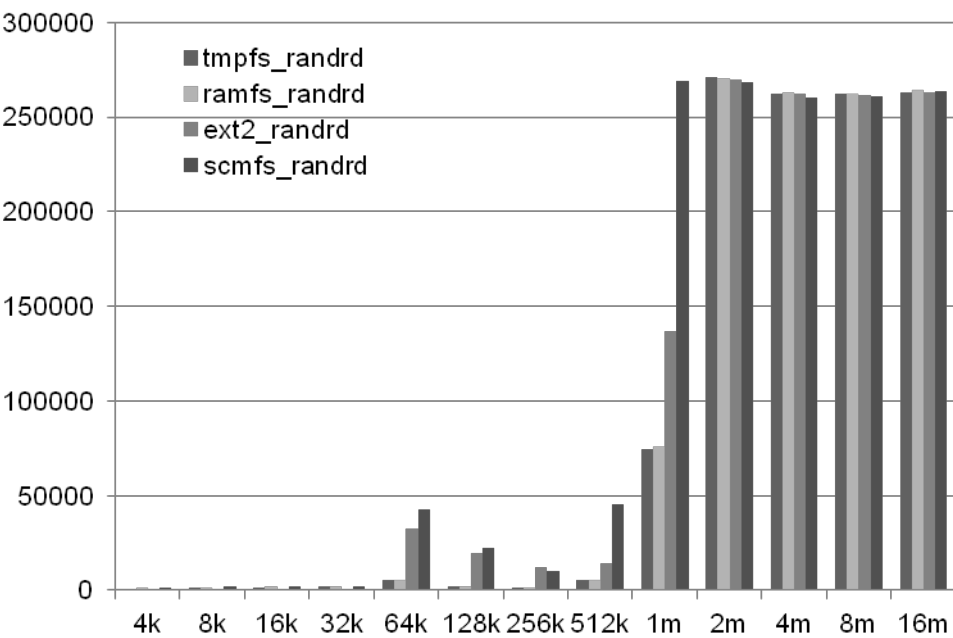
New FS on SCM



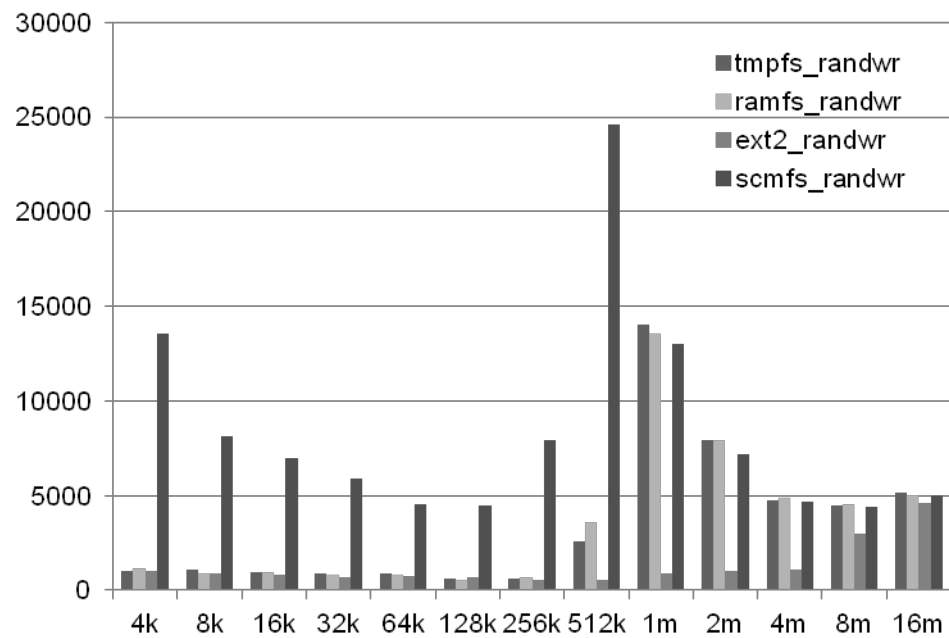
Wrap up

- Re-utilize Memory Management (MM) module in O/S to do block management
- Implement the file system in Virtual Address Space
- Keep all the files contiguous in Virtual Address Space
- **Cause more Data TLB Misses**

Data TLB Misses (IoZone Random)



Read



Write

SCMFS suffers from TLB misses.

Why higher TLB misses in SCMFS?

