Unioning of the Buffer Cache and Journaling Layers with Non-volatile Memory

Hyokyung Bahn (Ewha University)

Contents

- Reliability issues in storage systems
 - Consistency problem
 - Journaling techniques
- Consistency problem with non-volatile memory
 - Non-volatile memory technology overview
 - Data inconsistency with non-volatile memory
- Unioning of the Buffer cache and Journaling layers
 - In-place commit and system recovery of UBJ
 - Cache performance of UBJ
 - Performance evaluation
- Conclusion

A man working hard ...



A man working hard ...

A problem has been detected to your computer.	and windows has been sl	hut down to prevent damage
PAGE_FAULT_IN_NONPAGED_AREA	SYSTEM	
If this is the first time y restart your computer. If t these steps:	CRASH	or screen, 1, follow
Check to make sure any new If this is a new installati for any windows updates you	PLEASE WAIT	properly installed. software manufacturer
If problems continue, disak or software. Disable BIOS m If you need to use Safe Moc your computer, press F8 to select Safe Mode.		nstalled hardware ching or shadowing. components, restart Options, and then
Technical information:		
*** STOP: 0x0000050 (0x800	10205,0x00000001,0x88598	32A5,0x00000000)



 Sudden power failure may incur file system inconsistency in hierarchical memory systems



How to solve this problem?

Prevent data inconsistency by journaling techniques

- ext3, ext4, ReiserFS, XFS, btrFS
- Frequent commit increases write traffic to storage significantly, leading to the performance degradation



Non-volatile memory promises to replace DRAM in main memory

1. Scaling Limit of DRAM



2. Power consumption

As much as 40% of the total system energy is consumed by the main memory subsystem in a mid-range IBM eServer machine. (Querish, ISCA 2009)

Replacing DRAM with STT-RAM in data centers can reduce power by up to 75% (NVMW '10, Driskill)

3. Demand for fast memory access



As critical applications are becoming more data-centric, memory performance is fast becoming the key bottleneck

Non-volatile Memory Technology

	SRAM	DRAM	Disk	NAND	PCRAM	RRAM	MRAM
				Flash		(Memristor)	(STT-RAM)
Maturity	Product	Product	Product	Product	Advanced development	Early development	Advanced development
Cell Size	>100 F ²	6-8 F ²	(2/3) F ²	4-5 F ²	8-16 F ²	>5 F ²	37 F ²
Read	<10 ns	10-60 ns	8.5 ms	25 µs	48 ns	<10 ns	<10 ns
Latency							
Write	<10 ns	10-60 ns	9.5 ms	200 µs	40-150 ns	~10 ns	12.5 ns
Latency							
Energy per	>1 pJ	2 pJ	100-	10 nJ	100 pJ	2 pJ	0.02 pJ
bit access			1000 mJ				
Static Power	Yes	Yes	Yes	No	No	No	No
Endurance	>1015	>1015	>1015	10 ⁴	10 ⁸	105	>1015
Nonvolatility	No	No	Yes	Yes	Yes	Yes	Yes
	Cu	rrent Memor	y Technolog	gies	Emerg	ing NVM Techno	logies

Source: T. Perez, C. A. F. D Rose, Technical Report, PUCRS, 2010

Scalability [



High-performance

Non-volatile Memory Technology

	SRAM	DRAM	Disk	NAND Flash	PCRAM		RRAM Memristor	MRAM (STT-RAM)
Maturity	Product	Product	Product	Product	Advanced development		arly Jevelopment	Advanced development
Cell Size	>100 F ²	6-8 F ²	(2/3) F ²	4-5 F ²	8-16 F ²	Γ	>5 F ²	37 F ²
Read Latency	<10 ns	10-60 ns	8.5 ms	25 µs	48 ns		<10 ns	<10 ns
Write Latency	<10 ns	10-60 ns	9.5 ms	200 µs	40-150 ns		*10 ns	12.5 ns
Energy per bit access	>1 pJ	2 pJ	100- 1000 mJ	10 nJ	100 pJ		2 bì	0.02 pJ
Static Power	Yes	Yes	Yes	No	No		No	No
Endurance	>1015	>1015	>1015	10 ⁴	10 ⁸		105	>1015
Nonvolatility	No	No	Yes	Yes	Yes	Γ	/es	Yes
	Cu	rrent Memor	y Technolog	gies	Emerg	ir	g NVM Techr	ologies









(Optimistic expectations)





Unioning of Buffer cache and Journaling Layers (UBJ)

- Goal is providing data consistency without sacrificing performance
- Simply adopting non-volatile memory does not suffice
- Novel buffer cache architecture Journal Area called "UBJ"
- Subsume functions of caching and journaling by using a data block for dual purposes
- Make a journaling effect just by changing the status of cache block instead of storage writes



Workings of UBJ



Workings of UBJ



System recovery of UBJ



Cache performance of UBJ



Secondary storage

1. Original buffer cache



Secondary storage

2. Separate journaling

Buffer Cache+ Journal area

Secondary storage

3. UBJ

Cache performance of UBJ



- Prototype of UBJ on Linux 2.6.38
- Compare with ext4 in journal-mode
 - logs both data and metadata
- Intel Core i3-2100 CPU
 - 3.1GHz and 4GB of DDR2-800 memory
- Emulate non-volatile memory with DRAM
- Three benchmarks
 - Filebench, IOzone, Postmark

Filebench



Improve execution time and throughput by 30.7% and 59.8% on average

Iozone



Improve performance by 110% on average

Postmark



Improve performance by 109% on average

 Effectiveness of UBJ on performance as the commit period changes



The latency of ext4 becomes small as the commit period is longer The latency of UBJ is not sensitive to the commit period changes

Conclusion

- Novel non-volatile memory buffer cache architecture
- Subsumes the functions of caching and journaling
 - Buffer cache blocks $\leftarrow \rightarrow$ Journal logs
 - Notion of a frozen state
 - In-place Commit
 - Journal log block as well as a cache block
- Performance results
 - Implemented on Linux 2.6.38
 - Compared to ext4 in journal mode
 - Improve I/O performance by %76 and up to 240%