Scmfs works here, use small page size (4K)

RamDisk works here, use big page size (2M)

Memory Map (x86_64)

(=47 bits) user space
hole caused by [48:63] sign extension
(=47 bits) nvmalloc space
(=40 bits) guard hole
(=64 TB) direct mapping of all phys.
(=40 bits) hole
(=45 bits) vmalloc/ioremap space
(=40 bits) hole
(=40 bits) virtual memory map (1TB)
(=512 MB) kernel text mapping, from phys 0
(=1536 MB) module mapping space

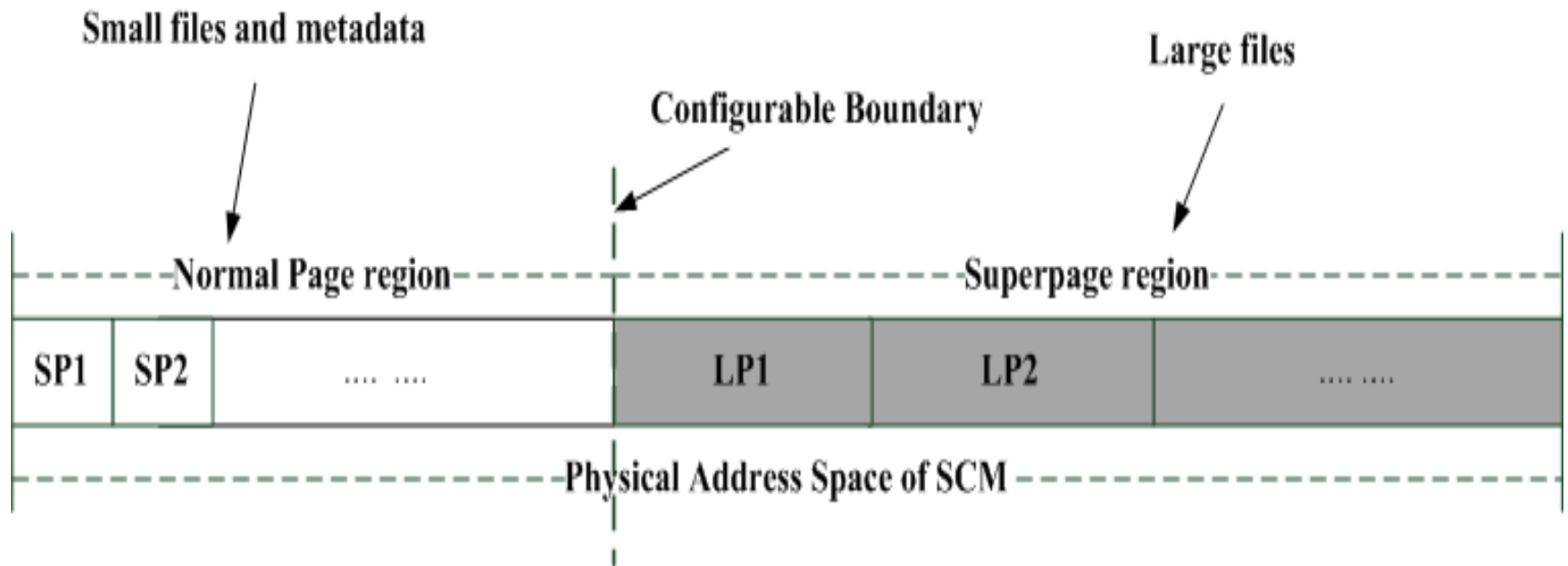
00000000000000000000 - 00000000000000000000	ffff000000000000 - ffff7fffffff
ffff800000000000 - ffff80fffffff	ffff880000000000 - ffffc7fffffff
ffffc80000000000 - ffffc8fffffff	ffffc90000000000 - ffffe8fffffff
ffffe90000000000 - ffffe9fffffff	ffffea0000000000 - ffffeafffffffff
ffffffffff80000000 - ffffffffafa0000000	fffffffffffa000000 - ffffffff00000000



How to reduce Data TLB Misses?

- Utilize two types of memory pages (i.e. 4K and 2M Bytes on x86_64)
- Use regular 4kb page for small files and metadata
- Pre-allocate superpages and allocate large files on them

Management of two types of memory pages

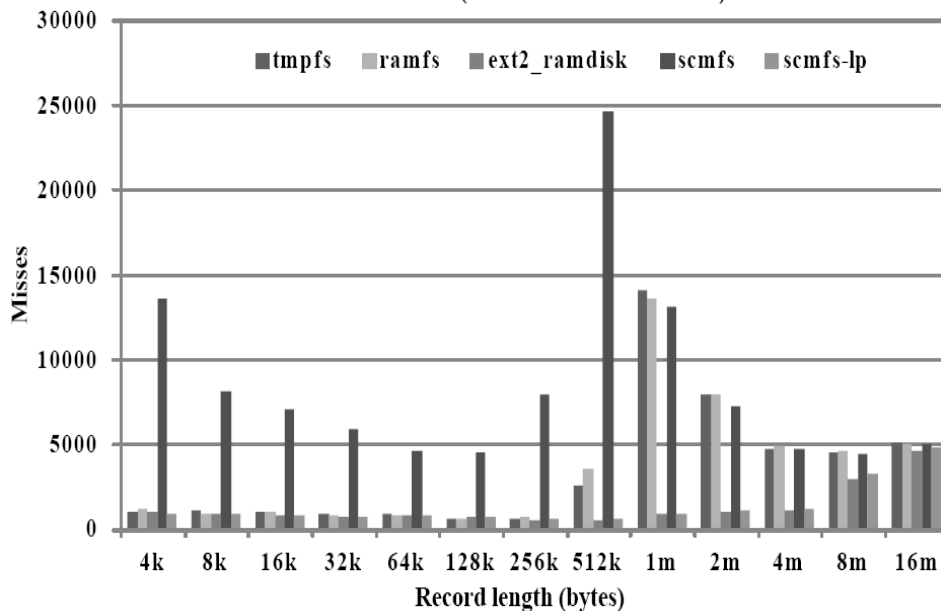


Evaluation

- Environment
 - 2.33GHz Intel Core2 Quad Processor Q8200
 - 8GB RAM, 4GB is used as SCM.
 - Linux 2.6.33
- Benchmarks
 - IoZone
 - Postmark
- Use Performance Counter to analyze the results

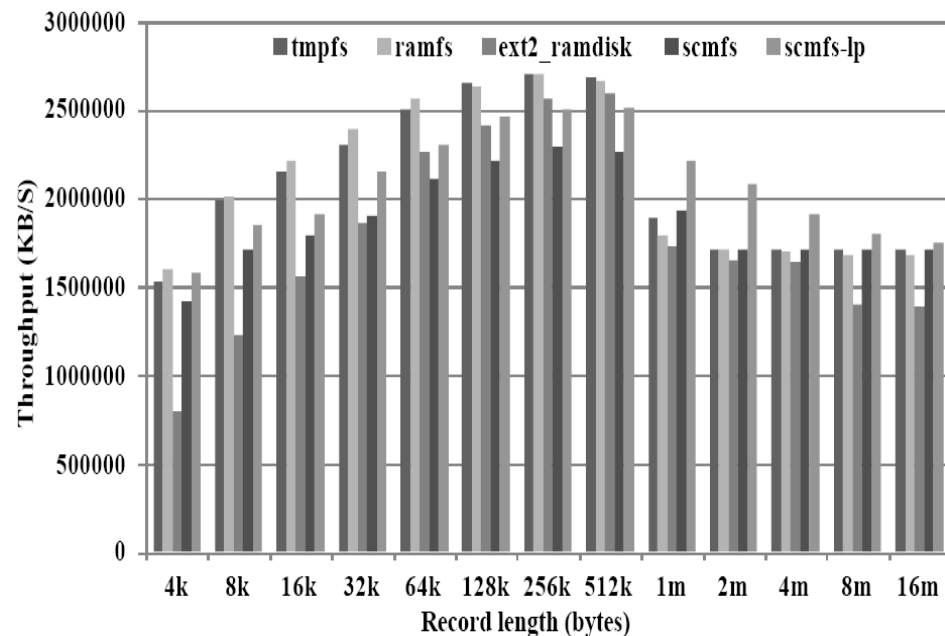
IOZONE Random Write

Data TLB misses (Random Write Workload)



Data TLB Misses

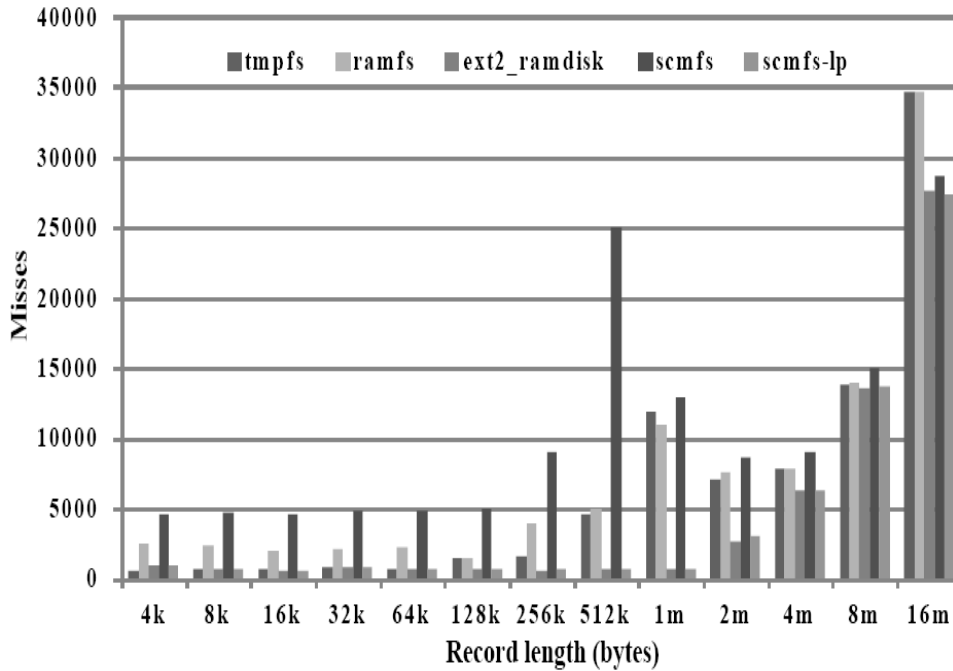
Random Write Performance



Throughput(kbytes/s)

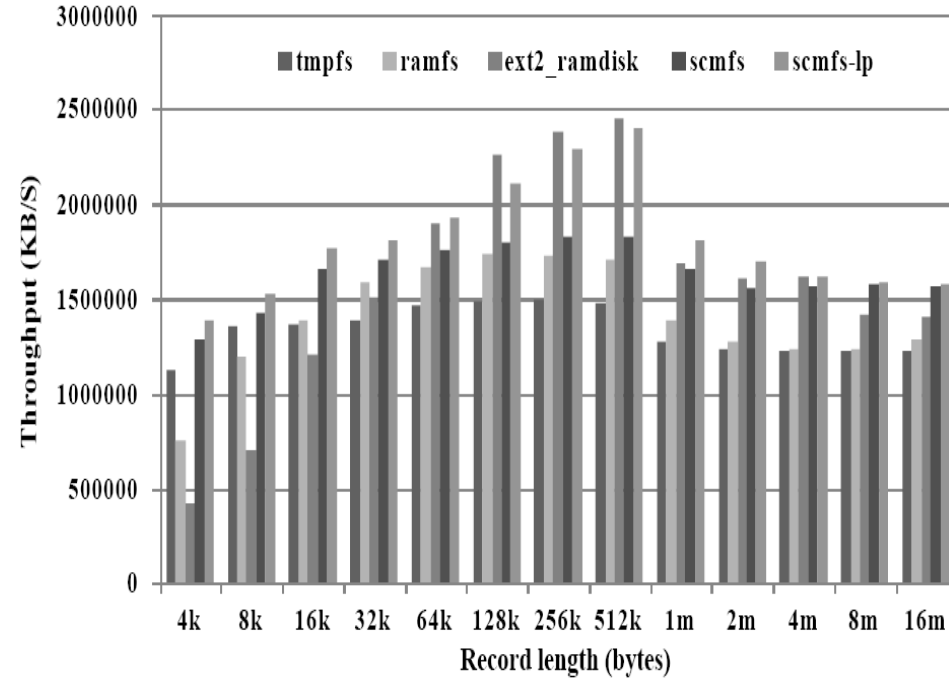
IOZONE Sequential Write

Data TLB misses (Sequential Write Workload)



Data TLB misses

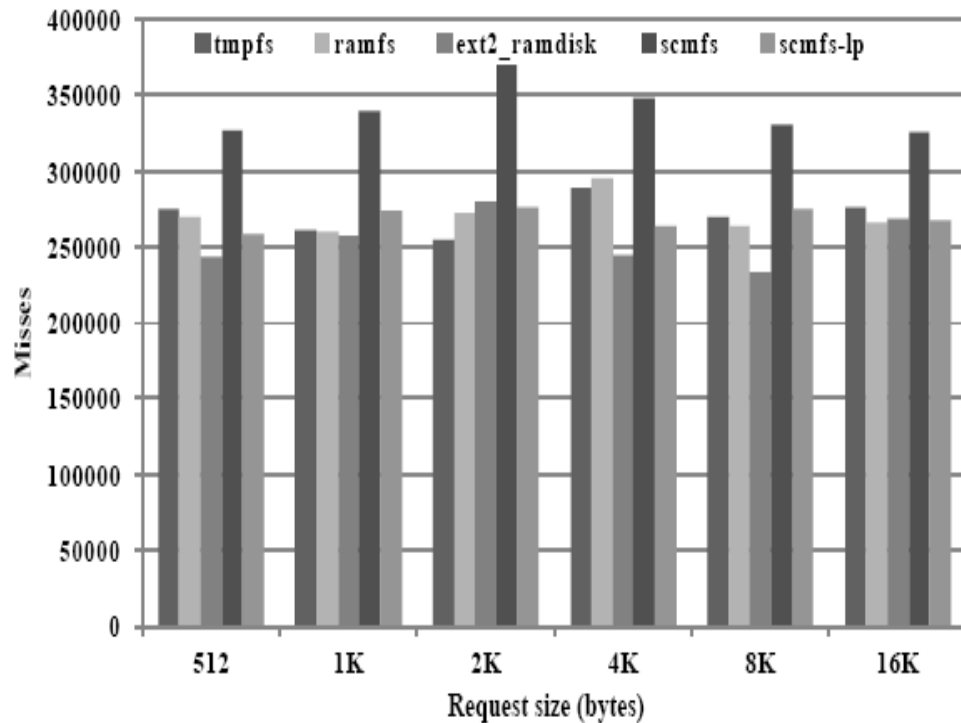
Sequential Write Performance



Throughput (KB/S)

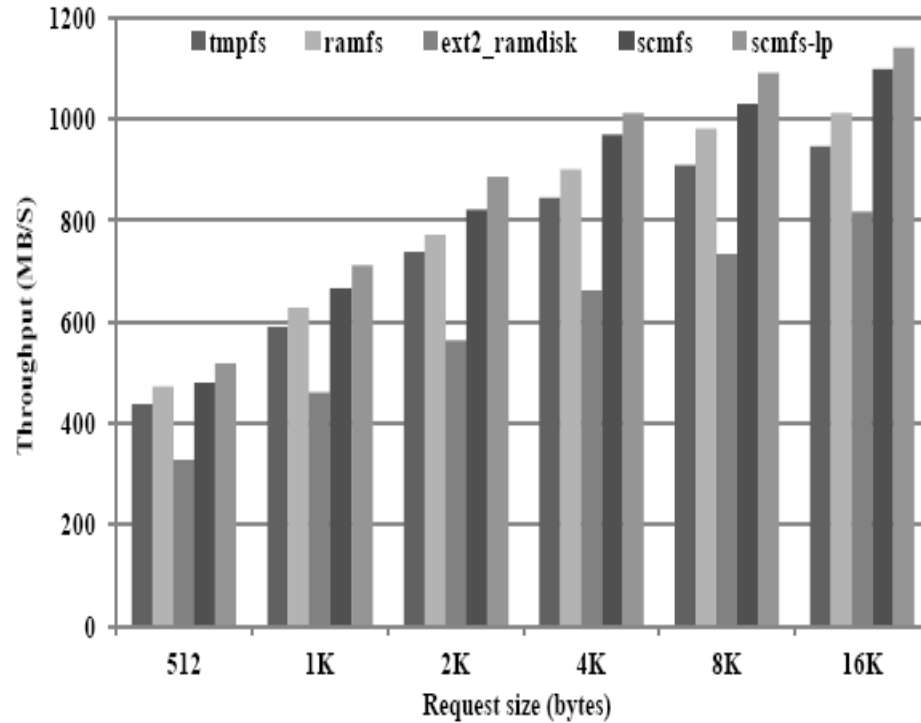
Postmark Write

Data TLB misses (Write Workload)



Data TLB misses

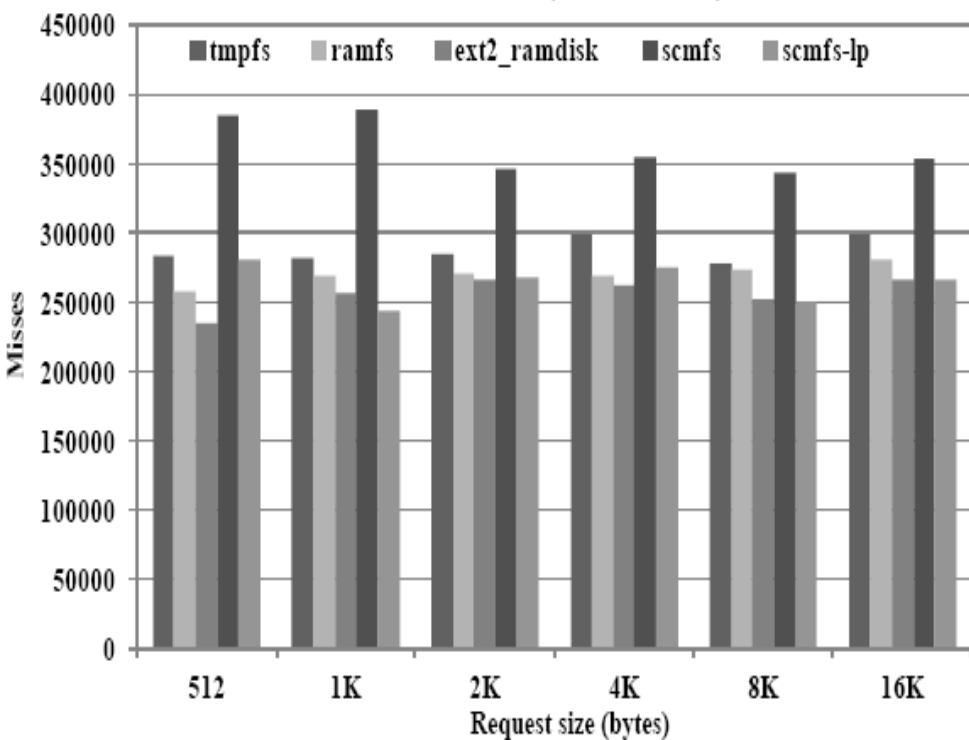
Write performance



Throughput (MB/S)

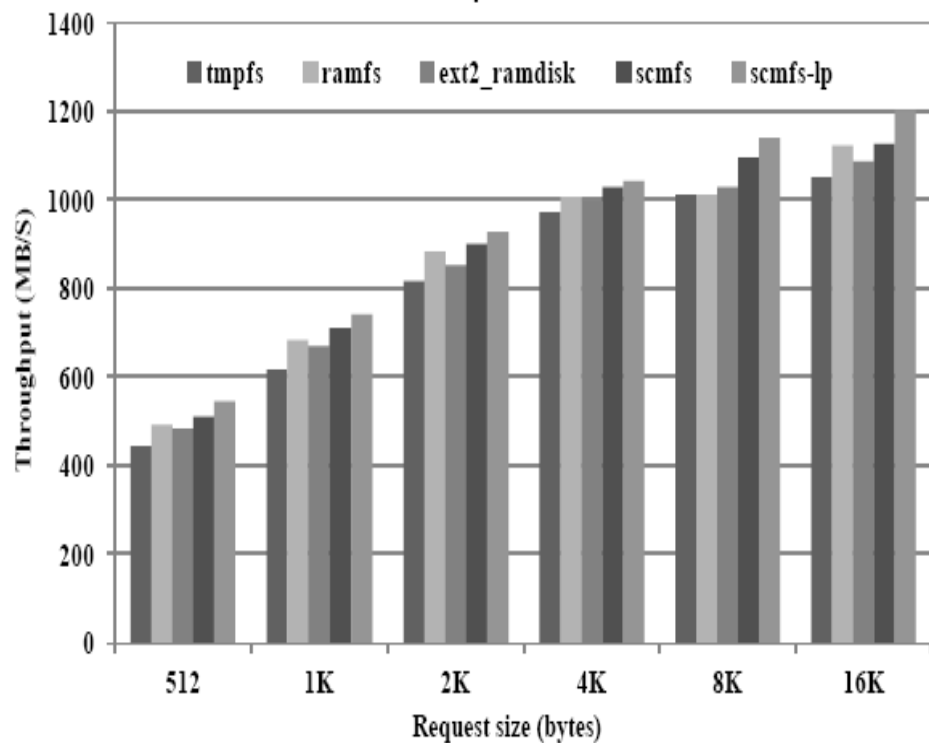
Postmark Read

Data TLB misses (Read Workload)



Data TLB misses

Read performance



Throughput (MB/S)

Conclusion

- We effectively reduced Data TLB misses by utilizing both regular and super pages.
- The FS's performance is further improved.
- Design of File System should adapt to the change of hardware hierarchy.
- Performance depends on more factors than ever.

Thanks